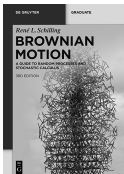


Werner Linde
Probability Theory

Also of Interest



Brownian Motion. A Guide to Random Processes and Stochastic Calculus

René L. Schilling, 2021

ISBN 978-3-11-074125-4, e-ISBN (PDF) 978-3-11-074127-8,
e-ISBN (EPUB) 978-3-11-074149-0



Probability Theory and Statistics with Real World Applications. Univariate and Multivariate Models Applications

Peter Zörnig, 2024

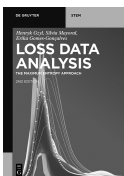
ISBN 978-3-11-133220-8, e-ISBN (PDF) 978-3-11-133227-7,
e-ISBN (EPUB) 978-3-11-133232-1



Bitcoin: A Game-Theoretic Analysis

Micah Warren, 2023

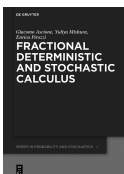
ISBN 978-3-11-077283-8, e-ISBN (PDF) 978-3-11-077284-5,
e-ISBN (EPUB) 978-3-11-077305-7



Loss Data Analysis. The Maximum Entropy Approach

Henryk Gzyl, Silvia Mayoral, Erika Gomes-Gonçalves, 2023

ISBN 978-3-11-104738-6, e-ISBN (PDF) 978-3-11-104818-5,
e-ISBN (EPUB) 978-3-11-104970-0



Fractional Deterministic and Stochastic Calculus

Giacomo Ascione, Yuliya Mishura, Enrica Pirozzi, 2023

ISBN 978-3-11-077981-3, e-ISBN (PDF) 978-3-11-078001-7,
e-ISBN (EPUB) 978-3-11-078022-2

Werner Linde

Probability Theory

A First Course in Probability Theory and Statistics

2nd edition

DE GRUYTER

Mathematics Subject Classification 2020

Primary: 60-01, 62-01; Secondary: 60A05

Author

Prof. Dr. Werner Linde
Friedrich-Schiller-Universität Jena
Fakultät für Mathematik & Informatik
Institut für Stochastik
D-07737 Jena
Germany
Werner.Linde@mathematik.uni-jena.de

ISBN 978-3-11-132484-5

e-ISBN (PDF) 978-3-11-132506-4

e-ISBN (EPUB) 978-3-11-132517-0

Library of Congress Control Number: 2024931814

Bibliographic information published by the Deutsche Nationalbibliothek

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available on the Internet at <http://dnb.dnb.de>.

© 2024 Walter de Gruyter GmbH, Berlin/Boston

Cover image: Density Diagram created by Werner Linde

Typesetting: VTeX UAB, Lithuania

Printing and binding: CPI books GmbH, Leck

www.degruyter.com

To my wife Karin

Preface

This book is intended as an introductory course for students in mathematics, physical sciences, engineering, or in other related fields. It is based on the experience of probability lectures taught during the past 25 years, where the spectrum reached from two-hour introductory courses, over Measure Theory and advanced probability classes, to such topics as Stochastic Processes and Mathematical Statistics. Until 2012 these lectures were delivered to students at the University of Jena (Germany), and since 2013 to those at the University of Delaware in Newark (USA).

The book is the completely revised version of the German edition “Stochastik für das Lehramt,” which appeared in 2014 at De Gruyter. At most universities in Germany, there exist special classes in Probability Theory for students who want to become teachers of mathematics in high schools. Besides basic facts about Probability Theory, these courses are also supposed to give an introduction into Mathematical Statistics. Thus, the original main intention for the German version was to write a book that helps those students understand Probability Theory better. But soon the book turned out to also be useful as introduction for students in other fields, e. g., in mathematics, physics, and so on. Thus we decided, in order to make the book applicable for a broader audience, to provide a translation in the English language.

During numerous years of teaching, I learned the following:

- Probabilistic questions are usually easy to formulate, generally have a tight relation to everyday problems, and therefore attract the interest of the audience. Every student knows the phenomena that occur when one rolls a die, plays cards, tosses a coin, or plays a lottery. Thus, an initial interest in Probability Theory exists.
- In contrast, after a short time many students have very serious difficulties with understanding the presented topics. Consequently, a common opinion among students is that Probability Theory is a very complicated topic, causing a lot of problems and troubles.

Surely there exist several reasons for the bad image of Probability Theory among students. But, as we believe, the most important one is as follows. In Probability Theory, the type of problems and questions considered, as well as the way of thinking, differs considerably from the problems, questions, and thinking in other fields of mathematics, i. e., from fields with which the students became acquainted before attending a probability course. For example, in Calculus a function has a well-described domain of definition; mostly it is defined by a concrete formula, has certain properties as continuity, differentiability, and so on. A function is something very concrete which can be made vivid by drawing its graph. In contrast, in Probability Theory functions are mostly investigated as random variables. They are defined on a completely unimportant, nonspecified sample space, and they generally do not possess a concrete formula for their definition. It may even happen that only the existence of a function (random variable) is known. The only property of a random variable which really matters is the distribution of its values.

This and many other similar techniques make the whole theory something mysterious and not completely comprehensible.

Considering this observation, we organized the book in a way that tries to make probabilistic problems more understandable and that puts the focus more onto explanations of the definitions, notations, and results. The tools we use to do this are examples; we present at least one before a new definition, in order to motivate it, followed by more examples after the definition to make it comprehensible. Here we act upon the maxim expressed by Einstein's quote:¹

Example isn't another way to teach, it is the only way to teach.

Presenting the basic results and methods in Probability Theory *without* using results, facts, and notations from Measure Theory is, in our opinion, as difficult as to square the circle. Either one restricts oneself to discrete probability measures and random variables or one has to be imprecise. There is no other choice! In some places, it is possible to avoid the use of measure-theoretic facts, such as the Lebesgue integral, or the existence of infinite product measures, and so on, but the price is high.² Of course, I also struggled with the problem of missing facts from Measure Theory while writing this book. Therefore, I tried to include some ideas and results about σ -fields, measures, and integrals, hoping that a few readers become interested and want to learn more about Measure Theory. For those, we refer to the books [Coh13, Dud02], or [Bil12] as good sources.

In this context, let us make some remark about the verification of the presented results. Whenever it was possible, we tried to prove the stated results. Times have changed; when I was a student, every theorem presented in a mathematical lecture was proved – really every one. Facts and results without proof were doubtful and soon forgotten. And a tricky and elegant proof is sometimes more impressive than the proven result (at least to us). Hopefully, some readers will like some of the proofs in this book as much as we did.

One of most used applications of Probability Theory is Mathematical Statistics. When I met former students of mine, I often asked them which kind of mathematics they are mainly using now in their daily work. The overwhelming majority of them answered that one of their main fields of mathematical work is statistical problems. Therefore, we decided to include an introductory chapter about Mathematical Statistics. Nowadays, due to the existence of good and fast statistical programs, it is very easy to analyze data, to evaluate confidence regions, or to test a given hypothesis. But

¹ See <http://www.alberteinstein.com/quotes/einsteinquotes.html>

² For example, several years ago, to avoid the use of the Lebesgue integral, I introduced the expected value of a random variable as a Riemann integral via its distribution function. This is mathematically correct, but at the end almost no students understood what the expected value really is. Try to prove that the expected value is linear using this approach!

do those who use these programs also always know what they are doing? Since we doubt that this is so, we stressed the focus in this chapter to the question of why the main statistical methods work and on what mathematical background they rest. We also investigate how precise statistical decisions are and what kinds of errors may occur.

The organization of this book differs a little bit from those in many other first-course books about Probability Theory. Having Measure Theory in the back of our minds causes us to think that probability measures are the most important ingredient of Probability Theory; random variables come in second. On the contrary, many other authors go exactly the other way. They start with random variables, and probability measures then occur as their distribution on their range spaces (mostly \mathbb{R}). In this case, a standard normal probability measure does not exist, only a standard normal distributed random variable. Both approaches have their advantages and disadvantages, but as we said, for us the probability measures are interesting in their own right, and therefore we start with them in Chapter 1, followed by random variables in Section 3.

The book also contains some facts and results that are more advanced and usually not part of an introductory course in Probability Theory. Such topics are, for example, the investigation of product measures, order statistics, and so on. We have assigned those more involved sections with a star. They may be skipped at a first reading without loss in the following chapters.

At the end of each chapter, one finds a collection of some problems related to the contents of the section. Here we restricted ourselves to a few problems in the actual task; the solutions of these problems are helpful to the understanding of the presented topics. The problems are mainly taken from our collection of homeworks and exams during the past years. For those who want to work with more problems we refer to many books, e. g., [GS01, Gha19, Pao06], or [Rss14], which contain a huge collection of probabilistic problems, ranging from easy to difficult, from natural to artificial, from interesting to boring.

Finally, I want to express my thanks to those who supported my work at the translation and revision of the present book. Many students at the University of Delaware helped me improve my English and correct wrong phrases and expressions. To mention all of them is impossible. But among them were a few students who read whole chapters and, without them, the book would have never been finished (or readable). In particular, I want to mention Emily Wagner and Spencer Walker. They both really did a great job. Many thanks! Let me also express my gratitude to Colleen McInerney, Rachel Austin, Daniel Atadan, and Quentin Dubroff, all students in Delaware and attending my classes for some time. They also read whole sections of the book and corrected my broken English. Finally, my thanks go to Professor Anne Leucht from the Technical University in Braunschweig (Germany); her field of work is Mathematical Statistics, and her hints and remarks about Chapter 8 in this book were important to me.

And last but not least, I want to thank the Department of Mathematical Sciences at the University of Delaware for the excellent working conditions after my retirement in Germany.

Newark, Delaware, June 6, 2016

Werner Linde

Changes in the second edition: The first edition of my textbook was well received by scholars and students alike, and I would like to thank all of them for their comments and positive criticisms.

There are a few changes to the second edition. For instance, I added more than 40 new examples, so that now the book contains about 280 of them. Among the new ones are some classical examples as, e. g., the “Boy or Girl Paradox,” the “Secretary Problem,” the “Two-Envelope Paradox,” or “Gambler’s Ruin,” which may be found in Sections 5.4 or 5.5, respectively. Other examples have been included for better understanding of the general, sometimes quite abstract, topics.

Section 8 about Mathematical Statistics contains now three tables which summarize the main tests and confidence regions for normal distributed populations. I hope that these résumés help to get a quick overview about the most used techniques in Mathematical Statistics. Furthermore, there is a new Section 8.6.4 about confidence regions for hypergeometric distributed samples, an important topic missing in the first edition.

A completely new ingredient in this edition are short summaries at the end of almost every section. Here I give a compressed overview about basic notions and results presented in the preceding section. The aim of these abstracts is to tell the reader what were the most important statements and what can possibly be omitted from the first reading.

I believe that graphical presentations of abstract mathematical statements are a very helpful aid for better understanding, not only for beginners. Therefore, I added more than 80 new figures, so that now the book contains more than 100 of them. The increase is mainly due to the fact that there exist now powerful tools as, e. g., TikZ for drawing convincing figures, tools which either did not yet exist or which I was not aware of when writing the first edition.

Besides I added several new problems, updated the list of references, and completed it by adding a few classical books about Measure Theory and Probability as, for example, [Hal14, Kal21] or [Par05].

Some minor misprints or incorrect arguments have been eliminated, a few parts were rewritten in order to make them, as I hope, clearer and better understandable.

Finally, I would like to thank my former student Frank Aurzada, TU Darmstadt, for the fruitful discussions about some newly added examples. Last but not least, I want to express my gratitude to Nadja Schedensack from De Gruyter for her advice and helpful remarks concerning layout and TEX problems.

Jena, Germany, January 8, 2024

Werner Linde

Contents

Preface — VII

1 Probabilities — 1

- 1.1 Probability spaces — **1**
- 1.1.1 Sample spaces — **1**
- 1.1.2 σ -fields of events* — **3**
- 1.1.3 Probability measures — **6**
- 1.2 Basic properties of probability measures — **10**
- 1.3 Discrete probability measures — **15**
- 1.4 Special discrete probability measures — **20**
 - 1.4.1 Dirac measure — **20**
 - 1.4.2 Uniform distribution on a finite set — **21**
 - 1.4.3 Binomial distribution — **27**
 - 1.4.4 Multinomial distribution — **30**
 - 1.4.5 Poisson distribution — **34**
 - 1.4.6 Hypergeometric distribution — **36**
 - 1.4.7 Geometric distribution — **41**
 - 1.4.8 Negative binomial distribution — **43**
- 1.5 Continuous probability measures — **48**
- 1.6 Special continuous distributions — **52**
 - 1.6.1 Uniform distribution on an interval — **52**
 - 1.6.2 Normal distribution — **55**
 - 1.6.3 Gamma distribution — **58**
 - 1.6.4 Exponential distribution — **61**
 - 1.6.5 Erlang distribution — **63**
 - 1.6.6 Chi-squared distribution — **64**
 - 1.6.7 Beta distribution — **65**
 - 1.6.8 Cauchy distribution — **68**
- 1.7 Distribution function — **69**
- 1.8 Multivariate continuous distributions — **81**
 - 1.8.1 Multivariate density functions — **81**
 - 1.8.2 Multivariate uniform distribution — **83**
- 1.9 Products of probability spaces* — **90**
 - 1.9.1 Product σ -fields and measures — **90**
 - 1.9.2 Product measures: discrete case — **95**
 - 1.9.3 Product measures: continuous case — **99**
- 1.10 Problems — **103**

2 Conditional probabilities and independence — 111

- 2.1 Conditional probabilities — **111**

2.2 Independence of events — **120**

2.3 Problems — **127**

3 Random variables and their distribution — 131

3.1 Transformation of random values — **131**

3.2 Probability distribution of a random variable — **133**

3.3 Special random variables — **143**

3.4 Random vectors — **145**

3.5 Joint and marginal distributions — **147**

3.5.1 Marginal distributions: discrete case — **149**

3.5.2 Marginal distributions: continuous case — **154**

3.6 Independence of random variables — **157**

3.6.1 Independence of discrete random variables — **161**

3.6.2 Independence of continuous random variables — **165**

3.7 Order statistics* — **169**

3.8 Problems — **176**

4 Operations on random variables — 180

4.1 Mappings of random variables — **180**

4.2 Linear transformations — **185**

4.3 Coin tossing versus uniform distribution — **189**

4.3.1 Binary fractions — **189**

4.3.2 Binary fractions of random numbers — **191**

4.3.3 Random numbers generated by coin tossing — **193**

4.4 Simulation of random variables — **196**

4.5 Addition of random variables — **202**

4.5.1 Sums of discrete random variables — **204**

4.5.2 Sums of continuous random variables — **207**

4.6 Sums of certain random variables — **210**

4.7 Products and quotients of random variables — **222**

4.7.1 Student's t -distribution — **226**

4.7.2 F-distribution — **228**

4.8 Problems — **231**

5 Expected value, variance, and covariance — 235

5.1 Expected value — **235**

5.1.1 Expected value of discrete random variables — **235**

5.1.2 Expected value of certain discrete random variables — **238**

5.1.3 Expected value of continuous random variables — **242**

5.1.4 Expected value of certain continuous random variables — **245**

5.1.5 Properties of the expected value — **249**

5.2 Variance — **257**

- 5.2.1 Higher moments of random variables — **257**
- 5.2.2 Variance of random variables — **261**
- 5.2.3 Variance of certain random variables — **263**
- 5.3 Covariance and correlation — **269**
- 5.3.1 Covariance — **269**
- 5.3.2 Correlation coefficient — **276**
- 5.4 Some paradoxes and examples — **279**
- 5.4.1 Boy or girl paradox — **279**
- 5.4.2 Randomly chosen entries — **282**
- 5.4.3 Secretary problem — **283**
- 5.4.4 Two-envelope paradox — **287**
- 5.5 Gambler's ruin — **293**
- 5.6 Problems — **304**

- 6 Normally distributed random vectors — 310**
- 6.1 Representation and density — **310**
- 6.2 Expected value and covariance matrix — **318**
- 6.3 Problems — **324**

- 7 Limit theorems — 327**
- 7.1 Laws of large numbers — **327**
- 7.1.1 Chebyshev's inequality — **327**
- 7.1.2 Infinite sequences of independent random variables* — **330**
- 7.1.3 Borel–Cantelli lemma* — **333**
- 7.1.4 Weak law of large numbers — **340**
- 7.1.5 Strong law of large numbers — **341**
- 7.2 Central limit theorem — **347**
- 7.3 Problems — **367**

- 8 Mathematical statistics — 370**
- 8.1 Statistical models — **370**
- 8.1.1 Nonparametric statistical models — **370**
- 8.1.2 Parametric statistical models — **374**
- 8.2 Statistical hypothesis testing — **376**
- 8.2.1 Hypotheses and tests — **376**
- 8.2.2 Power function and significance tests — **379**
- 8.3 Tests for binomial distributed populations — **386**
- 8.4 Tests for normally distributed populations — **391**
- 8.4.1 Fisher's theorem — **392**
- 8.4.2 Quantiles — **395**
- 8.4.3 Z-tests or Gauss tests — **400**
- 8.4.4 t-tests — **404**

8.4.5	χ^2 -tests for the variance —	406
8.4.6	Two-sample Z-tests —	408
8.4.7	Two-sample t-tests —	410
8.4.8	F-tests —	412
8.5	Point estimators —	414
8.5.1	Maximum likelihood estimation —	416
8.5.2	Unbiased estimators —	423
8.5.3	Risk function —	427
8.6	Confidence regions and intervals —	432
8.6.1	Construction of confidence regions —	432
8.6.2	Normally distributed samples —	435
8.6.3	Binomial distributed populations —	438
8.6.4	Hypergeometric distributed populations —	444
8.7	Problems —	447

A Appendix — 451

A.1	Notations —	451
A.2	Elements of set theory —	451
A.2.1	Set operations —	451
A.2.2	Preimages of sets —	454
A.2.3	Problems —	455
A.3	Combinatorics —	456
A.3.1	Binomial coefficients —	456
A.3.2	Drawing balls out of an urn —	461
A.3.3	Multinomial coefficients —	465
A.3.4	Problems —	467
A.4	Vectors and matrices —	468
A.5	Some analytic tools —	472

Bibliography — 479

Index — 481

1 Probabilities

1.1 Probability spaces

The basic concern of Probability Theory is to model experiments involving randomness, that is, experiments with nondetermined outcome, shortly called *random experiments*. The Russian mathematician A. N. Kolmogorov established the modern Probability Theory in 1933 by publishing his book (cf. [Kol33]) *Grundbegriffe der Wahrscheinlichkeitsrechnung*. In it, he postulated the following:

Random experiments are described by probability spaces $(\Omega, \mathcal{A}, \mathbb{P})$.



The triple $(\Omega, \mathcal{A}, \mathbb{P})$ comprises a *sample space* Ω , a σ -field \mathcal{A} of events, and a mapping \mathbb{P} from \mathcal{A} to $[0, 1]$, called *probability measure* or *probability distribution*.

Let us now explain the three different components of a probability space in detail. We start with the sample space.

1.1.1 Sample spaces

Definition 1.1.1. The **sample space** Ω is a nonempty set that contains (at least) all possible outcomes of the random experiment.

Remark 1.1.2. Due to mathematical reasons, sometimes it can be useful to choose Ω larger than necessary. It is only important that the sample space contains *all* possible results.

Example 1.1.3. When rolling a die one time, the natural choice for the sample space is $\Omega = \{1, \dots, 6\}$. However, it would also be possible to take $\Omega = \{1, 2, \dots\}$ or even $\Omega = \mathbb{R}$. In contrast, $\Omega = \{1, \dots, 5\}$ is not suitable for the description of the experiment.

Example 1.1.4. Roll a die until the number “6” shows up for the first time. Record the number of necessary rolls until the first appearance of “6”. The suitable sample space in this case is $\Omega = \{1, 2, \dots\}$. Any finite set $\{1, 2, \dots, N\}$ is not appropriate because, even if we choose N very large, we can never be 100 % sure that the first “6” really appears during the first N rolls.

Example 1.1.5. Suppose in a lottery 6 numbers out of 49 are chosen. If we record the 6 numbers in the way they are chosen, then a suitable sample space is

$$\Omega = \{(\omega_1, \dots, \omega_6) : 1 \leq \omega_i \leq 49, \omega_i \neq \omega_j, i \neq j\}.$$

But, in general, the chosen numbers are published ordered by their size. If this is so, as corresponding sample space we may choose

$$\Omega = \{(\omega_1, \dots, \omega_6) : 1 \leq \omega_1 < \dots < \omega_6 \leq 49\}.$$

Example 1.1.6. Say we have three urns, each containing white and black balls. Choose one urn at random and take out a ball. Register if the chosen ball is white or black. Letting $\Omega = \{w, b\}$ as sample space is not appropriate. It does not take into account which urn we had chosen. One suitable sample space would be $\Omega = \{(w, i), (b, i), i = 1, 2, 3\}$. Then, for instance, $(w, 1)$ occurs if we choose urn 1, take out a ball of this urn, and the chosen ball is a white one.

Example 1.1.7. A light bulb is switched on at time zero and burns for a certain period of time. At some random time $t > 0$, it burns out. To describe this experiment, we have to take into account all possible times $t > 0$. Therefore, a natural choice for the sample space in this case is $\Omega = (0, \infty)$, or, if we do not exclude that the bulb is defective from the very beginning, then $\Omega = [0, \infty)$.

Example 1.1.8. Customers arrive at the counter of a bank at certain random times. So, for example, the first customer shows up at time t_1 , the second at time $t_2 > t_1$, and so on. Then the sample space consists of infinite sequences $t_1 < t_2 < \dots$ of positive real numbers assuming that at least one customer enters the bank that day. But, in fact, because the number of customers per day is finite, one may also choose

$$\Omega = \{(t_1, \dots, t_n) : 0 < t_1 < \dots < t_n, n \in \mathbb{N}\} \cup \{0\},$$

where $\{0\}$ occurs if no customer arrives at that day.

Subsets of the sample space Ω are called **events**. In other words, the powerset $\mathcal{P}(\Omega)$ is the collection of all possible events. For example, when we roll a die once there are exactly $2^6 = 64$ possible events, as, for example,

$$\emptyset, \{1\}, \dots, \{6\}, \{1, 2\}, \dots, \{1, 6\}, \{2, 3\}, \dots, \{2, 6\}, \dots, \{1, 2, 3, 4, 5\}, \Omega.$$

Among all events, there are some of special interest, the so-called **elementary events**. These are events containing exactly one element. In Example 1.1.3, the elementary events are

$$\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \text{ and } \{6\}.$$

Remark 1.1.9. Never confuse the elementary events with the points that they contain. Look at Example 1.1.3. There we have $6 \in \Omega$ and for the generated elementary event holds $\{6\} \in \mathcal{P}(\Omega)$.

Let $A \subseteq \Omega$ be an event. After executing the random experiment, one observes a result $\omega \in \Omega$. Then two cases are possible:

1. The outcome ω belongs to A . In this case, we say that the event A **occurred**.
2. If ω is not in A , that is, if $\omega \in A^c$, then the event A did **not occur**.

Example 1.1.10. Roll a die once and let $A = \{2, 4\}$. Say the outcome was number “6”. Then A did not occur. But, if we obtained number “2,” then A occurred.

Example 1.1.11. Suppose we roll a die twice. The describing sample space consists of all 36 pairs of numbers from 1 to 6. If

$$A = \{(1, 1), \dots, (6, 6)\},$$

then A occurs if and only if the outcome of the first roll equals the that of the second roll.

Example 1.1.12. In Example 1.1.7, the occurrence of an event $A = [T, \infty)$ tells us that the light bulb burned out after time T or, in other words, at time T it was still shining.

Let us formulate some easy *rules* for the occurrence of events.

1. By the choice of the sample space, the event Ω always occurs. Therefore, Ω is also called the **certain** event.
2. The empty set never occurs. Thus it is called the **impossible** event.
3. An event A occurs if and only if the complementary event A^c does not, and vice versa, A does not occur if and only if A^c does.
4. If A and B are two events, then $A \cup B$ occurs if at least one of the two events occurs. Hereby we do not exclude that A and B may both occur.
5. The event $A \cap B$ occurs if and only if A and B both occur.

1.1.2 σ -fields of events*

The basic aim of Probability Theory is to assign to each event A a number $\mathbb{P}(A)$ in $[0, 1]$, which describes the likelihood of its occurrence. If the occurrence of an event A is very likely, then $\mathbb{P}(A)$ should be close to 1 while $\mathbb{P}(A)$ close to zero suggests that the appearance of A is very unlikely.¹ The mapping $A \mapsto \mathbb{P}(A)$ must possess certain natural properties. Unfortunately, by mathematical reason it is not always possible to assign to each event A a number $\mathbb{P}(A)$ such that $A \mapsto \mathbb{P}(A)$ has the desired properties. The solution is ingenious and one of the key observations in Kolmogorov’s approach: one chooses a subset $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ such that $\mathbb{P}(A)$ is only defined for $A \in \mathcal{A}$. If $A \notin \mathcal{A}$, then $\mathbb{P}(A)$ does not exist. Of course, \mathcal{A} should be chosen as large as possible and, moreover, at least “ordinary” sets should belong to \mathcal{A} .

¹ If the weather forecast predicts a 70 % chance of rain, you will surely take your umbrella with you. On the contrary, if the forecast is only 20 %, you will probably not do so. Why? In this case the probability of the occurrence of the event A , “it rains,” is significantly less likely.



In the case of “large” sample spaces, it is in general impossible to assign to each event a meaningful likelihood of its occurrence. Consequently, for “large” sample spaces, as, for example, \mathbb{R} or \mathbb{R}^n , the probability $\mathbb{P}(A)$ is defined only for certain special events A .

The collection \mathcal{A} of events for which $\mathbb{P}(A)$ is well defined has to satisfy some algebraic conditions. More precisely, the following properties are supposed.

Definition 1.1.13. A collection \mathcal{A} of subsets of Ω is called a σ -field if

- (1) $\emptyset \in \mathcal{A}$,
- (2) if $A \in \mathcal{A}$ then $A^c \in \mathcal{A}$, and
- (3) for countably many A_1, A_2, \dots in \mathcal{A} , it follows that $\bigcup_{j=1}^{\infty} A_j \in \mathcal{A}$.

Let us verify some easy properties of σ -fields.

Proposition 1.1.14. *Let \mathcal{A} be a σ -field of subsets of Ω . Then the following are valid:*

- (i) $\Omega \in \mathcal{A}$.
- (ii) If A_1, A_2, \dots, A_n are finitely many sets in \mathcal{A} , then $\bigcup_{j=1}^n A_j \in \mathcal{A}$.
- (iii) If A_1, A_2, \dots belong to \mathcal{A} , then so does $\bigcap_{j=1}^{\infty} A_j$.
- (iv) Whenever $A_1, \dots, A_n \in \mathcal{A}$, then $\bigcap_{j=1}^n A_j \in \mathcal{A}$.

Proof. Assertion (i) is a direct consequence of $\emptyset \in \mathcal{A}$ combined with property (2) of σ -fields.

To verify (ii), let A_1, \dots, A_n be in \mathcal{A} . Set $A_{n+1} = A_{n+2} = \dots = \emptyset$. Then for all $j = 1, 2, \dots$, we have $A_j \in \mathcal{A}$ and, by property (3) of σ -fields, also $\bigcup_{j=1}^{\infty} A_j \in \mathcal{A}$. But note that we have $\bigcup_{j=1}^{\infty} A_j = \bigcup_{j=1}^n A_j$, hence (ii) is valid.

To prove (iii), we first observe that $A_j \in \mathcal{A}$ yields $A_j^c \in \mathcal{A}$, hence $\bigcup_{j=1}^{\infty} A_j^c \in \mathcal{A}$. Another application of (2) implies $(\bigcup_{j=1}^{\infty} A_j^c)^c \in \mathcal{A}$. De Morgan’s rule asserts

$$\left(\bigcup_{j=1}^{\infty} A_j^c \right)^c = \bigcap_{j=1}^{\infty} A_j,$$

which completes the proof of (iii).

Assertion (iv) may be derived from an application of (ii) to the complementary sets, as we did in the proof of (iii). Or one can use the method in the proof of (ii), but this time we choose $A_{n+1} = A_{n+2} = \dots = \Omega$. \square

Corollary 1.1.15. *If sets A and B belong to a σ -field \mathcal{A} , then so do $A \cup B$, $A \cap B$, $A \setminus B$, and $A \Delta B$.*

The easiest examples of σ -fields are either $\mathcal{A} = \{\emptyset, \Omega\}$ or $\mathcal{A} = \mathcal{P}(\Omega)$. However, the former σ -field is much too small for applications while the latter is generally too big, at least if the sample space is uncountably infinite. We will shortly indicate how one constructs suitable σ -fields in the case of “large” sample spaces as, for example, \mathbb{R} or \mathbb{R}^n .

Proposition 1.1.16. *Let \mathcal{C} be an arbitrary nonempty collection of subsets of Ω . Then there is a σ -field \mathcal{A} possessing the following properties:*

1. *It holds that $\mathcal{C} \subseteq \mathcal{A}$ or, verbally, each set $C \in \mathcal{C}$ belongs to the σ -field \mathcal{A} .*
2. *The σ -field \mathcal{A} is the smallest one possessing this property. That is, whenever \mathcal{A}' is another σ -field with $\mathcal{C} \subseteq \mathcal{A}'$, then $\mathcal{A} \subseteq \mathcal{A}'$.*

Proof. Let Φ be the collection of all σ -fields \mathcal{A}' on Ω for which $\mathcal{C} \subseteq \mathcal{A}'$, that is,

$$\Phi := \{\mathcal{A}' \subseteq \mathcal{P}(\Omega) : \mathcal{C} \subseteq \mathcal{A}', \mathcal{A}' \text{ is a } \sigma\text{-field}\}.$$

The collection Φ is nonempty because it contains at least one element, namely the power set of Ω . Of course, $\mathcal{P}(\Omega)$ is a σ -field and $\mathcal{C} \subseteq \mathcal{P}(\Omega)$ trivially, hence $\mathcal{P}(\Omega) \in \Phi$.

Next define \mathcal{A} by

$$\mathcal{A} := \bigcap_{\mathcal{A}' \in \Phi} \mathcal{A}' = \{A \subseteq \Omega : A \in \mathcal{A}', \forall \mathcal{A}' \in \Phi\}.$$

It is not difficult to prove that \mathcal{A} is a σ -field with $\mathcal{C} \subseteq \mathcal{A}$. Indeed, if $C \in \mathcal{C}$, then $C \in \mathcal{A}'$ for all $\mathcal{A}' \in \Phi$, hence, by construction of \mathcal{A} , we get $C \in \mathcal{A}$.

Furthermore, \mathcal{A} is also the smallest σ -field containing \mathcal{C} . To see this, take an arbitrary σ -field $\tilde{\mathcal{A}}$ containing \mathcal{C} . Then $\tilde{\mathcal{A}} \in \Phi$, which implies $\mathcal{A} \subseteq \tilde{\mathcal{A}}$ because \mathcal{A} is the intersection over all σ -fields in Φ . This completes the proof. \square

Definition 1.1.17. Let \mathcal{C} be an arbitrary nonempty collection of subsets of Ω . The smallest σ -field containing \mathcal{C} is called the **σ -field generated by \mathcal{C}** . It is denoted by $\sigma(\mathcal{C})$.

Remark 1.1.18. The σ -field $\sigma(\mathcal{C})$ is characterized by the three following properties:

1. $\sigma(\mathcal{C})$ is a σ -field.
2. $\mathcal{C} \subseteq \sigma(\mathcal{C})$.
3. If $\mathcal{C} \subseteq \mathcal{A}'$ for some σ -field \mathcal{A}' , then $\sigma(\mathcal{C}) \subseteq \mathcal{A}'$.

Example 1.1.19. Let Ω be an arbitrary nonempty set. If $\mathcal{C} = \{C_1, \dots, C_n\}$ where

$$C_i \cap C_j = \emptyset \quad \text{if } i \neq j \quad \text{and} \quad \bigcup_{j=1}^n C_j = \Omega,$$

then $\sigma(\mathcal{C})$ consists of the empty set and any finite unions of the C_j s.

For example, if $\Omega = \{1, \dots, 6\}$ and $C_1 = \{1, 2\}$, $C_2 = \{3\}$, and $C_3 = \{4, 5, 6\}$, then the elements of the generated σ -field are the empty set \emptyset and C_1 , C_2 , C_3 , $C_1 \cup C_2$, $C_1 \cup C_3$, $C_2 \cup C_3$, and Ω .

Definition 1.1.20. Let $\mathcal{C} \subseteq \mathcal{P}(\mathbb{R})$ be the collection of all finite closed intervals in \mathbb{R} , that is,

$$\mathcal{C} = \{[a, b] : a < b, a, b \in \mathbb{R}\}.$$

The σ -field generated by \mathcal{C} is denoted by $\mathcal{B}(\mathbb{R})$ and is called the **Borel σ -field**. If $B \in \mathcal{B}(\mathbb{R})$, then it is said to be a **Borel set**.

Remark 1.1.21. By construction, every closed interval in \mathbb{R} is a Borel set. Furthermore, the properties of σ -fields also imply that complements of such intervals, their countable unions, and intersections are Borel sets. One might believe that all subsets of \mathbb{R} are Borel sets. This is not the case; for the construction of a non-Borel set, we refer to [Gha19, Example 1.21] or [Dud02, pages 105–108].

Remark 1.1.22. There exist many other systems of subsets in \mathbb{R} generating $\mathcal{B}(\mathbb{R})$. Let us only mention two of them:

$$\mathcal{C}_1 = \{(-\infty, b] : b \in \mathbb{R}\} \quad \text{or} \quad \mathcal{C}_2 = \{(a, \infty) : a \in \mathbb{R}\}.$$

Summary: There is a σ -field of subsets in \mathbb{R} (also in \mathbb{R}^n , as we will see later on), the collection of Borel sets, for which we may always define their probability of occurrence. All sets of interest are Borel sets. Thus, in fact, knowing probabilities only for those sets is necessary from a mathematical point of view, but for our purposes this is only a theoretical restriction.

1.1.3 Probability measures

The occurrence of an event in a random experiment is not completely haphazard. Although we are not able to predict the outcome of the next trial, the occurrence or nonoccurrence of an event follows certain rules. Some events are more likely to occur, others less. The degree of likelihood of an event A is described by a number $\mathbb{P}(A)$, called the probability of the occurrence of A (in short, probability of A). The most common scale for probabilities is $0 \leq \mathbb{P}(A) \leq 1$, where the larger $\mathbb{P}(A)$, the more likely A is to occur. One could also think of other scales as $0 \leq \mathbb{P}(A) \leq 100$. In fact, this is even quite often used; in this sense, a chance of 50 % equals a probability of $1/2$.

What does it mean that an event A has probability $\mathbb{P}(A)$? For example, what does it tell us that an event occurs with probability $1/2$? Does this mean a half-occurrence of A ? Surely not.

To answer this question, we have to assume that we execute an experiment not only once² but several, say n , times. Thereby we have to ensure that the conditions of

² It does not make sense to speak of the probability of an event that can be executed only once. For example, it is (mathematically) absurd to ask for the probability that the Eiffel Tower will be in Paris for yet another 100 years.

the experiment do not change and that the single results do not depend on each other. Let

$$a_n(A) := \text{Number of trials where } A \text{ occurs.}$$

The quantity $a_n(A)$ is called the **absolute frequency** of the occurrence of A in n trials. Observe that $a_n(A)$ is a random number with $0 \leq a_n(A) \leq n$. Next we set

$$r_n(A) := \frac{a_n(A)}{n} \quad (1.1)$$

and name it the **relative frequency** of the occurrence of A in n trials. This number is random as well, but now $0 \leq r_n(A) \leq 1$.

Example 1.1.23. At www.westlotto.de/lotto-6aus49 one finds a summary of the 6223 drawings in the German lottery between 10/09/1955 and 06/21/2023. Every time there were 6 numbers chosen out of 49. To describe the experiment, we take as sample space

$$\Omega = \{A \subset \{1, \dots, 49\} : |A| = 6\}.$$

Then we get $|\Omega| = \binom{49}{6}$. Given $1 \leq j \leq 49$, define the event \mathcal{S}_j as

$$\mathcal{S}_j = \{A \in \Omega : j \in A\}.$$

That is, \mathcal{S}_j occurs provided that the number j was among the chosen ones. Since $|\mathcal{S}_j| = \binom{48}{5}$, we obtain

$$\mathbb{P}(\mathcal{S}_j) = \frac{\binom{48}{5}}{\binom{49}{6}} = \frac{6}{49} \approx 0.1224.$$

In the above cited summary, one finds the absolute frequency of all events \mathcal{S}_j with $j = 1, \dots, 49$. For example, the frequencies \mathcal{S}_6 , \mathcal{S}_{20} , and \mathcal{S}_{45} equal 825, 719, and 699, respectively. That is, during the past $n = 6223$ drawings, the numbers 6, 20, and 45 appeared 825, 719, and 699 times, respectively. Thus, their relative frequencies are

$$r_n(\mathcal{S}_6) = \frac{825}{6223} \approx 0.1326, \quad r_n(\mathcal{S}_{20}) = \frac{719}{6223} \approx 0.1155, \quad \text{and}$$

$$r_n(\mathcal{S}_{45}) = \frac{699}{6223} \approx 0.1123.$$

One should compare these relative frequencies with the expected one (assumed that all numbers are equally likely), given by

$$\mathbb{P}(\mathcal{S}_j) = \frac{6}{49} = 0.1224, \quad j = 1, \dots, 49.$$

It is somehow intuitively clear³ that the relative frequencies of an event A converge to a (nonrandom) number as $n \rightarrow \infty$. And this limit is exactly the desired probability of the occurrence of the event A . Let us express this in a different way: say we execute an experiment n times for some large n . Then, on average, we will observe $n \cdot \mathbb{P}(A)$ occurrences of A . For example, when rolling a fair die many times, an even number will happen in approximately half the cases.

Or, in the setting of Example 1.1.23, on average each number from 1 to 49 should approximately show up $6223 \cdot \frac{6}{49} = 762$ times. Thus, the observed frequencies tell us that either the drawings were not fair (some numbers are more likely than others) or that the number $n = 6223$ of drawings is still too small.

Which natural properties of $A \mapsto \mathbb{P}(A)$ may be deduced from $r_n(A) \xrightarrow{n \rightarrow \infty} \mathbb{P}(A)$?

1. Since $0 \leq r_n(A) \leq 1$, we conclude $0 \leq \mathbb{P}(A) \leq 1$.
2. Because of $r_n(\Omega) = 1$ for each $n \geq 1$, we get $\mathbb{P}(\Omega) = 1$.
3. The property $r_n(\emptyset) = 0$ yields $\mathbb{P}(\emptyset) = 0$.
4. Let A and B be two disjoint events. Then $r_n(A \cup B) = r_n(A) + r_n(B)$, hence the limits should satisfy a similar relation, that is,

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) . \quad (1.2)$$

Definition 1.1.24. A mapping \mathbb{P} fulfilling eq. (1.2) for disjoint A and B is called **finitely additive**.

Remark 1.1.25. Applying eq. (1.2) successively leads to the following. If A_1, \dots, A_n are disjoint, then

$$\mathbb{P}\left(\bigcup_{j=1}^n A_j\right) = \sum_{j=1}^n \mathbb{P}(A_j) .$$

Finite additivity is a very useful property of probabilities, and in the case of finite sample spaces, it completely suffices to build a fruitful theory. But as soon as the sample space is infinite, it is too weak. To see this, let us come back to Example 1.1.4. Assume we want to evaluate the probability of the event $A = \{2, 4, 6, \dots\}$, that is, the first “6” appears in an even number of trials. Then we have to split A into (infinitely) many disjoint events $\{2\}, \{4\}, \dots$. The finite additivity of \mathbb{P} does not suffice to get $\mathbb{P}(A) = \mathbb{P}(\{2\}) + \mathbb{P}(\{4\}) + \dots$. In order to evaluate $\mathbb{P}(A)$ in this way, we need the following stronger property of \mathbb{P} .

Definition 1.1.26. A mapping \mathbb{P} is said to be σ -**additive** provided that for countably many disjoint A_1, A_2, \dots in Ω we get

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} \mathbb{P}(A_j) .$$

³ We will discuss this question more precisely in Section 7.1.

Let us summarize what we have until now: a mapping \mathbb{P} assigning each event its probability should possess the following natural properties:

1. For all A , one has $0 \leq \mathbb{P}(A) \leq 1$.
2. We have $\mathbb{P}(\emptyset) = 0$ and $\mathbb{P}(\Omega) = 1$.
3. The mapping \mathbb{P} is σ -additive.

Thus, given a sample space Ω , we look for a function \mathbb{P} defined on $\mathcal{P}(\Omega)$ satisfying the previous properties. But, as already mentioned, if Ω is uncountable, for example, if $\Omega = \mathbb{R}$, then only very special⁴ \mathbb{P} with these properties exist.

To overcome these difficulties, in such cases we have to restrict \mathbb{P} to a σ -field $\mathcal{A} \subseteq \mathcal{P}(\Omega)$.

Definition 1.1.27. Let Ω be a sample space and let \mathcal{A} be a σ -field of subsets of Ω . A function $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ is called a **probability measure** or **probability distribution** on (Ω, \mathcal{A}) if

1. $\mathbb{P}(\emptyset) = 0$ and $\mathbb{P}(\Omega) = 1$.
2. \mathbb{P} is σ -additive, that is, for each sequence of disjoint sets $A_j \in \mathcal{A}$, $j = 1, 2, \dots$, it follows that

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} \mathbb{P}(A_j). \quad (1.3)$$

Remark 1.1.28. Note that the left-hand side of eq. (1.3) is well defined. Indeed, since \mathcal{A} is a σ -field, $A_j \in \mathcal{A}$ implies $\bigcup_{j=1}^{\infty} A_j \in \mathcal{A}$ as well.

Now we are in a position to define probability spaces in the exact way.

Definition 1.1.29. A **probability space** is a triple $(\Omega, \mathcal{A}, \mathbb{P})$, where Ω is a sample space, \mathcal{A} denotes a σ -field consisting of subsets of Ω , and $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ is a probability measure.

Remark 1.1.30. Given $A \in \mathcal{A}$, the number $\mathbb{P}(A)$ describes its probability or, more precisely, its probability of occurrence. Subsets A of Ω with $A \notin \mathcal{A}$ do *not* possess a probability.

Let us demonstrate a simple example on how to construct a probability space for a given random experiment. Several other examples will follow soon.

Example 1.1.31. We ask for a probability space that describes rolling a fair die once. Of course, $\Omega = \{1, \dots, 6\}$ and $\mathcal{A} = \mathcal{P}(\Omega)$. The mapping $\mathbb{P} : \mathcal{P}(\Omega) \rightarrow [0, 1]$ is given by

$$\mathbb{P}(A) = \frac{|A|}{6}, \quad A \subseteq \{1, \dots, 6\}.$$

Recall that $|A|$ denotes the cardinality of the set A .

⁴ Discrete ones as we will investigate in Section 1.3.

Remark 1.1.32. Suppose we want to find a model for some concrete random experiment. How do we do this? In most cases, the sample space is immediately determined by the results we expect. If the question about Ω is settled, the choice of the σ -field depends on the size of the sample space. If Ω is finite or countably finite, then we may choose $\mathcal{A} = \mathcal{P}(\Omega)$. If $\Omega = \mathbb{R}$ or even \mathbb{R}^n , we take the corresponding Borel σ -fields. The challenging task is the determination of the probability measure \mathbb{P} . Here the following approaches are possible:

1. *Theoretical considerations* quite often lead to the determination of \mathbb{P} . For example, since the faces of a fair die are all equally likely, this already describes \mathbb{P} completely. Similar arguments can be used for certain games or also for lotteries.
2. If theoretical considerations are neither possible nor available then *statistically investigations* may help. This approach is based on the fact that the relative frequencies $r_n(A)$ converge to $\mathbb{P}(A)$. Thus, one executes n trials of the experiment and records the relative frequency of the occurrence of A . For example, one may question n randomly chosen persons or do n independent measurements of the same item. Then $r_n(A)$ may be used to approximate the value of $\mathbb{P}(A)$.
3. Sometimes also *subjective or experience-based* approaches can be used to find approximate probabilities. These may be erroneous, but they might give some hint for the correct distribution. For example, if a new product is on the market, the distribution of its lifetime is not yet known. At the beginning one uses data of an already existing similar product. After some time, data about the new product become available, the probabilities can be determined more accurately.

Summary: A probability space $(\Omega, \mathcal{A}, \mathbb{P})$ is a triple where Ω is a nonempty set, called sample space, \mathcal{A} denotes a σ -field of subsets of Ω , and, finally, $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ is a probability measure (or probability distribution). The probability measure \mathbb{P} possesses the following properties: it is normalized so that $\mathbb{P}(\emptyset) = 0$ and $\mathbb{P}(\Omega) = 1$, and it is σ -additive. Given an event $A \in \mathcal{A}$, the number $\mathbb{P}(A)$ describes the likelihood of its occurrence.

1.2 Basic properties of probability measures

Probability measures obey many useful properties. Let us summarize the most important ones in the next proposition.

Proposition 1.2.1. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Then the following are valid:*

- (1) *The measure \mathbb{P} is also finitely additive.*
- (2) *If $A, B \in \mathcal{A}$ satisfy $A \subseteq B$, then $\mathbb{P}(B \setminus A) = \mathbb{P}(B) - \mathbb{P}(A)$.*
- (3) *We have $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$ for $A \in \mathcal{A}$.*
- (4) *Probability measures are **monotone**, that is, if $A \subseteq B$ for some $A, B \in \mathcal{A}$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$.*

(5) Probability measures are **subadditive**, that is, for all (not necessarily disjoint) events⁵ $A_j \in \mathcal{A}$,

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) \leq \sum_{j=1}^{\infty} \mathbb{P}(A_j). \quad (1.4)$$

(6) Probability measures are **continuous from below**, that is, whenever $A_j \in \mathcal{A}$ satisfy $A_1 \subseteq A_2 \subseteq \dots$,

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) = \lim_{j \rightarrow \infty} \mathbb{P}(A_j).$$

(7) In a similar way, each probability measure is **continuous from above**: if $A_j \in \mathcal{A}$ satisfy $A_1 \supseteq A_2 \supseteq \dots$, then

$$\mathbb{P}\left(\bigcap_{j=1}^{\infty} A_j\right) = \lim_{j \rightarrow \infty} \mathbb{P}(A_j).$$

Proof. To prove (1), choose disjoint A_1, \dots, A_n in \mathcal{A} and set $A_{n+1} = A_{n+2} = \dots = \emptyset$. Then A_1, A_2, \dots are infinitely many disjoint events in \mathcal{A} , hence the σ -additivity of \mathbb{P} implies

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} \mathbb{P}(A_j).$$

Observe that $\bigcup_{j=1}^{\infty} A_j = \bigcup_{j=1}^n A_j$ and $\mathbb{P}(A_j) = 0$ if $j > n$, so the previous equation reduces to

$$\mathbb{P}\left(\bigcup_{j=1}^n A_j\right) = \sum_{j=1}^n \mathbb{P}(A_j),$$

and \mathbb{P} is finitely additive.

To prove (2), write $B = A \cup (B \setminus A)$ and observe that this is a disjoint decomposition of B . Hence, by the finite additivity of \mathbb{P} , we obtain

$$\mathbb{P}(B) = \mathbb{P}(A) + \mathbb{P}(B \setminus A).$$

Relocating $\mathbb{P}(A)$ to the left-hand side proves (2).

An application of (2) to Ω and A leads to

$$\mathbb{P}(A^c) = \mathbb{P}(\Omega \setminus A) = \mathbb{P}(\Omega) - \mathbb{P}(A) = 1 - \mathbb{P}(A),$$

which proves (3).

⁵ Estimate (1.4) is also known as Boole's inequality.

The monotonicity is an easy consequence of (2). Indeed,

$$\mathbb{P}(B) - \mathbb{P}(A) = \mathbb{P}(B \setminus A) \geq 0,$$

implying $\mathbb{P}(B) \geq \mathbb{P}(A)$.

To prove inequality (1.4), choose arbitrary A_1, A_2, \dots in \mathcal{A} . Set $B_1 := A_1$ and, if $j \geq 2$, then

$$B_j := A_j \setminus (A_1 \cup \dots \cup A_{j-1}).$$

Then B_1, B_2, \dots are disjoint subsets in \mathcal{A} with $\bigcup_{j=1}^{\infty} B_j = \bigcup_{j=1}^{\infty} A_j$. Furthermore, by the construction, $B_j \subseteq A_j$, hence $\mathbb{P}(B_j) \leq \mathbb{P}(A_j)$. An application of all these properties yields

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) = \mathbb{P}\left(\bigcup_{j=1}^{\infty} B_j\right) = \sum_{j=1}^{\infty} \mathbb{P}(B_j) \leq \sum_{j=1}^{\infty} \mathbb{P}(A_j).$$

Thus (5) is proved.

Let us turn now to the continuity from below. Choose A_1, A_2, \dots in \mathcal{A} satisfying $A_1 \subseteq A_2 \subseteq \dots$. With $A_0 := \emptyset$, set

$$B_k := A_k \setminus A_{k-1}, \quad k = 1, 2, \dots$$

The B_k s are disjoint and, moreover, $\bigcup_{k=1}^{\infty} B_k = \bigcup_{j=1}^{\infty} A_j$. Furthermore, because of $A_{k-1} \subseteq A_k$, from (2) we get $\mathbb{P}(B_k) = \mathbb{P}(A_k) - \mathbb{P}(A_{k-1})$. When putting this all together, it follows that

$$\begin{aligned} \mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) &= \mathbb{P}\left(\bigcup_{k=1}^{\infty} B_k\right) = \sum_{k=1}^{\infty} \mathbb{P}(B_k) = \lim_{j \rightarrow \infty} \sum_{k=1}^j \mathbb{P}(B_k) \\ &= \lim_{j \rightarrow \infty} \sum_{k=1}^j [\mathbb{P}(A_k) - \mathbb{P}(A_{k-1})] = \lim_{j \rightarrow \infty} [\mathbb{P}(A_j) - \mathbb{P}(A_0)] = \lim_{j \rightarrow \infty} \mathbb{P}(A_j), \end{aligned}$$

where we used $\mathbb{P}(A_0) = \mathbb{P}(\emptyset) = 0$. This proves the continuity from below.

Thus it remains to prove (7). For this, choose $A_j \in \mathcal{A}$ with $A_1 \supseteq A_2 \supseteq \dots$. Then the complementary sets satisfy $A_1^c \subseteq A_2^c \subseteq \dots$. The continuity from below lets us conclude that

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j^c\right) = \lim_{j \rightarrow \infty} \mathbb{P}(A_j^c) = \lim_{j \rightarrow \infty} [1 - \mathbb{P}(A_j)] = 1 - \lim_{j \rightarrow \infty} \mathbb{P}(A_j). \quad (1.5)$$

But

$$\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j^c\right) = 1 - \mathbb{P}\left(\left(\bigcup_{j=1}^{\infty} A_j^c\right)^c\right) = 1 - \mathbb{P}\left(\bigcap_{j=1}^{\infty} A_j\right),$$

and plugging this into eq. (1.5) gives

$$\mathbb{P}\left(\bigcap_{j=1}^{\infty} A_j\right) = \lim_{j \rightarrow \infty} \mathbb{P}(A_j),$$

as asserted. \square

Remark 1.2.2. Property (2) becomes false without the assumption $A \subseteq B$. But because of $B \setminus A = B \setminus (A \cap B)$ and $A \cap B \subseteq B$, we always have

$$\mathbb{P}(B \setminus A) = \mathbb{P}(B) - \mathbb{P}(A \cap B). \quad (1.6)$$

Another useful property of probability measures is as follows.

Proposition 1.2.3. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Then for all $A_1, A_2 \in \mathcal{A}$, it follows that*

$$\mathbb{P}(A_1 \cup A_2) = \mathbb{P}(A_1) + \mathbb{P}(A_2) - \mathbb{P}(A_1 \cap A_2). \quad (1.7)$$

Proof. Write the union of the two sets as

$$A_1 \cup A_2 = A_1 \cup [A_2 \setminus (A_1 \cap A_2)]$$

and note that the two sets on the right-hand side are disjoint. Because of $A_1 \cap A_2 \subseteq A_2$, property (2) of Proposition 1.2.1 applies and leads to

$$\mathbb{P}(A_1 \cup A_2) = \mathbb{P}(A_1) + \mathbb{P}(A_2 \setminus (A_1 \cap A_2)) = \mathbb{P}(A_1) + [\mathbb{P}(A_2) - \mathbb{P}(A_1 \cap A_2)].$$

This completes the proof. \square

Given $A_1, A_2, A_3 \in \mathcal{A}$, an application of the previous proposition to A_1 and $A_2 \cup A_3$ implies

$$\begin{aligned} \mathbb{P}(A_1 \cup A_2 \cup A_3) &= \mathbb{P}(A_1) + \mathbb{P}(A_2) + \mathbb{P}(A_3) \\ &\quad - [\mathbb{P}(A_1 \cap A_2) + \mathbb{P}(A_1 \cap A_3) + \mathbb{P}(A_2 \cap A_3)] \\ &\quad + \mathbb{P}(A_1 \cap A_2 \cap A_3). \end{aligned}$$

Another application of eq. (1.7) to the second and third terms in the right-hand sum proves the following result (compare Figure 1.1).

Proposition 1.2.4. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and let A_1, A_2 , and A_3 be in \mathcal{A} . Then*

$$\begin{aligned} \mathbb{P}(A_1 \cup A_2 \cup A_3) &= \mathbb{P}(A_1) + \mathbb{P}(A_2) + \mathbb{P}(A_3) \\ &\quad - [\mathbb{P}(A_1 \cap A_2) + \mathbb{P}(A_1 \cap A_3) + \mathbb{P}(A_2 \cap A_3)] \\ &\quad + \mathbb{P}(A_1 \cap A_2 \cap A_3). \end{aligned}$$

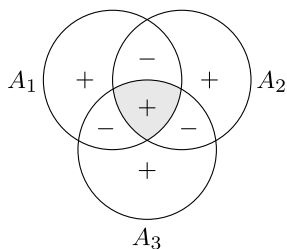


Figure 1.1: The inclusion–exclusion formula for three sets.

Remark 1.2.5. A generalization of Propositions 1.2.3 and 1.2.4 from 2 or 3 to an arbitrary number of sets can be found in Problem 1.7. It is the so-called **inclusion–exclusion formula**. For example, if $n = 4$, this formula says, given events A_1, A_2, A_3, A_4 , that

$$\mathbb{P}\left(\bigcup_{i=1}^4 A_i\right) = \sum_{i=1}^4 \mathbb{P}(A_i) - \sum_{1 \leq i < j \leq 4} \mathbb{P}(A_i \cap A_j) + \sum_{1 \leq i < j < k \leq 4} \mathbb{P}(A_i \cap A_j \cap A_k) - \mathbb{P}\left(\bigcap_{i=1}^4 A_i\right).$$

Let us explain two easy examples of how the properties of probability measures apply.

Example 1.2.6. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Suppose two events A and B in \mathcal{A} satisfy

$$\mathbb{P}(A) = 0.5, \quad \mathbb{P}(B) = 0.4, \quad \text{and} \quad \mathbb{P}(A \cap B) = 0.2.$$

Which probabilities do $A \cup B$, $A \setminus B$, $A^c \cup B^c$, and $A^c \cap B$ possess?

Answer: An application of Proposition 1.2.4 gives

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) = 0.4 + 0.5 - 0.2 = 0.7.$$

Furthermore, by eq. (1.6), one gets

$$\mathbb{P}(A \setminus B) = \mathbb{P}(A) - \mathbb{P}(A \cap B) = 0.5 - 0.2 = 0.3.$$

Finally, by De Morgan's rules and another application of eq. (1.6), we get

$$\mathbb{P}(A^c \cup B^c) = 1 - \mathbb{P}(A \cap B) = 0.8 \quad \text{and} \quad \mathbb{P}(A^c \cap B) = \mathbb{P}(B \setminus A) = \mathbb{P}(B) - \mathbb{P}(A \cap B) = 0.2.$$

In summary, say one has to take two exams A and B . The probability of passing exam A is 0.5, the probability of passing B equals 0.4, and to pass both it is 0.2. Then with probability 0.7 one passes at least one of the exams, with 0.3 exam A , but not B , with probability 0.8 one fails at least once, and, finally, the probability to pass B but not A is 0.2.

Example 1.2.7. Choose at random (all numbers are equally likely) a number from 1 to 1000. How likely is it that the chosen number is neither divisible by 3, nor by 5, nor by 7?

Answer: Let D be the set of numbers in $\{1, \dots, 1000\}$ which are neither divisible by 3, nor by 5, nor by 7. Then it follows that

$$D^c = A \cup B \cup C$$

where A consists of multiples of 3, B contains numbers divisible by 5, and C comprises multiples of 7. Moreover, the numbers in $A \cap B$ are multiples of 15, $A \cap C$ contains multiples of 21, and $B \cap C$ consists of multiples of 35. Finally, we note that $A \cap B \cap C$ includes only numbers divisible by 105. Easy calculations lead to

$$\begin{aligned} \mathbb{P}(A) &= 0.333, & \mathbb{P}(B) &= 0.2, & \mathbb{P}(C) &= 0.142, & \mathbb{P}(A \cap B) &= 0.066, \\ \mathbb{P}(A \cap C) &= 0.047, & \mathbb{P}(B \cap C) &= 0.028, & \text{and} & \mathbb{P}(A \cap B \cap C) &= 0.009. \end{aligned}$$

This implies

$$\begin{aligned} \mathbb{P}(D^c) &= \mathbb{P}(A) + \mathbb{P}(B) + \mathbb{P}(C) \\ &\quad - [\mathbb{P}(A \cap B) + \mathbb{P}(A \cap C) + \mathbb{P}(B \cap C)] + \mathbb{P}(A \cap B \cap C) \\ &= 0.333 + 0.2 + 0.142 - 0.066 - 0.047 - 0.028 + 0.009 = 0.543. \end{aligned}$$

So we finally arrive at

$$\mathbb{P}(D) = 1 - \mathbb{P}(D^c) = 1 - 0.543 = 0.457.$$

1.3 Discrete probability measures

We start with the investigation of *finite* sample spaces. They describe random experiments where only finitely many different results may occur, as, for example, rolling a die n times, tossing a coin finitely often, and so on. Suppose the sample space contains N different elements. Then we may enumerate these elements as follows:

$$\Omega = \{\omega_1, \dots, \omega_N\}.$$

As σ -field we choose $\mathcal{A} = \mathcal{P}(\Omega)$.

Given an arbitrary probability measure $\mathbb{P} : \mathcal{P}(\Omega) \rightarrow \mathbb{R}$, set

$$p_j := \mathbb{P}(\{\omega_j\}), \quad j = 1, \dots, N. \tag{1.8}$$

In this way we assign to each probability measure \mathbb{P} numbers p_1, \dots, p_N . Which properties do they possess? The answer to this question gives the following proposition.

Proposition 1.3.1. *If \mathbb{P} is a probability measure on $\mathcal{P}(\Omega)$, then the numbers p_j defined by eq. (1.8) satisfy*

$$0 \leq p_j \leq 1 \quad \text{and} \quad \sum_{j=1}^N p_j = 1. \quad (1.9)$$

Proof. The first property is an immediate consequence of having $\mathbb{P}(A) \geq 0$ for all $A \subseteq \Omega$. The second property of the p_j s follows from

$$1 = \mathbb{P}(\Omega) = \mathbb{P}\left(\bigcup_{j=1}^N \{\omega_j\}\right) = \sum_{j=1}^N \mathbb{P}(\{\omega_j\}) = \sum_{j=1}^N p_j. \quad \square$$

Conclusion. Each probability measure \mathbb{P} generates a sequence $(p_j)_{j=1}^N$ of real numbers satisfying the properties (1.9). Moreover, if $A \subseteq \Omega$, then we have

$$\mathbb{P}(A) = \sum_{\{j:\omega_j \in A\}} p_j. \quad (1.10)$$

In particular, the assignment $\mathbb{P} \rightarrow (p_j)_{j=1}^N$ is one-to-one.

Property (1.10) is an easy consequence of $A = \bigcup_{\{j:\omega_j \in A\}} \{\omega_j\}$. Furthermore, it tells us that \mathbb{P} is uniquely determined by the p_j s. Note that two probability measures \mathbb{P}_1 and \mathbb{P}_2 on (Ω, \mathcal{A}) coincide if $\mathbb{P}_1(A) = \mathbb{P}_2(A)$ for all $A \in \mathcal{A}$.

Now let us look at the reverse question. Suppose we are given an *arbitrary* sequence $(p_j)_{j=1}^N$ of real numbers satisfying the conditions (1.9).

Proposition 1.3.2. *Define \mathbb{P} on $\mathcal{P}(\Omega)$ by*

$$\mathbb{P}(A) = \sum_{\{j:\omega_j \in A\}} p_j. \quad (1.11)$$

Then \mathbb{P} is a probability measure satisfying $\mathbb{P}(\{\omega_j\}) = p_j$ for all $j \leq n$.

Proof. The map \mathbb{P} has values in $[0, 1]$ and $\mathbb{P}(\Omega) = 1$ by $\sum_{j=1}^N p_j = 1$. Since the summation over the empty set equals zero, $\mathbb{P}(\emptyset) = 0$.

Thus it remains to show that \mathbb{P} is σ -additive. Take disjoint subsets A_1, A_2, \dots of Ω . Since Ω is finite, there are at most finitely many of the A_j s which are nonempty. Say, for simplicity, these are the first n sets A_1, \dots, A_n . Then we get

$$\begin{aligned} \mathbb{P}\left(\bigcup_{k=1}^{\infty} A_k\right) &= \mathbb{P}\left(\bigcup_{k=1}^n A_k\right) = \sum_{\{j:\omega_j \in \bigcup_{k=1}^n A_k\}} p_j \\ &= \sum_{k=1}^n \sum_{\{j:\omega_j \in A_k\}} p_j = \sum_{k=1}^{\infty} \sum_{\{j:\omega_j \in A_k\}} p_j = \sum_{k=1}^{\infty} \mathbb{P}(A_k), \end{aligned}$$

hence \mathbb{P} is σ -additive.

By the construction, $\mathbb{P}(\{\omega_j\}) = p_j$, which completes the proof. \square

Summary: If $\Omega = \{\omega_1, \dots, \omega_N\}$, then probability measures \mathbb{P} on $\mathcal{P}(\Omega)$ can be identified with sequences $(p_j)_{j=1}^N$ satisfying the conditions (1.9).

$$\{\text{Probability measures } \mathbb{P} \text{ on } \mathcal{P}(\Omega)\} \iff \{\text{Sequences } (p_j)_{j=1}^N \text{ for which (1.9) hold}\}$$



Hereby the assignment from the left- to the right-hand side goes via $p_j = \mathbb{P}(\{\omega_j\})$ while in the other direction \mathbb{P} is given by eq. (1.11).

Example 1.3.3. Assume $\Omega = \{1, 2, 3\}$. Then each probability measure \mathbb{P} on $\mathcal{P}(\Omega)$ is uniquely determined by the three numbers $p_1 = \mathbb{P}(\{1\})$, $p_2 = \mathbb{P}(\{2\})$, and $p_3 = \mathbb{P}(\{3\})$. These numbers satisfy $p_1, p_2, p_3 \geq 0$ and $p_1 + p_2 + p_3 = 1$. Conversely, any three numbers p_1, p_2 , and p_3 with these properties generate a probability measure on $\mathcal{P}(\Omega)$ via (1.11). For example, if $A = \{1, 3\}$, then $\mathbb{P}(A) = p_1 + p_3$.

Next we treat *countably infinite* sample spaces, that is, $\Omega = \{\omega_1, \omega_2, \dots\}$. Also here we may take $\mathcal{P}(\Omega)$ as σ -field and, as in the case of finite sample spaces, given a probability measure \mathbb{P} on $\mathcal{P}(\Omega)$, we set

$$p_j := \mathbb{P}(\{\omega_j\}), \quad j = 1, 2, \dots$$

Then $(p_j)_{j=1}^\infty$ obeys the following properties:

$$p_j \geq 0 \quad \text{and} \quad \sum_{j=1}^{\infty} p_j = 1. \quad (1.12)$$

The proof is the same as in the finite case. The only difference is that here we have to use the σ -additivity of \mathbb{P} because this time $\Omega = \bigcup_{j=1}^{\infty} \{\omega_j\}$. By the same argument, it follows for $A \subseteq \Omega$ that

$$\mathbb{P}(A) = \sum_{\{j \geq 1: \omega_j \in A\}} p_j.$$

Hence, again the p_j s determine \mathbb{P} completely.

Conversely, let $(p_j)_{j=1}^\infty$ be an *arbitrary* sequence of real numbers with properties (1.12).

Proposition 1.3.4. *The mapping \mathbb{P} defined by*

$$\mathbb{P}(A) = \sum_{\{j \geq 1: \omega_j \in A\}} p_j \quad (1.13)$$

is a probability measure on $\mathcal{P}(\Omega)$ with $\mathbb{P}(\{\omega_j\}) = p_j$, $1 \leq j < \infty$.

Proof. The proof is analogous to that of Proposition 1.3.2 with one important exception. In the case $|\Omega| < \infty$, we used that there are at most finitely many disjoint nonempty subsets. This is no longer valid. Thus a different argument is needed.

Given disjoint subsets A_1, A_2, \dots , in Ω set

$$I_k = \{j \geq 1 : \omega_j \in A_k\}.$$

Then $I_k \cap I_l = \emptyset$ if $k \neq l$, thus,

$$\mathbb{P}(A_k) = \sum_{j \in I_k} p_j \quad \text{and} \quad \mathbb{P}\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{j \in I} p_j,$$

where $I = \bigcup_{k=1}^{\infty} I_k$. Use that $i \in I$ if and only if there is some $k \geq 1$ with $i \in I_k$ or, equivalently, with $\omega_i \in A_k$, that is, if and only if $\omega_i \in \bigcup_{k=1}^{\infty} A_k$.

Since $p_j \geq 0$, Remark A.5.6 applies and leads to

$$\mathbb{P}\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{j \in I} p_j = \sum_{k=1}^{\infty} \sum_{j \in I_k} p_j = \sum_{k=1}^{\infty} \mathbb{P}(A_k).$$

Thus \mathbb{P} is σ -additive.

The equality $\mathbb{P}(\{\omega_j\}) = p_j$, $1 \leq j < \infty$, is again a direct consequence of the definition of \mathbb{P} . \square

Summary: If $\Omega = \{\omega_1, \omega_2, \dots\}$, then probability measures \mathbb{P} on $\mathcal{P}(\Omega)$ can be identified with (infinite) sequences $(p_j)_{j=1}^{\infty}$ possessing the properties (1.12).



$$\{\text{Probability measures } \mathbb{P} \text{ on } \mathcal{P}(\Omega)\} \iff \{\text{Sequences } (p_j)_{j=1}^{\infty} \text{ satisfying (1.12)}\}$$

Again, the assignment from the left- to the right-hand side goes via $p_j = \mathbb{P}(\{\omega_j\})$ while the other direction rests upon eq. (1.13).

Example 1.3.5. For $\Omega = \mathbb{N}$ and $j \geq 1$, let $p_j = 2^{-j}$. These p_j s satisfy conditions (1.12) (check this!). The generated probability measure \mathbb{P} on $\mathcal{P}(\mathbb{N})$ is then given by

$$\mathbb{P}(A) := \sum_{j \in A} \frac{1}{2^j}.$$

For example, if $A = \{2, 4, 6, \dots\}$, then we get

$$\mathbb{P}(A) = \sum_{j \in A} \frac{1}{2^j} = \sum_{k=1}^{\infty} \frac{1}{2^{2k}} = \frac{1}{1 - 1/4} - 1 = \frac{1}{3}.$$

Or, if $B = \{N + 1, N + 2, \dots\}$ for a certain $N \in \mathbb{N}$, then one obtains

$$\mathbb{P}(B) = \sum_{j=N+1}^{\infty} \frac{1}{2^j} = 2^{-N-1} \sum_{j=0}^{\infty} \frac{1}{2^j} = 2^{-N-1} \cdot 2 = 2^{-N}.$$

Another way to evaluate the probability of B is

$$\mathbb{P}(B) = 1 - \mathbb{P}(B^c) = 1 - \sum_{j=1}^N \frac{1}{2^j} = 1 - \left[\frac{1 - 2^{-N-1}}{1 - \frac{1}{2}} - 1 \right] = 2^{-N}.$$

Example 1.3.6. Let $\Omega = \mathbb{Z} \setminus \{0\}$, that is, $\Omega = \{1, -1, 2, -2, \dots\}$. With $c > 0$ specified later on, assume

$$p_k = \frac{c}{k^2}, \quad k \in \Omega.$$

The number $c > 0$ has to be chosen so that the conditions (1.12) are satisfied, hence it has to fulfill

$$1 = c \sum_{k \in \mathbb{Z} \setminus \{0\}} \frac{1}{k^2} = 2c \sum_{k=1}^{\infty} \frac{1}{k^2}.$$

But, as is well known,⁶

$$\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6},$$

which implies $c = \frac{3}{\pi^2}$. Thus \mathbb{P} on $\mathcal{P}(\Omega)$ is uniquely described by

$$\mathbb{P}(\{k\}) = \frac{3}{\pi^2} \frac{1}{k^2}, \quad k \in \mathbb{Z} \setminus \{0\}.$$

For example, if $A = \mathbb{N}$, then

$$\mathbb{P}(A) = \frac{3}{\pi^2} \sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{3}{\pi^2} \frac{\pi^2}{6} = \frac{1}{2}.$$

Or if $A = \{2, 4, 6, \dots\}$, it follows that

$$\mathbb{P}(A) = \frac{3}{\pi^2} \sum_{k=1}^{\infty} \frac{1}{(2k)^2} = \frac{1}{4} \mathbb{P}(\mathbb{N}) = \frac{1}{8}.$$

For later purposes, we want to combine the two cases of finite and countably infinite sample spaces and thereby introduce a slight generalization.

⁶ We refer to [Mor16], where one can find an easy proof of this fact. The problem to compute the value of the sum is known as “Basel problem.” The first solution was found in 1734 by Leonhard Euler. Note that $\sum_{k \geq 1} 1/k^2 = \zeta(2)$ with Riemann’s ζ -function.

Let Ω be an arbitrary sample space. A probability measure \mathbb{P} is said to be **discrete** if there is an at most countably infinite set $D \subseteq \Omega$ (i. e., either D is finite or countably infinite) such that $\mathbb{P}(D) = 1$. Then for $A \subseteq \Omega$,

$$\mathbb{P}(A) = \mathbb{P}(A \cap D) = \sum_{\omega \in D} \mathbb{P}(\{\omega\}).$$

Since $\mathbb{P}(D^c) = 0$, this says that \mathbb{P} is **concentrated** on D . Of course, all previous results for a finite or countably infinite sample space carry over to this more general setting.



Discrete probability measures \mathbb{P} are concentrated on an at most countably infinite set D . They are uniquely determined by the values $\mathbb{P}(\{\omega\})$, where $\omega \in D$.

Of course, if the sample space is either finite or countably infinite, then *all* probability measures on this space are discrete. Nondiscrete probability measures will be introduced and investigated in Section 1.5.

Example 1.3.7. We once more model a single rolling of a die, but now we take as sample space $\Omega = \mathbb{R}$. Define $\mathbb{P}(\{\omega\}) = \frac{1}{6}$ if $\omega = 1, \dots, 6$ and $\mathbb{P}(\{\omega\}) = 0$ otherwise. If $D = \{1, \dots, 6\}$, then $\mathbb{P}(D) = 1$, hence \mathbb{P} is discrete. Given $A \subseteq \mathbb{R}$, it follows that

$$\mathbb{P}(A) = \frac{|A \cap D|}{6}.$$

For example, we have $\mathbb{P}([-2, 2]) = \frac{1}{3}$ and $\mathbb{P}([3, \infty)) = \frac{2}{3}$.

Another, maybe a little artificial, example is as follows.

Example 1.3.8. It is known that the set \mathbb{Q} of rational numbers is countably infinite. Hence $\mathbb{Q} = \{q_1, q_2, \dots\}$ with rational numbers q_k . Take any sequence $(p_k)_{k=1}^{\infty}$ of positive numbers with $\sum_{k=1}^{\infty} p_k = 1$. Then \mathbb{P} defined by

$$\mathbb{P}(A) = \sum_{\{k: q_k \in A\}} p_k, \quad A \subseteq \mathbb{R},$$

is a discrete probability measure on \mathbb{R} with $\mathbb{P}(\mathbb{Q}) = 1$. The problem with this probability measure is that it is completely impossible to evaluate $\mathbb{P}(B)$ for almost all $B \subseteq \mathbb{R}$, even if the p_k s are known.

1.4 Special discrete probability measures

1.4.1 Dirac measure

The simplest discrete probability measure is concentrated at a single point. That is, there exists an $\omega_0 \in \Omega$ such that $\mathbb{P}(\{\omega_0\}) = 1$. This probability measure is denoted by δ_{ω_0} .

Consequently, for each $A \in \mathcal{P}(\Omega)$, one has

$$\delta_{\omega_0}(A) = \begin{cases} 1 & \text{if } \omega_0 \in A, \\ 0 & \text{if } \omega_0 \notin A. \end{cases} \quad (1.14)$$

Definition 1.4.1. The probability measure δ_{ω_0} defined by eq. (1.14) is called the **Dirac measure** or **point measure** at ω_0 .

Which random experiment does $(\Omega, \mathcal{P}(\Omega), \delta_{\omega_0})$ model? It describes the experiment where, with probability one, the value ω_0 occurs. Thus, in fact it is a deterministic experiment, not random.

Dirac measures are useful tools to represent general discrete probability measures. Assume \mathbb{P} is concentrated on $D = \{\omega_1, \omega_2, \dots\}$ and let $p_j = \mathbb{P}(\{\omega_j\})$. Then we may write

$$\mathbb{P} = \sum_{j=1}^{\infty} p_j \delta_{\omega_j}. \quad (1.15)$$

Conversely, if a measure \mathbb{P} is represented as in eq. (1.15) with certain $\omega_j \in \Omega$ and numbers $p_j \geq 0$, $\sum_{j=1}^{\infty} p_j = 1$, then \mathbb{P} is discrete with $\mathbb{P}(D) = 1$, where $D = \{\omega_1, \omega_2, \dots\}$.

1.4.2 Uniform distribution on a finite set

The sample space is finite, say $\Omega = \{\omega_1, \dots, \omega_N\}$, and we assume that all elementary events are equally likely, that is,

$$\mathbb{P}(\{\omega_1\}) = \dots = \mathbb{P}(\{\omega_N\}).$$

A typical example is a fair die, where $\Omega = \{1, \dots, 6\}$.

Since $1 = \mathbb{P}(\Omega) = \sum_{j=1}^N \mathbb{P}(\{\omega_j\})$, we immediately get $\mathbb{P}(\{\omega_j\}) = 1/N$ for all $j \leq N$. If $A \subseteq \Omega$, an application of eq. (1.11) leads to

$$\mathbb{P}(A) = \frac{|A|}{N} = \frac{|A|}{|\Omega|}. \quad (1.16)$$

Definition 1.4.2. The probability measure \mathbb{P} defined by eq. (1.16) is called the **uniform distribution** or **Laplace distribution** on the finite set Ω .

The following formula may be helpful for remembrance. If \mathbb{P} is the uniform distribution on a finite sample space Ω , then

$$\mathbb{P}(A) = \frac{\text{Number of cases favorable for } A}{\text{Number of possible cases}}, \quad A \subseteq \Omega.$$



Example 1.4.3. In a lottery, 6 numbers are chosen out of 49 and each number appears only once. What is the probability that the chosen numbers are exactly the six marked on my lottery coupon?

Answer: Let us give two different approaches to answer this question.

Approach 1: We record the chosen numbers in the order they show up. As a sample space, we may take

$$\Omega := \{(\omega_1, \dots, \omega_6) : \omega_i \in \{1, \dots, 49\}, \omega_i \neq \omega_j \text{ if } i \neq j\}.$$

Then the number of possible cases is

$$|\Omega| = 49 \cdot 48 \cdot 47 \cdot 46 \cdot 45 \cdot 44 = \frac{49!}{43!}.$$

Let A be the event that the numbers on my lottery coupon appear. Which cardinality does A possess?

Say, for simplicity, that in our coupon the numbers 1, 2, ..., 6 are marked. Then it is favorable for A if these numbers appear in this order. But it is also favorable if (2, 1, 3, ..., 6) shows up, that is, any permutation of 1, ..., 6 is favorable. Hence $|A| = 6!$, which leads to⁷

$$\mathbb{P}(A) = \frac{6!}{49 \cdots 44} = \frac{1}{\binom{49}{6}} = 7.15112 \times 10^{-8}.$$

Approach 2: We assume that the chosen numbers are already ordered by their size (as they are published in a newspaper). In this case our sample space is

$$\Omega := \{(\omega_1, \dots, \omega_6) : 1 \leq \omega_1 < \cdots < \omega_6 \leq 49\},$$

and now

$$|\Omega| = \binom{49}{6}.$$

Why? Any set of six different numbers may be written exactly in one way in increasing order and thus choosing six ordered numbers is exactly the same as choosing a (nonordered) set of six numbers. And there are $\binom{49}{6}$ possibilities to choose six numbers. In this setting, we have $|A| = 1$, thus also here we get

$$\mathbb{P}(A) = \frac{1}{\binom{49}{6}}.$$

⁷ To get an impression about the size of this number, assume we buy lottery coupons with all possible choices of the six numbers. If each coupon is 0.5-mm thick, then all coupons together have a size of 6.992 km, which is about 4.3 miles. And in this row of 4.3 miles there exists exactly one coupon with the six numbers chosen in the lottery.

Example 1.4.4. A fair coin is labeled with “0” and “1”. Toss it n times and record the sequence of zeroes and ones in the order of their appearance. Thus,

$$\Omega := \{0, 1\}^n = \{(\omega_1, \dots, \omega_n) : \omega_i \in \{0, 1\}\},$$

and $|\Omega| = 2^n$. The coin is assumed to be fair, hence each sequence of zeroes and ones is equally likely. Therefore, whenever $A \subseteq \Omega$, one has

$$\mathbb{P}(A) = \frac{|A|}{2^n}.$$

Take, for example, the event A where for some fixed $i \leq n$ the i th toss equals “0”, that is,

$$A = \{(\omega_1, \dots, \omega_n) : \omega_i = 0\}.$$

Then $|A| = 2^{n-1}$ leads to the (not surprising) result

$$\mathbb{P}(A) = \frac{2^{n-1}}{2^n} = \frac{1}{2}.$$

Or let A occur if we observe for some given $k \leq n$ exactly k times the number “1”. Then $|A| = \binom{n}{k}$, and we get

$$\mathbb{P}(A) = \binom{n}{k} \cdot \frac{1}{2^n}.$$

Suppose $n \geq 2$. How likely is it that the first and the last toss coincide? There are 2^{n-2} possibilities that the first and the last toss are both 0 and also 2^{n-2} ways for both tosses to be 1. Hence, the probability of this event equals

$$\frac{2^{n-2} + 2^{n-2}}{2^n} = \frac{1}{2}.$$

Example 1.4.5. We have k particles that we distribute randomly into n boxes. All possible distributions of the particles are assumed to be equally likely. How do we get $\mathbb{P}(A)$ for a given event A ?

Answer: In this formulation, the question is not asked correctly because we did not fix when two distributions of particles coincide. Compare Figures 1.2 and 1.3 to understand why it is important whether or not the particles are anonymous or distinguishable.

Let us illustrate this problem in the case of two particles and two boxes. If the particles are *not distinguishable (anonymous)* then there are three different ways to distribute the particles into the two boxes. Thus, assuming that all distributions are equally likely, each elementary event has probability $1/3$.

On the other hand, if the particles are *distinguishable*, that is, they carry names, here 1 and 2, then there exist four different ways of distributing them, hence each elementary event has probability $1/4$.



Figure 1.2: Distributing two distinguishable particles into two boxes. Each event has probability $1/4$.

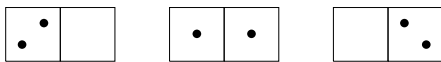


Figure 1.3: Distributing two anonymous particles into two boxes. Each event occurs with probability $1/3$.

Let us answer the above question in the two cases (distinguishable and anonymous) separately.

Distinguishable particles: Recall that we have k particles and n boxes. So we may enumerate the particles from 1 to k and each distribution of particles is uniquely described by a sequence (a_1, \dots, a_k) , where $a_j \in \{1, \dots, n\}$. For example, $a_1 = 3$ means that particle one is in box 3. Hence, a suitable sample space is

$$\Omega = \{(a_1, \dots, a_k) : 1 \leq a_i \leq n\}.$$

Since $|\Omega| = n^k$, for events $A \subseteq \Omega$, it follows that

$$\mathbb{P}(A) = \frac{|A|}{n^k}.$$

Anonymous particles: We record how many of the k particles are in box 1, how many are in box 2, and so on up to box n . Thus, as sample space we may choose

$$\Omega = \{(k_1, \dots, k_n) : 0 \leq k_j \leq k, k_1 + \dots + k_n = k\}.$$

The sequence (k_1, \dots, k_n) occurs if box 1 contains k_1 particles, box 2 contains k_2 , and so on. From the results in case 3 of Section A.3.2, we derive

$$|\Omega| = \binom{n+k-1}{k}.$$

Hence, if $A \subseteq \Omega$, then

$$\mathbb{P}(A) = \frac{|A|}{|\Omega|} = |A| \frac{k!(n-1)!}{(n+k-1)!}.$$

Summary: If we distribute k particles into n boxes and assume that all distributions are equally likely,⁸ then in the case of distinguishable or of anonymous particles for any set A of distributions,

$$\mathbb{P}(A) = \frac{|A|}{n^k} \quad \text{or} \quad \mathbb{P}(A) = |A| \frac{k!(n-1)!}{(n+k-1)!},$$

respectively.

Let us evaluate $\mathbb{P}(A)$ for some concrete event A in both cases. Suppose $k \leq n$ and select k of the n boxes. Set

$$A := \{\text{In each of the chosen } k \text{ boxes, there is exactly one particle}\}. \quad (1.17)$$

To simplify the notation, assume that the first k boxes have been chosen. The general case is treated in a similar way. Then in the “distinguishable case” the event A occurs if and only if for some permutation $\pi \in S_k$ the sequence $(\pi(1), \dots, \pi(k), 0, \dots, 0)$ appears. Thus $|A| = k!$ and

$$\mathbb{P}(A) = \frac{k!}{n^k}. \quad (1.18)$$

In the “anonymous case,” it follows that $|A| = 1$ (why?). Hence here we obtain

$$\mathbb{P}(A) = \frac{k!(n-1)!}{(n+k-1)!}. \quad (1.19)$$

Additional question: For $k \leq n$, define the event B by

$$B := \{\text{Each of the } n \text{ boxes contains at most 1 particle}\}.$$

Find $\mathbb{P}(B)$ in both cases.

Answer: The event B is the (disjoint) union of the following events: the k particles are distributed in a given collection of k boxes. The probability of this event was calculated in eqs. (1.18) and (1.19), respectively. Since there are $\binom{n}{k}$ possibilities to choose k boxes out of n , we get $\mathbb{P}(B) = \binom{n}{k} \mathbb{P}(A)$, with A as defined by (1.17), that is,

$$\begin{aligned} \mathbb{P}(B) &= \binom{n}{k} \cdot \frac{k!}{n^k} = \frac{n!}{(n-k)! n^k} \quad \text{and} \\ \mathbb{P}(B) &= \binom{n}{k} \cdot \frac{k!(n-1)!}{(n+k-1)!} = \frac{n!(n-1)!}{(n-k)!(n+k-1)!}, \end{aligned} \quad (1.20)$$

respectively.

⁸ Compare Example 1.4.20 and the following remark.

Example 1.4.6. Suppose we distribute 6 particles into 3 boxes such that all distributions are equally likely. What is the probability that each of the three boxes contains exactly two particles?

Answer: Let us first assume that the particles are distinguishable. If A is the event that each box contains 2 particles, then it follows that

$$|A| = \binom{6}{2, 2, 2} = \frac{6!}{2^3} = 90.$$

Consequently, we obtain

$$\mathbb{P}(A) = \frac{|A|}{3^6} = \frac{10}{81}.$$

Let us turn now to the case of indistinguishable particles. Here we have $|A| = 1$ (why?), hence in this case it follows that

$$\mathbb{P}(A) = \frac{6!(3-1)!}{(3+6-1)!} = \frac{2 \cdot 6!}{8!} = \frac{1}{28}.$$

Example 1.4.7. Let us check the validity of formula (1.20) by an easy test. In Example A.3.15, we evaluate the number of tiles of a domino in the following way. There are 7 boxes and we place two particles into these boxes. Hereby, the particles are anonymous and all distributions are equally likely. Then the event B defined by “at most one particle in each box” corresponds to tiles with different numbers of dots. Now choose one of the 28 tiles at random. Since there are 21 tiles with different numbers of dots, it follows that

$$\mathbb{P}(B) = \frac{21}{28} = \frac{3}{4}.$$

On the other hand, formula (1.20) leads to (recall that $n = 7$ and $k = 2$)

$$\mathbb{P}(B) = \frac{n!(n-1)!}{(n-k)!(n+k-1)!} = \frac{7! \cdot 6!}{5! \cdot 8!} = \frac{6}{8} = \frac{3}{4}.$$

So we see, in this case formula (1.20) gives the correct value.

Example 1.4.8. Suppose the sample space Ω is given by

$$\Omega = \{(k_1, \dots, k_6) : k_1 + \dots + k_6 = 4, k_j \in \mathbb{N}_0\}$$

and endowed with the uniform distribution. That is, all possible representations of the number 4 by six nonnegative integers are equally likely. So, for example, the occurrence of $(1, 2, 0, 0, 1, 0)$ is as likely as that of $(0, 0, 0, 0, 0, 4)$ or that of $(4, 0, 0, 0, 0, 0)$.

Questions:

(1) What is the cardinality of Ω ?

(2) How likely is it to observe those $(k_1, \dots, k_6) \in \Omega$ for which $0 \leq k_j \leq 1$?

Answers: An equivalent model is as follows: one has six boxes and places four anonymous particles into these boxes. In this way, k_j describes the number of particles in box j with $1 \leq j \leq 6$.

In the setting of case 3 of Section A.3.2, we have $k = 4$ and $n = 6$. Hence, the cardinality of the sample space equals

$$|\Omega| = \binom{n+k-1}{n-1} = \binom{9}{5} = 126.$$

To answer the second question, note that all numbers k_j satisfy $k_j \leq 1$ if and only if each of the six boxes contains at most one particle. According to eq. (1.20), the probability of this event is given by

$$\frac{n!(n-1)!}{(n-k)!(n+k-1)!} = \frac{6! \cdot 5!}{2! \cdot 9!} = \frac{5}{42} \approx 0.119048.$$

Another (more direct) way to obtain the last result is as follows: there are $\binom{6}{4}$ ways to write 4 as a sum of six numbers which are either 0 or 1. Hence, the desired probability equals

$$\frac{\binom{6}{4}}{|\Omega|} = \frac{15}{126} = \frac{5}{42}.$$

Summary: The uniform distribution \mathbb{P} on a sample space $\Omega = \{\omega_1, \dots, \omega_N\}$ is characterized by

$$\mathbb{P}(\{\omega_1\}) = \dots = \mathbb{P}(\{\omega_N\}) = \frac{1}{N}, \quad \text{or, equivalently, by}$$

$$\mathbb{P}(A) = \frac{|A|}{N} = \frac{\text{Number of cases favorable for } A}{\text{Number of possible cases}}, \quad A \subseteq \Omega.$$

1.4.3 Binomial distribution

The sample space is $\Omega = \{0, 1, \dots, n\}$ for some $n \geq 1$ and p is a real number with $0 \leq p \leq 1$.

Proposition 1.4.9. *There exists a unique probability measure $B_{n,p}$ on $\mathcal{P}(\Omega)$ satisfying*

$$B_{n,p}(\{k\}) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, \dots, n. \quad (1.21)$$

Proof. In order to use Proposition 1.3.2, we have to verify $B_{n,p}(\{k\}) \geq 0$ and $\sum_{k=0}^n B_{n,p}(\{k\}) = 1$. The first property is obvious because of $0 \leq p \leq 1$ and $0 \leq 1-p \leq 1$. To prove the second, we apply the binomial theorem (Proposition A.3.8) with $a = p$ and with $b = 1-p$. This leads to

$$\sum_{k=0}^n B_{n,p}(\{k\}) = \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} = (p + (1-p))^n = 1.$$

Hence the assertion follows by Proposition 1.3.2 with $p_k = B_{n,p}(\{k\})$, $k = 0, \dots, n$. \square

Compare Figure 1.4 for the values of $B_{n,p}(\{k\})$ in the case $n = 9$ and with, $p = 1/2$ and $p = 1/4$, respectively.

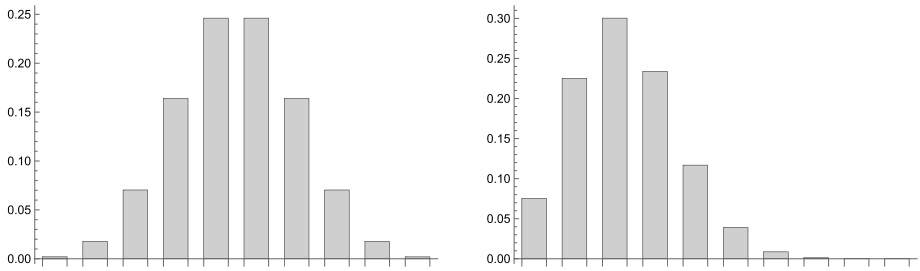


Figure 1.4: Probabilities $B_{n,p}(\{k\})$ with $n = 9$, $p = 1/2$, and $p = 1/4$, $k = 0, \dots, 9$.

Definition 1.4.10. The probability measure $B_{n,p}$ defined by eq. (1.21) is called the **binomial distribution** with parameters n and p .

Remark 1.4.11. Observe that $B_{n,p}$ acts as follows. If $A \subseteq \{0, \dots, n\}$, then

$$B_{n,p}(A) = \sum_{k \in A} \binom{n}{k} p^k (1-p)^{n-k}.$$

Furthermore, for $p = 1/2$, we get

$$B_{n,1/2}(\{k\}) = \binom{n}{k} \frac{1}{2^n}.$$

As we saw in Example 1.4.4, this probability describes the k -fold occurrence of “1” when tossing a fair coin n times.

Which random experiment describes the binomial distribution? To answer this question, let us first look at the case $n = 1$. Here we have $\Omega = \{0, 1\}$ with

$$B_{n,p}(\{0\}) = 1 - p \quad \text{and} \quad B_{n,p}(\{1\}) = p.$$

If we identify “0” with *failure* and “1” with *success*, then the binomial distribution describes an experiment where either success or failure may occur, and the success probability is p . Now we execute the same experiment n times and every time we may observe either failure or success. If we have k times success, then there are $\binom{n}{k}$ ways to obtain

these k successes during the n trials. The probability for success is p and for failure $1-p$. By the independence of the single trials, the probability for the sequence is $p^k(1-p)^{n-k}$. Multiplying this probability with the number of different positions of successes, we finally arrive at $\binom{n}{k}p^k(1-p)^{n-k}$, the value of $B_{n,p}(\{k\})$.

Example 1.4.12. An exam consists of 100 problems where each of the question may be answered either with “yes” or “no.” To pass the exam, at least 60 questions have to be answered correctly. Let p be the probability to answer a single question correctly. How big does p have to be in order to pass the exam with a probability greater than 75%?

Answer: The number p has to be chosen such that the following estimate is satisfied:

$$\sum_{k=60}^{100} \binom{100}{k} p^k (1-p)^{100-k} \geq 0.75.$$

Numerical calculations show that this is valid if and only if $p \geq 0.62739$.

Example 1.4.13. In an auditorium there are N students. Find the probability that at least two of them have their birthday on April 1.

Answer: We do not take leap years into account and assume that there are no twins among the students. Finally, we make the (probably unrealistic) assumption that all days of a year are equally likely as birthdays. Say success occurs if a student has birthday on April 1. Under the above assumptions, the success probability is $1/365$. Hence the number of students having birthday on April 1 is binomially distributed with parameters N and $p = 1/365$. We ask for the probability of $A = \{2, 3, \dots, N\}$. This may be evaluated by

$$\begin{aligned} \sum_{k=2}^N \binom{N}{k} \left(\frac{1}{365}\right)^k \left(\frac{364}{365}\right)^{N-k} &= 1 - B_{N,1/365}(\{0\}) - B_{N,1/365}(\{1\}) \\ &= 1 - \left(\frac{364}{365}\right)^N - \frac{N}{365} \left(\frac{364}{365}\right)^{N-1}. \end{aligned}$$

For example, $N = 500$ this probability is approximately 0.397895.

Example 1.4.14. In an urn there are 40 white balls and 60 black ones. One chooses balls one after another from the urn with replacement. How often does one have to choose balls in order to observe with probability greater than 0.5 at least 10 white balls?

Answer: The success probability is $p = 40/100 = 2/5$. So we ask for the minimal number $n \geq 10$ for which

$$B_{n,p}(\{10, 11, \dots, n\}) \geq 0.5 \quad \Leftrightarrow \quad \sum_{k=10}^n \binom{n}{k} \left(\frac{2}{5}\right)^k \left(\frac{3}{5}\right)^{n-k} \geq 0.5.$$

Numerical calculations give the following probabilities:

n	$B_{n,p}(\{10, 11, \dots, n\})$
20	0.244663
21	0.308558
22	0.375648
23	0.443771
24	0.510920

So we see that $n = 24$ is the minimal number of trials to obtain at least 10 white balls with probability greater than or equal to $1/2$. Of course, if we increase the number of trials then it becomes more and more likely to get at least 10 white balls. For example, if one takes out 40 balls, then the probability for at least 10 white balls is about 0.9845.

Summary: The binomial distribution describes the following setting. One executes n times independently the same experiment where each time either success or failure may appear. The success probability is p . Then $B_{n,p}(\{k\})$ is the probability to observe exactly k successes or, equivalently, $n - k$ failures.

1.4.4 Multinomial distribution

Given natural numbers n and m , the sample space for the multinomial distribution is⁹

$$\Omega := \{(k_1, \dots, k_m) \in \mathbb{N}_0^m : k_1 + \dots + k_m = n\}.$$

With certain nonnegative real numbers p_1, \dots, p_m satisfying $p_1 + \dots + p_m = 1$, set

$$\mathbb{P}(\{(k_1, \dots, k_m)\}) := \binom{n}{k_1, \dots, k_m} p_1^{k_1} \dots p_m^{k_m}, \quad (k_1, \dots, k_m) \in \Omega. \quad (1.22)$$

Recall that the multinomial coefficients appearing in eq. (1.22) were defined in eq. (A.16) as

$$\binom{n}{k_1, \dots, k_m} = \frac{n!}{k_1! \dots k_m!}.$$

The next result shows that eq. (1.22) defines a probability measure.

Proposition 1.4.15. *There is a unique probability measure \mathbb{P} on $\mathcal{P}(\Omega)$ such that (1.22) holds for all $(k_1, \dots, k_m) \in \Omega$.*

⁹ By case 3 in Section A.3.2, the cardinality of Ω is $\binom{n+m-1}{n}$.

Proof. An application of the multinomial theorem (Proposition A.3.20) implies

$$\begin{aligned} \sum_{(k_1, \dots, k_m) \in \Omega} \mathbb{P}(\{(k_1, \dots, k_m)\}) &= \sum_{\substack{k_1 + \dots + k_m = n \\ k_i \geq 0}} \binom{n}{k_1, \dots, k_m} p_1^{k_1} \cdots p_m^{k_m} \\ &= (p_1 + \cdots + p_m)^n = 1^n = 1. \end{aligned}$$

Since $\mathbb{P}(\{(k_1, \dots, k_m)\}) \geq 0$, the assertion follows by Proposition 1.3.2. \square

In view of the preceding proposition, the following definition is justified.

Definition 1.4.16. The probability measure \mathbb{P} defined by eq. (1.22) is called the **multinomial distribution** with parameters n , m , and p_1, \dots, p_m .

Remark 1.4.17. Sometimes it is useful to regard the multinomial distribution on the larger sample space $\Omega = \mathbb{N}_0^m$. In this case we have to modify eq. (1.22) slightly as follows:

$$\mathbb{P}(\{(k_1, \dots, k_m)\}) = \begin{cases} \binom{n}{k_1, \dots, k_m} p_1^{k_1} \cdots p_m^{k_m} & \text{if } k_1 + \cdots + k_m = n, \\ 0 & \text{if } k_1 + \cdots + k_m \neq n. \end{cases}$$

Which random experiment does the multinomial distribution describe? To answer this question, let us recall the model for the binomial distribution. In an urn there are balls of two different colors, say white and red. The proportion of the white balls is p , hence $1 - p$ is the proportion of the red ones. If we choose n balls with replacement, then $B_{n,p}(\{k\})$ is the probability to observe exactly k white balls.

What happens if in the urn there are balls of more than two different colors, say of m , and the proportions of the colored balls are p_1, \dots, p_m with $p_1 + \cdots + p_m = 1$?

As in the model for the binomial distribution, we choose n balls with replacement. Given integers $k_j \geq 0$, one asks now for the probability of the following event: balls of the first color showed up k_1 times, those of the second k_2 times, and so on. Of course, this probability is zero whenever $k_1 + \cdots + k_m \neq n$. But if the sum is n , then $p_1^{k_1} \cdots p_m^{k_m}$ is the probability for k_j balls of color j in some fixed order. There are $\binom{n}{k_1, \dots, k_m}$ ways to order the balls without changing the frequency of the colors. Thus the desired probability equals $\binom{n}{k_1, \dots, k_m} p_1^{k_1} \cdots p_m^{k_m}$.

Remark 1.4.18. If $m = 2$, then $p_2 = 1 - p_1$, as well as $\binom{n}{k_1, k_2} = \binom{n}{k_1, n-k_1} = \binom{n}{k_1}$. Consequently, in this case the multinomial distribution coincides with the binomial distribution B_{n,p_1} .

Example 1.4.19. Suppose in an urn there are 3 white, 5 red, and 4 black balls. Choose 9 balls with replacement. How likely is it to observe three balls of each color.

Answer: The success probabilities are $p_1 = 3/12$, $p_2 = 5/12$, and $p_3 = 4/12$. Hence, the desired probability equals

$$\binom{9}{3, 3, 3} \left(\frac{3}{12}\right)^3 \left(\frac{5}{12}\right)^3 \left(\frac{4}{12}\right)^3 = \frac{9! \cdot 60^3}{6^3 \cdot 12^9} \approx 0.0703286.$$

Additional question: How likely is it to observe no white ball among the 9 chosen ones?

Answer: An approach via the multinomial distributions would be possible. Then one has to sum over all possible choices of red and black balls. But the problem can be handled in a more direct way. Say success occurs if the chosen ball is white and failure if this is not so. Then the success probability is $3/12 = 1/4$. Thus, we ask for a total of 9 failures which has the probability $(3/4)^9 \approx 0.0750847$.

Example 1.4.20. Suppose we have m boxes B_1, \dots, B_m and n particles that we place successively into these boxes. Thereby p_j is the probability to place a single particle into box B_j . What is the probability that after distributing all n particles there are k_1 particles in the first box, k_2 in the second, and all the way up to k_m in the last one?

Answer: This probability is given by formula (1.22), that is,

$$\mathbb{P}\{k_1 \text{ particles are in } B_1, \dots, k_m \text{ particles are in } B_m\} = \binom{n}{k_1, \dots, k_m} p_1^{k_1} \cdots p_m^{k_m}.$$

For example, if we place 5 particles into 3 boxes with probabilities $1/2$, $1/3$, and $1/6$, respectively, then

$$\mathbb{P}\{k_1 \text{ are in } B_1, k_2 \text{ are in } B_2, k_3 \text{ are in } B_3\} = \binom{5}{k_1, k_2, k_3} \left(\frac{1}{2}\right)^{k_1} \left(\frac{1}{3}\right)^{k_2} \left(\frac{1}{6}\right)^{k_3}$$

provided that $k_1 + k_2 + k_3 = 5$.

Remark 1.4.21. In the case that all m possible different outcomes of an experiment are equally likely, that is, we have

$$p_1 = \cdots = p_m = \frac{1}{m},$$

then

$$\mathbb{P}(\{(k_1, \dots, k_m)\}) = \binom{n}{k_1, \dots, k_m} \frac{1}{m^n}, \quad k_1 + \cdots + k_m = n. \quad (1.23)$$

Example 1.4.22. Roll a fair die 12 times. How likely is it that each of the six possible results appears exactly twice?

Answer: An equivalent formulation of the problem is as follows: there are 6 boxes and 12 particles. Place these 12 particles randomly into the 6 boxes so that all boxes are equally likely. Let the event A occur if in each of the six boxes there are two particles or, equivalently, if each of the numbers from 1 to 6 appears twice. Then we get

$$\mathbb{P}(A) = \binom{12}{2, 2, 2, 2, 2, 2} \left(\frac{1}{6}\right)^{12} = \frac{12!}{2^6 6^{12}} \approx 0.00343829.$$

Example 1.4.23. Suppose that in Example 1.4.20 all m boxes are chosen with probability $1/m$. Then, if $n \leq m$, one may ask for the probability that each of the first n boxes

B_1, \dots, B_n contains exactly one particle. By eq. (1.23), it follows that

$$\mathbb{P}(\{\underbrace{(1, \dots, 1)}_n, 0, \dots, 0\}) = \binom{n}{\underbrace{1, \dots, 1}_n, 0, \dots, 0} \frac{1}{m^n} = \frac{n!}{m^n}. \quad (1.24)$$

Remark 1.4.24. From a different point of view, we investigated the last problem already in Example 1.4.5. But why do we get in eq. (1.24) the same answer as in the case of distinguishable particles although the n distributed ones are anonymous?

Answer: The crucial point is that we assumed in the anonymous case that all partitions of the particles are equally likely. And this is not valid when distributing the particles successively. To see this, assume $n = m = 2$. Then there exist three different ways to distribute the particles, but they have different probabilities:

$$\mathbb{P}(\{(0, 2)\}) = \mathbb{P}(\{(2, 0)\}) = \frac{1}{4} \quad \text{while} \quad \mathbb{P}(\{(1, 1)\}) = \frac{1}{2}.$$

Thus, although the distributed particles are not distinguishable, they get names due to the successive distribution (first particle, second particle, etc.).

Example 1.4.25. Six people randomly enter a train with three coaches. Each person chooses his wagon independently of the others and all coaches are equally likely to be chosen. Find the probability that there are two people in each coach.

Answer: We have $m = 3$, $n = 6$, and $p_1 = p_2 = p_3 = \frac{1}{3}$. Hence the probability we are looking for is

$$\mathbb{P}(\{(2, 2, 2)\}) = \binom{6}{2, 2, 2} \frac{1}{3^6} = \frac{6!}{2! 2! 2!} \frac{1}{3^6} = \frac{10}{81} = 0.12345679.$$

Check how this result is related to that presented in Example 1.4.6.

Example 1.4.26. In a country 40 % of the cars are gray, 20 % are black, and 10 % are red. The remaining cars have different colors. Now we observe at random 10 cars. What is the probability to see two gray cars, four black, and one red?

Answer: By assumption $m = 4$ (gray, black, red, and others), $p_1 = 2/5$, $p_2 = 1/5$, $p_3 = 1/10$, and $p_4 = 3/10$. Thus the probability of the vector $(2, 4, 1, 3)$ is given by

$$\begin{aligned} \binom{10}{2, 4, 1, 3} \left(\frac{2}{5}\right)^2 \left(\frac{1}{5}\right)^4 \left(\frac{1}{10}\right)^1 \left(\frac{3}{10}\right)^3 &= \frac{10!}{2! 4! 1! 3!} \cdot \frac{2^2}{5^2} \cdot \frac{1}{5^4} \cdot \frac{1}{10} \cdot \frac{3^3}{10^3} \\ &= 0.00870912. \end{aligned}$$

Summary: The multinomial distribution is a generalization of the binomial distribution from two (failure or success) to $m \geq 2$ possible results. In each single experiment, m different results may occur (e. g., m different colors of balls) and each time the j th result, $1 \leq j \leq m$, shows up with probability p_j . If one executes the experiment n times, then the multinomial distribution describes the probability of the following event: the

first result occurs k_1 times, the second k_2 times, and so on. Here k_1, \dots, k_m are some nonnegative integers with $k_1 + \dots + k_m = n$.

1.4.5 Poisson distribution

The sample space for this distribution is $\mathbb{N}_0 = \{0, 1, 2, \dots\}$. Furthermore, $\lambda > 0$ is a given parameter.

Proposition 1.4.27. *There exists a unique probability measure Pois_λ on $\mathcal{P}(\mathbb{N}_0)$ such that*

$$\text{Pois}_\lambda(\{k\}) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k \in \mathbb{N}_0. \quad (1.25)$$

Proof. Because of $e^{-\lambda} > 0$, $\text{Pois}_\lambda(\{k\}) > 0$ follows. Thus it suffices to verify

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = 1.$$

But this is a direct consequence of

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{-\lambda} e^\lambda = 1. \quad \square$$

Definition 1.4.28. The probability measure Pois_λ on $\mathcal{P}(\mathbb{N}_0)$ satisfying Eq. (1.25) is called the **Poisson distribution** with parameter $\lambda > 0$.

Compare Figure 1.5 for the values of the Poisson distribution in the case $\lambda = 5$. Most likely are the events $k = 4$ and $k = 5$.

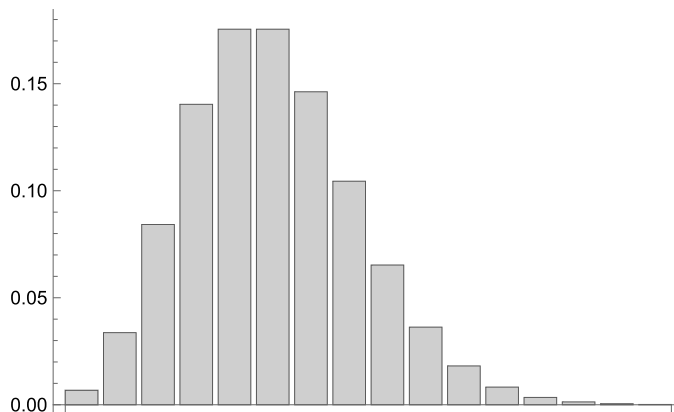


Figure 1.5: The values $\text{Pois}_5(\{k\})$, $k = 0, \dots, 15$.

The Poisson distribution describes experiments where the number of trials is big, but the single success probability is small. More precisely, the following limit theorem holds.

Proposition 1.4.29 (Poisson's limit theorem). *Let $(p_n)_{n=1}^{\infty}$ be a sequence of numbers with $0 < p_n \leq 1$ and*

$$\lim_{n \rightarrow \infty} n p_n = \lambda$$

for some $\lambda > 0$. Then for all $k \in \mathbb{N}_0$,

$$\lim_{n \rightarrow \infty} B_{n,p_n}(\{k\}) = \text{Pois}_{\lambda}(\{k\}).$$

Proof. Write

$$\begin{aligned} B_{n,p_n}(\{k\}) &= \binom{n}{k} p_n^k (1-p_n)^{n-k} \\ &= \frac{1}{k!} \frac{n(n-1) \cdots (n-k+1)}{n^k} (n p_n)^k (1-p_n)^n (1-p_n)^{-k} \\ &= \frac{1}{k!} \left[\frac{n}{n} \cdot \frac{n-1}{n} \cdots \frac{n-k+1}{n} \right] (n p_n)^k (1-p_n)^n (1-p_n)^{-k}, \end{aligned}$$

and investigate the behavior of the different parts of the last equation separately. Each fraction in the left-hand brackets tends to 1, hence the whole factor in brackets tends to 1. By assumption, we have $n p_n \rightarrow \lambda$, thus, $\lim_{n \rightarrow \infty} (n p_n)^k = \lambda^k$. Moreover, because of $n p_n \rightarrow \lambda$ with $\lambda > 0$, we get $p_n \rightarrow 0$, which implies $\lim_{n \rightarrow \infty} (1-p_n)^{-k} = 1$.

Thus, it remains to determine the behavior of $(1-p_n)^n$ as $n \rightarrow \infty$. Proposition A.5.1 asserts that if a sequence of real numbers $(x_n)_{n \geq 1}$ converges to $x \in \mathbb{R}$, then

$$\lim_{n \rightarrow \infty} \left(1 + \frac{x_n}{n} \right)^n = e^x.$$

Setting $x_n := -n p_n$, by assumption $x_n \rightarrow -\lambda$, hence

$$\lim_{n \rightarrow \infty} (1-p_n)^n = \lim_{n \rightarrow \infty} \left(1 + \frac{x_n}{n} \right)^n = e^{-\lambda}.$$

If we combine all the different parts, then this completes the proof due to

$$\lim_{n \rightarrow \infty} B_{n,p_n}(\{k\}) = \frac{1}{k!} \lambda^k e^{-\lambda} = \text{Pois}_{\lambda}(\{k\}). \quad \square$$

The previous theorem allows two *conclusions*:

(1) Whenever n is large and p is small, without hesitation one may replace $B_{n,p}$ by Pois_{λ} , where $\lambda = n \cdot p$. In this way, one avoids the (sometimes) difficult evaluation of the binomial coefficients.

Example 1.4.30. In Example 1.4.13, we found the probability that among N students there are at least two having their birthday on April 1. We then used the binomial distribution with parameters N and $p = 1/365$. Hence the approximating Poisson distribution has parameter $\lambda = N/365$ and the corresponding probability is given by

$$\text{Pois}_\lambda(\{2, 3, \dots\}) = 1 - (1 + \lambda)e^{-\lambda} = 1 - \left(1 + \frac{N}{365}\right)e^{-N/365}.$$

If again $N = 500$, hence $\lambda = 500/365$, the approximate probability equals 0.397719. Compare this value with the “precise” probability 0.397895 obtained in Example 1.4.13.

(2) Poisson’s limit theorem explains why the Poisson distribution describes experiments with many trials and small success probability. For example, if we look for a model for the number of car accidents per year, then the Poisson distribution is a good choice. There are many cars, but the probability¹⁰ that a single driver is involved in an accident is quite small.

Similarly, the Poisson distribution is used to model the number of customers entering some shop, to describe the number of phone calls arriving at a call center, or the number of daily accesses to a website. This is due to the fact that there are many potential customers but with small probability a single one enters the shop. In the same way, there are many people possessing a phone, but the probability that a single one calls a certain center is very small.

Later on we will investigate other examples where the Poisson distribution appears in a natural way.

Summary: The Poisson distribution Pois_λ occurs as the limit of the binomial distribution $B_{n,p}$ in the following sense: the number n of trials tends to infinity while at the same time the success probability p becomes smaller and smaller. In other words, for large n and small success probability p , it follows that $B_{n,p}(A) \approx \text{Pois}_\lambda(A)$ with $\lambda = np$.

1.4.6 Hypergeometric distribution

Among N delivered machines M are defective. One chooses n of the N machines randomly and checks them. What is the probability to observe m defective machines in the sample of size n ?

First note that there are $\binom{N}{n}$ ways to choose n machines for checking. In order to observe m defective ones, these have to be taken from the M defective. The remaining $n-m$ machines are nondefective, hence they must be chosen from the $N-M$ nondefective ones. There are $\binom{M}{m}$ ways to take the defective machines and $\binom{N-M}{n-m}$ possibilities for the nondefective ones.

¹⁰ To call it a “success” probability in this case is perhaps not quite appropriate.

Thus the following approach describes this experiment:

$$H_{N,M,n}(\{m\}) := \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad 0 \leq m \leq n. \quad (1.26)$$

Recall that in Section A.3.1 we agreed that $\binom{n}{k} = 0$ whenever $k > n$. This turns out be useful in the definition of $H_{N,M,n}$. For example, if $m > M$, then the probability to observe m defective machines is, of course, zero.

We want to prove now that eq. (1.26) defines a probability measure.

Proposition 1.4.31. *There exists a unique probability measure $H_{N,M,n}$ on the powerset of $\{0, \dots, n\}$ satisfying eq. (1.26).*

Proof. Vandermonde's identity (cf. Proposition A.3.9) asserts that for all k, m , and n in \mathbb{N}_0 ,

$$\sum_{j=0}^k \binom{n}{j} \binom{m}{k-j} = \binom{n+m}{k}. \quad (1.27)$$

Now replace n by M , next m by $N - M$, then k by n , and, finally, j by m . Doing so, eq. (1.27) leads to

$$\sum_{m=0}^n \binom{M}{m} \binom{N-M}{n-m} = \binom{N}{n}.$$

But this implies

$$\sum_{m=0}^n H_{N,M,n}(\{m\}) = \frac{1}{\binom{N}{n}} \cdot \sum_{m=0}^n \binom{M}{m} \binom{N-M}{n-m} = \frac{1}{\binom{N}{n}} \cdot \binom{N}{n} = 1.$$

Clearly, $H_{N,M,n}(\{m\}) \geq 0$, which completes the proof by virtue of Proposition 1.3.2. \square

Definition 1.4.32. The probability measure $H_{N,M,n}$ defined by eq. (1.26) is called the **hypergeometric distribution** with parameters N, M , and n .

Example 1.4.33. A retailer gets a delivery of 100 machines; 10 of them are defective. He chooses 8 machines at random and tests them. Find the probability that 2 or more of the tested machines are defective.

Answer: The desired probability is

$$\sum_{m=2}^8 \frac{\binom{10}{m} \binom{90}{8-m}}{\binom{100}{8}} = 0.18195.$$

See Figure 1.6 for another example of a hypergeometric distribution. There are 40 defective machines among 100 and one checks a sample of size 15.

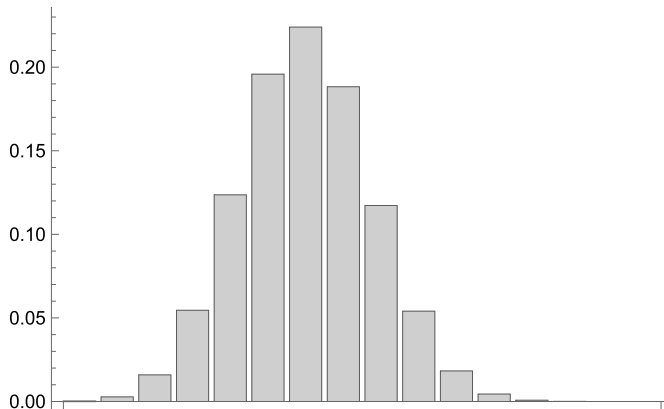


Figure 1.6: Probabilities $H_{100,40,15}(\{m\})$ with $m = 0, \dots, 15$. The maximal value is about 0.224 attained at $m = 6$. That is, if there are 40 defective items among 100, most likely in a sample of 15 there are 6 defective ones.

Remark 1.4.34. In the daily practice, the opposite question is more important. The size N of the delivery is known and, of course, also the size of the tested sample. The number M of defective machines is unknown. Now suppose we observed m defective machines among the n tested. Does this (random) number m lead to some information about the number M of defective machines in the delivery? We will investigate this problem in Proposition 8.5.16.

Example 1.4.35. In a pond there are 200 fish. One day the owner of the pond catches 20 fish, marks them, and puts them back into the pond. After a while the owner catches once more 20 fish. Find the probability that among these fish there is exactly one marked.

Answer: We have $N = 200$, $M = 20$, and $n = 20$. Hence the desired probability is

$$H_{200,20,20}(\{1\}) = \frac{\binom{20}{1}\binom{180}{19}}{\binom{200}{20}} = 0.26967.$$

Remark 1.4.36. The previous example is not very realistic because in general the number N of fish is unknown. Known are M and n , the (random) number m was observed. Also here one may ask whether the knowledge of m leads to some information about N . This question will be investigated later in Proposition 8.5.18.

Example 1.4.37. In a lottery, 6 numbers are chosen randomly out of 49. Suppose we bought a lottery coupon with six numbers. What is the probability that exactly k , $k = 0, \dots, 6$, of our numbers appear in the drawing?

Answer: There are $n = 6$ numbers randomly chosen out of $N = 49$. Among the 49 numbers, $M = 6$ are “defective.” These are the six numbers on our coupon, and we ask

for the probability that k of the “defective” are among the chosen six. The question is answered by the hypergeometric distribution $H_{49,6,6}$, that is, the probability of k correct numbers on our coupon is given by

$$H_{49,6,6}(\{k\}) = \frac{\binom{6}{k} \binom{43}{6-k}}{\binom{49}{6}}, \quad k = 0, \dots, 6.$$

The numerical values of these probabilities for $k = 0, \dots, 6$ are

k	$H_{49,6,6}(\{k\})$
0	0.435965
1	0.413019
2	0.132378
3	0.0176504
4	0.00096862
5	0.0000184499
6	$7.15112 \cdot 10^{-8}$

Remark 1.4.38. Another model for the hypergeometric distribution is as follows: in an urn there are N balls, M of them are white, the remaining $N - M$ are red. Choose n balls out of the urn *without replacing* the chosen ones. Then $H_{N,M,n}(\{m\})$ is the probability to observe m white balls among the n chosen.

If we do the same experiment, but now *replacing* the chosen balls, then this is described by the binomial distribution. The success probability for a white ball is $p = M/N$, hence now the probability for m white balls is given by

$$B_{n,M/N}(\{m\}) = \binom{n}{m} \left(\frac{M}{N}\right)^m \left(1 - \frac{M}{N}\right)^{n-m}.$$

It is intuitively clear that for large N and M (and comparably small n) the difference between both models (replacing and nonreplacing) is insignificant. Imagine there are 10^6 white and also 10^6 red balls in an urn. When choosing two balls, it does not matter a lot whether the first ball was replaced or not.

The next proposition makes the previous observation more precise.

Proposition 1.4.39. *If $0 \leq m \leq n$ and $0 \leq p \leq 1$, then*

$$\lim_{\substack{N, M \rightarrow \infty \\ M/N \rightarrow p}} H_{N,M,n}(\{m\}) = B_{n,p}(\{m\}).$$

Proof. Suppose first $0 < p < 1$. Then the definition of the hypergeometric distribution yields

$$\begin{aligned}
 \lim_{\substack{N, M \rightarrow \infty \\ M/N \rightarrow p}} H_{N, M, n}(\{m\}) &= \lim_{\substack{N, M \rightarrow \infty \\ M/N \rightarrow p}} \frac{\frac{M \cdots (M-m+1)}{m!} \frac{(N-M) \cdots (N-M-(n-m)+1)}{(n-m)!}}{\frac{N(N-1) \cdots (N-n+1)}{n!}} \\
 &= \lim_{\substack{N, M \rightarrow \infty \\ M/N \rightarrow p}} \binom{n}{m} \frac{[\frac{M}{N} \cdots (\frac{M-m+1}{N})][(\frac{1}{N}) \cdots (\frac{1}{N})]}{(1 - \frac{1}{N}) \cdots (1 - \frac{n+1}{N})} \\
 &= \binom{n}{m} p^m (1-p)^{n-m} = B_{n,p}(\{m\}).
 \end{aligned} \tag{1.28}$$

Note that if either $m = 0$ or $m = n$, then the first or the second brackets in the second line of eq. (1.28) become 1, thus they do not appear.

The cases $p = 0$ and $p = 1$ have to be treated separately. For example, if $p = 0$, the fraction in the second line of eq. (1.28) converges to zero provided that $m \geq 1$. If $m = 0$, then

$$\lim_{\substack{N, M \rightarrow \infty \\ M/N \rightarrow 0}} \frac{(1 - \frac{M}{N}) \cdots (1 - \frac{M-n+1}{N})}{(1 - \frac{1}{N}) \cdots (1 - \frac{n+1}{N})} = 1 = B_{n,0}(\{0\}).$$

The case $p = 1$ is treated similarly. Hence, the proposition is also valid in the border cases. □

Example 1.4.40. Suppose there are $N = 200$ balls in an urn. In the first table, there are $M = 80$ white balls. Choosing $n = 10$ balls with or without replacement, we get the following numerical values. Note that $p = M/N = 2/5$. In the second table, we execute the same experiment, but now there are 100 white balls among 200.

m	$H_{200,80,10}(\{m\})$	$B_{10,0.4}(\{m\})$	m	$H_{200,100,10}(\{m\})$	$B_{10,0.5}(\{m\})$
1	0.0372601	0.0403108	1	0.00847281	0.00976563
2	0.118268	0.120932	2	0.0410287	0.0439453
3	0.217696	0.214991	3	0.115292	0.117188
4	0.257321	0.250823	4	0.2082	0.205078
5	0.204067	0.200658	5	0.25247	0.246094
6	0.10995	0.111477	6	0.2082	0.205078
7	0.0397376	0.0424673	7	0.115292	0.117188
8	0.00921879	0.0106168	8	0.0410287	0.0439453
9	0.0012395	0.00157286	9	0.00847281	0.00976563

Let us shortly analyze the table on the right. Here 50 % of the balls are white. Drawing 10 balls, it is more likely to observe 4, 5, or 6 white balls in the case of nonreplacement, while it is the other way around in the remaining cases 1, 2, 3, 7, 8, and 9. Try to find a heuristic explanation for this phenomenon.

Summary: The hypergeometric distribution may be viewed as a counterpart to the binomial distribution in the following sense: in an urn there are M white balls and $N - M$ black. Choosing randomly n balls *without replacement*,

$$H_{N,M,n}(\{m\}) := \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}$$

is the probability to observe m white balls. In contrast, the binomial distribution $B_{n,p}$ with $p = M/N$ applies in the case of *replacing* the chosen balls. Of course, in this case the number N of balls does not matter, only the proportion of white balls is of interest.

1.4.7 Geometric distribution

At first glance, the model for the geometric distribution looks as that for the binomial distribution. In each single trial, we may observe “0” or “1”, that is, failure or success. Again the success probability is a fixed number p . While in the case of the binomial distribution we executed a fixed number of trials, now this number is random. More precisely, we execute the experiment until we observe success for the first time. Recorded is the number of necessary trials until this first success shows up. Or, in other words, a number $k \geq 1$ occurs if and only if the first $k - 1$ trials were all failures and the k th one is a success, that is, we observe the sequence $(\underbrace{0, \dots, 0}_{k-1}, 1)$. Since failure appears with probability $1 - p$ and success shows up with probability p , the following approach is plausible:

$$G_p(\{k\}) := p(1-p)^{k-1}, \quad k \in \mathbb{N}. \quad (1.29)$$

Proposition 1.4.41. *If $0 < p < 1$, then (1.29) defines a probability measure on $\mathcal{P}(\mathbb{N})$.*

Proof. Because of $p(1-p)^{k-1} > 0$, it suffices to verify $\sum_{k=1}^{\infty} G_p(\{k\}) = 1$. Using the formula for the sum of a geometric series, this follows directly from

$$\sum_{k=1}^{\infty} p(1-p)^{k-1} = p \sum_{k=0}^{\infty} (1-p)^k = p \frac{1}{1-(1-p)} = 1.$$

Observe that by assumption $1 - p < 1$, thus the formula for the geometric series applies. \square

Definition 1.4.42. The probability measure G_p on $\mathcal{P}(\mathbb{N})$ defined by Eq. (1.29) is called the **geometric distribution** with parameter p .

If $p = 0$, then success will never show up, thus, G_p is not a probability measure. On the other hand, for $p = 1$, success appears with probability one in the first trial, that is, $G_p = \delta_1$. Therefore, this case is of no interest.

Remark 1.4.43. Some authors define the geometric distribution in a slightly different way. They ask for the probability for the first success in the $(k+1)$ th trial. This is described by

$$\tilde{G}_p(\{k\}) := p(1-p)^k, \quad k \in \mathbb{N}_0.$$

To our opinion, the shift by 1 is a little bit confusing. Therefore, we decided to define the geometric distribution as we did in eq. (1.29).

Example 1.4.44. Given a number $n \in \mathbb{N}$, let $A_n = \{k \in \mathbb{N} : k > n\}$. Find $G_p(A_n)$.

Answer: We answer this question by two different approaches.

At first, we remark that A_n occurs if and only if the first success shows up strictly after n trials or, equivalently, if and only if the first n trials were all failures. But this event has probability $B_{n,p}(\{0\}) = (1-p)^n$, hence $G_p(A_n) = (1-p)^n$.

In the second approach, we use eq. (1.29) directly and obtain

$$\begin{aligned} G_p(A_n) &= \sum_{k=n+1}^{\infty} G_p(\{k\}) = p \sum_{k=n+1}^{\infty} (1-p)^{k-1} = p(1-p)^n \sum_{k=0}^{\infty} (1-p)^k \\ &= p(1-p)^n \frac{1}{1-(1-p)} = (1-p)^n. \end{aligned}$$

Example 1.4.45. Roll a die until number “6” occurs for the first time. What is the probability that this happens in roll k ?

Answer: The success probability is $1/6$, hence the probability of the first occurrence of “6” in the k th trial is $(1/6)(5/6)^{k-1}$.

k	$G_{1/6}(\{k\})$
1	0.166667
2	0.138889
3	0.115741
⋮	⋮
12	0.022431
13	0.018693

Example 1.4.46. Roll a die until the first “6” shows up. What is the probability that this happens at an even number of trials?

Answer: The first “6” has to appear in the second, or fourth, or sixth, and so on, trial. Hence, the probability of this event is

$$\sum_{k=1}^{\infty} G_{1/6}(\{2k\}) = \frac{1}{6} \sum_{k=1}^{\infty} (5/6)^{2k-1} = \frac{5}{36} \sum_{k=1}^{\infty} (5/6)^{2k-2} = \frac{5}{36} \frac{1}{1-(5/6)^2} = \frac{5}{11}.$$

Example 1.4.47. Play a series of games where p is the chance of winning. Whenever you put x dollars into the pool you get back $2x$ dollars if you win. If you lose, then the x dollars are lost.

Apply the following strategy. After losing, double the amount in the pool in the next game. Say you start with \$1 and lose, then next time put \$2 into the pool, then \$4, and so on until you win for the first time. As easily seen, in the k th game, the stakes are 2^{k-1} dollars.

Suppose for some $k \geq 1$ you lost $k-1$ games and won the k th one. How much money did you lose? If $k = 1$, then you lost nothing, while for $k \geq 2$ you spent

$$1 + 2 + 4 + \cdots + 2^{k-2} = 2^{k-1} - 1$$

dollars. Note that $2^{k-1} - 1 = 0$ if $k = 1$, hence for all $k \geq 1$ the total loss is $2^{k-1} - 1$ dollars.

On the other hand, if you win the k th game, you gain 2^{k-1} dollars. Consequently, no matter what the results are, you will always win $2^{k-1} - (2^{k-1} - 1) = 1$ dollar.¹¹

Let X be the amount of money needed to follow this strategy. One needs 1 dollar to play the first game, $1 + 2 = 3$ dollars to play the second, until

$$1 + 2 + 4 + \cdots + 2^{k-1} = 2^k - 1$$

to play the k th game. Thus, in this case we have $X = 2^k - 1$ and

$$\mathbb{P}\{X = 2^k - 1\} = \mathbb{P}\{\text{first win in game } k\} = p(1-p)^{k-1}, \quad k = 1, 2, \dots$$

In particular, if $p = 1/2$ then this probability equals 2^{-k} . For example, if one starts the game with $127 = 2^7 - 1$ dollars in the pocket, then one goes bankrupt if the first success appears after game 7. The probability for this equals $\sum_{k=8}^{\infty} 2^{-k} = 2^{-7} = 0.0078125$.

1.4.8 Negative binomial distribution

The geometric distribution describes the probability for having the first success in trial k . Given a fixed $n \geq 1$, we ask now for the probability that in trial k success appears not for the first but for the n th time. Of course, this question makes only sense if $k \geq n$. But how to determine this probability for those k ?

Thus, take $k \geq n$ and suppose we had success in trial k . When is this the n th one? This is the case if and only if we had $n-1$ successes during the first $k-1$ trials or, equivalently, $k-n$ failures. There exist $\binom{k-1}{k-n}$ possibilities to distribute the $k-n$ failures among the first $k-1$ trials. Furthermore, the probability for n successes is p^n and for $k-n$ failures

¹¹ Starting the first game with x dollars, one will always win x dollars no matter what happens.

it is $(1-p)^{k-n}$, hence the probability for observing the n th success in trial k is given by

$$B_{n,p}^-(\{k\}) := \binom{k-1}{k-n} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots \quad (1.30)$$

We still have to verify that there is a probability measure satisfying eq. (1.30).

Proposition 1.4.48. *By*

$$B_{n,p}^-(\{k\}) = \binom{k-1}{k-n} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots,$$

a probability measure $B_{n,p}^-$ on $\mathcal{P}(\{n, n+1, \dots\})$ is defined.

Proof. Of course, $B_{n,p}^-(\{k\}) \geq 0$. Hence it remains to show

$$\sum_{k=n}^{\infty} B_{n,p}^-(\{k\}) = 1 \quad \text{or, equivalently,} \quad \sum_{k=0}^{\infty} B_{n,p}^-(\{k+n\}) = 1. \quad (1.31)$$

Because of Proposition A.3.11, we get

$$\begin{aligned} B_{n,p}^-(\{k+n\}) &= \binom{n+k-1}{k} p^n (1-p)^k \\ &= \binom{-n}{k} p^n (-1)^k (1-p)^k = \binom{-n}{k} p^n (p-1)^k, \end{aligned} \quad (1.32)$$

where the generalized binomial coefficient is defined in eq. (A.14) as

$$\binom{-n}{k} = \frac{-n(-n-1)\cdots(-n-k+1)}{k!}.$$

In Proposition A.5.2, we proved for $|x| < 1$ that

$$\sum_{k=0}^{\infty} \binom{-n}{k} x^k = \frac{1}{(1+x)^n}. \quad (1.33)$$

Note that $0 < p < 1$, hence eq. (1.33) applies with $x = p-1$ and leads to

$$\sum_{k=0}^{\infty} \binom{-n}{k} (p-1)^k = \frac{1}{p^n}. \quad (1.34)$$

Combining eqs. (1.32) and (1.34) implies

$$\sum_{k=0}^{\infty} B_{n,p}^-(\{k+n\}) = p^n \sum_{k=0}^{\infty} \binom{-n}{k} (p-1)^k = p^n \frac{1}{p^n} = 1,$$

thus the equations in (1.31) are valid and this completes the proof. \square

Definition 1.4.49. The probability measure $B_{n,p}^-$ with

$$B_{n,p}^-(\{k\}) := \binom{k-1}{k-n} p^n (1-p)^{k-n} = \binom{k-1}{n-1} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots$$

is called the **negative binomial distribution** with parameters $n \geq 1$ and $p \in (0, 1)$. Of course, $B_{1,p}^- = G_p$.

Remark 1.4.50. We saw in eq. (1.32) that

$$B_{n,p}^-(\{k+n\}) = \binom{n+k-1}{k} p^n (1-p)^k = \binom{-n}{k} p^n (p-1)^k. \quad (1.35)$$

Alternatively, one may define the negative binomial distribution also via Eq. (1.35). Then it describes the event that the n th success appears in trial $n+k$. The advantage of this approach is that now $k \in \mathbb{N}_0$, that is, the restriction $k \geq n$ is no longer needed. Its disadvantage is that we are interested in what happens in trial k , not in trial $k+n$.

Example 1.4.51. Roll a die successively. Determine the probability that in the 20th trial number “6” appears for the fourth time.

Answer: We have $p = 1/6$, $n = 4$, and $k = 20$. Therefore, the probability for this event is given by

$$B_{4,1/6}^-(\{20\}) = \binom{19}{16} \left(\frac{1}{6}\right)^4 \left(\frac{5}{6}\right)^{16} = 0.0404407.$$

Let us ask now for the probability that the fourth success appears (strictly) before trial 21. This probability is given by

$$\sum_{k=4}^{20} \binom{k-1}{3} \left(\frac{1}{6}\right)^4 \left(\frac{5}{6}\right)^{k-4} = 0.433454.$$

Example 1.4.52. There are two urns, say U_0 and U_1 , each containing N balls. Choose one of the two urns at random and take out a ball. Hereby U_0 is chosen with probability $1-p$, hence U_1 with probability p . Repeat the procedure until we choose the last (the N th) ball out of one of the urns. What is the probability that there are m balls left in the other urn, where $m = 1, \dots, N$?

Answer: For $m = 1, \dots, N$ let A_m be the event that there are still m balls in one of the urns when choosing the last ball out of the other. Then A_m splits into the disjoint events $A_m = A_m^0 \cup A_m^1$, where

- A_m^0 occurs if we take the last ball out of U_0 and U_1 contains m balls, and
- A_m^1 occurs choosing the N th ball out of U_1 with m balls remaining in U_0 .

Let us start with evaluating the probability of A_m^1 . Say success occurs if we choose urn U_1 . Thus, if we take out the last ball of urn U_1 , then success occurred for the N th time. On the other hand, if there are still m balls in U_0 , then failure had occurred $N-m$ times.

Consequently, there are still m balls left in urn U_0 if and only if the N th success shows up in trial $N + (N - m) = 2N - m$. Therefore, we get

$$\mathbb{P}(A_m^1) = B_{N,p}^-(\{2N - m\}) = \binom{2N - m - 1}{N - m} p^N (1 - p)^{N - m}. \quad (1.36)$$

The probability of A_m^0 may be derived from that of A_m^1 by interchanging p and $1 - p$ (success occurs now with probability $1 - p$). This yields

$$\mathbb{P}(A_m^0) = B_{N,1-p}^-(\{2N - m\}) = \binom{2N - m - 1}{N - m} p^{N - m} (1 - p)^N. \quad (1.37)$$

Adding eqs. (1.36) and (1.37) leads to

$$\mathbb{P}(A_m) = \binom{2N - m - 1}{N - m} [p^N (1 - p)^{N - m} + p^{N - m} (1 - p)^N],$$

for $m = 1, \dots, N$.

If $p = 1/2$, that is, both urns are equally likely, the previous formula simplifies to

$$\mathbb{P}(A_m) = \binom{2N - m - 1}{N - m} 2^{-2N + m + 1} = \binom{2N - m - 1}{N - 1} 2^{-2N + m + 1}. \quad (1.38)$$

Remark 1.4.53. The case $p = 1/2$ in the previous problem is known as *Banach's matchbox problem*. In each of two matchboxes, there are N matches. One chooses randomly a matchbox (both boxes are equally likely) and takes out a match. What is the probability that there are still m matches left in the other box when taking the last match out of one of the boxes? The answer is given by eq. (1.38).

Remark 1.4.54. There exists a slightly different version of Banach's matchbox problem. Here one asks for the probability of the event \tilde{A}_m which occurs provided that there are m matches left in one of the boxes when choosing for the first time an empty one. Note that in this setting also $m = 0$ makes sense. To answer this modified problem, one has to ask in the previous calculations for the $(N + 1)$ th success instead of the N th one. In doing so, one obtains

$$\mathbb{P}(\tilde{A}_m) = \binom{2N - m}{N - m} 2^{-2N + m} = \binom{2N - m}{N} 2^{-2N + m}, \quad m = 0, \dots, N.$$

See Figure 1.7 for the values of these probabilities in the case $N = 20$.

What is more likely at the moment when choosing for the first time an empty box: the unchosen box contains one match or it contains none?

The answer is that both events are equally likely. This follows from

$$\begin{aligned} \mathbb{P}(\tilde{A}_0) &= \binom{2N}{N} 2^{-2N} = \frac{(2N)!}{(N!)^2} 2^{-2N} = \frac{2N}{N} \frac{(2N - 1)!}{N!(N - 1)!} 2^{-2N} \\ &= \binom{2N - 1}{N - 1} 2^{-2N + 1} = \mathbb{P}(\tilde{A}_1). \end{aligned}$$

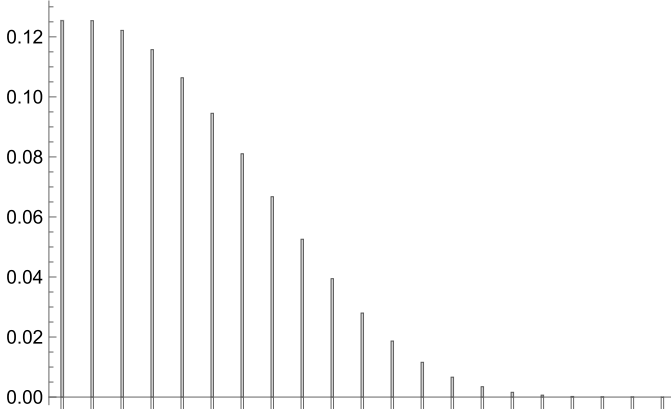


Figure 1.7: Probabilities for $0 \leq m \leq 20$ matches left in one box when choosing for the first time an empty one. At the beginning, both boxes contained $N = 20$ matches.

Example 1.4.55. We continue Example 1.4.52 with $0 < p < 1$ and ask the following question: What is the probability that U_1 becomes empty before U_0 ?

Answer: This happens if and only if U_0 is nonempty when choosing U_1 for the N th time, that is, when in U_0 there are m balls left for some $m = 1, \dots, N$. Because of eq. (1.36), this probability is given by

$$\begin{aligned} \sum_{m=1}^N \mathbb{P}(A_m^1) &= p^N \sum_{m=1}^N \binom{2N-m-1}{N-m} (1-p)^{N-m} \\ &= p^N \sum_{k=0}^{N-1} \binom{N+k-1}{k} (1-p)^k. \end{aligned} \quad (1.39)$$

Remark 1.4.56. Formula (1.39) leads to an interesting (known) property of the binomial coefficients. Since $\sum_{m=1}^N \mathbb{P}(A_m) = 1$, by eqs. (1.39) and (1.37), we obtain

$$\sum_{k=0}^{N-1} \binom{N+k-1}{k} [p^N (1-p)^k + (1-p)^N p^k] = 1$$

or, setting $n = N - 1$, to

$$\sum_{k=0}^n \binom{n+k}{k} [p^{n+1} (1-p)^k + (1-p)^{n+1} p^k] = 1.$$

In particular, if $p = 1/2$, this yields

$$\sum_{k=0}^n \binom{n+k}{k} \frac{1}{2^k} = 2^n, \quad n = 0, 1, 2, \dots$$

Summary: One executes independently arbitrarily many experiments where either success or failure may occur. Hereby, the success probability equals $0 < p < 1$. The probability to observe the n th success in trial $k \geq n$ is given by the negative binomial distribution $B_{n,p}^-$ defined by

$$B_{n,p}^-(\{k\}) := \binom{k-1}{k-n} p^n (1-p)^{k-n} = \binom{k-1}{n-1} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots$$

Equivalently, when one asks for the n th success in trial $n + \ell$, then

$$B_{n,p}^-(\{n + \ell\}) = \binom{-n}{\ell} p^n (p-1)^\ell, \quad \ell = 0, 1, 2, \dots$$

If $n = 1$, that is, one looks for the first success, then $G_p = B_{1,p}^-$ is called the geometric distribution, and

$$G_p(\{k\}) = p(1-p)^{k-1}, \quad k = 1, 2, \dots$$

1.5 Continuous probability measures

Discrete probability measures are inappropriate for the description of random experiments where uncountably many different results may appear. Typical examples of such experiments are the lifetime of an item, the duration of a phone call, the measuring result of workpiece, and so on.

Discrete probability measures are concentrated on a finite or countably infinite set of points. An extension to larger sets is impossible. For example, there is no¹² probability measure \mathbb{P} on $[0, 1]$ with $\mathbb{P}(\{t\}) > 0$ for $t \in [0, 1]$.

Consequently, in order to describe random experiments with “many” possible different outcomes, another approach is needed. To explain this “new” approach, let us shortly recall how we evaluated $\mathbb{P}(A)$ in the discrete case. If Ω is either finite or countably infinite and if $p(\omega) = \mathbb{P}(\{\omega\})$, $\omega \in \Omega$, then with this $p: \Omega \rightarrow \mathbb{R}$ we have

$$\mathbb{P}(A) = \sum_{\omega \in A} p(\omega). \tag{1.40}$$

If the sample space is \mathbb{R} or \mathbb{R}^n , then such a representation is no longer possible. Indeed, if \mathbb{P} is not discrete, then, we will have $p(\omega) = 0$ for all possible observations ω . Therefore, the sum in eq. (1.40) has to be replaced by an integral over a more general function. We start with introducing functions p , which may be used for representing $\mathbb{P}(A)$ via an integral.

¹² Compare with Problem 1.38.

Definition 1.5.1. A Riemann integrable function $p : \mathbb{R} \rightarrow \mathbb{R}$ is called a **probability density function**, or simply a **density function**, if

$$p(x) \geq 0, \quad x \in \mathbb{R}, \quad \text{and} \quad \int_{-\infty}^{\infty} p(x) \, dx = 1. \quad (1.41)$$

Remark 1.5.2. Let us formulate more precisely the second condition for p in the previous definition. For all finite intervals $[a, b]$ in \mathbb{R} , the function p is Riemann integrable on $[a, b]$ and, moreover,

$$\lim_{\substack{a \rightarrow -\infty \\ b \rightarrow \infty}} \int_a^b p(x) \, dx = 1.$$

The density functions we will use later on are either continuous or piecewise continuous, that is, they are compositions of finitely many continuous functions. These functions are Riemann integrable, hence in this case it remains to verify the two conditions (1.41).

Example 1.5.3. Define p on \mathbb{R} by $p(x) = 0$ if $x < 0$ and by $p(x) = e^{-x}$ if $x \geq 0$. Then p is piecewise continuous, $p(x) \geq 0$ if $x \in \mathbb{R}$, and it satisfies

$$\int_{-\infty}^{\infty} p(x) \, dx = \lim_{b \rightarrow \infty} \int_0^b e^{-x} \, dx = \lim_{b \rightarrow \infty} [-e^{-x}]_0^b = 1 - \lim_{b \rightarrow \infty} e^{-b} = 1.$$

Hence, p is a density function.

Definition 1.5.4. Let p be a probability density function. Given a finite interval $[a, b]$, its probability (of occurrence) is defined by

$$\mathbb{P}([a, b]) := \int_a^b p(x) \, dx.$$

A graphical presentation of the previous definition is as follows. As visualized in Figure 1.8, the probability $\mathbb{P}([a, b])$ is the area under the graph of the density p , taken from a to b .

Let us illustrate Definition 1.5.4 with the density function regarded in Example 1.5.3. Then

$$\mathbb{P}([a, b]) = \int_a^b e^{-x} \, dx = [-e^{-x}]_a^b = e^{-a} - e^{-b}$$

whenever $0 \leq a < b < \infty$. On the other hand, if $a < b < 0$, then $\mathbb{P}([a, b]) = 0$ while for $a < 0 \leq b$ the probability of $[a, b]$ is calculated by

$$\mathbb{P}([a, b]) = \mathbb{P}([0, b]) = 1 - e^{-b}.$$

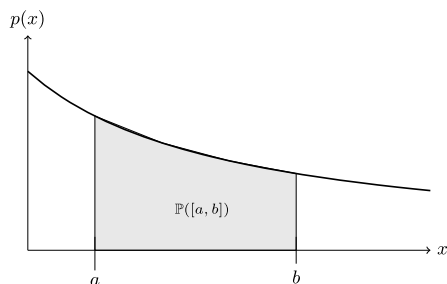


Figure 1.8: The size of the gray shaded area defines the probability of the interval $[a, b]$.

Remark 1.5.5. Definition 1.5.4 of the probability measure \mathbb{P} does not fit into the scheme presented in Section 1.1.3. Why? Probability measures are defined on σ -fields. But the collection of finite intervals in \mathbb{R} is not a σ -field. It is neither closed under taking complements nor is the union of intervals in general again an interval. Furthermore, it is far from being clear in which sense \mathbb{P} should be σ -additive.

The next result justifies the approach in Definition 1.5.4. Its proof rests upon an extension theorem in Measure Theory (cf. [Bau01, Coh13] or [Dud02]).

Proposition 1.5.6. *Let $\mathcal{B}(\mathbb{R})$ be the σ -field of Borel sets introduced in Definition 1.1.20. Then for each density function p , there exists a unique probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R})$ such that*

$$\mathbb{P}([a, b]) = \int_a^b p(x) \, dx \quad \text{for all } a < b. \quad (1.42)$$

Definition 1.5.7. A probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R})$ is said to be **continuous** provided that there exists a density function p such that for $a < b$,

$$\mathbb{P}([a, b]) = \int_a^b p(x) \, dx. \quad (1.43)$$

The function p is called a **density function**, or simply **density**, of \mathbb{P} .

Remark 1.5.8. The mathematically correct name would be “absolutely continuous”. But since we do not treat so-called “singularly continuous” probability measures, there is no need to distinguish between them, and we may shorten the notation to “continuous.”

Remark 1.5.9. Note that changing the density function at finitely many points does not change the generated probability measure. For instance, if we define $p(x) = 0$ if $x \leq 0$ and $p(x) = e^{-x}$ if $x > 0$, then this density function is different from that in Example 1.5.3 but, of course, generates the same probability measure.

Moreover, observe that eq. (1.43) is valid for all $a < b$ if and only if for each $t \in \mathbb{R}$,

$$\mathbb{P}((-\infty, t]) = \int_{-\infty}^t p(x) dx. \quad (1.44)$$

Consequently, \mathbb{P} is continuous if and only if there is a density p with eq. (1.44) for $t \in \mathbb{R}$.

Proposition 1.5.10. *Let $\mathbb{P} : \mathcal{B}(\mathbb{R}) \rightarrow [0, 1]$ be a continuous probability measure with density p . Then the following are valid:*

1. $\mathbb{P}(\mathbb{R}) = 1$.
2. For each $t \in \mathbb{R}$, one has $\mathbb{P}(\{t\}) = 0$. More generally, if $A \subseteq \mathbb{R}$ is either finite or countably infinite, then $\mathbb{P}(A) = 0$.
3. For all $a < b$, we have

$$\mathbb{P}((a, b)) = \mathbb{P}([a, b)) = \mathbb{P}([a, b]) = \mathbb{P}(]a, b]) = \int_a^b p(x) dx.$$

Proof. Let us start with proving $\mathbb{P}(\mathbb{R}) = 1$. For $n \geq 1$, set $A_n := [-n, n]$ and note that $A_1 \subseteq A_2 \subseteq \dots$, as well as $\bigcup_{n=1}^{\infty} A_n = \mathbb{R}$. Thus we may use that \mathbb{P} is continuous from below and, by the properties of the density p , we obtain

$$\mathbb{P}(\mathbb{R}) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \lim_{n \rightarrow \infty} \int_{-n}^n p(x) dx = \int_{-\infty}^{\infty} p(x) dx = 1.$$

To verify the second property, fix $t \in \mathbb{R}$ and define for each $n \geq 1$ the intervals B_n by $B_n := [t, t + \frac{1}{n}]$. Now we have $B_1 \supseteq B_2 \supseteq \dots$ and $\bigcap_{n=1}^{\infty} B_n = \{t\}$. Use this time the continuity from above. Then we get

$$\mathbb{P}(\{t\}) = \lim_{n \rightarrow \infty} \mathbb{P}(B_n) = \lim_{n \rightarrow \infty} \int_t^{t+\frac{1}{n}} p(x) dx = 0.$$

If $A = \{t_1, t_2, \dots\}$, then the σ -additivity of \mathbb{P} , together with $\mathbb{P}(\{t_j\}) = 0$, gives

$$\mathbb{P}(A) = \sum_{j=1}^{\infty} \mathbb{P}(\{t_j\}) = 0,$$

as asserted.

The third property is an immediate consequence of the second. Observe

$$[a, b] = (a, b) \cup \{a\} \cup \{b\},$$

hence $\mathbb{P}([a, b]) = \mathbb{P}((a, b)) + \mathbb{P}(\{a\}) + \mathbb{P}(\{b\})$, proving (1.43) due to $\mathbb{P}(\{a\}) = \mathbb{P}(\{b\}) = 0$. \square

Remark 1.5.11. Say a set $C \subseteq \mathbb{R}$ can be represented as $C = \bigcup_{j=1}^{\infty} I_j$ with disjoint (open or half-open or closed) intervals I_j , then

$$\mathbb{P}(C) = \sum_{j=1}^{\infty} \int_{I_j} p(x) \, dx := \int_C p(x) \, dx.$$

More generally, if a set B may be written as $B = \bigcap_{n=1}^{\infty} C_n$ where the C_n s are unions of disjoint intervals and satisfy $C_1 \supseteq C_2 \supseteq \dots$, then

$$\mathbb{P}(B) = \lim_{n \rightarrow \infty} \mathbb{P}(C_n).$$

In this way, one may evaluate $\mathbb{P}(B)$ for a large class of subsets $B \subseteq \mathbb{R}$.

Summary: There are two completely different types of probability measures:

$$\left[\mathbb{P} \text{ discrete} \Leftrightarrow \mathbb{P}([a, b]) = \sum_{a \leq \omega \leq b} \mathbb{P}(\{\omega\}) \right] \text{ and } \left[\mathbb{P} \text{ continuous} \Leftrightarrow \mathbb{P}([a, b]) = \int_a^b p(x) \, dx \right].$$

1.6 Special continuous distributions

1.6.1 Uniform distribution on an interval

Let $I = [\alpha, \beta]$ be a finite interval of real numbers. Define a function $p : \mathbb{R} \rightarrow \mathbb{R}$ by

$$p(x) := \begin{cases} \frac{1}{\beta - \alpha} & \text{if } x \in [\alpha, \beta], \\ 0 & \text{if } x \notin [\alpha, \beta]. \end{cases} \quad (1.45)$$

Proposition 1.6.1. *The mapping p defined by eq. (1.45) is a probability density function.*

Proof. Note that p is piecewise continuous, hence Riemann integrable. Moreover, $p(x) \geq 0$ for $x \in \mathbb{R}$ and

$$\int_{-\infty}^{\infty} p(x) \, dx = \int_{\alpha}^{\beta} \frac{1}{\beta - \alpha} \, dx = \frac{1}{\beta - \alpha} (\beta - \alpha) = 1.$$

Consequently, p is a probability density. □

Definition 1.6.2. The probability measure \mathbb{P} generated by the density in eq. (1.45) is called the **uniform distribution** on the interval $I = [\alpha, \beta]$.

How is $\mathbb{P}([a, b])$ evaluated for some interval $[a, b]$? Let us first treat the case that $[a, b] \subseteq I$. Then

$$\mathbb{P}([a, b]) = \int_a^b \frac{1}{\beta - \alpha} dx = \frac{b - a}{\beta - \alpha} = \frac{\text{Length of } [a, b]}{\text{Length of } [\alpha, \beta]}. \quad (1.46)$$

This explains why \mathbb{P} is called the “uniform distribution.” The probability of an interval $[a, b] \subseteq I$ depends only on its length, not on its position. Shifting $[a, b]$ inside I does not change its probability of occurrence (compare Figure 1.9).

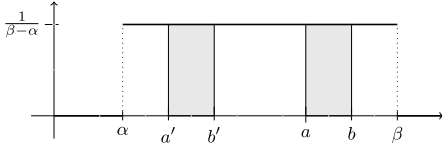


Figure 1.9: If \mathbb{P} is the uniform distribution on $[a, \beta]$, then $\mathbb{P}([a, b]) = \mathbb{P}([a', b'])$.

If $[a, b]$ is arbitrary, not necessarily contained in I , then $\mathbb{P}([a, b])$ can be easily calculated by

$$\mathbb{P}([a, b]) = \mathbb{P}([a, b] \cap I).$$

Example 1.6.3. Let \mathbb{P} be the uniform distribution on $[0, 1]$. Which probabilities do $[-1, 0.5]$, $[0, 0.25] \cup [0.75, 1]$, $(-\infty, t]$ if $t \in \mathbb{R}$, and $A \subseteq \mathbb{R}$, where $A = \bigcup_{n=1}^{\infty} [\frac{1}{2^{n+1/2}}, \frac{1}{2^n}]$, have?

Answer: The first two intervals have probability $\frac{1}{2}$. If $t \in \mathbb{R}$, then

$$\mathbb{P}((-\infty, t]) = \begin{cases} 0 & \text{if } t < 0, \\ t & \text{if } 0 \leq t \leq 1, \\ 1 & \text{if } t > 1. \end{cases}$$

Finally, observe that the intervals $[\frac{1}{2^{n+1/2}}, \frac{1}{2^n}]$ are disjoint subsets of $[0, 1]$. Hence we get

$$\mathbb{P}(A) = \sum_{n=1}^{\infty} \left[\frac{1}{2^n} - \frac{1}{2^{n+1/2}} \right] = (1 - 2^{-1/2}) \sum_{n=1}^{\infty} \frac{1}{2^n} = 1 - 2^{-1/2}.$$

Example 1.6.4. A stick of length L is randomly broken into two pieces. Find the probability that the size of one piece is at least twice that of the other.

Answer: This event happens if and only if the point at which the stick is broken is either in $[0, L/3]$ or in $[2L/3, L]$. Assuming that the point at which the stick is broken is uniformly distributed on $[0, L]$, the desired probability is $\frac{2}{3}$. Another way to get this result is as follows. The size of each piece is less than twice as that of the other if the point at

which the stick is broken is in $[L/3, 2L/3]$. Hence, the probability of the complementary event is $1/3$, leading again to $2/3$ for the desired probability.

Example 1.6.5. Choose at random a real number uniformly distributed in $[-1, 1]$. How likely is it that its square is less than $1/4$?

Answer: A number $x \in [-1, 1]$ satisfies $x^2 \leq 1/4$ if and only if $-1/2 \leq x \leq 1/2$. Hence, if A is the event that the square is less than $1/4$, then

$$\mathbb{P}(A) = \frac{\text{Length of } [-\frac{1}{2}, \frac{1}{2}]}{\text{Length of } [-1, 1]} = \frac{1}{2}.$$

Example 1.6.6. Let $C_0 := [0, 1]$. Extract from C_0 the interval $(\frac{1}{3}, \frac{2}{3})$, thus there remains $C_1 = [0, \frac{1}{3}] \cup [\frac{2}{3}, 1]$. To construct C_2 , extract from C_1 the two middle intervals $(\frac{1}{9}, \frac{2}{9})$ and $(\frac{7}{9}, \frac{8}{9})$, hence $C_2 = [0, \frac{1}{9}] \cup [\frac{2}{9}, \frac{1}{3}] \cup [\frac{2}{3}, \frac{7}{9}] \cup [\frac{8}{9}, 1]$.

Suppose that through this method we already got the set C_n which is a union of 2^n disjoint closed intervals of length 3^{-n} . In order to construct C_{n+1} , split each of the 2^n intervals into three intervals of length 3^{-n-1} and erase the middle one. In this way, we get C_{n+1} , which consists of 2^{n+1} disjoint intervals of length 3^{-n-1} . Finally, one defines

$$C = \bigcap_{n=1}^{\infty} C_n.$$

The set C is known as the **Cantor set**. Let \mathbb{P} be the uniform distribution on $[0, 1]$. Which value does $\mathbb{P}(C)$ have?

Answer: First observe that $C_0 \supset C_1 \supset C_2 \supset \dots$, hence, using that \mathbb{P} is continuous from above, it follows that

$$\mathbb{P}(C) = \lim_{n \rightarrow \infty} \mathbb{P}(C_n). \quad (1.47)$$

The set C_n is a disjoint union of 2^n intervals of length 3^{-n} . Consequently, it follows that $\mathbb{P}(C_n) = \frac{2^n}{3^n}$, which by eq. (1.47) implies $\mathbb{P}(C) = 0$.

One might conjecture that $C = \emptyset$. On the contrary, C is even uncountably infinite. To see this, we have to make the construction of the Cantor set a little bit more precise.

Given $n \geq 1$, let

$$A_n = \{\alpha = (\alpha_1, \dots, \alpha_n) : \alpha_1, \dots, \alpha_{n-1} \in \{0, 2\}, \alpha_n = 1\}.$$

If $\alpha = (\alpha_1, \dots, \alpha_n) \in A_n$, set $x_\alpha = \sum_{k=1}^n \frac{\alpha_k}{3^k}$ and $I_\alpha = (x_\alpha, x_\alpha + \frac{1}{3^n})$. In this notation,

$$I_{(1)} = \left(\frac{1}{3}, \frac{2}{3}\right), \quad I_{(0,1)} = \left(\frac{1}{9}, \frac{2}{9}\right), \quad I_{(2,1)} = \left(\frac{7}{9}, \frac{8}{9}\right), \quad \text{and} \quad I_{(0,0,1)} = \left(\frac{1}{27}, \frac{2}{27}\right).$$

Then, if $C_0 = [0, 1]$, for $n \geq 1$ we have

$$C_n = C_{n-1} \setminus \bigcup_{\alpha \in A_n} I_\alpha, \quad \text{hence } C = [0, 1] \setminus \bigcup_{n=1}^{\infty} \bigcup_{\alpha \in A_n} I_\alpha.$$

Take now any sequence x_1, x_2, \dots with $x_k \in \{0, 2\}$ and set $x = \sum_{k=1}^{\infty} \frac{x_k}{3^k}$. Then x cannot belong to any I_α because otherwise at least one of the x_k s should satisfy $x_k = 1$. Thus $x \in C$, and the number of x that may be represented by x_k s with $x_k \in \{0, 2\}$ is uncountably infinite.

Summary: The uniform distribution \mathbb{P} on an interval $I = [a, \beta]$ is characterized by the following property: if $[a, b] \subseteq I$, then

$$\mathbb{P}([a, b]) = \int_a^b \frac{1}{\beta - a} dx = \frac{b - a}{\beta - a} = \frac{\text{Length of } [a, b]}{\text{Length of } [a, \beta]}.$$

1.6.2 Normal distribution

This section is devoted to the most important probability measure, the normal distribution. Before we can introduce it, we need the following result.

Proposition 1.6.7. *We have*

$$\int_{-\infty}^{\infty} e^{-x^2/2} dx = \sqrt{2\pi}. \quad (1.48)$$

Proof. Set

$$a := \int_{-\infty}^{\infty} e^{-x^2/2} dx$$

and note that $a > 0$. Then we get

$$a^2 = \left(\int_{-\infty}^{\infty} e^{-x^2/2} dx \right) \left(\int_{-\infty}^{\infty} e^{-y^2/2} dy \right) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy.$$

Change the variables in the right-hand double integral as follows: $x := r \cos \theta$ and $y := r \sin \theta$, where $0 < r < \infty$ and $0 \leq \theta < 2\pi$. Observe that

$$dx dy = |D(r, \theta)| dr d\theta$$

with

$$D(r, \theta) = \det \begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{pmatrix} = \det \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix} = r \cos^2 \theta + r \sin^2 \theta = r.$$

Using $x^2 + y^2 = r^2 \cos^2 \theta + r^2 \sin^2 \theta = r^2$, this change of variables leads to

$$a^2 = \int_0^{2\pi} \int_0^\infty r e^{-r^2/2} dr d\theta = \int_0^{2\pi} [-e^{-r^2/2}]_0^\infty d\theta = 2\pi,$$

which, due to $a > 0$, implies $a = \sqrt{2\pi}$. This completes the proof. \square

Given $\mu \in \mathbb{R}$ and $\sigma > 0$, let

$$p_{\mu,\sigma}(x) := \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}, \quad x \in \mathbb{R}. \quad (1.49)$$

Remark 1.6.8. The graph of the function $x \mapsto p_{\mu,\sigma}(x)$ is “bell-shaped,” has its maximal value $1/\sqrt{2\pi}\sigma$ at $x = \mu$, attains only positive values, and is symmetric around $x = \mu$. If $x \rightarrow \infty$ or $x \rightarrow -\infty$, then $p_{\mu,\sigma}(x)$ tends to zero very rapidly. The bigger the $\sigma > 0$, the flatter the graph of $p_{\mu,\sigma}$. Nevertheless, as we will show in the next proposition, the area under the graph of $p_{\mu,\sigma}$ always equals 1, no matter how big or small $\sigma > 0$ is (compare also Figure 1.10).

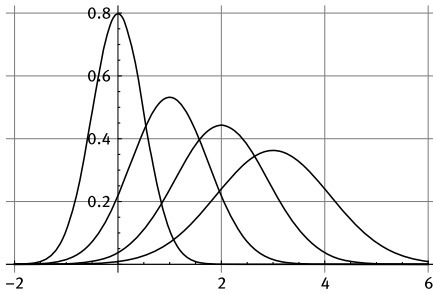


Figure 1.10: The function $p_{\mu,\sigma}$ with parameters $\mu = 0, 1, 2, 3$ and $\sigma = 0.5, 0.75, 0.9, 1.1$.

Proposition 1.6.9. If $\mu \in \mathbb{R}$ and $\sigma > 0$, then $p_{\mu,\sigma}$ is a probability density function.

Proof. We have to verify

$$\int_{-\infty}^{\infty} p_{\mu,\sigma}(x) dx = 1, \quad \text{or} \quad \int_{-\infty}^{\infty} e^{(x-\mu)^2/2\sigma^2} dx = \sqrt{2\pi} \sigma.$$

Setting $u := (x - \mu)/\sigma$, it follows that $dx = \sigma du$, hence Proposition 1.6.7 leads to

$$\int_{-\infty}^{\infty} e^{(x-\mu)^2/2\sigma^2} dx = \sigma \int_{-\infty}^{\infty} e^{-u^2/2} du = \sigma \sqrt{2\pi}.$$

This completes the proof. \square

Definition 1.6.10. The probability measure generated by $p_{\mu,\sigma}$ is called the **normal distribution** with expected value μ and variance σ^2 . It is denoted by $\mathcal{N}(\mu, \sigma^2)$, that is, for all $a < b$,

$$\mathcal{N}(\mu, \sigma^2)([a, b]) = \frac{1}{\sqrt{2\pi}\sigma} \int_a^b e^{-(x-\mu)^2/2\sigma^2} dx .$$

Remark 1.6.11. At this moment, the numbers $\mu \in \mathbb{R}$ and $\sigma > 0$ are nothing else than parameters. Why they are called “expected value” and “variance” will become clear after we introduce these notations in Section 5.

Definition 1.6.12. The probability measure $\mathcal{N}(0, 1)$ is called the **standard normal distribution**. It is given by

$$\mathcal{N}(0, 1)([a, b]) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx .$$

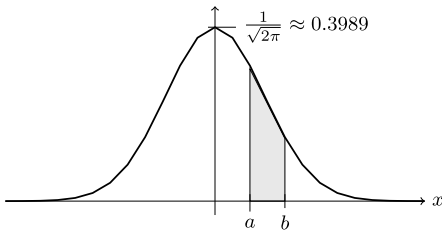


Figure 1.11: The area of the gray shaded region coincides with $\mathcal{N}(0, 1)([a, b])$.

For example, we have

$$\mathcal{N}(0, 1)([-1, 1]) = \frac{1}{\sqrt{2\pi}} \int_{-1}^1 e^{-x^2/2} dx = 0.682689, \quad \text{or}$$

$$\mathcal{N}(0, 1)([2, 4]) = \frac{1}{\sqrt{2\pi}} \int_2^4 e^{-x^2/2} dx = 0.0227185 .$$

Summary: The normal distribution with expected value $\mu \in \mathbb{R}$ and variance $\sigma^2 > 0$ acts as

$$\mathcal{N}(\mu, \sigma^2)([a, b]) = \frac{1}{\sqrt{2\pi}\sigma} \int_a^b e^{-(x-\mu)^2/2\sigma^2} dx .$$

The probability measure $\mathcal{N}(0, 1)$ is called the standard normal distribution.

1.6.3 Gamma distribution

Euler's **gamma function** is a mapping from $(0, \infty)$ to \mathbb{R} defined by

$$\Gamma(x) := \int_0^{\infty} s^{x-1} e^{-s} ds, \quad x > 0.$$

The graph of the gamma function for small entries is drawn in Figure 1.12.

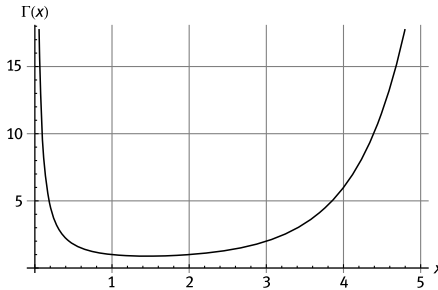


Figure 1.12: The graph of the gamma function.

Let us summarize the main properties of the gamma function.

Proposition 1.6.13.

1. *The gamma function maps $(0, \infty)$ continuously to $(0, \infty)$ and possesses continuous derivatives of any order.*
2. *If $x > 0$, then*

$$\Gamma(x + 1) = x \Gamma(x). \tag{1.50}$$

3. *For $n \in \mathbb{N}$, it follows that $\Gamma(n) = (n - 1)!$. In particular,*

$$\Gamma(1) = \Gamma(2) = 1 \quad \text{and} \quad \Gamma(3) = 2.$$

4. *One has $\Gamma(1/2) = \sqrt{\pi}$, which implies*

$$\Gamma\left(n + \frac{1}{2}\right) = \frac{(2n - 1)!!}{2^n} \sqrt{\pi}, \quad n = 1, 2, \dots \tag{1.51}$$

Here the double factorial is defined by $(2n - 1)!! = (2n - 1)(2n - 3) \cdots 3 \cdot 1$.

Proof. For the proof of the continuity and differentiability, we refer to [Art64].

The proof of eq. (1.50) is carried out by integration by parts as follows:

$$\Gamma(x + 1) = \int_0^{\infty} s^x e^{-s} ds = [-s^x e^{-s}]_0^{\infty} + \int_0^{\infty} x s^{x-1} e^{-s} ds = x \Gamma(x).$$

Note that $s^x e^{-s} = 0$ if $s = 0$ or $s \rightarrow \infty$.

From

$$\Gamma(1) = \int_0^{\infty} e^{-s} ds = 1$$

and eq. (1.50), it follows, as claimed, that

$$\Gamma(n) = (n-1)\Gamma(n-1) = (n-1)(n-2)\Gamma(n-2) = \dots = (n-1) \cdot \dots \cdot 1 \cdot \Gamma(1) = (n-1)!$$

To prove the fourth assertion, we use Proposition 1.6.7. Because of

$$\sqrt{2\pi} = \int_{-\infty}^{\infty} e^{-t^2/2} dt = 2 \int_0^{\infty} e^{-t^2/2} dt,$$

it follows that

$$\int_0^{\infty} e^{-t^2/2} dt = \sqrt{\frac{\pi}{2}}. \quad (1.52)$$

Substituting $s = t^2/2$, thus $ds = t dt$, by eq. (1.52), the integral for $\Gamma(1/2)$ transforms into

$$\Gamma(1/2) = \int_0^{\infty} s^{-1/2} e^{-s} ds = \int_0^{\infty} \frac{\sqrt{2}}{t} e^{-t^2/2} t dt = \sqrt{2} \int_0^{\infty} e^{-t^2/2} dt = \sqrt{\pi}.$$

Formula (1.51) follows by a repeated application of eq. (1.50). Indeed, then we get

$$\begin{aligned} \Gamma\left(n + \frac{1}{2}\right) &= \frac{2n-1}{2} \cdot \Gamma\left(n - \frac{1}{2}\right) \\ &= \frac{2n-1}{2} \cdot \frac{2n-3}{2} \cdot \Gamma\left(n - \frac{3}{2}\right) = \dots = \frac{(2n-1)!!}{2^n} \Gamma\left(\frac{1}{2}\right). \end{aligned}$$

This completes the proof. □

If $x \rightarrow \infty$, then $\Gamma(x)$ increases very rapidly. More precisely, the following is valid (cf. [Art64]):

Proposition 1.6.14 (Stirling's formula for the gamma function). *For $x > 0$, there exists a number $\theta \in (0, 1)$ such that*

$$\Gamma(x) = \sqrt{\frac{2\pi}{x}} \left(\frac{x}{e}\right)^x e^{\theta/12x}. \quad (1.53)$$

In view of

$$n! = \Gamma(n+1) = n\Gamma(n),$$

formula (1.53) leads to¹³ the following.

Corollary 1.6.15 (Stirling's formula for n -factorial). *For each natural number n , there is some $\theta \in (0, 1)$ depending on n such that*

$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\theta/12n}. \quad (1.54)$$

In particular, we have

$$\lim_{n \rightarrow \infty} \frac{e^n}{n^{n+1/2}} n! = \sqrt{2\pi}.$$

Our next aim is to introduce a continuous probability measure with density tightly related to the Γ -function. Given $\alpha, \beta > 0$, define $p_{\alpha, \beta}$ from \mathbb{R} to \mathbb{R} by (see Figure 1.13 for some examples)

$$p_{\alpha, \beta}(x) := \begin{cases} 0 & \text{if } x \leq 0, \\ \frac{1}{\alpha^\beta \Gamma(\beta)} x^{\beta-1} e^{-x/\alpha} & \text{if } x > 0. \end{cases} \quad (1.55)$$

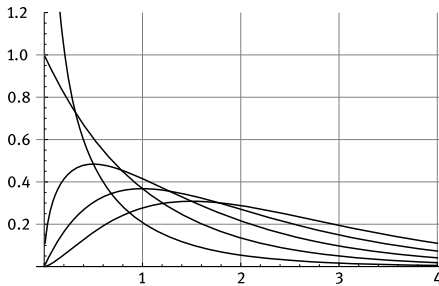


Figure 1.13: The functions $p_{1, \beta}$ with $\beta = 0.5, 1, 1.5, 2$ and 2.5 from top left to bottom left.

Proposition 1.6.16. *For all $\alpha, \beta > 0$, the function $p_{\alpha, \beta}$ in eq. (1.55) is a probability density.*

Proof. Of course, $p_{\alpha, \beta}(x) \geq 0$. Thus it remains to verify

$$\int_{-\infty}^{\infty} p_{\alpha, \beta}(x) dx = 1. \quad (1.56)$$

¹³ cf. also [Spi08, Chapter 27, Problem 19].

By the definition of $p_{\alpha,\beta}$, we have

$$\int_{-\infty}^{\infty} p_{\alpha,\beta}(x) dx = \frac{1}{\alpha^\beta \Gamma(\beta)} \int_0^{\infty} x^{\beta-1} e^{-x/\alpha} dx.$$

Substituting in the right-hand integral $u := x/\alpha$, thus $dx = \alpha du$, the right-hand side becomes

$$\frac{1}{\Gamma(\beta)} \int_0^{\infty} u^{\beta-1} e^{-u} du = \frac{1}{\Gamma(\beta)} \Gamma(\beta) = 1.$$

Hence eq. (1.56) is valid, and $p_{\alpha,\beta}$ is a probability density function. \square

Definition 1.6.17. The probability measure $\Gamma_{\alpha,\beta}$ with density function $p_{\alpha,\beta}$ is called the **gamma distribution** with positive parameters α and β . For all $0 \leq a < b < \infty$,

$$\Gamma_{\alpha,\beta}([a, b]) = \frac{1}{\alpha^\beta \Gamma(\beta)} \int_a^b x^{\beta-1} e^{-x/\alpha} dx. \quad (1.57)$$

Remark 1.6.18. Since $p_{\alpha,\beta}(x) = 0$ for $x \leq 0$, it follows that $\Gamma_{\alpha,\beta}((-\infty, 0]) = 0$. Hence, if $a < b$ are arbitrary, then

$$\Gamma_{\alpha,\beta}([a, b]) = \Gamma_{\alpha,\beta}([0, \infty) \cap [a, b]).$$

Remark 1.6.19. If $\beta \notin \mathbb{N}$, then the integral in eq. (1.57) cannot be expressed by elementary functions. Only numerical evaluations are possible.

1.6.4 Exponential distribution

An important special gamma distribution is the exponential distribution. This probability measure is defined as follows.

Definition 1.6.20. For $\lambda > 0$, let $E_\lambda := \Gamma_{\lambda^{-1}, 1}$ be the **exponential distribution** with parameter $\lambda > 0$.

Remark 1.6.21. The probability density function p_λ of E_λ is given by

$$p_\lambda(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ \lambda e^{-\lambda x} & \text{if } x > 0. \end{cases}$$

Consequently, if $0 \leq a < b < \infty$, then the probability of $[a, b]$ can be evaluated by

$$E_{\lambda}([a, b]) = \lambda \int_a^b e^{-\lambda x} dx = e^{-\lambda a} - e^{-\lambda b}.$$

Moreover,

$$E_{\lambda}([t, \infty)) = e^{-\lambda t}, \quad t \geq 0.$$

See Figure 1.14 for certain graphs of densities generating exponential distributions.

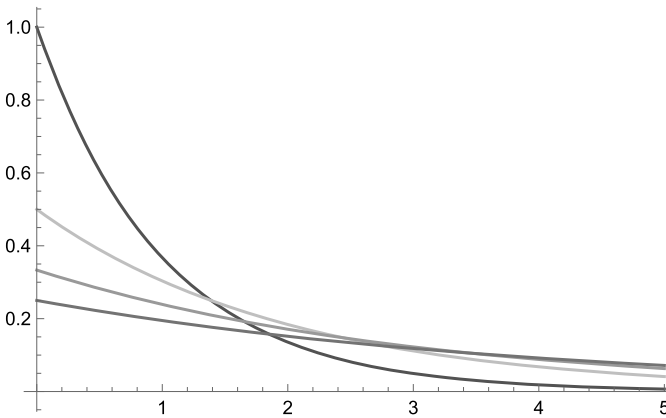


Figure 1.14: The densities of E_{λ} with $\lambda = 1, 1/2, 1/3$, and $1/4$.

Remark 1.6.22. The exponential distribution plays an important role for the description of lifetimes. For instance, it is used to determine the probability that the lifetime of a component part or the duration of a phone call exceeds a certain time $T > 0$. Furthermore, it is applied to describe the time between the arrivals of customers at a counter or in a shop.

Example 1.6.23. Suppose that the duration of phone calls is exponentially distributed with parameter $\lambda = 0.1$. What is the probability that a call lasts less than two time units? Or what is the probability that it lasts between one and two units? Or more than five units?

Answer: These probabilities are evaluated by

$$\begin{aligned} E_{0.1}([0, 2]) &= 1 - e^{-0.2} = 0.181269, \\ E_{0.1}([1, 2]) &= e^{-0.1} - e^{-0.2} = 0.08611, \\ E_{0.1}([5, \infty)) &= e^{-0.5} = 0.60653. \end{aligned}$$

1.6.5 Erlang distribution

Another important class of gamma distributions is that of Erlang distributions defined as follows.

Definition 1.6.24. For $\lambda > 0$ and $n \in \mathbb{N}$, let $E_{\lambda,n} := \Gamma_{\lambda^{-1},n}$. This probability measure is called the **Erlang distribution** with parameters λ and n .

Remark 1.6.25. The density $p_{\lambda,n}$ of the Erlang distribution is (see also Figure 1.15)

$$p_{\lambda,n}(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ \frac{\lambda^n}{(n-1)!} x^{n-1} e^{-\lambda x} & \text{if } x > 0. \end{cases}$$

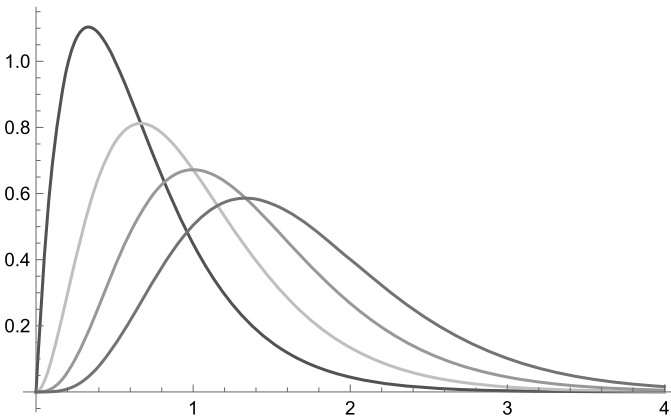


Figure 1.15: The densities of the Erlang distribution with $\lambda = 3$ and $n = 2, 3, 4$, and 5 .

Of course, $E_{\lambda,1} = E_\lambda$. Thus, the Erlang distribution may be viewed as generalized exponential distribution.

An important property of the Erlang distribution is as follows.

Proposition 1.6.26. *If $t > 0$, then*

$$E_{\lambda,n}([t, \infty)) = \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t}.$$

Proof. We have to show that for $t > 0$,

$$\int_t^\infty p_{\lambda,n}(x) dx = \frac{\lambda^n}{(n-1)!} \int_t^\infty x^{n-1} e^{-\lambda x} dx = \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t}. \quad (1.58)$$

This is done by induction over n .

If $n = 1$, then eq. (1.58) is valid due to

$$\int_t^{\infty} p_{\lambda,1}(x) \, dx = \int_t^{\infty} \lambda e^{-\lambda x} \, dx = e^{-\lambda t}.$$

Suppose now eq. (1.58) is proven for some $n \geq 1$. Next, we have to show that it is also valid for $n + 1$. Thus, we know

$$\frac{\lambda^n}{(n-1)!} \int_t^{\infty} x^{n-1} e^{-\lambda x} \, dx = \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t} \quad (1.59)$$

and want

$$\frac{\lambda^{n+1}}{n!} \int_t^{\infty} x^n e^{-\lambda x} \, dx = \sum_{j=0}^n \frac{(\lambda t)^j}{j!} e^{-\lambda t}. \quad (1.60)$$

Let us integrate the integral in eq. (1.60) by parts as follows. Set $u := x^n$, hence $u' = n x^{n-1}$, and $v' = e^{-\lambda x}$, thus $v = -\lambda^{-1} e^{-\lambda x}$. Doing so and using eq. (1.59), the left-hand side of eq. (1.60) becomes

$$\begin{aligned} \frac{\lambda^{n+1}}{n!} \int_t^{\infty} x^n e^{-\lambda x} \, dx &= \left[-\frac{\lambda^n}{n!} x^n e^{-\lambda x} \right]_t^{\infty} + \frac{\lambda^n}{(n-1)!} \int_t^{\infty} x^{n-1} e^{-\lambda x} \, dx \\ &= \frac{(\lambda t)^n}{n!} e^{-\lambda t} + \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t} = \sum_{j=0}^n \frac{(\lambda t)^j}{j!} e^{-\lambda t}. \end{aligned}$$

This proves eq. (1.60) and, consequently, eq. (1.58) is valid for all $n \geq 1$. \square

1.6.6 Chi-squared distribution

Another important class of gamma distributions is that of χ^2 -distributions. These probability measures play a crucial role in Mathematical Statistics (cf. Chapter 8).

Definition 1.6.27. For $n \geq 1$, let

$$\chi_n^2 := \Gamma_{2,n/2}.$$

This probability measure is called the χ^2 -distribution with n degrees of freedom.

Remark 1.6.28. At the moment, the integer $n \geq 1$ in Definition 1.6.27 is only a parameter. The term “degree of freedom” will become clear when we apply the χ^2 -distribution to statistical problems.

Remark 1.6.29. The density p of a χ_n^2 -distribution is given by (compare Figure 1.16)

$$p(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ \frac{x^{n/2-1} e^{-x/2}}{2^{n/2} \Gamma(n/2)} & \text{if } x > 0, \end{cases}$$

i. e., if $0 \leq a < b$, then

$$\chi_n^2([a, b]) = \frac{1}{2^{n/2} \Gamma(n/2)} \int_a^b x^{n/2-1} e^{-x/2} dx .$$

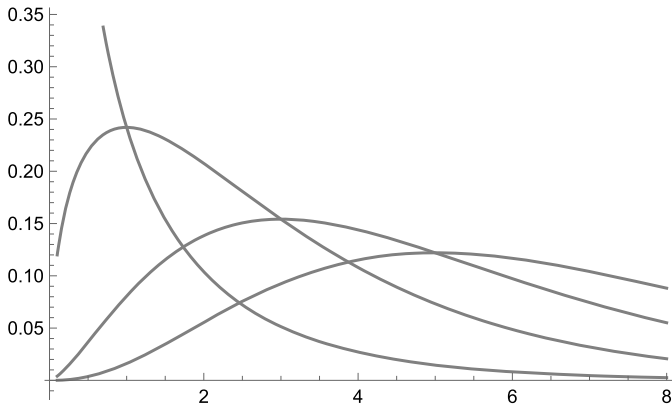


Figure 1.16: Density functions of χ_n^2 -distributions, $n = 1, 3, 5$, and $n = 7$.

Summary: For each $a, \beta > 0$, the probability measure $\Gamma_{a,\beta}$ is given by

$$\Gamma_{a,\beta}([a, b]) = \frac{1}{a^\beta \Gamma(\beta)} \int_a^b x^{\beta-1} e^{-x/a} dx, \quad 0 \leq a < b < \infty .$$

Of special interest are the exponential distribution $E_\lambda = \Gamma_{\lambda^{-1}, 1}$, the Erlang distribution defined by $E_{\lambda, n} = \Gamma_{\lambda^{-1}, n}$, and the chi-squared distribution $\chi_n^2 = \Gamma_{2, n/2}$. Here $\lambda > 0$ and $n \in \mathbb{N}$.

1.6.7 Beta distribution

Tightly connected with the gamma function is Euler's **beta function** B . It maps $(0, \infty)^2$ to \mathbb{R} and is defined by

$$B(x, y) := \int_0^1 s^{x-1} (1-s)^{y-1} ds, \quad x, y > 0. \quad (1.61)$$

The link between the gamma and beta functions is the following important identity:

$$B(x, y) = \frac{\Gamma(x) \cdot \Gamma(y)}{\Gamma(x+y)}, \quad x, y > 0. \quad (1.62)$$

For a proof of eq. (1.62), we refer to Problem 1.34.

Let us summarize further properties of the beta function. They are either easy to prove or follow via eq. (1.62) by the properties of the gamma function.

Proposition 1.6.30.

1. The beta function is continuous on $(0, \infty) \times (0, \infty)$ with values in $(0, \infty)$.
2. For $x, y > 0$, one has $B(x, y) = B(y, x)$.
3. If $x, y > 0$, then

$$B(x+1, y) = \frac{x}{x+y} B(x, y). \quad (1.63)$$

4. For $x > 0$, one has $B(x, 1) = 1/x$.
5. If $m, n \geq 1$ are integers, then

$$B(m, n) = \frac{(m-1)! (n-1)!}{(m+n-1)!}.$$

6. It holds

$$B\left(\frac{1}{2}, \frac{1}{2}\right) = \pi. \quad (1.64)$$

Definition 1.6.31. Let $\alpha, \beta > 0$. The probability measure $\mathcal{B}_{\alpha, \beta}$ defined by

$$\mathcal{B}_{\alpha, \beta}([a, b]) := \frac{1}{B(\alpha, \beta)} \int_a^b x^{\alpha-1} (1-x)^{\beta-1} dx, \quad 0 \leq a < b \leq 1,$$

is called the **beta distribution** with parameters α and β .

That is, the density function $q_{\alpha, \beta}$ of $\mathcal{B}_{\alpha, \beta}$ is given by

$$q_{\alpha, \beta}(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} & \text{if } 0 < x < 1, \\ 0 & \text{otherwise.} \end{cases}$$

Compare Figure 1.17 for the densities of beta distributions with certain pairs of parameters $\alpha, \beta > 0$. Note that the densities are bounded provided that $\alpha, \beta \geq 1$ and unbounded whenever one of the parameters is less than 1.

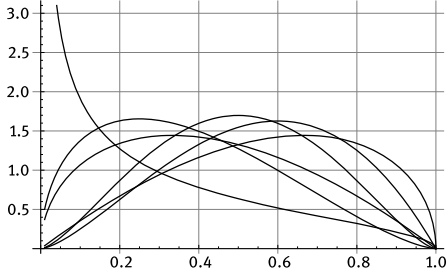


Figure 1.17: Density functions of the beta distribution with parameters $(0.5, 1.5)$, $(1.5, 2.5)$, $(2.5, 2)$, $(1.5, 2)$, $(2, 1.5)$, and $(2.5, 2.5)$.

Remark 1.6.32. It is easy to see that $q_{\alpha,\beta}$ is a density function. Of course, it is nonnegative and, moreover, we have

$$\int_{-\infty}^{\infty} q_{\alpha,\beta}(x) dx = \frac{1}{B(\alpha,\beta)} \int_0^1 x^{\alpha-1}(1-x)^{\beta-1} dx = \frac{B(\alpha,\beta)}{B(\alpha,\beta)} = 1.$$

Furthermore, since $q_{\alpha,\beta}(x) = 0$ if $x \notin [0, 1]$, the probability measure $\mathcal{B}_{\alpha,\beta}$ is concentrated on $[0, 1]$, that is, $\mathcal{B}_{\alpha,\beta}([0, 1]) = 1$ or, equivalently, $\mathcal{B}_{\alpha,\beta}(\mathbb{R} \setminus [0, 1]) = 0$.

Example 1.6.33. Choose independently n numbers x_1, \dots, x_n in $[0, 1]$ according to the uniform distribution. Ordering these numbers by their size, we get $x_1^* \leq \dots \leq x_n^*$. In Example 3.7.11, we will show that the k th largest number x_k^* is $\mathcal{B}_{k,n-k+1}$ -distributed. In other words, if $0 \leq a < b \leq 1$, then

$$\mathbb{P}\{a \leq x_k^* \leq b\} = \mathcal{B}_{k,n-k+1}([a, b]) = \frac{n!}{(k-1)!(n-k)!} \int_a^b x^{k-1}(1-x)^{n-k} dx.$$

Among the beta distributions $\mathcal{B}_{\alpha,\beta}$, the one with $\alpha = \beta = \frac{1}{2}$ is of special interest. It plays an important role in the investigation of the symmetric random walk as well as of the Brownian motion (see Proposition 5.5.19).

Proposition 1.6.34. The density q of the beta distribution $\mathcal{B}_{1/2,1/2}$ is given by

$$q(x) = \begin{cases} \frac{1}{\pi} \frac{1}{\sqrt{x(1-x)}} & \text{if } 0 < x < 1, \\ 0 & \text{otherwise.} \end{cases} \quad (1.65)$$

Furthermore, if $0 \leq a < b \leq 1$, then it follows that

$$\mathcal{B}_{1/2,1/2}([a, b]) = \frac{2}{\pi} [\arcsin(\sqrt{b}) - \arcsin(\sqrt{a})]. \quad (1.66)$$

Proof. Using eq. (1.64), representation (1.65) of the density easily follows from

$$q(x) = \frac{1}{B(1/2, 1/2)} x^{\frac{1}{2}-1} (1-x)^{\frac{1}{2}-1} = \frac{1}{\pi} \frac{1}{\sqrt{x(1-x)}}, \quad 0 < x < 1.$$

Assertion (1.66) is a consequence of

$$\frac{2}{\pi} \frac{d}{dx} \arcsin(\sqrt{x}) = \frac{2}{\pi} \frac{1}{\sqrt{1-(\sqrt{x})^2}} \cdot \frac{1}{2\sqrt{x}} = \frac{1}{\pi} \frac{1}{\sqrt{x(1-x)}}, \quad 0 < x < 1,$$

in view of the fundamental theorem of Calculus. See Figure 1.18 for the graph of the density q . \square

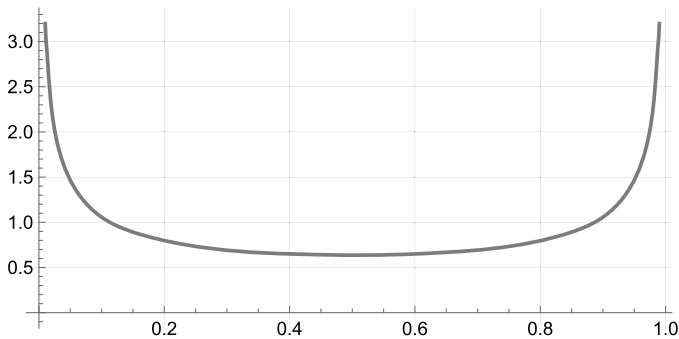


Figure 1.18: The density of the arcsine distribution.

Definition 1.6.35. The probability measure $\mathcal{B}_{1/2,1/2}$ is called the **arcsine distribution**. Its density is given by eq. (1.65), and for any $0 < a \leq 1$ it follows that

$$\mathcal{B}_{1/2,1/2}([0, a]) = \frac{2}{\pi} \arcsin(\sqrt{a}). \quad (1.67)$$

1.6.8 Cauchy distribution

We start with the following statement.

Proposition 1.6.36. *The function p defined by*

$$p(x) = \frac{1}{\pi} \cdot \frac{1}{1+x^2}, \quad x \in \mathbb{R}, \quad (1.68)$$

is a probability density.

Proof. Of course, $p(x) > 0$ for $x \in \mathbb{R}$. Let us now investigate $\int_{-\infty}^{\infty} p(x) dx$. Because of $\lim_{b \rightarrow \infty} \arctan(b) = \pi/2$ and $\lim_{a \rightarrow -\infty} \arctan(a) = -\pi/2$, it follows that

$$\int_{-\infty}^{\infty} p(x) dx = \frac{1}{\pi} \lim_{a \rightarrow -\infty} \lim_{b \rightarrow \infty} \int_a^b \frac{1}{1+x^2} dx = \frac{1}{\pi} \lim_{a \rightarrow -\infty} \lim_{b \rightarrow \infty} [\arctan x]_a^b = 1.$$

Thus, as asserted, p is a probability density. □

Definition 1.6.37. The probability measure \mathbb{P} with density p from Eq. (1.68) is called the **Cauchy distribution**. It is characterized by

$$\mathbb{P}([a, b]) = \frac{1}{\pi} \int_a^b \frac{1}{1+x^2} dx = \frac{1}{\pi} [\arctan(b) - \arctan(a)], \quad -\infty < a < b < \infty.$$

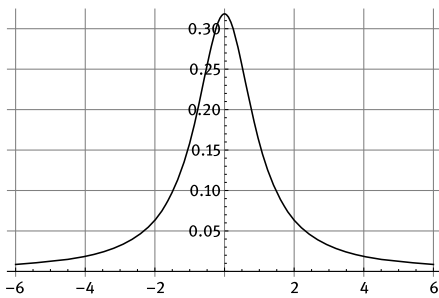


Figure 1.19: The density function of the Cauchy distribution.

Remark 1.6.38. Comparing the density function of the Cauchy distribution in Figure 1.19 with that of the standard normal distribution in Figure 1.11, at a first glance both functions look very similar. But, in fact they differ significantly. While the density of the standard normal distribution tends to 0 exponentially as $x \rightarrow \infty$, the convergence of the density of the Cauchy distribution is only of quadratic order. That is, events lying far away from zero possess an essentially greater probability for the Cauchy distribution as they do in the case of the normal. This property of the Cauchy distribution is sometimes expressed by saying that it is a distribution possessing “heavy tails.”

1.7 Distribution function

In this section we always assume that the sample space is \mathbb{R} , even if the random experiment has only finitely or countably infinitely many different outcomes. For example,

rolling a die once is modeled by $(\mathbb{R}, \mathcal{P}(\mathbb{R}), \mathbb{P})$, where $\mathbb{P}(\{k\}) = 1/6$, $k = 1, \dots, 6$, and $\mathbb{P}(\{x\}) = 0$ whenever $x \notin \{1, \dots, 6\}$.

Thus, let \mathbb{P} be a probability measure either defined on $\mathcal{B}(\mathbb{R})$ (continuous case) or on $\mathcal{P}(\mathbb{R})$ (discrete case).

Definition 1.7.1. The function $F : \mathbb{R} \rightarrow [0, 1]$ defined by

$$F(t) := \mathbb{P}((-\infty, t]), \quad t \in \mathbb{R}, \quad (1.69)$$

is called the **(cumulative) distribution function** of \mathbb{P} .

Remark 1.7.2. To shorten the notation, we will mostly call F in (1.69) the “distribution function” instead of, as is often done in the literature, the “cumulative distribution function,” abbreviated CDF.

Remark 1.7.3. If \mathbb{P} is discrete, that is, $\mathbb{P}(D) = 1$ for some $D = \{x_1, x_2, \dots\}$, then its distribution function can be evaluated by

$$F(t) = \sum_{x_j \leq t} \mathbb{P}(\{x_j\}) = \sum_{x_j \leq t} p_j,$$

where $p_j = \mathbb{P}(\{x_j\})$, while for a continuous \mathbb{P} with probability density p ,

$$F(t) = \int_{-\infty}^t p(x) \, dx.$$

Example 1.7.4. Let \mathbb{P} be the uniform distribution on $\{1, \dots, 6\}$. Then

$$F(t) = \begin{cases} 0 & \text{if } t < 1, \\ \frac{k}{6} & \text{if } k \leq t < k+1, \quad k \in \{1, \dots, 5\}, \\ 1 & \text{if } t \geq 6. \end{cases}$$

(See Figure 1.20.)

Example 1.7.5. The distribution function of the binomial distribution $B_{n,p}$ is given by

$$F(t) = \sum_{0 \leq k \leq t} \binom{n}{k} p^k (1-p)^{n-k}, \quad 0 \leq t < \infty,$$

and $F(t) = 0$ if $t < 0$. (See Figure 1.21.)

Example 1.7.6. The distribution function of the exponential distribution E_λ is (see Figure 1.22 for an example)

$$F(t) = \begin{cases} 0 & \text{if } t < 0, \\ 1 - e^{-\lambda t} & \text{if } t \geq 0. \end{cases}$$

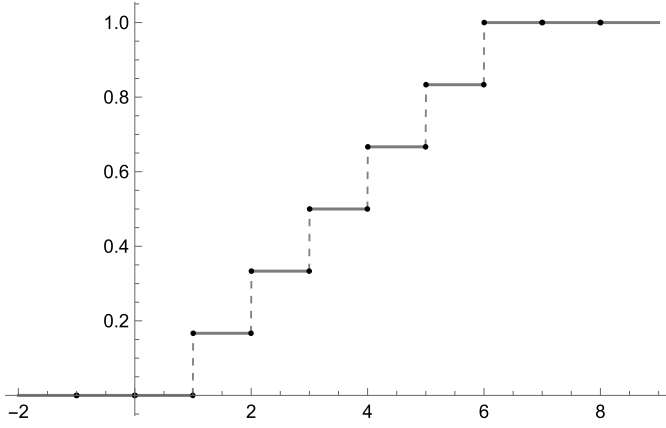


Figure 1.20: Distribution function of the uniform distribution on $\{1, \dots, 6\}$.

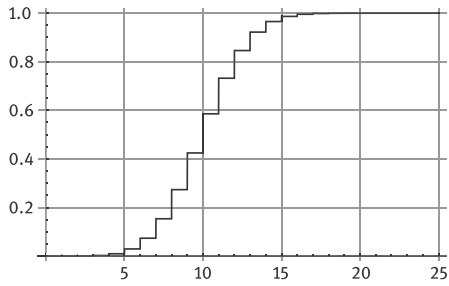


Figure 1.21: Distribution function of the binomial distribution $B_{25,0.4}$.

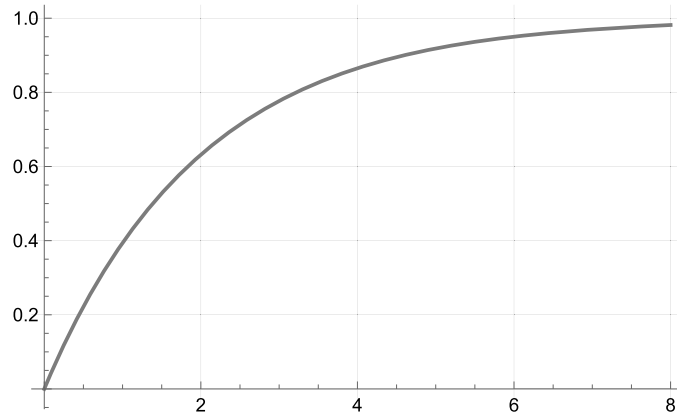


Figure 1.22: Distribution function of $E_{0.5}$.

Example 1.7.7. How does the distribution function F of the Erlang distribution $E_{\lambda,n}$ look like? In view of Proposition 1.6.26, it follows that $F(t) = 0$ if $t < 0$ and

$$F(t) = 1 - \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t}, \quad 0 \leq t < \infty.$$

In the case $n = 1$, we rediscover the function stated in Example 1.7.6. Compare Figure 1.23 for certain distribution functions of the Erlang distribution.

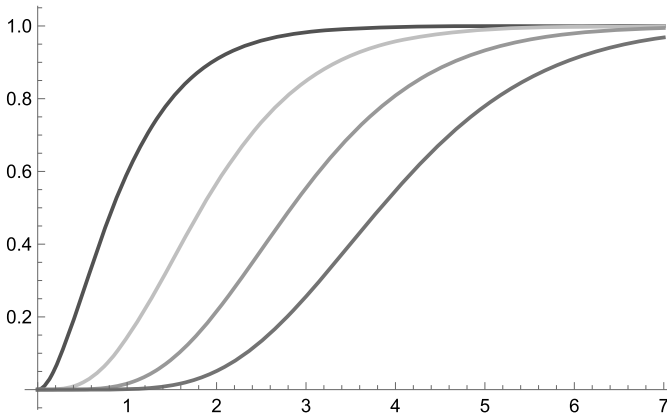


Figure 1.23: Distribution functions of $E_{2,n}$, where from the left to the right $n = 2, 4, 6, 8$.

Example 1.7.8. The distribution function of the standard normal distribution is denoted¹⁴ by Φ , therefore also called **Gaussian Φ -function** (see Figure 1.24),

$$\Phi(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx, \quad t \in \mathbb{R}. \quad (1.70)$$

Remark 1.7.9. The Gaussian Φ -function is tightly related to the **Gaussian error function** defined by (compare Figure 1.25)

$$\operatorname{erf}(t) = \frac{2}{\sqrt{\pi}} \int_0^t e^{-x^2} dx, \quad t \in \mathbb{R}.$$

Observe that $\operatorname{erf}(-t) = -\operatorname{erf}(t)$.

¹⁴ Sometimes also denoted as “norm(·).”

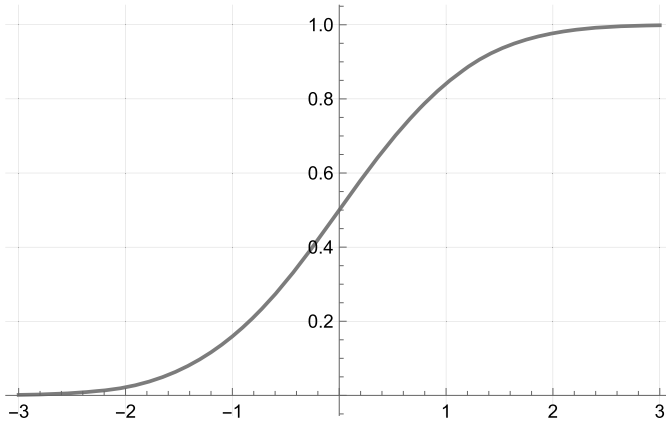


Figure 1.24: Distribution function of the standard normal distribution (Φ -function).

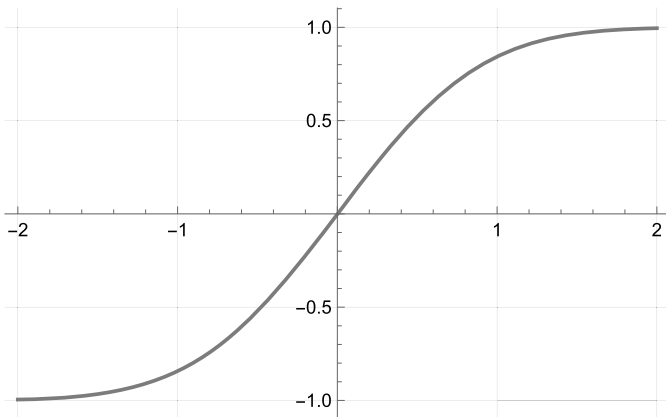


Figure 1.25: The Gaussian error function $t \mapsto \operatorname{erf}(t)$.

The link between the Φ and the error function is

$$\Phi(t) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{t}{\sqrt{2}} \right) \right] \quad \text{and} \quad \operatorname{erf}(t) = 2\Phi(\sqrt{2}t) - 1, \quad t \in \mathbb{R}. \quad (1.71)$$

Example 1.7.10. Let \mathbb{P} be the uniform distribution on the interval $[\alpha, \beta]$. Then its distribution function (cf. Figure 1.26) is

$$F(t) = \begin{cases} 0 & \text{if } t < \alpha, \\ \frac{t-\alpha}{\beta-\alpha} & \text{if } \alpha \leq t \leq \beta, \\ 1 & \text{if } t > \beta. \end{cases}$$

In particular, for the uniform distribution on $[0, 1]$, one obtains

$$F(t) = \begin{cases} 0 & \text{if } t < 0, \\ t & \text{if } 0 \leq t \leq 1, \\ 1 & \text{if } t > 1. \end{cases}$$

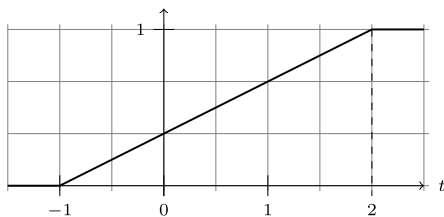


Figure 1.26: Distribution function of the uniform distribution on $[-1, 2]$.

Example 1.7.11. In view of eq. (1.67), the distribution function of the arcsine distribution (see Definition 1.6.35) is given by (compare Figure 1.27)

$$F(t) = \begin{cases} 0 & \text{if } -\infty < t < 0, \\ \frac{2}{\pi} \arcsin(\sqrt{t}) & \text{if } 0 \leq t \leq 1, \\ 1 & \text{if } 1 < t < \infty. \end{cases}$$

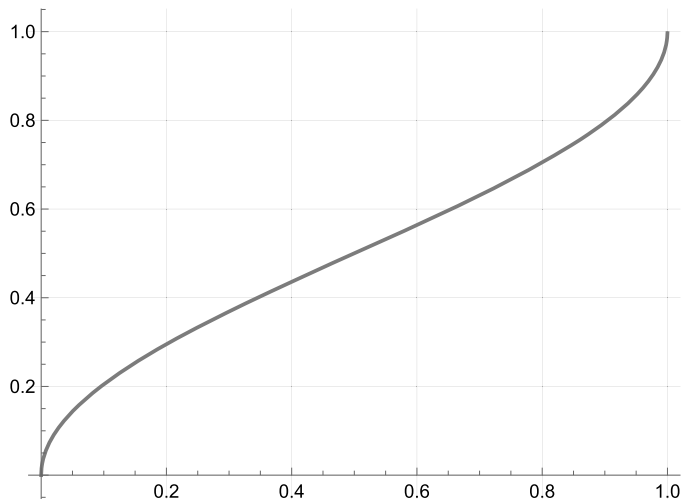


Figure 1.27: The distribution function of the arcsine distribution.

Example 1.7.12. The distribution function of the Cauchy distribution is (compare Figure 1.28)

$$F(t) = \frac{1}{\pi} \int_{-\infty}^t \frac{1}{1+x^2} dx = \frac{\arctan x}{\pi} + \frac{1}{2}, \quad t \in \mathbb{R}.$$

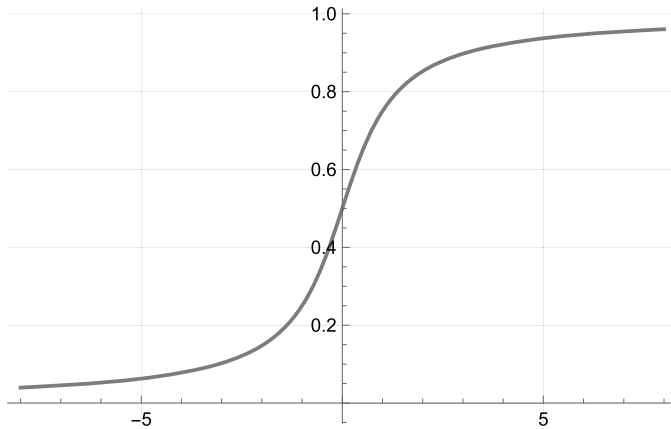


Figure 1.28: The distribution function of the Cauchy distribution.

The next proposition lists the main properties of distribution functions.

Proposition 1.7.13. *Let F be the distribution function of a probability measure \mathbb{P} on \mathbb{R} , discrete or continuous. Then F possesses the following properties:*

- (1) *Function F is nondecreasing.*
- (2) *It holds*

$$F(-\infty) = \lim_{t \rightarrow -\infty} F(t) = 0 \quad \text{and} \quad F(\infty) = \lim_{t \rightarrow \infty} F(t) = 1.$$

- (3) *Function F is continuous from the right.*

Proof. Suppose $s < t$. This implies $(-\infty, s] \subset (-\infty, t]$, hence, since \mathbb{P} is monotone, we obtain

$$F(s) = \mathbb{P}((-\infty, s]) \leq \mathbb{P}((-\infty, t]) = F(t).$$

Thus F is nondecreasing.

Take any sequence $(t_n)_{n \geq 1}$ that decreases monotonically to $-\infty$. Set $A_n := (-\infty, t_n]$. Then $A_1 \supseteq A_2 \supseteq \dots$, as well as $\bigcap_{n=1}^{\infty} A_n = \emptyset$. Since \mathbb{P} is continuous from above, it fol-

lows that

$$\lim_{n \rightarrow \infty} F(t_n) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \mathbb{P}(\emptyset) = 0.$$

This being true for any sequence $(t_n)_{n \geq 1}$ tending to $-\infty$ implies $F(-\infty) = 0$.

The proof of $F(\infty) = 1$ is very similar. In this case, let $(t_n)_{n \geq 1}$ be a sequence which increases monotonically to ∞ . If as before $A_n := (-\infty, t_n]$, this time we get $A_1 \subseteq A_2 \subseteq \dots$ and $\bigcup_{n=1}^{\infty} A_n = \mathbb{R}$. By the continuity of \mathbb{P} from below, now we obtain

$$\lim_{n \rightarrow \infty} F(t_n) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \mathbb{P}(\mathbb{R}) = 1.$$

Again, since the t_n s were arbitrary, $F(\infty) = 1$.

Thus it remains to prove that F is continuous from the right. To do this, we take $t \in \mathbb{R}$ and a decreasing sequence $(t_n)_{n \geq 1}$ tending to t . We have to show that if $n \rightarrow \infty$, then $F(t_n) \rightarrow F(t)$.

As before set $A_n := (-\infty, t_n]$. Again $A_1 \supseteq A_2 \supseteq \dots$, but now $\bigcap_{n=1}^{\infty} A_n = (-\infty, t]$. Another application of the continuity from above yields

$$F(t) = \mathbb{P}((-\infty, t]) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \lim_{n \rightarrow \infty} F(t_n).$$

This is valid for each $t \in \mathbb{R}$, hence F is continuous from the right. \square

Properties (1), (2), and (3) in Proposition 1.7.13 characterize distribution functions. More precisely, the following result is true. Its proof is based on an extension theorem in Measure Theory (cf. [Bau01]). Therefore, we can show here only its main ideas.

Proposition 1.7.14. *Let $F : \mathbb{R} \rightarrow \mathbb{R}$ be an arbitrary function possessing the properties stated in Proposition 1.7.13. Then there exists a unique probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R})$ such that*

$$F(t) = \mathbb{P}((-\infty, t]), \quad t \in \mathbb{R}.$$

Idea of the proof: If $a < b$, set

$$\mathbb{P}_0((a, b]) := F(b) - F(a).$$

In this way, we get a mapping \mathbb{P}_0 defined on the collection of all half-open intervals $\{(a, b] : a < b\}$. The key point is to verify that \mathbb{P}_0 can be uniquely extended to a probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R})$. One way to do this is to introduce a so-called outer measure \mathbb{P}^* defined on $\mathcal{P}(\mathbb{R})$ by

$$\mathbb{P}^*(B) := \inf \left\{ \sum_{i=1}^{\infty} \mathbb{P}_0((a_i, b_i]) : B \subseteq \bigcup_{i=1}^{\infty} (a_i, b_i] \right\}.$$

Generally, this outer measure is not σ -additive. Therefore, one restricts \mathbb{P}^* to $\mathcal{B}(\mathbb{R})$. If \mathbb{P} denotes this restriction, the most difficult part of the proof is to verify that \mathbb{P} is σ -additive. After this has been done, by the construction, \mathbb{P} is the probability measure possessing distribution function F .

The uniqueness of \mathbb{P} follows by a general uniqueness theorem for probability measures (see Theorem 5.7 in [Sch17] for a proof) asserting the following:

Let \mathbb{P}_1 and \mathbb{P}_2 be two probability measures on (Ω, \mathcal{A}) and let $\mathcal{E} \subseteq \mathcal{A}$ be a collection of events closed under taking intersections and generating \mathcal{A} . If $\mathbb{P}_1(E) = \mathbb{P}_2(E)$ for all $E \in \mathcal{E}$, then $\mathbb{P}_1 = \mathbb{P}_2$. In our case $\mathcal{E} = \{(-\infty, t] : t \in \mathbb{R}\}$ and $\mathcal{A} = \mathcal{B}(\mathbb{R})$.

Conclusion. If the outcomes of a random experiment are real numbers, then this experiment can also be described by a function $F : \mathbb{R} \rightarrow \mathbb{R}$ possessing the properties in Proposition 1.7.13. Then $F(t)$ is the probability to observe a result that is less than or equal to t .

Let us state further properties of distribution functions.

Proposition 1.7.15. *If F is the distribution function of a probability measure \mathbb{P} , then for all $a < b$,*

$$F(b) - F(a) = \mathbb{P}((a, b]).$$

Proof. Observing that $(-\infty, a] \subseteq (-\infty, b]$, this is an immediate consequence of

$$F(b) - F(a) = \mathbb{P}((-\infty, b]) - \mathbb{P}((-\infty, a]) = \mathbb{P}((-\infty, b] \setminus (-\infty, a]) = \mathbb{P}((a, b]). \quad \square$$

Since F is nondecreasing and bounded, for each $t \in \mathbb{R}$ the left-hand limit

$$F(t-0) := \lim_{\substack{s \rightarrow t \\ s < t}} F(s)$$

exists and, moreover, $F(t-0) \leq F(t)$. Furthermore, by the right continuity of F , one has $F(t-0) = F(t)$ if and only if F is continuous at the point t .

If this is not so, then $h = F(t) - F(t-0) > 0$, that is, F possesses at $t \in \mathbb{R}$ a **jump** of height $h > 0$. This height is directly connected with the value of $\mathbb{P}(\{t\})$.

Proposition 1.7.16. *The distribution function F of a probability measure \mathbb{P} has a jump of height $h \geq 0$ at $t \in \mathbb{R}$ if and only if $\mathbb{P}(\{t\}) = h$.*

Proof. Let $(t_n)_{n \geq 1}$ be a sequence of real numbers increasing monotonically to t . Then, using that \mathbb{P} is continuous from above, it follows that

$$h = F(t) - F(t-0) = \lim_{n \rightarrow \infty} [F(t) - F(t_n)] = \lim_{n \rightarrow \infty} \mathbb{P}((t_n, t]) = \mathbb{P}(\{t\}).$$

Observe that $\bigcap_{n=1}^{\infty} (t_n, t] = \{t\}$. This proves the assertion. See Figure 1.29 for a visualization of the result. \square

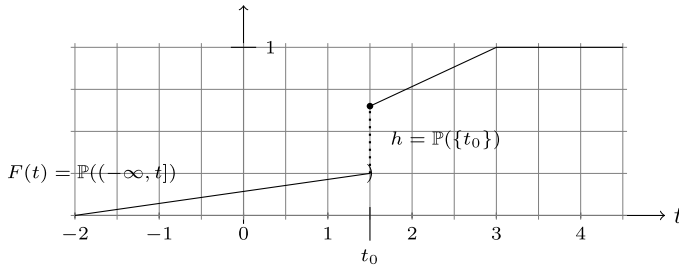


Figure 1.29: The height h of the jump of F at t_0 coincides with the probability $\mathbb{P}(\{t_0\})$.

Corollary 1.7.17. *The function F is continuous at $t \in \mathbb{R}$ if and only if $\mathbb{P}(\{t\}) = 0$.*

Proof. Since F is nondecreasing and continuous from the right, it is continuous at some point $t \in \mathbb{R}$ if and only if $F(t - 0) = F(t)$. But, in view of Proposition 1.7.16, this happens if and only if $\mathbb{P}(\{t\}) = 0$. \square

Example 1.7.18. Suppose the function F is defined by

$$F(t) = \begin{cases} 0 & \text{if } t < -1, \\ 1/3 & \text{if } -1 \leq t < 0, \\ 1/2 & \text{if } 0 \leq t < 1, \\ 2/3 & \text{if } 1 \leq t < 2, \\ 1 & \text{if } t \geq 2. \end{cases}$$

Then F fulfills the assumptions of Proposition 1.7.14. Hence there is a probability measure \mathbb{P} with $F(t) = \mathbb{P}((-\infty, t])$. What does \mathbb{P} look like?

Answer: The function F has jumps at $-1, 0, 1,$ and 2 with heights $1/3, 1/6, 1/6,$ and $1/3,$ respectively. Therefore,

$$\mathbb{P}(\{-1\}) = 1/3, \quad \mathbb{P}(\{0\}) = 1/6, \quad \mathbb{P}(\{1\}) = 1/6, \quad \text{and} \quad \mathbb{P}(\{2\}) = 1/3,$$

hence \mathbb{P} is the discrete probability measure concentrated on $D = \{-1, 0, 1, 2\}$ with $\mathbb{P}(\{t\})$, $t \in D$, given above.

Suppose now that \mathbb{P} is continuous with density function p . Recall that then

$$F(t) = \mathbb{P}((-\infty, t]) = \int_{-\infty}^t p(x) \, dx, \quad t \in \mathbb{R}. \quad (1.72)$$

In particular, since F is a function of the upper bound in an integral, it is continuous.

Next we investigate the question whether we may evaluate the density p knowing the distribution function F .

Proposition 1.7.19. *Suppose p is continuous at some $t \in \mathbb{R}$. Then F is differentiable at t with*

$$F'(t) = \frac{d}{dt}F(t) = p(t).$$

Proof. This follows immediately by an application of the fundamental theorem of calculus to representation (1.72) of F . \square

Remark 1.7.20. Let F be the distribution function of a probability measure \mathbb{P} . If F is continuous, then $\mathbb{P}(\{t\}) = 0$ for all $t \in \mathbb{R}$. But does this also imply that \mathbb{P} is continuous, that is, that \mathbb{P} has a density? The answer is negative. There exist probability measures \mathbb{P} on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ with a continuous distribution function but without possessing a density. Such probability measures are called **singularly continuous**.

To get an impression of how such probability measures look, let us shortly sketch the construction of an example. Let C be the Cantor set introduced in Example 1.6.6. The basic idea is to transfer the uniform distribution on $[0, 1]$ to a probability measure \mathbb{P} with $\mathbb{P}(C) = 1$. The transformation is done by the function f defined as follows. If $x \in [0, 1]$ is represented as $x = \sum_{k=1}^{\infty} \frac{x_k}{2^k}$ with $x_k \in \{0, 1\}$, then $f(x) = \sum_{k=1}^{\infty} \frac{2x_k}{3^k}$. Note that f maps $[0, 1]$ into C . If $\tilde{\mathbb{P}}$ denotes the uniform distribution on $[0, 1]$, define the probability measure \mathbb{P} by

$$\mathbb{P}(B) = \tilde{\mathbb{P}}\{x \in [0, 1] : f(x) \in B\}.$$

Then for all $t \in \mathbb{R}$, we have $\mathbb{P}(\{t\}) = 0$, but since $\mathbb{P}(C) = 1$, \mathbb{P} cannot have a density. Indeed, such a density should vanish outside C . But, as we saw, the probability of C with respect to the uniform distribution is zero. Hence the only possible density would be $p(t) = 0$, $t \in \mathbb{R}$. This contradiction shows that \mathbb{P} is not continuous in our sense.

Assuming a little bit more than the continuity of F , the corresponding probability measure possesses a density (cf. [Coh13]).

Proposition 1.7.21. *Let F be the distribution function of a probability measure \mathbb{P} . If F is continuous and continuously differentiable with the exception of at most finitely many points, then \mathbb{P} is continuous. That is, there is a density function p such that*

$$F(t) = \mathbb{P}((-\infty, t]) = \int_{-\infty}^t p(x) dx, \quad t \in \mathbb{R}.$$

Remark 1.7.22. Proposition 1.7.19 implies $p(t) = F'(t)$ for those t where $F'(t)$ exists. If F is not differentiable at t , define $p(t)$ arbitrarily, for example, $p(t) = 0$.

Example 1.7.23. For some $\alpha, \beta > 0$, define F by

$$F(t) = \begin{cases} 0 & \text{if } t \leq 0, \\ 1 - e^{-\alpha t^\beta} & \text{if } t > 0. \end{cases}$$

It is easy to see that this function satisfies the conditions of Proposition 1.7.13. Moreover, it is continuous and continuously differentiable on $\mathbb{R} \setminus \{0\}$. By Proposition 1.7.21, the corresponding probability measure \mathbb{P} is continuous and, since

$$F'(t) = \begin{cases} 0 & \text{if } t < 0, \\ \alpha \beta t^{\beta-1} e^{-\alpha t^\beta} & \text{if } t > 0, \end{cases}$$

a suitable density function is $p(t) = F'(t)$, $t \neq 0$, and $p(0) = 0$.

Example 1.7.24. For some $\alpha > 0$, let

$$F_\alpha(t) = \begin{cases} 0 & \text{if } t \leq 1, \\ 1 - t^{-\alpha} & \text{if } t > 1. \end{cases}$$

Then F_α is continuous, nondecreasing with $F_\alpha(-\infty) = 0$ and $F_\alpha(\infty) = 1$. Hence, it is a distribution function of a probability measure \mathbb{P}_α . Moreover, at each point $t \neq 1$, the function F_α is continuously differentiable with derivative

$$p_\alpha(t) = \begin{cases} 0 & \text{if } t < 1, \\ \alpha t^{-\alpha-1} & \text{if } t > 1. \end{cases}$$

So we see that p_α is a density of \mathbb{P}_α where we may define, for example, $p_\alpha(1) = 0$.

How does one get $\mathbb{P}_\alpha([a, b])$ for some $1 \leq a < b < \infty$? There is no need to evaluate the integral $\int_a^b p_\alpha(x) dx$. The much easier way is to use

$$\mathbb{P}_\alpha([a, b]) = F_\alpha(b) - F_\alpha(a) = \frac{1}{a^\alpha} - \frac{1}{b^\alpha}.$$

Summary: Let \mathbb{P} be a probability measure (discrete or continuous) on the Borel sets of \mathbb{R} . Then its (cumulative) distribution function F is defined by

$$F(t) = \mathbb{P}\{x \in \mathbb{R} : x \leq t\}, \quad t \in \mathbb{R}.$$

Distribution functions are characterized by the three properties stated in Proposition 1.7.13. Two probability measures coincide if and only if they possess the same distribution function.

1.8 Multivariate continuous distributions

1.8.1 Multivariate density functions

In this section we suppose that $\Omega = \mathbb{R}^n$. A subset $Q \subset \mathbb{R}^n$ is called a (closed, n -dimensional) **box**¹⁵ provided that for some real numbers $a_i < b_i$, $1 \leq i \leq n$,

$$Q = \{(x_1, \dots, x_n) \in \mathbb{R}^n : a_i \leq x_i \leq b_i, 1 \leq i \leq n\}. \quad (1.73)$$

Definition 1.8.1. A Riemann integrable function $p : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be an n -dimensional **probability density function**, or simply n -dimensional **density function**, if $p(x) \geq 0$ for $x \in \mathbb{R}^n$ and, furthermore,

$$\int_{\mathbb{R}^n} p(x) \, dx := \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} p(x_1, \dots, x_n) \, dx_n \cdots dx_1 = 1.$$

Suppose a box Q is represented with certain $a_i < b_i$ as in eq. (1.73). Then we set

$$\mathbb{P}(Q) = \int_Q p(x) \, dx = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} p(x_1, \dots, x_n) \, dx_n \cdots dx_1. \quad (1.74)$$

In analogy to Definition 1.1.20, we introduce now the Borel σ -field $\mathcal{B}(\mathbb{R}^n)$.

Definition 1.8.2. Let \mathcal{C} be the collection of all boxes in \mathbb{R}^n . Then $\sigma(\mathcal{C}) := \mathcal{B}(\mathbb{R}^n)$ denotes the **Borel σ -field**. Recall that the existence of $\sigma(\mathcal{C})$ was proven in Proposition 1.1.16. In other words, $\mathcal{B}(\mathbb{R}^n)$ is the smallest σ -field containing all (closed) boxes in \mathbb{R}^n . Sets in $\mathcal{B}(\mathbb{R}^n)$ are called (n -dimensional) **Borel sets**.

Remark 1.8.3. As in the univariate case, there exist several other collections of subsets in \mathbb{R}^n generating $\mathcal{B}(\mathbb{R}^n)$. For example, one may choose the collection of open boxes or the sets which may be written as

$$(-\infty, t_1] \times \cdots \times (-\infty, t_n], \quad t_1, \dots, t_n \in \mathbb{R}.$$

With the previous notations, the following multivariate extension theorem is valid. Compare with Proposition 1.5.6 for the univariate case.

Proposition 1.8.4. *Let \mathbb{P} be defined on boxes by eq. (1.74). Then \mathbb{P} admits a unique extension to a probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R}^n)$.*

Definition 1.8.5. A probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R}^n)$ is called **continuous** provided that there exists a probability density $p : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\mathbb{P}(Q) = \int_Q p(x) \, dx$ for all boxes $Q \subseteq \mathbb{R}^n$. The function p is said to be the **density function**, or simply **density**, of \mathbb{P} .

¹⁵ Also called “hyperrectangle.”

Remark 1.8.6. It is easy to see that the validity of eq. (1.74) for all boxes is equivalent to the following. If $t_j \in \mathbb{R}$ and $B_{t_1, \dots, t_n} := (-\infty, t_1] \times \dots \times (-\infty, t_n]$, then

$$\mathbb{P}(B_{t_1, \dots, t_n}) = \int_{B_{t_1, \dots, t_n}} p(x) \, dx = \int_{-\infty}^{t_1} \dots \int_{-\infty}^{t_n} p(x_1, \dots, x_n) \, dx_n \dots dx_1. \quad (1.75)$$

Thus \mathbb{P} is continuous if and only if eq. (1.75) is satisfied for all $t_j \in \mathbb{R}$.

Let us first give an example of a multivariate probability density function.

Example 1.8.7. Consider $p : \mathbb{R}^3 \rightarrow \mathbb{R}$ defined by

$$p(x_1, x_2, x_3) = \begin{cases} 48 x_1 x_2 x_3 & \text{if } 0 \leq x_1 \leq x_2 \leq x_3 \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Of course, $p(x) \geq 0$ for $x \in \mathbb{R}^3$. Moreover,

$$\begin{aligned} \int_{\mathbb{R}^3} p(x) \, dx &= 48 \int_0^1 \int_0^{x_3} \int_0^{x_2} x_1 x_2 x_3 \, dx_1 dx_2 dx_3 \\ &= 48 \int_0^1 \int_0^{x_3} \frac{x_3 x_2^3}{2} \, dx_2 dx_3 = 48 \int_0^1 \frac{x_3^5}{8} \, dx_3 = 1, \end{aligned}$$

hence it is a density function on \mathbb{R}^3 . For example, if \mathbb{P} is the generated probability measure, then

$$\mathbb{P}([0, 1/2]^3) = 48 \int_0^{1/2} \int_0^{x_3} \int_0^{x_2} x_1 x_2 x_3 \, dx_1 dx_2 dx_3 = \frac{1}{2^6} = \frac{1}{64}.$$

There exists a quite general approach to construct multivariate density functions.

Proposition 1.8.8. Let p_1, \dots, p_n be (univariate) density functions. Define a function $p : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$p(x) = p_1(x_1) \cdots p_n(x_n), \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n. \quad (1.76)$$

Then p is a multivariate distribution density.

Proof. Of course, we have $p(x) \geq 0$ for all $x \in \mathbb{R}^n$. Moreover, an application of Proposition A.5.5 (note that all p_j are nonnegative) implies that

$$\int_{\mathbb{R}^n} p(x) \, dx = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} [p_1(x_1) \cdots p_n(x_n)] \, dx_n \cdots dx_1$$

$$= \left(\int_{-\infty}^{\infty} p_1(x_1) dx_1 \right) \cdots \left(\int_{-\infty}^{\infty} p_n(x_n) dx_n \right) = 1 \cdots 1 = 1.$$

Thus, p is a density as asserted. \square

Remark 1.8.9. Observe that not all multivariate densities p may be represented as a product of univariate ones as stated in eq. (1.76). As can be seen easily, the function p in Example 1.8.7 may not be written as a product of three univariate densities. We will investigate this and related problems more thoroughly in Section 1.9.3.

1.8.2 Multivariate uniform distribution

Our next aim is the introduction and investigation of a special multivariate distribution, the *uniform distribution* on a set K in \mathbb{R}^n . To do so, we remember how we defined the uniform distribution on an interval I in \mathbb{R} . Its density p is given by

$$p(s) = \begin{cases} \frac{1}{|I|} & \text{if } s \in I, \\ 0 & \text{if } s \notin I. \end{cases}$$

Here $|I|$ denotes the length of the interval I . Let now $K \subset \mathbb{R}^n$ be bounded. In order to introduce a similar density for the uniform distribution on K , the length of the underlying set has to be replaced by the n -dimensional volume, which we will denote by $\text{vol}_n(K)$. But how is this volume defined? To answer this question, let us first investigate a box Q represented as in eq. (1.73). It is immediately clear that its n -dimensional volume is evaluated by

$$\text{vol}_n(Q) = \prod_{i=1}^n (b_i - a_i).$$

If $n = 1$, then Q is an interval and its one-dimensional volume is nothing else as its length. For $n = 2$ the box Q is the rectangle $[a_1, b_1] \times [a_2, b_2]$ and

$$\text{vol}_2(Q) = (b_1 - a_1)(b_2 - a_2)$$

coincides with the area of Q . If $n = 3$, then $\text{vol}_3(Q)$ is the ordinary volume of bodies in \mathbb{R}^3 .

For arbitrary $K \subset \mathbb{R}^n$, the definition of its volume $\text{vol}_n(K)$ is more involved. Let us shortly sketch one way how this can be done. Setting

$$\text{vol}_n(K) := \inf \left\{ \sum_{j=1}^{\infty} \text{vol}_n(Q_j) : K \subseteq \bigcup_{j=1}^{\infty} Q_j, Q_j \text{ box} \right\}, \quad (1.77)$$

at least for Borel sets $K \subseteq \mathbb{R}^n$, a suitable volume is defined. Compare Figure 1.30 to get an impression how the two-dimensional volume of an ellipse is evaluated.

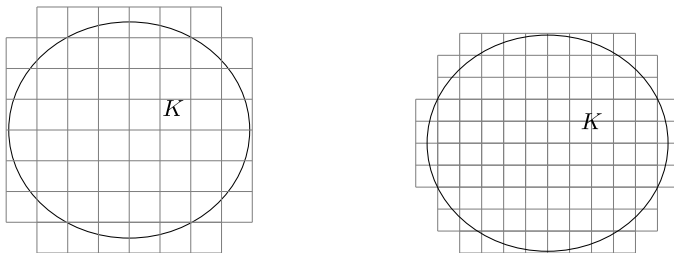


Figure 1.30: An ellipse K covered by boxes, here squares. The smaller the covering boxes, the better the approximation of $\text{vol}_2(K)$.

In the case of “ordinary” sets such as balls, ellipsoids, or similar bodies, this approach leads to the known values. The reason is the basic formula

$$\text{vol}_n(K) = \int \cdots \int_K \mathbf{1} \, dx_n \cdots dx_1 \quad (1.78)$$

valid for Borel sets $K \subseteq \mathbb{R}^n$. For example, if K is the square in \mathbb{R}^2 (see Figure 1.31) with corner points $(1, 0)$, $(0, 1)$, $(-1, 0)$, and $(0, -1)$, then

$$\begin{aligned} \text{vol}_2(K) &= \iint_K \mathbf{1} \, dx_2 dx_1 = \int_{-1}^0 \int_{-x_1-1}^{1+x_1} dx_2 dx_1 + \int_0^1 \int_{x_1-1}^{1-x_1} dx_2 dx_1 \\ &= 2 \int_{-1}^0 (x_1 + 1) dx_1 + 2 \int_0^1 (1 - x_1) dx_1 \\ &= 2 \left[\frac{x_1^2}{2} + x_1 \right]_{-1}^0 + 2 \left[x_1 - \frac{x_1^2}{2} \right]_0^1 = 2. \end{aligned}$$

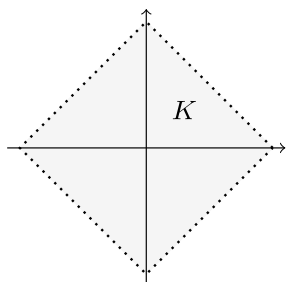


Figure 1.31: The square $K \subseteq \mathbb{R}^2$ with corner points $(1, 0)$, $(0, 1)$, $(-1, 0)$, and $(0, -1)$.

Remark 1.8.10. The n -dimensional volume shares several important properties:

1. It is invariant under shifts. That is, if

$$K + z = \{x + z : x \in K\}$$

denotes the shift of K by $z \in \mathbb{R}^n$, then it follows that

$$\text{vol}_n(K + z) = \text{vol}_n(K), \quad z \in \mathbb{R}^n.$$

2. Unitary transformations of a set do not change its volume. More precisely, if $U : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a unitary transformation (cf. (A.24) for the definition), then

$$\text{vol}_n(U(K)) = \text{vol}_n(K).$$

3. The n -dimensional volume is homogeneous of power n under dilation. Thus, given $\alpha > 0$, let $\alpha K = \{\alpha x : x \in K\}$ be the stretched (if $\alpha > 1$) or shrunk (if $\alpha < 1$) set K . Then it follows that

$$\text{vol}_n(\alpha K) = \alpha^n \text{vol}_n(K).$$

Example 1.8.11. Let $K_n(r)$ be the n -dimensional ball of radius $r > 0$, that is,

$$K_n(r) = \{x \in \mathbb{R}^n : |x| \leq r\} = \{(x_1, \dots, x_n) \in \mathbb{R}^n : x_1^2 + \dots + x_n^2 \leq r^2\}.$$

If

$$V_n(r) := \text{vol}_n(K_n(r)), \quad r > 0,$$

denotes the n -dimensional volume of this ball, property (3) in Remark 1.8.10 implies $V_n(r) = V_n \cdot r^n$, where $V_n = V_n(1)$. But for $K_n = K_n(1)$, eq. (1.78) gives

$$\begin{aligned} V_n &= \int_{K_n} \cdots \int \mathbf{1} \, dx_n \cdots dx_1 = \int_{-1}^1 \left[\int_{\{x_2^2 + \dots + x_n^2 \leq 1 - x_1^2\}} \cdots \int \mathbf{1} \, dx_n \cdots dx_2 \right] dx_1 \\ &= \int_{-1}^1 V_{n-1}(\sqrt{1 - x_1^2}) \, dx_1 = \int_{-1}^1 V_{n-1}(\sqrt{1 - s^2}) \, ds. \end{aligned}$$

Hence, due to $V_{n-1}(r) = r^{n-1} V_{n-1}(1) = r^{n-1} V_{n-1}$, we obtain

$$V_n = V_{n-1} \cdot \int_{-1}^1 (1 - s^2)^{(n-1)/2} \, ds = 2 V_{n-1} \cdot \int_0^1 (1 - s^2)^{(n-1)/2} \, ds.$$

The change of the variables $s = y^{1/2}$, thus $ds = \frac{1}{2}y^{-1/2} dy$, yields

$$\begin{aligned} V_n &= V_{n-1} \cdot \int_0^1 y^{-1/2} (1-y)^{(n-1)/2} dy = V_{n-1} B\left(\frac{1}{2}, \frac{n+1}{2}\right) \\ &= \sqrt{\pi} V_{n-1} \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2} + 1)}. \end{aligned}$$

Hereby we used eq. (1.62), as well as $\Gamma(1/2) = \sqrt{\pi}$. Starting with $V_1 = 2$, a recursive application of the last formula finally leads to

$$\text{vol}_n(K_n(r)) = V_n(r) = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)} r^n = \frac{2\pi^{n/2}}{n\Gamma(\frac{n}{2})} r^n, \quad r > 0.$$

If we distinguish between even and odd dimensions, properties of the Γ function imply

$$V_{2k}(r) = \frac{\pi^k}{k!} r^{2k} \quad \text{and} \quad V_{2k+1}(r) = \frac{2^{k+1}\pi^k}{(2k+1)!!} r^{2k+1},$$

where $(2k+1)!! = 1 \cdot 3 \cdot 5 \cdots (2k-1)(2k+1)$. The first volumes are (see Figure 1.32)

$$V_1 = 2, \quad V_2 = \pi, \quad V_3 = \frac{4\pi}{3}, \quad V_4 = \frac{\pi^2}{2}, \quad V_5 = \frac{8\pi^2}{15}, \quad V_6 = \frac{\pi^3}{6}, \quad V_7 = \frac{16\pi^3}{105}, \quad V_8 = \frac{\pi^4}{24}.$$

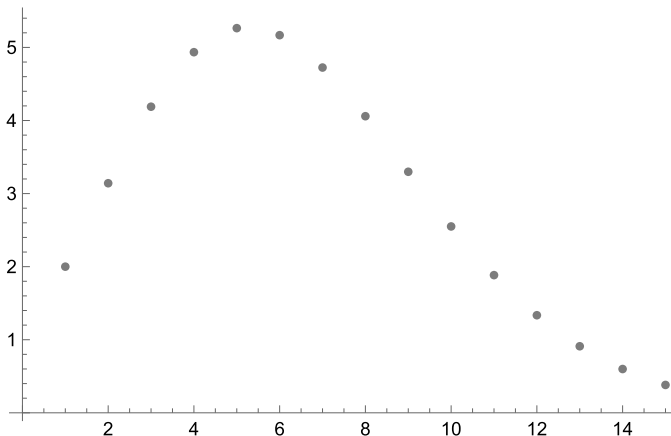


Figure 1.32: The volume of the unit ball as a function of the dimension. It attains its maximal value in dimension 5.

After the question about the volume is settled, we are now in the position to introduce the uniform distribution on bounded Borel sets in \mathbb{R}^n . Thus let $K \subseteq \mathbb{R}^n$ be a bounded Borel set in \mathbb{R}^n with volume $\text{vol}_n(K) > 0$. Define $p : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$p(x) := \begin{cases} \frac{1}{\text{vol}_n(K)} & \text{if } x \in K, \\ 0 & \text{if } x \notin K. \end{cases} \quad (1.79)$$

Proposition 1.8.12. *The function p defined by eq. (1.79) is an (n -dimensional) probability density function.*

Proof. By virtue of eq. (1.78), one has

$$\begin{aligned} \int_{\mathbb{R}^n} p(x) \, dx &= \int_K \frac{1}{\text{vol}_n(K)} \, dx = \frac{1}{\text{vol}_n(K)} \int_K \cdots \int_K \mathbf{1} \, dx_n \cdots dx_1 \\ &= \frac{\text{vol}_n(K)}{\text{vol}_n(K)} = 1. \end{aligned}$$

Since $p(x) \geq 0$ if $x \in \mathbb{R}^n$, as asserted, p is a probability density function. \square

Definition 1.8.13. The probability measure \mathbb{P} on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ with density p given by eq. (1.79) is said to be the (multivariate) **uniform distribution on K** .

Let \mathbb{P} be the uniform distribution on K . How do we get $\mathbb{P}(B)$ for a Borel set B ? Let us first assume $B \subseteq K$. Then

$$\mathbb{P}(B) = \int_B p(x) \, dx = \frac{1}{\text{vol}_n(K)} \int_B \cdots \int_B \mathbf{1} \, dx_n \cdots dx_1 = \frac{\text{vol}_n(B)}{\text{vol}_n(K)}.$$

If $B \subseteq \mathbb{R}^n$ is arbitrary, that is, B is not necessarily a subset of K , from $\mathbb{P}(B) = \mathbb{P}(B \cap K)$ it follows that

$$\mathbb{P}(B) = \frac{\text{vol}_n(B \cap K)}{\text{vol}_n(K)}.$$

If $n = 1$ and K is an interval, the latter formula coincides with eq. (1.46).

Example 1.8.14. Two friends agree to meet in a restaurant between 1 and 2 pm. Both friends go to the restaurant randomly during this hour. After they arrive, they wait 20 minutes each. What is the probability that they meet each other?

Answer: Let t_1 be the moment when the first of the two friends enters the restaurant, while t_2 is the arrival time of the second one. They arrive independently of each other, thus we may assume that the point $t := (t_1, t_2)$ is uniformly distributed in the square $Q := [1, 2]^2$. Observing that 20 minutes are a third of an hour, they meet each other if and only if $1 \leq t_1, t_2 \leq 2$ and $|t_1 - t_2| \leq 1/3$.

Setting $B := \{(t_1, t_2) \in \mathbb{R}^2 : |t_1 - t_2| \leq 1/3\}$, it is easy to see that $\text{vol}_2(B \cap Q) = 5/9$ (see Figure 1.33). Hence, if \mathbb{P} is the uniform distribution on Q , because of $\text{vol}_2(Q) = 1$, it follows that $\mathbb{P}(B) = 5/9$. Therefore, the probability that the friends meet is $5/9$.

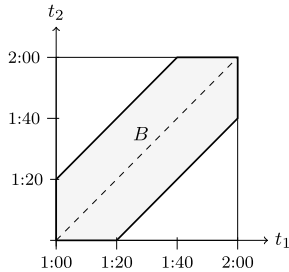


Figure 1.33: The two friends meet if (t_1, t_2) belongs to the gray region.

Example 1.8.15. Place n particles independently of each other into a ball K_R of radius $R > 0$ according to the uniform distribution. Let K_r be a smaller ball of radius $r > 0$ contained in K_R . Find the probability that exactly k of the n particles are inside K_r for some $k = 0, \dots, n$.

Answer: In the first step, we determine the probability that a single particle is in K_r . Since we assumed that the particles are uniformly distributed in K_R , this probability equals

$$p := \frac{\text{vol}_3(K_r)}{\text{vol}_3(K_R)} = \frac{(4/3)\pi r^3}{(4/3)\pi R^3} = \left(\frac{r}{R}\right)^3.$$

For each of the n particles, this p is the “success” probability to be inside K_r , hence the number of particles in K_r is $B_{n,p}$ -distributed with $p = (r/R)^3$. Thus,

$$\mathbb{P}\{k \text{ particles in } K_r\} = B_{n,p}(\{k\}) = \binom{n}{k} \left(\frac{r}{R}\right)^{3k} \left(\frac{R-r}{R}\right)^{3(n-k)}, \quad k = 0, \dots, n.$$

If the number n of particles is big and r is much smaller than R , then the number of particles in K_r is approximately Pois_λ distributed, where $\lambda = np = \frac{nr^3}{R^3}$. In other words,

$$\mathbb{P}\{k \text{ particles in } K_r\} \approx \frac{1}{k!} \left(\frac{nr^3}{R^3}\right)^k e^{-nr^3/R^3}.$$

Example 1.8.16 (Buffon’s needle test). Take a needle of length $a < 1$ and throw it randomly on a lined sheet of paper. Say the distance between two lines on the paper is 1. Find the probability that the needle cuts a line.

Answer: What is random in this experiment? Choose the two lines such that between them the midpoint of the needle lies. Let $x \in [0, 1]$ be the distance of the midpoint of the needle to the lower line. Furthermore, denote by $\theta \in [-\pi/2, \pi/2]$ the angle of the needle to a line perpendicular to the lines on the paper. For example, if $\theta = 0$, then the needle is perpendicular to the lines on the paper while for $\theta = \pm\pi/2$ it lies parallel.

Hence, to throw a needle randomly is equivalent to choosing a point (θ, x) uniformly distributed in $K = [-\pi/2, \pi/2] \times [0, 1]$.

The needle cuts the lower line (compare Figure 1.34) if and only if $\frac{a}{2} \cos \theta \geq x$, and it cuts the upper line provided that $\frac{a}{2} \cos \theta \geq 1 - x$. If the set A is as in Figure 1.35 defined by

$$A = \left\{ (\theta, x) \in [-\pi/2, \pi/2] \times [0, 1] : x \leq \frac{a}{2} \cos \theta \text{ or } 1 - x \leq \frac{a}{2} \cos \theta \right\},$$

then we get

$$\mathbb{P}\{\text{the needle cuts a line}\} = \mathbb{P}(A) = \frac{\text{vol}_2(A)}{\text{vol}_2(K)} = \frac{\text{vol}_2(A)}{\pi}.$$

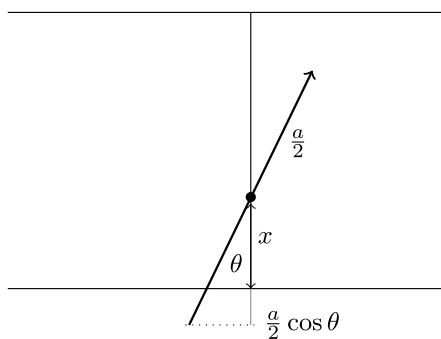


Figure 1.34: A needle of length $a < 1$ hits the lower line. Here $0 \leq x \leq 1$ denotes the distance of its midpoint to the lower line and θ is its angle to the perpendicular line.

But it follows that

$$\text{vol}_2(A) = 2 \int_{-\pi/2}^{\pi/2} \frac{a}{2} \cos \theta \, d\theta = 2a,$$

hence

$$\mathbb{P}(A) = \frac{2a}{\pi}.$$

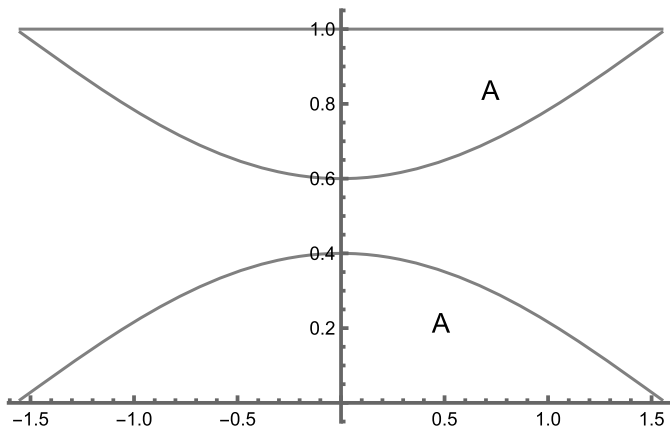


Figure 1.35: The set A where the needle cuts the lower or upper line, respectively. Since $a < 1$, the lower and upper parts of A do not overlap.

Remark 1.8.17. Suppose we throw the same needle n times. Let r_n be the relative frequency of the occurrence of A , that is,

$$r_n = \frac{\text{Number of throws where the needle cuts a line}}{n}.$$

As mentioned in Section 1.1.3, if $n \rightarrow \infty$, then r_n approaches $\mathbb{P}(A) = \frac{2a}{\pi}$. Thus for large n , we have $r_n \approx \frac{2a}{\pi}$ or, equivalently, $\pi \approx \frac{2a}{r_n}$. Consequently, throwing the needle sufficiently often, $\frac{2a}{r_n}$ should be close¹⁶ to π .

Summary: The uniform distribution \mathbb{P} on a bounded Borel set $K \subseteq \mathbb{R}^n$ with $\text{vol}_n(K) > 0$ is defined by

$$\mathbb{P}(B) = \frac{\text{vol}_n(B \cap K)}{\text{vol}_n(K)}, \quad B \in \mathcal{B}(\mathbb{R}^n).$$

1.9 Products of probability spaces*

1.9.1 Product σ -fields and measures

Suppose we execute n (maybe different) random experiments so that the outcomes do not depend on each other. In order to describe these n experiments, two different approaches are possible. Firstly, we record each single result separately, that is, we have n

¹⁶ For example, in 1901 Mr. Lazzarini threw a needle of length $a = 5/6$ exactly 3408 times. In 1808 of the cases, the needle cut a line, leading to 3.1416 as an approximate value of π .

(maybe different) probability spaces $(\Omega_1, \mathcal{A}_1, \mathbb{P}_1)$ to $(\Omega_n, \mathcal{A}_n, \mathbb{P}_n)$ modeling the outcomes of the first up to the n th experiment.

A second possible approach is that we combine the n experiments into a single one. Thus, instead of n different outcomes ω_1 to ω_n , we observe now a single vector $\omega = (\omega_1, \dots, \omega_n)$. The sample space in this approach is given by $\Omega = \Omega_1 \times \dots \times \Omega_n$.

Example 1.9.1. When rolling a die n times, the outcome is a series of n numbers ω_1 to ω_n , each in $\{1, \dots, 6\}$. Now, imagine we have a die with 6^n equally likely faces. On these faces, all possible sequences of length n with entries from $\{1, \dots, 6\}$ are written. Roll this die **once**. The first experiment may be described by n probability spaces, one for each roll. The second experiment involves only one probability space. Nevertheless, both experiments lead to the same result, a random sequence of numbers from 1 to 6.

It is intuitively clear that both approaches to this experiment (rolling a die n times) are equivalent; they differ only by the point of view. But how to come from one model to the other? One direction is immediately clear. If the random result is a vector $\omega = (\omega_1, \dots, \omega_n)$, then its coordinates may be taken as the results of the single experiments.¹⁷ But how about the other direction? That is, we are given n probability spaces $(\Omega_1, \mathcal{A}_1, \mathbb{P}_1), \dots, (\Omega_n, \mathcal{A}_n, \mathbb{P}_n)$ and have to construct a model for the joint execution of these experiments.

Of course, the “new” sample space is

$$\Omega = \Omega_1 \times \dots \times \Omega_n, \quad (1.80)$$

but what are \mathcal{A} and \mathbb{P} ? We start with the construction of the product σ -field.

Definition 1.9.2. Let \mathcal{A}_j be σ -fields on Ω_j , $1 \leq j \leq n$. Set $\Omega = \Omega_1 \times \dots \times \Omega_n$. Then

$$\mathcal{A} = \sigma\{A_1 \times \dots \times A_n : A_j \in \mathcal{A}_j\}$$

is called the **product σ -field** of \mathcal{A}_1 to \mathcal{A}_n , denoted by $\mathcal{A} = \mathcal{A}_1 \otimes \dots \otimes \mathcal{A}_n$.

Remark 1.9.3. In other words, \mathcal{A} is the smallest σ -field containing measurable rectangle sets, that is, sets of the form $A_1 \times \dots \times A_n$ with $A_j \in \mathcal{A}_j$, $1 \leq j \leq n$.

It is easy to see that $\mathcal{P}(\Omega_1) \otimes \dots \otimes \mathcal{P}(\Omega_n) = \mathcal{P}(\Omega)$. A more complicated example is as follows.

Proposition 1.9.4. Suppose $\Omega_1 = \dots = \Omega_n = \mathbb{R}$, hence $\Omega = \mathbb{R}^n$. Then the σ -field $\mathcal{B}(\mathbb{R}^n)$ of Borel sets in \mathbb{R}^n is the n -fold product of the σ -fields $\mathcal{B}(\mathbb{R})$ of Borel sets in \mathbb{R} , that is,

¹⁷ Of course, one still has to verify that the distribution of the coordinates is the same as in the single experiments. But before we can do this, we need a probability measure describing the distribution of the vectors (cf. Proposition 1.9.10).

$$\mathcal{B}(\mathbb{R}^n) = \underbrace{\mathcal{B}(\mathbb{R}) \otimes \cdots \otimes \mathcal{B}(\mathbb{R})}_{n \text{ times}}.$$

Proof. We only give a sketch of the proof. Let Q be a box as in eq. (1.73). Then we have $Q = A_1 \times \cdots \times A_n$, where the A_j s are intervals, hence in $\mathcal{B}(\mathbb{R})$. By the construction of the product σ -field, it follows that $Q \in \mathcal{B}(\mathbb{R}) \otimes \cdots \otimes \mathcal{B}(\mathbb{R})$. But $\mathcal{B}(\mathbb{R}^n)$ is the smallest σ -field containing all boxes, which lets us conclude

$$\mathcal{B}(\mathbb{R}^n) \subseteq \mathcal{B}(\mathbb{R}) \otimes \cdots \otimes \mathcal{B}(\mathbb{R}).$$

The inclusion in the other direction may be proved as follows: fix $a_2 < b_2$ to $a_n < b_n$ and let

$$\mathcal{C}_1 = \{C \in \mathcal{B}(\mathbb{R}) : C \times [a_2, b_2] \times \cdots \times [a_n, b_n] \in \mathcal{B}(\mathbb{R}^n)\}.$$

It is not difficult to prove that \mathcal{C}_1 is a σ -field. If $C = [a_1, b_1]$, then

$$C \times [a_2, b_2] \times \cdots \times [a_n, b_n]$$

is a box, thus in $\mathcal{B}(\mathbb{R}^n)$. Consequently, \mathcal{C}_1 contains closed intervals, hence, since $\mathcal{B}(\mathbb{R})$ is the smallest σ -field with this property, it follows $\mathcal{C}_1 = \mathcal{B}(\mathbb{R})$. This tells us that for all $B_1 \in \mathcal{B}(\mathbb{R})$ and all $a_j < b_j$,

$$B_1 \times [a_2, b_2] \times \cdots \times [a_n, b_n] \in \mathcal{B}(\mathbb{R}^n).$$

In a next step, fix $B_1 \in \mathcal{B}(\mathbb{R})$ and $a_3 < b_3$ up to $a_n < b_n$, and set

$$\mathcal{C}_2 = \{C \in \mathcal{B}(\mathbb{R}) : B_1 \times C \times [a_3, b_3] \times \cdots \times [a_n, b_n] \in \mathcal{B}(\mathbb{R}^n)\}.$$

By the same arguments as before, but now using the first step, we get $\mathcal{C}_2 = \mathcal{B}(\mathbb{R})$, that is,

$$B_1 \times B_2 \times [a_3, b_3] \times \cdots \times [a_n, b_n] \in \mathcal{B}(\mathbb{R}^n)$$

for all $B_1, B_2 \in \mathcal{B}(\mathbb{R})$ and $a_j < b_j$.

Iterating further, we finally obtain that for all $B_j \in \mathcal{B}(\mathbb{R})$ it follows that

$$B_1 \times \cdots \times B_n \in \mathcal{B}(\mathbb{R}^n).$$

Since $\mathcal{B}(\mathbb{R}) \otimes \cdots \otimes \mathcal{B}(\mathbb{R}) = \sigma\{B_1 \times \cdots \times B_n : B_j \in \mathcal{B}(\mathbb{R})\}$ is the smallest σ -field containing sets $B_1 \times \cdots \times B_n$, this implies

$$\mathcal{B}(\mathbb{R}) \otimes \cdots \otimes \mathcal{B}(\mathbb{R}) \subseteq \mathcal{B}(\mathbb{R}^n)$$

and completes the proof. □

Let us now turn to the probability measure \mathbb{P} on (Ω, \mathcal{A}) that describes the combined experiment.

Definition 1.9.5. Let $(\Omega_1, \mathcal{A}_1, \mathbb{P}_1)$ to $(\Omega_n, \mathcal{A}_n, \mathbb{P}_n)$ be n probability spaces. Define Ω by eq. (1.80) and endow it with the product σ -field $\mathcal{A} = \mathcal{A}_1 \otimes \cdots \otimes \mathcal{A}_n$. A probability measure \mathbb{P} on (Ω, \mathcal{A}) is called the **product measure** of $\mathbb{P}_1, \dots, \mathbb{P}_n$ if

$$\mathbb{P}(A_1 \times \cdots \times A_n) = \mathbb{P}_1(A_1) \cdots \mathbb{P}_n(A_n) \quad \text{for all } A_j \in \mathcal{A}_j. \quad (1.81)$$

We write $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$ and if $\mathbb{P}_1 = \cdots = \mathbb{P}_n = \mathbb{P}_0$ set

$$\mathbb{P}_0^{\otimes n} := \underbrace{\mathbb{P}_0 \otimes \cdots \otimes \mathbb{P}_0}_{n \text{ times}}.$$

It is not clear at all whether product measures exist and, if this is so, whether condition (1.81) determines them uniquely. The next result shows that the answer to both questions is affirmative. Unfortunately, the proof is too complicated to be presented here. The idea is quite similar to that used in the introduction of volumes in eq. (1.77). The boxes appearing there have to be replaced by rectangle sets $A_1 \times \cdots \times A_n$ with $A_j \in \mathcal{A}_j$ and the volume of the boxes by $\mathbb{P}_1(A_1)$ to $\mathbb{P}_n(A_n)$, respectively. We refer to [Dur19, Section 1.7], [Kle20] or [Coh13], for a detailed proof for the existence (and uniqueness) of product measures.

Proposition 1.9.6. Let $(\Omega_1, \mathcal{A}_1, \mathbb{P}_1), \dots, (\Omega_n, \mathcal{A}_n, \mathbb{P}_n)$ be probability spaces. Define Ω by eq. (1.80) and let \mathcal{A} be the product σ -field of the \mathcal{A}_j s. Then there is a unique probability measure \mathbb{P} on (Ω, \mathcal{A}) satisfying

$$\mathbb{P}(A_1 \times \cdots \times A_n) = \mathbb{P}_1(A_1) \cdots \mathbb{P}_n(A_n) \quad \text{for all } A_j \in \mathcal{A}_j. \quad (1.82)$$

Hence, the product measure $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$ always exists and is uniquely determined by (1.82).

Remark 1.9.7. While the product measure of rectangle sets can be evaluated directly by eq. (1.82), it is more complicated to determine the probability for arbitrary non-rectangle sets. Compare Figures 1.36 and 1.37.

Example 1.9.8. Let $\Omega_1 = \Omega_2 = \{1, \dots, 6\}$ be endowed with the uniform distributions \mathbb{P}_1 and \mathbb{P}_2 . That is, $\mathbb{P}_1(\{\omega\}) = \mathbb{P}_2(\{\omega\}) = \frac{1}{6}$ for all $\omega = 1, \dots, 6$. Then

$$\Omega = \Omega_1 \times \Omega_2 = \{(\omega_1, \omega_2) : \omega_1, \omega_2 \in \{1, \dots, 6\}\},$$

and we get

$$(\mathbb{P}_1 \otimes \mathbb{P}_2)(A_1 \times A_2) = \mathbb{P}_1(A_1) \cdot \mathbb{P}_2(A_2) = \frac{|A_1|}{6} \cdot \frac{|A_2|}{6} = \frac{|A_1 \times A_2|}{36}$$

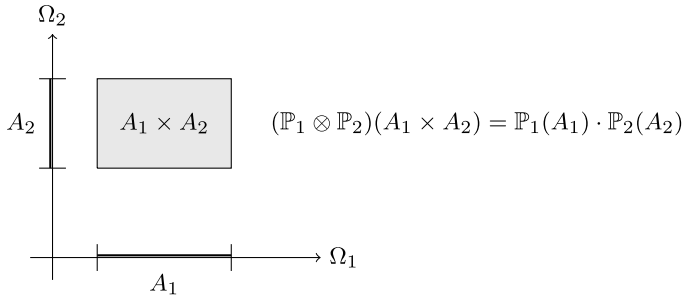


Figure 1.36: The product measure $\mathbb{P}_1 \otimes \mathbb{P}_2$ applied to rectangle sets in $\Omega_1 \times \Omega_2$.

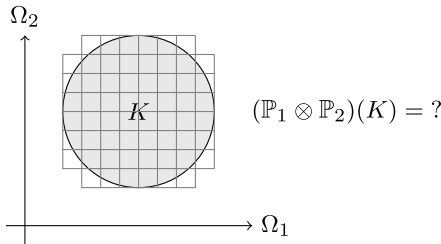


Figure 1.37: How to evaluate $(\mathbb{P}_1 \otimes \mathbb{P}_2)(K)$? Cover it in an optimal way by unions of rectangles.

for all $A_1 \subseteq \Omega_1$ and $A_2 \subseteq \Omega_2$. So we see that

$$(\mathbb{P}_1 \otimes \mathbb{P}_2)(B) = \frac{|B|}{36} = \mathbb{P}(B)$$

whenever $B = A_1 \times A_2$ is a rectangle set and where \mathbb{P} denotes the uniform distribution on $\Omega = \{1, \dots, 6\}^2$. Since $\mathbb{P}_1 \otimes \mathbb{P}_2$ is uniquely determined by its values at rectangle sets, it follows that the product measure is nothing else as the uniform distribution on Ω . For example, this implies that

$$(\mathbb{P}_1 \otimes \mathbb{P}_2)(\{(1, 1), (2, 2), \dots, (6, 6)\}) = \frac{6}{36} = \frac{1}{6}.$$

Note that $\{(1, 1), \dots, (6, 6)\}$ is no rectangle set, hence in this case, formula (1.82) does not apply.

Corollary 1.9.9. *Let $\mathbb{P}_1, \dots, \mathbb{P}_n$ be probability measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Then there is a unique probability measure \mathbb{P} on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ such that*

$$\mathbb{P}(B_1 \times \dots \times B_n) = \mathbb{P}_1(B_1) \cdots \mathbb{P}_n(B_n) \quad \text{for all } B_j \in \mathcal{B}(\mathbb{R}).$$

Proof. The proof is a direct consequence of Propositions 1.9.6 and 1.9.4. Indeed, take $\mathbb{P} = \mathbb{P}_1 \otimes \dots \otimes \mathbb{P}_n$ and observe that $\mathcal{B}(\mathbb{R}^n) = \mathcal{B}(\mathbb{R}) \otimes \dots \otimes \mathcal{B}(\mathbb{R})$. □

Let us shortly come back to the question asked at the beginning of this section. Suppose we observe a vector $\omega = (\omega_1, \dots, \omega_n)$. How are the coordinates distributed?

Proposition 1.9.10. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be the product probability space of $(\Omega_1, \mathcal{A}_1, \mathbb{P}_1)$ to $(\Omega_n, \mathcal{A}_n, \mathbb{P}_n)$. If $j \leq n$ and $A \in \mathcal{A}_j$, then*

$$\mathbb{P}\{(\omega_1, \dots, \omega_n) \in \Omega : \omega_j \in A\} = \mathbb{P}_j(A).$$

Proof. Observe that

$$\{(\omega_1, \dots, \omega_n) \in \Omega : \omega_j \in A\} = \Omega_1 \times \dots \times \Omega_{j-1} \times A \times \Omega_{j+1} \times \dots \times \Omega_n,$$

thus eq. (1.81) implies (see Figure 1.38)

$$\begin{aligned} &\mathbb{P}\{(\omega_1, \dots, \omega_n) \in \Omega : \omega_j \in A\} \\ &= \mathbb{P}_1(\Omega_1) \cdots \mathbb{P}_{j-1}(\Omega_{j-1}) \cdot \mathbb{P}_j(A) \cdot \mathbb{P}_{j+1}(\Omega_{j+1}) \cdots \mathbb{P}_n(\Omega_n) = \mathbb{P}_j(A), \end{aligned}$$

as asserted. □

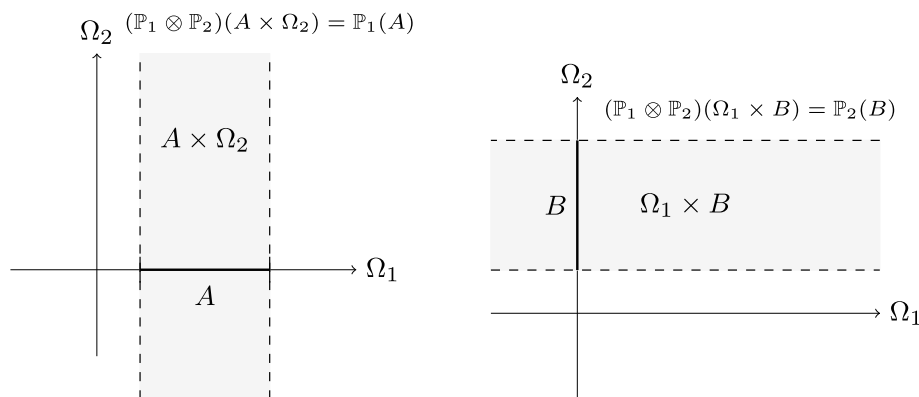


Figure 1.38: The product measure $\mathbb{P}_1 \otimes \mathbb{P}_2$ applied to the cylindrical shaped sets $A \times \Omega_2$ and $\Omega_1 \times B$.

How do we get product measures in concrete cases? We answer this question for discrete and continuous probability measures separately.

1.9.2 Product measures: discrete case

Let Ω_1 to Ω_n be either finite or countably infinite sets. Given probability measures \mathbb{P}_j defined on $\mathcal{P}(\Omega_j)$, $1 \leq j \leq n$, the following result characterizes the product measure of the \mathbb{P}_j s.

Proposition 1.9.11. *Probability measure \mathbb{P} is the product measure of $\mathbb{P}_1, \dots, \mathbb{P}_n$ if and only if*

$$\mathbb{P}(\{\omega\}) = P_1(\{\omega_1\}) \cdots P_n(\{\omega_n\}) \quad \text{for all } \omega = (\omega_1, \dots, \omega_n) \in \Omega. \quad (1.83)$$

Proof. One direction is easy. Indeed, if $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$, given $\omega = (\omega_1, \dots, \omega_n) \in \Omega$ set $A = \{\omega\}$ and $A_j = \{\omega_j\}$. Then $A = A_1 \times \cdots \times A_n$, hence

$$\mathbb{P}(\{\omega\}) = \mathbb{P}(A) = \mathbb{P}_1(A_1) \cdots \mathbb{P}_n(A_n) = P_1(\{\omega_1\}) \cdots P_n(\{\omega_n\}),$$

proving eq. (1.83).

To verify the other implication, let \mathbb{P} be a probability measure on $(\Omega, \mathcal{P}(\Omega))$ satisfying (1.83). We have to show that \mathbb{P} fulfills eq. (1.81). Thus choose arbitrary $A_j \subseteq \Omega_j$ and set $A = A_1 \times \cdots \times A_n$. By applying eq. (1.83), it follows that

$$\begin{aligned} \mathbb{P}(A) &= \sum_{\omega \in A} \mathbb{P}(\{\omega\}) = \sum_{(\omega_1, \dots, \omega_n) \in A} \mathbb{P}(\{(\omega_1, \dots, \omega_n)\}) \\ &= \sum_{\omega_1 \in A_1, \dots, \omega_n \in A_n} P_1(\{\omega_1\}) \cdots P_n(\{\omega_n\}) \\ &= \sum_{\omega_1 \in A_1} P_1(\{\omega_1\}) \cdots \sum_{\omega_n \in A_n} P_n(\{\omega_n\}) = P_1(A_1) \cdots P_n(A_n). \end{aligned}$$

This being true for all $A_j \subseteq \Omega_j$ shows that $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$, and the proof is complete. \square

Summary: In the discrete case, the product measure is characterized as follows. Given $A \subseteq \Omega$,

$$(\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n)(A) = \sum_{(\omega_1, \dots, \omega_n) \in A} P_1(\{\omega_1\}) \cdots P_n(\{\omega_n\}).$$

Example 1.9.12. Suppose two players, say U and V , each simultaneously toss a biased coin. On both coins, “0” (failure) appears with probability $1 - p$ and “1” (success) with probability p . Whoever gets the first “1” wins.

A pair $(k, l) \in \mathbb{N}^2$ occurs if player U has his first success in trial k and player V in trial l . Each single experiment is described by the geometric distribution G_p , hence the model for the combined experiment is $(\mathbb{N}^2, \mathcal{P}(\mathbb{N}^2), G_p^{\otimes 2})$. Here

$$G_p^{\otimes 2}(A) = \sum_{(k,l) \in A} G_p(\{k\})G_p(\{l\}) = \sum_{(k,l) \in A} p^2(1-p)^{k+l-2}, \quad A \subseteq \mathbb{N}^2.$$

If $A = \{(k, k) : k \geq 1\}$, that is, the game ends in a draw, then

$$G_p^{\otimes 2}(A) = p^2 \sum_{k=1}^{\infty} (1-p)^{2k-2} = \frac{p^2}{1 - (1-p)^2} = \frac{p}{2-p}.$$

Consequently, with probability

$$1 - \frac{p}{2-p} = \frac{2-2p}{2-p},$$

either U or V wins. Since both players have the same chance to win, it follows that U wins with probability $\frac{1-p}{2-p}$ and the same is true for V .

If we analyze the result, then we see that we have three possible outcomes of the game: U wins or V wins, each with probability $\frac{1-p}{2-p}$, or the game ends in a draw with probability $\frac{p}{2-p}$. In the case of a fair coin, that is, if $p = 1/2$, these three outcomes each occur with probability $1/3$. In the general setting, the following is true: the bigger the success probability p , the more likely the game ends in a draw. On the contrary, if p is near zero, then with great probability there will be a winner of the game.

Example 1.9.13. Toss a biased coin n times. Say the coin is labeled with “0” and “1” and $p \in [0, 1]$ is the probability of the occurrence of “1.” Recording each single result separately, the describing probability spaces are $(\{0, 1\}, \mathcal{P}(\{0, 1\}), \mathbb{P}_j)$, $1 \leq j \leq n$, with $\mathbb{P}_j(\{1\}) = p$. Which probability space does the combined result describe?

Answer: Of course, the sample space is $\{0, 1\}^n$ endowed with σ -field $\mathcal{P}(\Omega)$. Let $\omega = (\omega_1, \dots, \omega_n)$ be an arbitrary vector in Ω . Then by Proposition 1.9.11, the product measure \mathbb{P} of the \mathbb{P}_j s is characterized by

$$\mathbb{P}(\{\omega\}) = \mathbb{P}_1(\{\omega_1\}) \cdots \mathbb{P}_n(\{\omega_n\}) = p^k (1-p)^{n-k}$$

where $k = |\{j \leq n : \omega_j = 1\}| = \sum_{j=1}^n \omega_j$.

For example, tossing the coin five times, the sequence $(0, 0, 1, 1, 0)$ occurs with probability $p^2(1-p)^3$.

In the general case, let A be the set of sequences possessing exactly k times the number “1.” Then

$$\begin{aligned} (\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n)(A) &= \sum_{(\omega_1, \dots, \omega_n) \in A} \mathbb{P}_1(\{\omega_1\}) \cdots \mathbb{P}_n(\{\omega_n\}) \\ &= |A| \cdot p^k (1-p)^{n-k} = \binom{n}{k} p^k (1-p)^{n-k}. \end{aligned}$$

This is another justification for the model described by the binomial distribution $B_{n,p}$.

Example 1.9.14. Suppose we have two urns, both containing the same proportion of white balls. Choose from each urn n balls with replacement. How likely is it to take out the same number of white balls from each of the two urns?

Answer: The experiment is described by the product measure $B_{n,p}^{\otimes 2}$ where p denotes the proportion of white balls in each urn. Hence, the probability to observe k white balls from the first urn and ℓ from the second one equals

$$\binom{n}{k} \binom{n}{\ell} p^{k+\ell} (1-p)^{2n-k-\ell}, \quad 0 \leq k, \ell \leq n.$$

Consequently, if A is the event that one chooses the same number of white balls, then

$$B_{n,p}^{\otimes 2}(A) = \sum_{k=0}^n \binom{n}{k}^2 p^{2k} (1-p)^{2n-2k}. \quad (1.84)$$

Look at Figure 1.39 to get an impression how the probability of the event A depends on the success probability p . For example, if either $p = 0$ or $p = 1$, then, of course, we will for sure see the same number of white balls out of the two urns. Argue why we get the same probability for A when we replace p by $1 - p$.

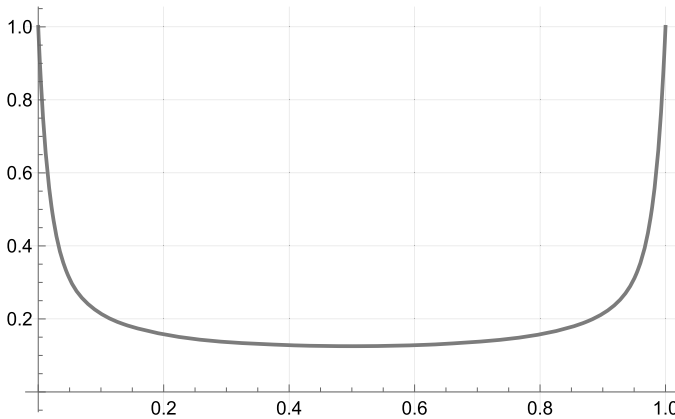


Figure 1.39: The probability to observe the same number of white balls, when choosing 20 balls, depending on the proportion p . The minimal value is 0.125371 attained at $p = 0.5$.

If $p = 1/2$, then eq. (1.84) reduces to (use formula (A.19))

$$\frac{1}{2^{2n}} \sum_{k=0}^n \binom{n}{k}^2 = \frac{1}{2^{2n}} \binom{2n}{n}.$$

Remark 1.9.15. Stirling's formula (Corollary 1.6.15) implies

$$\frac{1}{2} \frac{1}{\sqrt{\pi n}} < \frac{1}{2^{2n}} \binom{2n}{n} < \frac{1}{\sqrt{\pi n}}, \quad n = 1, 2, \dots$$

Hence, if in both urns the number of white balls equals that of black, then for large n the probability to observe from both urns the same number of white balls is approximately $1/\sqrt{\pi n}$.

1.9.3 Product measures: continuous case

Here we assume $\Omega_1 = \cdots = \Omega_n = \mathbb{R}$, hence the product sample space is $\Omega = \mathbb{R}^n$. Furthermore, each $\Omega_j = \mathbb{R}$ is endowed with the Borel σ -field. Because of Proposition 1.9.4, the product σ -field on $\Omega = \mathbb{R}^n$ is given by $\mathcal{B}(\mathbb{R}^n)$.

The next proposition characterizes the product measure of continuous probability measures.

Proposition 1.9.16. *Let $\mathbb{P}_1, \dots, \mathbb{P}_n$ be probability measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ with respective density functions p_1, \dots, p_n , that is,*

$$\mathbb{P}_j([a, b]) = \int_a^b p_j(x) dx, \quad 1 \leq j \leq n.$$

Define $p : \mathbb{R}^n \rightarrow [0, \infty)$ by

$$p(x) = p_1(x_1) \cdots p_n(x_n), \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n. \quad (1.85)$$

Then the product measure $\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$ is continuous with (n -dimensional) density p defined by (1.85). In other words, for each Borel set $A \subseteq \mathbb{R}^n$,

$$(\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n)(A) = \underbrace{\int_A \cdots \int}_A p_1(x_1) \cdots p_n(x_n) dx_n \cdots dx_1 = \int_A p(x) dx.$$

Proof. First note that p is a density of the product measure $\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$ if

$$(\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n)(Q) = \int_Q p(x) dx$$

for all boxes $Q = [a_1, b_1] \times \cdots \times [a_n, b_n]$. But this is an immediate consequence of

$$\begin{aligned} \int_Q p(x) dx &= \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} p_1(x_1) \cdots p_n(x_n) dx_n \cdots dx_1 \\ &= \left(\int_{a_1}^{b_1} p_1(x_1) dx_1 \right) \cdots \left(\int_{a_n}^{b_n} p_n(x_n) dx_n \right) = \mathbb{P}_1([a_1, b_1]) \cdots \mathbb{P}_n([a_n, b_n]) \\ &= (\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n)([a_1, b_1] \times \cdots \times [a_n, b_n]) = (\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n)(Q). \end{aligned}$$

This completes the proof. \square

Because of its importance, let us explain through several examples how Proposition 1.9.16 applies. Further applications, for example, the characterization of independent random variables, will follow in Sections 3 and 8.

Example 1.9.17. Let the probability measures \mathbb{P}_j , $1 \leq j \leq n$, be uniform distributions on $[\alpha_j, \beta_j]$. Thus

$$p_j(x) = \begin{cases} \frac{1}{\beta_j - \alpha_j} & \text{if } \alpha_j \leq x \leq \beta_j, \\ 0 & \text{otherwise,} \end{cases}$$

henceforth, if $x = (x_1, \dots, x_n)$, then

$$p(x) = p_1(x_1) \cdots p_n(x_n) = \begin{cases} \frac{1}{\prod_{j \leq n} (\beta_j - \alpha_j)} & \text{if } x \in K, \\ 0 & \text{otherwise.} \end{cases}$$

Here $K \subseteq \mathbb{R}^n$ is the box $[\alpha_1, \beta_1] \times \cdots \times [\alpha_n, \beta_n]$. Since $\prod_{j \leq n} (\beta_j - \alpha_j) = \text{vol}_n(K)$, it follows that the product measure $\mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$ is nothing else as the (n -dimensional) uniform distribution on K as introduced in¹⁸ Definition 1.8.13.

Summary: The product measure of n uniform distributions on intervals $[\alpha_j, \beta_j]$ is the uniform distribution on the box $[\alpha_1, \beta_1] \times \cdots \times [\alpha_n, \beta_n]$.

Example 1.9.18. Let \mathbb{P} be the uniform distribution on the unit sphere

$$K = \{(x_1, x_2) \in \mathbb{R}^2 : x_1^2 + x_2^2 \leq 1\}.$$

That is,

$$\mathbb{P}(A) = \frac{\text{vol}_2(A \cap K)}{\pi}, \quad A \in \mathcal{B}(\mathbb{R}^2).$$

The density p of \mathbb{P} is given by

$$p(x_1, x_2) = \begin{cases} \frac{1}{\pi} & \text{if } x_1^2 + x_2^2 \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Then there are **no** functions p_1, p_2 on \mathbb{R} for which

$$p(x_1, x_2) = p_1(x_1) \cdot p_2(x_2), \quad x_1, x_2 \in \mathbb{R}.$$

Indeed, if such a representation were to exist, then $p_2(x_2) = 0$ if $x_1^2 > 1 - x_2^2$ for all $-1 \leq x_1 \leq 1$. Hence, $p_2(x_2) = 0$ for all x_2 . This contradicts the representation of the density p . Consequently, \mathbb{P} is **not** a product measure.

¹⁸ This result was already used in Example 1.8.14. Indeed, the arrival times t_1 and t_2 were described by the uniform distributions on $[1, 2]$, thus the pair $t = (t_1, t_2)$ is distributed according to the product measure, which is the uniform distribution on $[1, 2] \times [1, 2]$. Similarly, in Example 1.8.16, we applied that the pair (θ, x) is uniformly distributed on $[-\pi/2, \pi/2] \times [0, 1]$.

Example 1.9.19. Assume now $\mathbb{P}_1 = \dots = \mathbb{P}_n = E_\lambda$, that is, we want to describe the product of n exponential distributions with parameter $\lambda > 0$. Since $p_j(s) = \lambda e^{-\lambda s}$ if $s \geq 0$ and $p_j(s) = 0$ if $s < 0$, their product measure $E_\lambda^{\otimes n}$ possesses the density

$$p(s_1, \dots, s_n) = \begin{cases} \lambda^n e^{-\lambda(s_1 + \dots + s_n)} & \text{if } s_1, \dots, s_n \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Which random experiment does $E_\lambda^{\otimes n}$ describe? Suppose we have n light bulbs of the same type with lifetime distributed according to E_λ . Switch on all n bulbs at once and record the times t_1, \dots, t_n where the first bulb, the second, and so on, burns out. If $t = (t_1, \dots, t_n) \in \mathbb{R}^n$ denotes the generated vector of these times, then for Borel sets $A \subseteq [0, \infty)^n$,

$$\mathbb{P}\{t \in A\} = E_\lambda^{\otimes n}(A) = \lambda^n \int_A e^{-\lambda(s_1 + \dots + s_n)} ds_1 \cdots ds_n.$$

For example, if we want to compute the probability of

$$A := \{(t_1, \dots, t_n) : 0 \leq t_1 \leq \dots \leq t_n\},$$

that is, that the second bulb burns longer than the first, the third longer than the second, and so on, then

$$E_\lambda^{\otimes n}(A) = \lambda^n \int_0^\infty e^{-\lambda s_n} \int_0^{s_n} e^{-\lambda s_{n-1}} \int_0^{s_{n-1}} \cdots \int_0^{s_3} e^{-\lambda s_2} \int_0^{s_2} e^{-\lambda s_1} ds_1 \cdots ds_n.$$

Iterative integration leads to $E_\lambda^{\otimes n}(A) = 1/n!$. This is more or less obvious by the following observation. Each ordering of the failure times is equally likely. And since there are $n!$ different ways to order these times, each ordering has probability $1/n!$. In particular, this is true for the ordering $t_1 \leq \dots \leq t_n$.

Next we ask how likely is it that all n bulbs still burn at time $T > 0$. That is, we ask for the probability $E_\lambda^{\otimes n}(B)$ where $B = [T, \infty)^n$. The properties of product measures imply that

$$E_\lambda^{\otimes n}(B) = \underbrace{E_\lambda([T, \infty)) \cdots E_\lambda([T, \infty))}_n = e^{n\lambda T}.$$

Consequently, the probability that at least one of the n bulbs burns out before time $T > 0$ equals $1 - e^{-n\lambda T}$. In other words, if we say that the system of n bulbs becomes defective if at least one bulb burns out, then the lifetime of this system is exponentially distributed with parameter $n\lambda$.

Next we give another example of a product measure that will play a crucial role in Sections 6 and 8.

Example 1.9.20. Let $\mathbb{P}_1, \dots, \mathbb{P}_n$ be standard normal distributions. The corresponding densities are

$$p_j(x_j) = \frac{1}{\sqrt{2\pi}} e^{-x_j^2/2}, \quad 1 \leq j \leq n.$$

Thus, by eq. (1.85), the density p of their product $\mathcal{N}(0, 1)^{\otimes n}$ coincides with

$$p(x) = \frac{1}{(2\pi)^{n/2}} e^{-\sum_{j=1}^n x_j^2/2} = \frac{1}{(2\pi)^{n/2}} e^{-|x|^2/2},$$

where $|x| = (\sum_{j=1}^n x_j^2)^{1/2}$ denotes the Euclidean distance of the vector x to 0 (see Section A.4).

See Figure 1.40 for the visualization of the density p in the case $n = 2$.

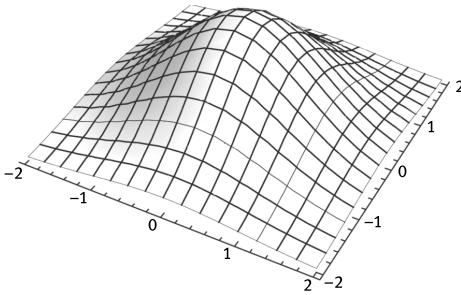


Figure 1.40: The density of the two-dimensional standard normal distribution.

Definition 1.9.21. The probability measure $\mathcal{N}(0, 1)^{\otimes n}$ on $\mathcal{B}(\mathbb{R}^n)$ is called the **n -dimensional, or multivariate, standard normal distribution**. It is described by

$$\mathcal{N}(0, 1)^{\otimes n}(B) = \frac{1}{(2\pi)^{n/2}} \int_B e^{-|x|^2/2} dx.$$

For example, if $K \subseteq \mathbb{R}^2$ denotes a circle of radius 1, then this leads to the following integral:

$$\begin{aligned} \mathcal{N}(0, 1)^{\otimes 2}(K) &= \frac{1}{2\pi} \iint_K e^{-x_1^2/2} e^{-x_2^2/2} dx_1 dx_2 \\ &= \frac{1}{\pi} \int_{-1}^1 e^{-x_1^2/2} \int_0^{\sqrt{1-x_1^2}} e^{-x_2^2/2} dx_2 dx_1 \\ &= \frac{2}{\pi} \int_0^1 e^{-x_1^2/2} \int_0^{\sqrt{1-x_1^2}} e^{-x_2^2/2} dx_2 dx_1 \approx 0.393. \end{aligned}$$

Example 1.9.22. Finally, we describe the n -fold product measure of general normal distributions $\mathcal{N}(\mu_j, \sigma_j^2)$ with $\mu_j \in \mathbb{R}$ and $\sigma_j^2 > 0$. The densities are

$$p_j(x_j) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-(x_j - \mu_j)^2/2\sigma_j^2}, \quad 1 \leq j \leq n.$$

Hence, the product measure $\mathcal{N}(\mu_1, \sigma_1^2) \otimes \cdots \otimes \mathcal{N}(\mu_n, \sigma_n^2)$ possesses the density

$$p(x) = \frac{1}{(2\pi)^{n/2}\sigma_1 \cdots \sigma_n} \exp\left(-\sum_{j=1}^n \frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right), \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n.$$

In particular, if $\mathcal{N}(\mu_1, \sigma_1^2) = \cdots = \mathcal{N}(\mu_n, \sigma_n^2) = \mathcal{N}(\mu, \sigma^2)$ by

$$\sum_{j=1}^n \frac{(x_j - \mu_j)^2}{2\sigma_j^2} = \frac{|x - \bar{\mu}|^2}{2\sigma^2},$$

where $\bar{\mu} = (\mu, \dots, \mu)$, the n -fold product measure of $\mathcal{N}(\mu, \sigma^2)$ acts as follows:

$$\mathcal{N}(\mu, \sigma^2)^{\otimes n}(B) = \frac{1}{(2\pi)^{n/2}\sigma^n} \int_B e^{-|x - \bar{\mu}|^2/2\sigma^2} dx, \quad B \in \mathcal{B}(\mathbb{R}^n). \quad (1.86)$$

Summary: Let $(\Omega_1, \mathcal{A}_1, \mathbb{P}_1)$ up to $(\Omega_n, \mathcal{A}_n, \mathbb{P}_n)$ be probability spaces. Then there exists a unique probability measure $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$ (the product measure) on $\Omega = \Omega_1 \times \cdots \times \Omega_n$ such that

$$\mathbb{P}(A_1 \times \cdots \times A_n) = \mathbb{P}_1(A_1) \cdots \mathbb{P}_n(A_n), \quad A_j \in \mathcal{A}_j.$$

In the discrete case, the product measure \mathbb{P} is described by

$$\mathbb{P}(A) = \sum_{(\omega_1, \dots, \omega_n) \in A} \mathbb{P}_1(\{\omega_1\}) \cdots \mathbb{P}_n(\{\omega_n\}), \quad A \subseteq \Omega,$$

while for continuous \mathbb{P}_j with densities p_j the product measure is given by

$$\mathbb{P}(B) = \int \cdots \int_B p_1(x_1) \cdots p_n(x_n) dx_n \cdots dx_1, \quad B \in \mathcal{B}(\mathbb{R}^n).$$

1.10 Problems

Problem 1.1. Let A , B , and C be three events in a sample space Ω . Express the following events in terms of these sets:

- Only A occurs.
- A and B occur, but C does not.
- At least one of the three events occurs.

- At least two of the events occur.
- At most one of the three events occurs.
- None of the events occur.
- Not more than two of the events occur.
- Exactly two of the events occur.

Problem 1.2. Suppose an urn contains black and white balls. Successively one draws n balls out of the urn. The event A_j occurs if the ball drawn in the j th trial is white. Hereby $1 \leq j \leq n$. Express the following events B_1, \dots, B_4 in terms of the A_j s:

$$\begin{aligned} B_1 &= \{\text{All drawn balls are white}\}, \\ B_2 &= \{\text{At least one of the balls is white}\}, \\ B_3 &= \{\text{Exactly one of the drawn balls is white}\}, \\ B_4 &= \{\text{All } n \text{ balls possess the same color}\}. \end{aligned}$$

Determine the cardinalities $|B_j|, j = 1, \dots, 4$.

Problem 1.3. Argue why for every $t \in \mathbb{R}$ the elementary event $\{t\}$ is a Borel set. What is wrong in the following argument? Since for any set $B \subseteq \mathbb{R}$ one has

$$B = \bigcup_{t \in B} \{t\},$$

every subset B of \mathbb{R} is union of Borel sets. Because $\mathcal{B}(\mathbb{R})$ is a σ -field, the union of Borel sets is a Borel set as well. Thus, any $B \subseteq \mathbb{R}$ is a Borel set.

Problem 1.4. Suppose \mathbb{P} is a σ -additive mapping from a σ -field \mathcal{A} to $[0, 1]$ with $\mathbb{P}(\Omega) = 1$. Show that then necessarily $\mathbb{P}(\emptyset) = 0$. Consequently, whenever a σ -additive mapping \mathbb{P} satisfies $\mathbb{P}(\Omega) = 1$, then it is a probability measure.

Problem 1.5. Let \mathbb{P} be a probability measure on (Ω, \mathcal{A}) . Given $A, B \in \mathcal{A}$, show that

$$\mathbb{P}(A \Delta B) = \mathbb{P}(A) + \mathbb{P}(B) - 2\mathbb{P}(A \cap B).$$

Problem 1.6. The events A and B possess the probabilities $\mathbb{P}(A) = 1/3$ and $\mathbb{P}(B) = 1/4$. Moreover, we know that $\mathbb{P}(A \cap B) = 1/6$. Compute $\mathbb{P}(A^c)$, $\mathbb{P}(A^c \cup B)$, $\mathbb{P}(A \cup B^c)$, $\mathbb{P}(A \cap B^c)$, $\mathbb{P}(A \Delta B)$, and $\mathbb{P}(A^c \cup B^c)$.

Problem 1.7 (Inclusion–exclusion formula). Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and let $A_1, \dots, A_n \in \mathcal{A}$ be some (not necessarily disjoint) events. Prove that

$$\mathbb{P}\left(\bigcup_{j=1}^n A_j\right) = \sum_{k=1}^n (-1)^{k+1} \sum_{1 \leq j_1 < \dots < j_k \leq n} \mathbb{P}(A_{j_1} \cap \dots \cap A_{j_k}).$$

Hint: One way to prove this is by induction over n , thereby using Proposition 1.2.3.

Problem 1.8. Use Problem 1.7 to investigate the following question: The numbers from 1 to n are ordered randomly. All orderings are equally likely. What is the probability that there exists an integer $m \leq n$ so that m is at position m of the ordering? Determine the limit of this probability as $n \rightarrow \infty$.

Another version of this problem is as follows. Suppose n persons attend a Christmas party. Each of the n participants brings a present with him. These presents are collected, mixed, and then randomly distributed among the guests. Compute the probability that at least one of the participants gets his own present.

Problem 1.9. Suppose there are N balls in an urn; k are white, l are red, and m are black. Thus, $k + l + m = N$. Choose n balls out of the urn. Find a formula for the probability that among the n chosen balls are those of all three colors. Investigate this problem if

1. the chosen ball is always replaced and
2. if $n \leq N$ and the balls are not replaced.

Hint: If A is the event that all three colors appear then compute $\mathbb{P}(A^c)$. To this end, write $A^c = A_1 \cup A_2 \cup A_3$ with suitable A_j s and apply Proposition 1.2.4.

Problem 1.10. As in Example 1.4.19, choose 9 balls with replacement out of an urn containing 3 white, 5 red and 4 black balls. How likely is it that among the 9 chosen balls are those of all three colors?

Problem 1.11. Suppose events A and B occur both with probability $1/2$. Prove that then

$$\mathbb{P}(A \cup B) = \mathbb{P}(A^c \cup B^c). \quad (1.87)$$

Does (1.87) remain valid assuming $\mathbb{P}(A) + \mathbb{P}(B) = 1$ instead of $\mathbb{P}(A) = \mathbb{P}(B) = \frac{1}{2}$?

Problem 1.12. Three men and three women sit down randomly on six chairs in a row. Find the probability that the three men and the three women sit side by side. What is the probability that next to each woman sits a man (to the right or to the left)?

Problem 1.13. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Prove the following: Whenever events A_1, A_2, \dots in \mathcal{A} satisfy $\mathbb{P}(A_1) = \mathbb{P}(A_2) = \dots = 1$, then this implies

$$\mathbb{P}\left(\bigcap_{j=1}^{\infty} A_j\right) = 1.$$

Problem 1.14 (Paradox of Chevalier de Méré). Chevalier de Méré mentioned that when rolling three fair indistinguishable dice there are 6 different possibilities for obtaining either 11 or 12 as the sum. Thus he concluded that both events (sum equals 11 or sum equals 12) should be equally likely. But experiments showed that this is not the case. Why was he wrong and what are the correct probabilities for both events?

Problem 1.15. A man has forgotten an important phone number. He only remembers that the seven-digit number contained three times “1,” and “4” and “6” twice each. He dials the seven numbers in random order. Find the probability that he dialed the correct one.

Problem 1.16. In an urn there are n black and m red balls. One successively draws *all* $n + m$ balls (without replacement). What is the probability that the ball chosen last is red?

Problem 1.17. A man has in his pocket n keys to open a door. Only one of the keys fits. He tries the keys one after another until he has chosen the correct one. Given an integer k , compute the probability that the correct key is the one chosen in the k th trial.

Evaluate this probability in each of the two following cases:

- The man always discards wrong keys.
- The man does not discard them, that is, he puts back wrong keys.

Problem 1.18 (Monty Hall problem). At the end of a quiz, the winner has the choice between three doors, say A , B , and C . Behind two of the doors there is a goat, behind the third one is a car. His prize is what is behind the chosen door.

Say the winner has chosen door A . Then the quizmaster (who knows what is behind each of the three doors) opens one of the two remaining doors (in our case either door B or door C) and shows that there is a goat behind it. After that the quizmaster asks the candidate whether or not he wants to revise his decision, that is, for example, if B was opened, to switch from A to C , or if he furthermore chooses door A .

Find the probabilities to win the car in both cases (switching or nonswitching).

Problem 1.19. In a lecture room there are N students. Evaluate the probability that at least two of the students were born on the same day of a year (day and month of their births are the same, but not necessarily the year). Hereby disregard leap years and assume that all days in a year are equally likely. How big must N be in order that this probability is greater than $1/2$?

Hint: Try to evaluate the probability of the complementary event.

Problem 1.20. In an urn there are balls labeled from 0 to 6 so that all numbers are equally likely. Choose successively and with replacement three balls. Find the probability that the three observed numbers sum up to 6.

Problem 1.21. As in Example 1.4.25, six persons enter independently of each other a train with three coaches. How likely is it that no coach remains empty? Find the probability that there are exactly four persons in one of the three coaches.

Problem 1.22. When sending messages from A to B , on average 3% are transmitted falsely. Suppose 300 messages are sent. What is the probability that at least three messages are transmitted falsely? Evaluate the exact probability by using the binomial

distribution as well as the approximate probability by using the Poisson distribution. Compute the probability (exact and approximate one) that all messages arrive correctly.

Problem 1.23 (C. Huygens, 1657). How often does one have to roll two fair dice in order to observe the sum 12 with a probability greater than or equal to $1/2$?

Problem 1.24. Suppose you are given 11 tiles labeled with letters. One tile is labeled with “M,” two with “P” four tiles are labeled with “I,” and, finally, also four with “S.” Order the tiles randomly in a row so that all orders are equally likely. Find the probability to end up with the word “MISSISSIPPI.”

Problem 1.25. The number of accidents in a city per week is assumed to be Poisson distributed with parameter 5. Find the probability that next week there will be either two or three accidents. How likely is it that there will be no accidents?

Problem 1.26. In a room there are 12 men and 8 women. One randomly chooses 5 of the 20 persons. Given $k \in \{0, \dots, 5\}$, what is the probability that among the five chosen are exactly k women? How likely is it that among the five persons are more women than men?

Problem 1.27. Two players A and B take turns rolling a die. The first to roll a “6” wins. Player A starts. Find the probability that A wins. Suppose now there is a third player C and the order of rolling the die is given by $ABCABCA \dots$. Find each player’s probability of winning.

Problem 1.28. Two players, say A and B , toss a biased coin where “head” appears with probability $0 < p < 1$. Whoever gets the first “head” wins. Player A starts, then B tosses twice, then again A once, B twice, and so on. Determine the number p for which the game is fair, that is, the probability that A (or B) wins is $1/2$.

Problem 1.29. In an urn there are 50 white and 200 red balls.

- (1) Take out 10 balls *with* replacement. What is the probability to observe four white balls? Give the exact value via the binomial distribution as well as the approximate one using the related Poisson distribution.
- (2) Next choose 10 balls *without* replacement. What is the probability to get four white balls in this case?
- (3) The number of balls in the urn is as above. But now we choose the balls with replacement until for the first time a white ball shows up. Find the probability of the following events:
 - (a) The first white ball shows up in the fourth trial.
 - (b) The first white ball appears strictly after the third trial.
 - (c) The first white ball is observed in an even number of trials, that is, in the second, or in the fourth, or in the sixth, and so on, trial.

Problem 1.30. Place successively and independently four particles into five boxes. Thereby each box is equally likely. Find the probabilities of the following events:

- $A := \{\text{Each box contains at most one particle}\}$ and
- $B := \{\text{All 4 particles are in the same box}\}$.

Problem 1.31. Four students did not attend at an exam because they were on vacation and drove home too late. Their excuse for missing the test was that they had a flat tire on their way back. The professor tells them: “no problem, you can make up your exam today”. He puts the four students in separate rooms and gives each a sheet of paper with exactly one question: “Which of the four tires was flat?” How likely is it that the four students gave the same answer?

Problem 1.32. Investigate the following generalization of Example 1.4.52: in urn U_0 there are M balls and in urn U_1 there are N balls for some $N, M \geq 1$. Choose U_0 with probability $1 - p$ and U_1 with probability p , and take out a ball from the chosen urn. Given $1 \leq m \leq M$, find the probability that there are m balls left in U_0 when choosing the last ball out of U_1 . How do these probabilities change when $1 \leq m \leq N$, and we assume that there are m balls in U_1 when choosing the last ball from U_0 ?

Problem 1.33. Let $n \in \mathbb{N}$. Use properties of the gamma function to evaluate the following integrals:

$$\int_0^{\infty} x^{2n} e^{-x^2/2} dx \quad \text{and} \quad \int_0^{\infty} x^{2n+1} e^{-x^2/2} dx.$$

Problem 1.34. Prove formula (1.62) that relates the beta and gamma functions.

Hint: Start with

$$\Gamma(x)\Gamma(y) = \int_0^{\infty} \int_0^{\infty} u^{x-1} v^{y-1} e^{-u-v} du dv$$

and change the variables as follows: $u = f(z, t) = zt$ and $v = g(z, t) = z(1 - t)$, where $0 \leq z < \infty$ and $0 \leq t \leq 1$.

Problem 1.35. Prove that for integers n and k with $0 \leq k \leq n$,

$$\binom{n}{k} = \frac{1}{(n+1)B(n-k+1, k+1)}$$

where $B(\cdot, \cdot)$ denotes Euler’s beta function (cf. formula (1.61)).

Problem 1.36. Write $x \in [0, 1)$ as finite or infinite decimal fraction $x = 0.x_1x_2\dots$ with $x_j \in \{0, \dots, 9\}$. Fix some $m \in \{0, \dots, 9\}$ and set

$$A_j = \{x \in [0, 1) : x_j = m\}.$$

That is, A_j contains those real numbers for which the j th digit in its decimal expansion equals m . For example, if $m = 4$, then $x = 0.2534114\dots$ belongs to A_4 and A_7 . Let \mathbb{P} be the uniform distribution on $[0, 1]$. Evaluate

$$\mathbb{P}(A_j) \quad \text{as well as} \quad \mathbb{P}\left(\bigcap_{j=1}^{\infty} A_j\right).$$

Problem 1.37. If $\gamma > 0$ and $x_0 \in \mathbb{R}$, set

$$f_{x_0, \gamma}(x) = \frac{1}{\pi} \left[\frac{\gamma}{(x - x_0)^2 + \gamma^2} \right], \quad x \in \mathbb{R}.$$

Show that $f_{x_0, \gamma}$ is a probability density. Why is the generated probability measure a generalization of the Cauchy distribution as introduced in Definition 1.6.37? Compute the distribution function of the probability measure with density $f_{x_0, \gamma}$.

Problem 1.38. Let $F : \mathbb{R} \rightarrow [0, 1]$ be the distribution function of a probability measure. Show that F possesses at most countably many points of discontinuity. Conclude from this and Proposition 1.7.16 the following: If \mathbb{P} is a probability measure on $\mathcal{B}(\mathbb{R})$, then there are at most countably infinitely many $t \in \mathbb{R}$ such that $\mathbb{P}(\{t\}) > 0$.

Problem 1.39. Let Φ be the distribution function of the standard normal distribution introduced in eq. (1.70). Show the following properties of Φ :

1. $\Phi(0) = \frac{1}{2}$.
2. For $t \in \mathbb{R}$, one has $\Phi(-t) = 1 - \Phi(t)$.
3. If $a > 0$, then

$$\mathcal{N}(0, 1)([-a, a]) = 2\Phi(a) - 1.$$

4. Prove formulas (1.71), that is,

$$\Phi(t) = \frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{t}{\sqrt{2}}\right) \right] \quad \text{and} \quad \operatorname{erf}(t) = 2\Phi(\sqrt{2}t) - 1, \quad t \in \mathbb{R}.$$

5. Compute

$$\lim_{t \rightarrow \infty} \frac{1 - \Phi(t)}{t^{-1}e^{-t^2/2}}.$$

Problem 1.40 (Bertrand paradox). Consider an equilateral triangle inscribed in a circle of radius $r > 0$. Suppose a chord of the circle is chosen at random. What is the probability that the chord is longer than a side of the triangle?

In this form, the problem allows different answers. Why? Because we did not define in which way the random chord is chosen.

1. The “random endpoints” method: Choose independently two uniformly distributed random points on the circumference of the circle and draw the chord joining them.

2. The “random radius” method: Choose a radius of the circle, that is, choose a random angle in $[0, 2\pi]$, choose independently a point on the radius according to the uniform distribution on $[0, r]$, and construct the chord through this point and perpendicular to the radius.
3. The “random midpoint” method: Choose a point within the circle according to the uniform distribution on the circle and construct a chord with the chosen point as its midpoint.

Answer the above question about the length of the chord in each of the three cases.

Problem 1.41. A stick of length $L > 0$ is randomly broken into three pieces. Hereby we assume that both break points are uniformly distributed on $[0, L]$ and independent of each other. What is the probability that these three parts piece together to form a triangle?

2 Conditional probabilities and independence

2.1 Conditional probabilities

In order to motivate the definition of conditional probabilities, let us start with the following easy example.

Example 2.1.1. Roll a fair die twice. The probability of the event “sum of both rolls equals 5” is $1/9$. Suppose now we were told that the first roll was an even number. Does this additional information make the event “sum equals 5” more likely? Or does it even diminish the probability of its occurrence? To answer this question, we apply the so-called technique of “restricting the sample space.” Since we know that the event $B = \{\text{First roll is even}\}$ had occurred, we may rule out elements in B^c and restrict our sample space. Choose B as the new sample space. Its cardinality is 18. Moreover, under this condition, an event A occurs if and only if $A \cap B$ does. Hence, the “new” probability of A under condition B , written $\mathbb{P}(A|B)$, is given by

$$\mathbb{P}(A|B) = \frac{|A \cap B|}{|B|} = \frac{|A \cap B|}{18}. \quad (2.1)$$

In the question above, we asked for $\mathbb{P}(A|B)$, and have

$$A = \{\text{Sum of both rolls equals 5}\} = \{(1, 4), (2, 3), (3, 2), (4, 1)\}.$$

Since $A \cap B = \{(2, 3), (4, 1)\}$, we obtain $\mathbb{P}(A|B) = 2/18 = 1/9$. Consequently, in this case, condition B does not change the probability of the occurrence of A .

Define now A as a set of pairs adding to 6. Then $\mathbb{P}(A) = 5/36$, while the conditional probability remains $1/9$. Note that now $A \cap B = \{(2, 4), (4, 2)\}$. Thus, in this case, condition B makes the occurrence of A less likely.

Before we state the definition of conditional probabilities in the general case, let us rewrite eq. (2.1) as follows:

$$\mathbb{P}(A|B) = \frac{|A \cap B|}{|B|} = \frac{|A \cap B|/36}{|B|/36} = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}. \quad (2.2)$$

Equation (2.2) gives us a hint how to introduce conditional probabilities in the general setting.

Definition 2.1.2. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Given events $A, B \in \mathcal{A}$ with $\mathbb{P}(B) > 0$, the **probability of A under condition B** is defined by

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}. \quad (2.3)$$

Remark 2.1.3. If we know the values of $\mathbb{P}(A \cap B)$ and $\mathbb{P}(B)$, then formula (2.3) allows us to evaluate $\mathbb{P}(A|B)$. Sometimes it happens that we know the values of $\mathbb{P}(B)$ and $\mathbb{P}(A|B)$ and want to calculate $\mathbb{P}(A \cap B)$. In order to do this, we rewrite eq. (2.3) as

$$\mathbb{P}(A \cap B) = \mathbb{P}(B) \mathbb{P}(A|B). \quad (2.4)$$

In this way, we get the desired value of $\mathbb{P}(A \cap B)$. Formula (2.4) is called the **law of multiplication**.

The next two examples show how this law applies.

Example 2.1.4. In an urn there are two white and two black balls. Choose two balls without replacing the first. We want to evaluate the probability of occurrence of a black ball in the first draw *and* of a white in the second. Let us first find a suitable mathematical model that describes this experiment. The sample space is given by $\Omega = \{(b, b), (b, w), (w, b), (w, w)\}$, and we consider the events

$$\begin{aligned} A &= \{\text{Second ball is white}\} = \{(b, w), (w, w)\} \quad \text{and} \\ B &= \{\text{First ball is black}\} = \{(b, b), (b, w)\}. \end{aligned}$$

The event of interest is then $A \cap B = \{(b, w)\}$.

Which probabilities can be directly determined? Of course, the probability of occurrence of B equals $1/2$ because the number of white and black balls is the same. Furthermore, if B had occurred, then in the urn two white balls and one black ball have remained. Under this condition, event A occurs with probability $2/3$, that is, $\mathbb{P}(A|B) = 2/3$. Using eq. (2.4), we obtain

$$\mathbb{P}(\{(b, w)\}) = \mathbb{P}(A \cap B) = \mathbb{P}(B) \cdot \mathbb{P}(A|B) = \frac{1}{2} \cdot \frac{2}{3} = \frac{1}{3}.$$

Example 2.1.5. Among three indistinguishable coins, there are two fair and one biased. Tossing the biased coin, “heads” appears with probability $1/3$, hence “tails” appears with probability $2/3$. We choose at random one of the three coins and toss it. Find the probability to observe “tails” at the biased coin.

To solve this problem, let us first mention that the sample space $\Omega = \{H, T\}$ is not adequate to describe that experiment. Why? Because the event $\{H\}$ may have different probabilities depending on the occurrence using a biased or a fair coin. We have to distinguish between the appearance of “heads” or “tails” for the different types of coin. Hence, an adequate choice of the sample space is

$$\Omega := \{(H, B), (T, B), (H, F), (T, F)\}.$$

Here, B stands for the biased and F for the fair coin. The event of interest is $\{(T, B)\}$. Set

$$T := \{(T, B), (T, F)\} \quad \text{and} \quad B := \{(H, B), (T, B)\}.$$

Then T occurs if “tails” appears regardless of the type of the coin while B occurs if we have chosen the biased coin. Of course, it follows that $\{(T, B)\} = T \cap B$. Since only one of the three coins is biased, we have $\mathbb{P}(B) = 1/3$. By assumption, $\mathbb{P}(T|B) = 2/3$, hence an application of eq. (2.4) leads to

$$\mathbb{P}(\{(T, B)\}) = \mathbb{P}(B) \mathbb{P}(T|B) = \frac{1}{3} \cdot \frac{2}{3} = \frac{2}{9}.$$

Next, we present two examples where formula (2.3) applies directly.

Example 2.1.6. Roll a die twice. One already knows that the first number is not “6”. What is the probability that the sum of both rolls is greater than or equal to “10”?

Answer: The model for this experiment is $\Omega = \{1, \dots, 6\}^2$ endowed with the uniform distribution \mathbb{P} on $\mathcal{P}(\Omega)$. The event $B := \{\text{First result is not “6”}\}$ contains 30 elements, namely

$$\{(1, 1), \dots, (5, 1), \dots, (1, 6), \dots, (5, 6)\},$$

and if A consists of pairs with the sum equal to or larger than 10, then

$$A = \{(4, 6), (5, 6), (6, 6), (5, 5), (6, 5), (6, 4)\}, \quad \text{hence } A \cap B = \{(4, 6), (5, 6), (5, 5)\}.$$

Therefore, it follows that

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{3/36}{30/36} = \frac{1}{10}.$$

In the case that all elementary events are equally likely, there exists a more direct way to evaluate $\mathbb{P}(A|B)$. We reduce the sample space as we already did in Example 2.1.1.

Proposition 2.1.7 (Reduction of the sample space). *Suppose the sample space Ω is finite and let \mathbb{P} be the uniform distribution on $\mathcal{P}(\Omega)$. Then for all events A and a nonempty B in Ω , we have*

$$\mathbb{P}(A|B) = \frac{|A \cap B|}{|B|}. \tag{2.5}$$

Proof. This easily follows from

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{|A \cap B|/|\Omega|}{|B|/|\Omega|} = \frac{|A \cap B|}{|B|}. \quad \square$$

Example 2.1.8. We want to investigate Example 2.1.6 once more, this time using formula (2.5) directly. Since $|A \cap B| = 3$ and $|B| = 30$, we get as before

$$\mathbb{P}(A|B) = \frac{|A \cap B|}{|B|} = \frac{3}{30} = \frac{1}{10}.$$

Remark 2.1.9. It is important to state that Proposition 2.1.7 becomes false for general probabilities \mathbb{P} on $\mathcal{P}(\Omega)$. Formula (2.5) is only valid in the case that \mathbb{P} is the *uniform distribution* on $\mathcal{P}(\Omega)$.

Example 2.1.10. Toss a fair coin 5 times. How likely is it to observe 3 times “1” under the condition that the first toss was a “1”?

Answer: The event A occurs if among the five tosses there are three with “1,” while B occurs provided the first toss is “1.” Then $|B| = 2^4 = 16$ while $|A \cap B| = \binom{4}{2} = 6$. Since by assumption all sequences of zeroes and ones are equally likely, Proposition 2.1.7 applies and leads to

$$\mathbb{P}(A|B) = \frac{|A \cap B|}{|B|} = \frac{6}{16} = \frac{3}{8}.$$

Suppose now that the coin is no longer fair. Say “1” occurs with probability p and “0” with probability $1 - p$. Then we can no longer evaluate the conditional probability by reducing the sample space. In this case one has to apply directly the definition of the conditional probabilities and obtains

$$\mathbb{P}(B) = p, \quad \mathbb{P}(A \cap B) = \binom{4}{2} p^3 (1-p)^2 \quad \Rightarrow \quad \mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = 6 p^2 (1-p)^2.$$

Example 2.1.11. The duration of a telephone call is exponentially distributed with parameter $\lambda > 0$. Find the probability that a call does not last more than 5 minutes provided it already lasted 2 minutes.

Solution: Let A be the event that the call does not last more than 5 minutes, that is, $A = [0, 5]$. We know it already lasted 2 minutes, hence event $B = [2, \infty)$ has occurred. Thus, under condition B , it follows that

$$E_\lambda(A|B) = \frac{E_\lambda(A \cap B)}{E_\lambda(B)} = \frac{E_\lambda([2, 5])}{E_\lambda([2, \infty))} = \frac{e^{-2\lambda} - e^{-5\lambda}}{e^{-2\lambda}} = 1 - e^{-3\lambda}.$$

Note the interesting fact that this conditional probability equals $E_\lambda([0, 3])$. What does this tell us? It says that the probability that a call lasts no more than another 3 minutes is independent of the fact that it has already lasted 2 minutes. This means that the duration of a call has not “become older.” Independent of the fact that it has already lasted 2 minutes, the probability for talking no more than another 3 minutes remains the same.

Let us come back to the general case. Fix an event $B \in \mathcal{A}$ with $\mathbb{P}(B) > 0$. Then

$$A \mapsto \mathbb{P}(A|B), \quad A \in \mathcal{A},$$

is a well-defined mapping from \mathcal{A} to $[0, 1]$. Its main properties are summarized in the next proposition.

Proposition 2.1.12. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be an arbitrary probability space. For each $B \in \mathcal{A}$ with $\mathbb{P}(B) > 0$, the mapping $A \mapsto \mathbb{P}(A|B)$ is a probability measure on \mathcal{A} . It is concentrated on B ,*

that is,

$$\mathbb{P}(B|B) = 1 \quad \text{or, equivalently,} \quad \mathbb{P}(B^c|B) = 0.$$

Proof. Of course, one has

$$\mathbb{P}(\emptyset|B) = \mathbb{P}(\emptyset \cap B)/\mathbb{P}(B) = 0 \quad \text{and} \quad \mathbb{P}(\Omega|B) = \mathbb{P}(\Omega \cap B)/\mathbb{P}(B) = \mathbb{P}(B)/\mathbb{P}(B) = 1.$$

Thus, it remains to prove that $\mathbb{P}(\cdot|B)$ is σ -additive. To this end, choose disjoint A_1, A_2, \dots in \mathcal{A} . Then also $A_1 \cap B, A_2 \cap B, \dots$ are disjoint and, using the σ -additivity of \mathbb{P} leads to

$$\begin{aligned} \mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j|B\right) &= \frac{\mathbb{P}(\left[\bigcup_{j=1}^{\infty} A_j\right] \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(\bigcup_{j=1}^{\infty} (A_j \cap B))}{\mathbb{P}(B)} \\ &= \frac{\sum_{j=1}^{\infty} \mathbb{P}(A_j \cap B)}{\mathbb{P}(B)} = \sum_{j=1}^{\infty} \frac{\mathbb{P}(A_j \cap B)}{\mathbb{P}(B)} = \sum_{j=1}^{\infty} \mathbb{P}(A_j|B). \end{aligned}$$

Consequently, as asserted, $\mathbb{P}(\cdot|B)$ is a probability. Since the identity $\mathbb{P}(B|B) = 1$ is obvious, this ends the proof. \square

Definition 2.1.13. The mapping $\mathbb{P}(\cdot|B)$ is called the **conditional probability** or also **conditional distribution** (under condition B).

Remark 2.1.14. The main advantage of Proposition 2.1.12 is that it implies that conditional probabilities share all the properties of “ordinary” probability measures. For example, it holds that

$$\mathbb{P}(A_2 \setminus A_1|B) = \mathbb{P}(A_2|B) - \mathbb{P}(A_1|B) \quad \text{provided that } A_1 \subseteq A_2,$$

or

$$\mathbb{P}(A_1 \cup A_2|B) = \mathbb{P}(A_1|B) + \mathbb{P}(A_2|B) - \mathbb{P}(A_1 \cap A_2|B).$$

But note, there do not exist similar rules for $\mathbb{P}(A|B)$ independent of the event B and with A fixed.

We come now to the so-called law of total probability. It allows us to evaluate the probability of an event A knowing only its conditional probabilities $\mathbb{P}(A|B_j)$ for certain $B_j \in \mathcal{A}$. More precisely, the following is valid.

Proposition 2.1.15 (Law of total probability). *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and let B_1, \dots, B_n in \mathcal{A} be disjoint with $\mathbb{P}(B_j) > 0$ and $\bigcup_{j=1}^n B_j = \Omega$. Then for each $A \in \mathcal{A}$, one has*

$$\mathbb{P}(A) = \sum_{j=1}^n \mathbb{P}(B_j) \mathbb{P}(A|B_j). \tag{2.6}$$

Proof. Let us start with the investigation of the right-hand side of eq. (2.6). By the definition of the conditional probability, this expression may be rewritten as

$$\sum_{j=1}^n \mathbb{P}(B_j) \mathbb{P}(A|B_j) = \sum_{j=1}^n \mathbb{P}(B_j) \frac{\mathbb{P}(A \cap B_j)}{\mathbb{P}(B_j)} = \sum_{j=1}^n \mathbb{P}(A \cap B_j). \quad (2.7)$$

The sets B_1, \dots, B_n are disjoint, hence so are $A \cap B_1, \dots, A \cap B_n$. Thus, the finite additivity of \mathbb{P} implies

$$\sum_{j=1}^n \mathbb{P}(A \cap B_j) = \mathbb{P}\left(\bigcup_{j=1}^n (A \cap B_j)\right) = \mathbb{P}\left(\left(\bigcup_{j=1}^n B_j\right) \cap A\right) = \mathbb{P}(\Omega \cap A) = \mathbb{P}(A).$$

Together with eq. (2.7), this proves eq. (2.6). \square

A first example illustrates how the law of total probability applies.

Example 2.1.16. Suppose we have n urns, each containing a certain (maybe different) number of white and black balls. Choose urn U_j with probability $\mathbb{P}(U_j)$, $1 \leq j \leq n$, and take out one ball of the chosen urn. Let W occur if the chosen ball is white. Then the law of total probability asserts that

$$\mathbb{P}(W) = \mathbb{P}(U_1)\mathbb{P}(W|U_1) + \dots + \mathbb{P}(U_n)\mathbb{P}(W|U_n),$$

where $\mathbb{P}(W|U_j)$ is the proportion of white balls in urn U_j , $1 \leq j \leq n$. In particular, if all urns are equally likely, then one gets

$$\mathbb{P}(W) = \frac{1}{n} \sum_{j=1}^n \mathbb{P}(W|U_j).$$

Example 2.1.17. A fair coin is tossed four times. Suppose we observe exactly k “heads” for some $k = 0, \dots, 4$. According to the observed k , we take k dice and roll them. Find the probability of the event $A = \{\text{Number “6” does not show up}\}$. Note that $k = 0$ means that we do not roll a die, hence in this case “6” cannot appear.

Solution: As sample space, we choose $\Omega = \{(k, Y), (k, N) : k = 0, \dots, 4\}$, where (k, Y) means that we rolled k dice and at least on one of them we got a “6”. In the same way, (k, N) stands for k dice and no “6”. Let $N = \{(0, N), \dots, (4, N)\}$ and let $B_k = \{(k, Y), (k, N)\}$, $k = 0, \dots, 4$. Then B_k occurs if we observed k “heads.” The conditional probabilities equal

$$\mathbb{P}(N|B_0) = 1, \quad \mathbb{P}(N|B_1) = 5/6, \quad \dots, \quad \mathbb{P}(N|B_4) = (5/6)^4,$$

while

$$\mathbb{P}(B_k) = \binom{4}{k} \frac{1}{2^4}, \quad k = 0, \dots, 4.$$

The events B_0, \dots, B_4 satisfy the assumptions of Proposition 2.1.15, thus Eq. (2.6) applies and leads to

$$\mathbb{P}(A) = \frac{1}{2^4} \sum_{k=0}^4 \binom{4}{k} (5/6)^k = \frac{1}{2^4} \left(\frac{5}{6} + 1 \right)^4 = \left(\frac{11}{12} \right)^4 = 0.706066743.$$

Example 2.1.18. Three different machines, M_1 , M_2 and M_3 , produce light bulbs. In a single day, M_1 produces 500 bulbs, M_2 yields 200, and M_3 provides 100. The quality of the produced bulbs depends on the machines: Among the light bulbs produced by M_1 , 5 % are defective; among those from M_2 , 10 % are defective; and only 2 % are defective from M_3 . At the end of a day, a controller chooses 1 of the 800 produced light bulbs at random and tests it. Determine the probability that the checked bulb is defective.

Solution: The probabilities that the checked bulb was produced by M_1 , M_2 , or M_3 are $5/8$, $1/4$, and $1/8$, respectively. The conditional probabilities for choosing a defective bulb produced by M_1 , M_2 or M_3 were given as $1/20$, $1/10$, and $1/50$, respectively. If D is the event that the tested bulb was defective, then the law of total probability yields

$$\mathbb{P}(D) = \frac{5}{8} \cdot \frac{1}{20} + \frac{1}{4} \cdot \frac{1}{10} + \frac{1}{8} \cdot \frac{1}{50} = \frac{47}{800} = 0.05875.$$

Let us look at Example 2.1.18 from a different point of view. When choosing a light bulb out of the 800 produced, there were certain fixed probabilities of whether it was produced by M_1 , M_2 , or M_3 , namely $5/8$, $1/4$, and $1/8$. These are the probabilities *before* checking a bulb. Therefore, they are called *a priori* probabilities. After checking a bulb, we obtained additional information that it was defective. Does this additional information change the probabilities which of the machines M_1 , M_2 , and M_3 produced it? More precisely, if the D above occurs when the tested bulb is defective, then we now ask for the conditional probabilities $\mathbb{P}(M_1|D)$, $\mathbb{P}(M_2|D)$, and $\mathbb{P}(M_3|D)$. To understand that these probabilities may differ considerably from the *a priori* probabilities, imagine that, for example, M_1 produces almost no defective bulbs. Then it will be very unlikely that the tested bulb has been produced by M_1 , although $\mathbb{P}(M_1)$ may be big.

Because $\mathbb{P}(M_1|D)$, $\mathbb{P}(M_2|D)$, and $\mathbb{P}(M_3|D)$ are the probabilities *after* executing the random experiment (choosing and testing the bulb), they are called *a posteriori* probabilities.

Let us now introduce the exact and general definition of *a priori* and *a posteriori* probabilities.

Definition 2.1.19. Suppose there is a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and there are disjoint events $B_1, \dots, B_n \in \mathcal{A}$ satisfying $\Omega = \bigcup_{j=1}^n B_j$. Then we call $\mathbb{P}(B_1), \dots, \mathbb{P}(B_n)$ the **a priori** probabilities of B_1, \dots, B_n . Let $A \in \mathcal{A}$ with $\mathbb{P}(A) > 0$ be given. Then the conditional probabilities $\mathbb{P}(B_1|A), \dots, \mathbb{P}(B_n|A)$ are said to be the **a posteriori** probabilities, that is, those after the occurrence of A .

To calculate the *a posteriori* probabilities, the next proposition turns out to be very useful.

Proposition 2.1.20 (Bayes' formula). *Suppose we are given disjoint events B_1 to B_n satisfying $\bigcup_{j=1}^n B_j = \Omega$ and $\mathbb{P}(B_j) > 0$. Let A be an event with $\mathbb{P}(A) > 0$. Then for each $j \leq n$ the following equation holds:*

$$\mathbb{P}(B_j|A) = \frac{\mathbb{P}(B_j) \mathbb{P}(A|B_j)}{\sum_{i=1}^n \mathbb{P}(B_i) \mathbb{P}(A|B_i)}. \quad (2.8)$$

Proof. Proposition 2.1.15 implies

$$\sum_{i=1}^n \mathbb{P}(B_i) \mathbb{P}(A|B_i) = \mathbb{P}(A).$$

Hence, the right-hand side of eq. (2.8) may also be written as

$$\frac{\mathbb{P}(B_j) \mathbb{P}(A|B_j)}{\mathbb{P}(A)} = \frac{\mathbb{P}(B_j) \frac{\mathbb{P}(A \cap B_j)}{\mathbb{P}(B_j)}}{\mathbb{P}(A)} = \frac{\mathbb{P}(A \cap B_j)}{\mathbb{P}(A)} = \mathbb{P}(B_j|A),$$

and the proposition is proven. \square

Remark 2.1.21. In the case $\mathbb{P}(A)$ is already known, Bayes' formula simplifies to

$$\mathbb{P}(B_j|A) = \frac{\mathbb{P}(B_j) \mathbb{P}(A|B_j)}{\mathbb{P}(A)}, \quad j = 1, \dots, n. \quad (2.9)$$

Remark 2.1.22. Let us treat the special case of two sets partitioning Ω . If $B_1 = B$, then necessarily $B_2 = B^c$, hence $\Omega = B \cup B^c$. Then formula (2.8) looks as follows:

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(B) \mathbb{P}(A|B)}{\mathbb{P}(B) \mathbb{P}(A|B) + \mathbb{P}(B^c) \mathbb{P}(A|B^c)} \quad (2.10)$$

and

$$\mathbb{P}(B^c|A) = \frac{\mathbb{P}(B^c) \mathbb{P}(A|B^c)}{\mathbb{P}(B) \mathbb{P}(A|B) + \mathbb{P}(B^c) \mathbb{P}(A|B^c)}. \quad (2.11)$$

Again, if the probability of A is known, the denominators in Eqs. (2.10) and (2.11) may be replaced by $\mathbb{P}(A)$.

Example 2.1.23. Let us use Bayes' formula to calculate the *a posteriori* probabilities in Example 2.1.18. Recall that D occurs if the tested bulb is defective. We already know $\mathbb{P}(D) = 47/800$, hence we may apply eq. (2.9). Doing so, we get

$$\begin{aligned} \mathbb{P}(M_1|D) &= \frac{\mathbb{P}(M_1) \mathbb{P}(D|M_1)}{\mathbb{P}(D)} = \frac{5/8 \cdot 1/20}{47/800} = 25/47, \\ \mathbb{P}(M_2|D) &= \frac{\mathbb{P}(M_2) \mathbb{P}(D|M_2)}{\mathbb{P}(D)} = \frac{1/4 \cdot 1/10}{47/800} = 20/47, \\ \mathbb{P}(M_3|D) &= \frac{\mathbb{P}(M_3) \mathbb{P}(D|M_3)}{\mathbb{P}(D)} = \frac{1/8 \cdot 1/50}{47/800} = 2/47. \end{aligned}$$

By assignment of the problem, the *a priori* probabilities were given by $\mathbb{P}(M_1) = 5/8$, $\mathbb{P}(M_2) = 1/4$, and $\mathbb{P}(M_3) = 1/8$. In the case that the tested light bulb was defective, these probabilities changed to $25/47$, $20/47$, and $2/47$. This tells us that it becomes less likely that the tested bulb was produced by M_1 or M_3 ; their probabilities diminish by 0.0930851 and 0.0824468 , respectively. On the other hand, the probability of M_2 increases by 0.175532 .

Finally, note that Proposition 2.1.12 implies that the sum of the *a posteriori* probabilities has to be 1. Because of $25/47 + 20/47 + 2/47 = 1$, this is true in that example.

Example 2.1.24. In order to figure out whether or not a person suffers from a certain disease, say disease X , a test is assumed to give a clue. If the tested person is sick, then the test is positive in 96 % of cases. If the person is well, then with 94 % accuracy the test will be negative. Furthermore, it is known that 0.4 % of the population suffers from the disease X .

Now a person, chosen at random, is tested. Suppose the result was positive. Find the probability that this person really suffers from X .

Solution: As sample space, we may choose

$$\Omega = \{(X, p), (X, n), (X^c, p), (X^c, n)\},$$

where, for example, (X, n) means that the person suffers from X and the test was negative. Set $A := \{(X, p), (X^c, p)\}$. Then A occurs if and only if the test turned out to be positive. Furthermore, event $B := \{(X, p), (X, n)\}$ occurs in the case that the tested person suffers from X . Known are

$$\mathbb{P}(A|B) = 0.96, \quad \mathbb{P}(A|B^c) = 0.06, \quad \text{and} \quad \mathbb{P}(B) = 0.004, \quad \text{hence} \quad \mathbb{P}(B^c) = 0.996.$$

Therefore, by eq. (2.10), the probability we asked for can be calculated as follows:

$$\begin{aligned} \mathbb{P}(B|A) &= \frac{\mathbb{P}(B)\mathbb{P}(A|B)}{\mathbb{P}(B)\mathbb{P}(A|B) + \mathbb{P}(B^c)\mathbb{P}(A|B^c)} \\ &= \frac{0.004 \cdot 0.96}{0.004 \cdot 0.96 + 0.996 \cdot 0.06} = \frac{0.00384}{0.0636} = 0.0603774. \end{aligned}$$

This tells us that it is quite unlikely that a randomly chosen person with a positive test is really sick. The chance for this being true is only about 6 %.

Summary: Given two events A and B in a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with $\mathbb{P}(B) > 0$, the probability of A under the condition of the occurrence of B is given by

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \quad \text{or, equivalently, by} \quad \mathbb{P}(A \cap B) = \mathbb{P}(B) \mathbb{P}(A|B).$$

The basic properties of the conditional probability are summarized in the “Law of total probability” and in “Bayes’ formula” (Propositions 2.1.15 and 2.1.20).

2.2 Independence of events

What does it mean that two events are independent or, more precisely, that they occur independently of each other? To get an idea, let us look at the following example.

Example 2.2.1. Roll a fair die twice. Event B occurs if the first number is even while event A consists of all pairs (x_1, x_2) , where $x_1 = 5$ or $x_1 = 6$. It is intuitively clear that these two events occur independently of each other. But how to express this mathematically? To answer this question, think about the probability of A under the condition B . The fact whether or not B has occurred has no influence on the occurrence of A . For the occurrence or nonoccurrence of A , it is completely insignificant what happened in the first roll. Mathematically, this means that $\mathbb{P}(A|B) = \mathbb{P}(A)$. Let us check whether this is true in this concrete case. Indeed, it holds that $\mathbb{P}(A) = 1/3$, as well as

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{6/36}{1/2} = 1/3.$$

The previous example suggests that the independence of A of B could be described by

$$\mathbb{P}(A) = \mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}. \quad (2.12)$$

But formula (2.12) has a disadvantage, namely we have to assume $\mathbb{P}(B) > 0$ to ensure that $\mathbb{P}(A|B)$ exists. To overcome this problem, rewrite eq. (2.12) as

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B). \quad (2.13)$$

In this form, we may take eq. (2.13) as the basis for the definition of independence.

Definition 2.2.2. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Two events A and B in \mathcal{A} are said to be (stochastically) **independent** provided

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \cdot \mathbb{P}(B). \quad (2.14)$$

In the case that eq. (2.14) does not hold, the events A and B are called (stochastically) **dependent**.

Remark 2.2.3. In the sequel, we use the notations “independent” and “dependent” without adding the word “stochastically.” Since we will not use other versions of independence, there should be no confusion.

Example 2.2.4. A fair die is rolled twice. Event A occurs if the first roll is either “1” or “2” while B occurs if the sum of both rolls equals 7. Are A and B independent?

Answer: We have $\mathbb{P}(A) = 1/3$, $\mathbb{P}(B) = 1/6$, and $\mathbb{P}(A \cap B) = 2/36 = 1/18$. Hence, we get $\mathbb{P}(A \cap B) = \mathbb{P}(A) \cdot \mathbb{P}(B)$ and so A and B are independent.

Question: Are A and B also independent if A is as before and B is defined as a set of pairs with sum 4?

Example 2.2.5. In an urn, there are $n \geq 2$ white balls and also n black balls. One chooses two balls without replacing the first. Let A be the event that the second ball is black while B occurs if the first ball was white. Are A and B independent?

Answer: The probability of B equals $1/2$. To calculate $\mathbb{P}(A)$, we use Proposition 2.1.15. Then we get

$$\mathbb{P}(A) = \mathbb{P}(B)\mathbb{P}(A|B) + \mathbb{P}(B^c)\mathbb{P}(A|B^c) = \frac{1}{2} \cdot \frac{n}{2n-1} + \frac{1}{2} \cdot \frac{n-1}{2n-1} = \frac{1}{2},$$

hence $\mathbb{P}(A) \cdot \mathbb{P}(B) = 1/4$.

On the other hand, we have

$$\mathbb{P}(A \cap B) = \mathbb{P}(B)\mathbb{P}(A|B) = \frac{1}{2} \cdot \frac{n}{2n-1} = \frac{n}{4n-2} \neq \frac{1}{4}.$$

Consequently, A and B are dependent.

Remark 2.2.6. Note that if $n \rightarrow \infty$, then

$$\mathbb{P}(A \cap B) = \frac{n}{4n-2} \rightarrow \frac{1}{4} = \mathbb{P}(A) \mathbb{P}(B).$$

This tells us the following: if n is big, then A and B are “almost” independent or, equivalently, the degree of dependence between A and B is very small. This question will be investigated more thoroughly in Chapter 5 when a measure for the degree of dependence is available.

Next, we prove some properties of independent events.

Proposition 2.2.7. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space.*

1. *For any $A \in \mathcal{A}$, the events A and \emptyset , as well as A and Ω , are independent.¹*
2. *If A and B are independent, then so are A and B^c , as well as A^c and B^c .*

Proof. We have

$$\mathbb{P}(A \cap \emptyset) = \mathbb{P}(\emptyset) = 0 = \mathbb{P}(A) \cdot 0 = \mathbb{P}(A) \cdot \mathbb{P}(\emptyset),$$

hence A and \emptyset are independent.

In the same way, the independence of A and Ω follows from

$$\mathbb{P}(A \cap \Omega) = \mathbb{P}(A) = \mathbb{P}(A) \cdot 1 = \mathbb{P}(A) \cdot \mathbb{P}(\Omega).$$

To prove the second part, assume that A and B are independent. Our aim is to show that A and B^c are independent as well. We know that

¹ For a more general result, see Problem 2.15.

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$$

and want to show that

$$\mathbb{P}(A \cap B^c) = \mathbb{P}(A) \mathbb{P}(B^c).$$

Let us start with the right-hand side of the latter equation. Using the independence of A and B and the fact $A \cap B \subseteq B$, it follows that

$$\begin{aligned} \mathbb{P}(A) \mathbb{P}(B^c) &= \mathbb{P}(A)(1 - \mathbb{P}(B)) = \mathbb{P}(A) - \mathbb{P}(A) \cdot \mathbb{P}(B) \\ &= \mathbb{P}(A) - \mathbb{P}(A \cap B) = \mathbb{P}(A \setminus (A \cap B)). \end{aligned} \quad (2.15)$$

Since $A \setminus (A \cap B) = A \setminus B = A \cap B^c$, using eq. (2.15), we derive

$$\mathbb{P}(A) \cdot \mathbb{P}(B^c) = \mathbb{P}(A \cap B^c).$$

Consequently, as asserted, A and B^c are independent.

If A and B are independent, then so are B and A , and as seen above, so are B and A^c . Another application of the first step, this time with A^c and B , shows that also A^c and B^c are independent. This completes the proof. \square

Suppose we are given n events A_1, \dots, A_n in \mathcal{A} . We want to figure out when they are independent. A first possible approach could be as follows.

Definition 2.2.8. Events A_1, \dots, A_n are said to be **pairwise independent** if, whenever $i \neq j$,

$$\mathbb{P}(A_i \cap A_j) = \mathbb{P}(A_i) \cdot \mathbb{P}(A_j).$$

In other words, for all $1 \leq i < j \leq n$ the events A_i and A_j are independent.

Unfortunately, for many purposes, the property of pairwise independence is too weak. For example, as we will see next, in general it does not imply the important equation

$$\mathbb{P}(A_1 \cap \dots \cap A_n) = \mathbb{P}(A_1) \cdot \dots \cdot \mathbb{P}(A_n). \quad (2.16)$$

Example 2.2.9. Roll a die twice and define events A_1, A_2 , and A_3 as follows:

$$A_1 := \{2, 4, 6\} \times \{1, \dots, 6\},$$

$$A_2 := \{1, \dots, 6\} \times \{1, 3, 5\},$$

$$A_3 := \{2, 4, 6\} \times \{1, 3, 5\} \cup \{1, 3, 5\} \times \{2, 4, 6\}.$$

Verbally this says that A_1 occurs if the first roll is even, A_2 occurs if the second one is odd, and A_3 occurs if either the first number is odd while the second is even or vice versa.

Direct calculations give $\mathbb{P}(A_1) = \mathbb{P}(A_2) = \mathbb{P}(A_3) = 1/2$, as well as

$$\mathbb{P}(A_1 \cap A_2) = \mathbb{P}(A_1 \cap A_3) = \mathbb{P}(A_2 \cap A_3) = \frac{1}{4}.$$

Hence, A_1 , A_2 , and A_3 are pairwise independent.

Since

$$A_1 \cap A_2 \cap A_3 = A_1 \cap A_2,$$

it follows that

$$\mathbb{P}(A_1 \cap A_2 \cap A_3) = \mathbb{P}(A_1 \cap A_2) = \frac{1}{4} \neq \frac{1}{8} = \mathbb{P}(A_1) \cdot \mathbb{P}(A_2) \cdot \mathbb{P}(A_3).$$

So, we found three pairwise independent events for which eq. (2.16) is not valid.

After mentioning that pairwise independence of A_1, \dots, A_n does not imply

$$\mathbb{P}(A_1 \cap \dots \cap A_n) = \mathbb{P}(A_1) \cdots \mathbb{P}(A_n), \quad (2.17)$$

it makes sense to ask whether or not pairwise independence can be derived from eq. (2.17). The next example shows that, in general, this is also not true.

Example 2.2.10. Let $\Omega = \{1, \dots, 12\}$ be endowed with the uniform distribution \mathbb{P} , that is, for any $A \subseteq \Omega$ we have $\mathbb{P}(A) = |A|/12$. Define events A_1, A_2 , and A_3 as $A_1 := \{1, \dots, 9\}$, $A_2 := \{6, 7, 8, 9\}$, and $A_3 := \{9, 10, 11, 12\}$. Direct calculations give

$$\mathbb{P}(A_1) = \frac{9}{12} = \frac{3}{4}, \quad \mathbb{P}(A_2) = \frac{4}{12} = \frac{1}{3}, \quad \text{and} \quad \mathbb{P}(A_3) = \frac{4}{12} = \frac{1}{3}.$$

Moreover, we have

$$\mathbb{P}(A_1 \cap A_2 \cap A_3) = \mathbb{P}(\{9\}) = \frac{1}{12} = \frac{3}{4} \cdot \frac{1}{3} \cdot \frac{1}{3} = \mathbb{P}(A_1) \cdot \mathbb{P}(A_2) \cdot \mathbb{P}(A_3),$$

hence eq. (2.17) is valid. But, because of

$$\mathbb{P}(A_1 \cap A_2) = \mathbb{P}(A_2) = \frac{1}{3} \neq \frac{1}{4} = \mathbb{P}(A_1) \cdot \mathbb{P}(A_2),$$

the events A_1, A_2 , and A_3 are **not** pairwise independent.

Remark 2.2.11. Summing up, Examples 2.2.9 and 2.2.10 show that neither pairwise independence nor eq. (2.17) are suitable to define the independence of more than two events. Why? On the one hand, independence should yield eq. (2.17) and, on the other hand, whenever A_1, \dots, A_n are independent, then so should be any subcollection of them. In particular, independence should imply pairwise independence.

A reasonable definition of independence of n events is as follows.

Definition 2.2.12. The events A_1, \dots, A_n are said to be **independent** provided that for each subset $I \subseteq \{1, \dots, n\}$ we have

$$\mathbb{P}\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} \mathbb{P}(A_i). \quad (2.18)$$

Remark 2.2.13. Of course, it suffices that eq. (2.18) is valid for sets $I \subseteq \{1, \dots, n\}$ satisfying $|I| \geq 2$. Indeed, if $|I| = 1$, then eq. (2.18) holds trivially.

Remark 2.2.14. Another way to introduce independence is as follows: For all $m \geq 2$ and $1 \leq i_1 < \dots < i_m \leq n$, it follows that

$$\mathbb{P}(A_{i_1} \cap \dots \cap A_{i_m}) = \mathbb{P}(A_{i_1}) \cdots \mathbb{P}(A_{i_m}).$$

Identify I with $\{i_1, \dots, i_m\}$ to see that both definitions are equivalent.

At a first glance, the previous Definition 2.2.12 looks complicated; in fact, it is not. To see this, let us once more investigate the case $n = 3$. Here exist exactly four different subsets $I \subseteq \{1, 2, 3\}$ with $|I| \geq 2$. These are $I = \{1, 2\}$, $I = \{1, 3\}$, $I = \{2, 3\}$, and $I = \{1, 2, 3\}$. Consequently, three events A_1, A_2 , and A_3 are independent if and only if the four following conditions hold **at once**:

$$\begin{aligned} \mathbb{P}(A_1 \cap A_2) &= \mathbb{P}(A_1) \cdot \mathbb{P}(A_2), \\ \mathbb{P}(A_1 \cap A_3) &= \mathbb{P}(A_1) \cdot \mathbb{P}(A_3), \\ \mathbb{P}(A_2 \cap A_3) &= \mathbb{P}(A_2) \cdot \mathbb{P}(A_3), \quad \text{as well as} \\ \mathbb{P}(A_1 \cap A_2 \cap A_3) &= \mathbb{P}(A_1) \cdot \mathbb{P}(A_2) \cdot \mathbb{P}(A_3). \end{aligned}$$

Examples 2.2.9 and 2.2.10 show that all four equations are really necessary. None of them is a consequence of the other three.

The independence of n events provides the following properties:

Proposition 2.2.15.

1. Let A_1, \dots, A_n be independent. For any $J \subseteq \{1, \dots, n\}$, the events $\{A_j : j \in J\}$ are independent as well. In particular, independence implies pairwise independence.
2. For each permutation π of $\{1, \dots, n\}$, the independence of A_1, \dots, A_n implies that of $A_{\pi(1)}, \dots, A_{\pi(n)}$.
3. Suppose for each $1 \leq j \leq n$ either $B_j = A_j$ or $B_j = A_j^c$ holds. Then the independence of A_1, \dots, A_n implies that of B_1, \dots, B_n .

Proof. The first two properties are an immediate consequence of the definition of independence.

2 For example, in the case $n = 3$ with A_1, A_2, A_3 also A_3, A_2, A_1 and A_2, A_3, A_1 are independent.

To prove the third assertion, reorder A_1, \dots, A_n such that³ $B_1 = A_1^c$. In the first step, we show that A_1^c, A_2, \dots, A_n are independent as well, that is, we have $B_1 = A_1^c, B_2 = A_2$, and so on. Given $I \subseteq \{1, \dots, n\}$, it has to hold that

$$\mathbb{P}\left(\bigcap_{i \in I} B_i\right) = \prod_{i \in I} \mathbb{P}(B_i).$$

In the case $1 \notin I$, this follows by the independence of A_1, \dots, A_n . If $1 \in I$, we apply Proposition 2.2.7 with⁴ A_1 and $C = \bigcap_{i \in I \setminus \{1\}} A_i = \bigcap_{i \in I \setminus \{1\}} B_i$. Then $A_1^c = B_1$ and C are independent as well. Hence, by the independence of A_2, \dots, A_n , we get

$$\mathbb{P}\left(\bigcap_{i \in I} B_i\right) = \mathbb{P}(B_1 \cap C) = \mathbb{P}(B_1) \cdot \mathbb{P}(C) = \mathbb{P}(B_1) \cdot \prod_{i \in I \setminus \{1\}} \mathbb{P}(B_i) = \prod_{i \in I} \mathbb{P}(B_i).$$

The general case then follows by reordering the A_j s and by an iterative application of the first step. This is exactly the procedure we did in the proof of Proposition 2.2.7 when verifying the independence of A^c and B^c for independent A and B . \square

The next two examples show how independence of more than two events appears in a natural way.

Example 2.2.16. Toss a fair coin n times. Let us assume that the coin is labeled with “0” and “1.” Choose a fixed sequence $(a_j)_{j=1}^n$ of numbers in $\{0, 1\}$ and suppose that the event A_j occurs if in the j th trial a_j comes up.

We claim now that A_1, \dots, A_n are independent. To verify this, choose a subset $I \subseteq \{1, \dots, n\}$ with $|I| = k$ for some $k = 2, \dots, n$. The cardinality of $\bigcap_{i \in I} A_i$ equals 2^{n-k} . Why? At k positions, the values of the tosses are fixed; at $n - k$ positions, they still may be either “0” or “1”. Consequently,

$$\mathbb{P}\left(\bigcap_{i \in I} A_i\right) = \frac{2^{n-k}}{2^n} = 2^{-k}. \quad (2.19)$$

The same argument as before gives $|A_j| = 2^{n-1}$, hence $\mathbb{P}(A_j) = 1/2, 1 \leq j \leq n$. Consequently, it follows that

$$\prod_{i \in I} \mathbb{P}(A_i) = \left(\frac{1}{2}\right)^{|I|} = 2^{-k}. \quad (2.20)$$

Combining Eqs. (2.19) and (2.20) gives

³ If all $B_j = A_j$, there is nothing to prove.

⁴ Why are A_1 and C independent? Give a short proof.

$$\mathbb{P}\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} \mathbb{P}(A_i),$$

and since I was arbitrary, the sets A_1, \dots, A_n are independent.

Remark 2.2.17. Even the simple Example 2.2.16 shows that it might be rather complicated to verify the independence of n given events. For example, if we modify the previous example by taking a biased coin, then the A_j s remain independent, but the proof becomes more complicated.

Example 2.2.18. A machine consists of n components. These components break down with certain probabilities p_1, \dots, p_n . Moreover, we assume that they break down independently of each other. Find the probability that a chosen machine stops working. Before answering this question, we have to determine the conditions.

Case 1: The machine stops working provided at least one component breaks down.

Let M be the event that the machine stops working. If $j \leq n$, assume A_j occurs if component j breaks down. By assumption, $\mathbb{P}(A_j) = p_j$. Since

$$M = \bigcup_{j=1}^n A_j,$$

by the independence⁵ it follows that

$$\mathbb{P}(M) = 1 - \mathbb{P}(M^c) = 1 - \mathbb{P}\left(\bigcap_{j=1}^n A_j^c\right) = 1 - \prod_{j=1}^n \mathbb{P}(A_j^c) = 1 - \prod_{j=1}^n (1 - p_j). \quad (2.21)$$

Case 2: The machine stops working provided all n components break down.

Using the same notation as in case 1, we now have

$$M = \bigcap_{j=1}^n A_j.$$

Hence, by the independence we obtain

$$\mathbb{P}(M) = \mathbb{P}\left(\bigcap_{j=1}^n A_j\right) = \prod_{j=1}^n p_j. \quad (2.22)$$

Remark 2.2.19. Formula (2.21) tells us the following: If among the n components there is one of bad quality, say component j_0 , then p_{j_0} is close to one; hence, $1 - p_{j_0}$ is close to zero, and so is $\prod_{j=1}^n (1 - p_j)$. Because of eq. (2.21), $\mathbb{P}(M)$ is large, and so the machine breaks down with a large probability.

⁵ In fact, we also have to use Proposition 2.2.15.

In the second case, the conclusion is as follows: if among the n components there is one of high quality, say component j_0 , then p_{j_0} is small and so is $\prod_{j=1}^n p_j$. By eq. (2.22), $\mathbb{P}(M)$ is also small, hence it is very unlikely that the machine stops working.

Summary: Two events A and B in a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ are said to be (stochastically) independent if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \cdot \mathbb{P}(B).$$

In general, events A_1, \dots, A_n are (stochastically) independent provided that for all $2 \leq m \leq n$ and all choices of indices $1 \leq i_1 < \dots < i_m \leq n$ one has

$$\mathbb{P}(A_{i_1} \cap \dots \cap A_{i_m}) = \mathbb{P}(A_{i_1}) \cdots \mathbb{P}(A_{i_m}).$$

In particular, then the A_j s are pairwise independent and $\mathbb{P}(A_1 \cap \dots \cap A_n) = \mathbb{P}(A_1) \cdots \mathbb{P}(A_n)$.

2.3 Problems

Problem 2.1 (Willem Jacob's Gravesande, 1736). On a ship there are 84 passengers from Belgium, 12 from England, and 4 from Germany.

- One passenger leaves the ship. How likely is it that he or she is
 - Belgian,
 - German,
 - Belgian or German?
- Two people leave the ship. How likely is it that at least one of them is Belgian?

Problem 2.2. The chance to win a certain game is 50%. One plays six games. Find the probability to win exactly four games. Evaluate the probability of this event under the condition to win at least two games. Suppose one had won exactly one of the two first games. Which probability does the event "winning 4 games" have under this condition?

Problem 2.3. Toss a fair coin six times. Define events A and B as follows:

$$A = \{\text{"Heads" appears exactly 3 times}\},$$

$$B = \{\text{The first and the second toss are "heads"}\}.$$

Evaluate $\mathbb{P}(A)$, $\mathbb{P}(A|B)$, and $\mathbb{P}(A|B^c)$.

Problem 2.4. Let A and B be as in Problem 1.30, that is, A occurs if each box contains at most one particle while B occurs if all four particles are in the same box.

Find now $\mathbb{P}(A|C)$ and $\mathbb{P}(B|C)$ with $C = \{\text{The first box remains empty}\}$.

Problem 2.5. Justify why Propositions 2.1.15 and 2.1.20 (Law of total probability and Bayes' formula) remain valid for *infinitely many* disjoint sets B_1, B_2, \dots satisfying $\mathbb{P}(B_j) > 0$ and $\bigcup_{j=1}^{\infty} B_j = \Omega$.

Prove that Proposition 2.1.15 also holds without assuming $\bigcup_{j=1}^n B_j = \Omega$. But then we have to suppose $A \subseteq \bigcup_{j=1}^n B_j$.

Problem 2.6. To go to work, a man can either use a train, a bus, or his car. He chooses the train 50 %, the bus 30 %, and the car 20 % of workdays. If he takes the train, he arrives on time with probability 0.95. By bus, he is on time with probability 0.8, and by car with probability 0.7.

- Evaluate the probability that the man is at work on time.
- How big is this probability given the man does *not* use the car?
- Assume the man arrived at work on time. What are then the probabilities that he came by train, bus, or car?

Problem 2.7. Let U_1 , U_2 , and U_3 be three urns containing five balls each. Urn U_1 contains four white balls and one black ball, U_2 has three white balls and two black balls and, finally, U_3 contains two white balls and three black balls. Choose one urn at random (each urn is equally likely) and, without replacing the first ball, take two balls out of the chosen urn.

- Give a suitable sample space for this random experiment.
- Find the probability to observe two balls of different color.
- Assume the chosen balls were of different color. What are the probabilities that the balls were taken out of U_1 , U_2 , or U_3 ?

Problem 2.8. Suppose we have three indistinguishable dice. Two of them are fair, the remaining one is biased. For the latter, the number “6” appears with probability $1/5$ while all other numbers have probability $4/25$. We choose at random one of the dice and roll it.

- Find a suitable sample space for the description of this experiment.
- Give the probability of occurrence of $\{1\}$ to $\{6\}$ in that experiment.
- Suppose we have observed the number “2” on the chosen die. Find the probability that this die was the biased one.

Problem 2.9 (P. S. Laplace, 1774). Two urns U_1 and U_2 contain w_1 white and b_1 black balls and w_2 white and b_2 black balls, respectively. Choose one of the two urns at random and take out n balls without replacement. Among the n chosen balls, w are white and b are black. Find a formula for the probability that we took off the balls from U_1 ?

Find the numerical value for the likelihood of U_1 in the case that there are 8 white and 7 black balls in U_1 , 5 white and 15 black balls in U_2 , and that we chose 6 balls, with 4 of them white, hence 2 black.

Problem 2.10 (P. S. Laplace, 1786). In an urn there are three balls which are known to be either white or black. Choosing n balls with replacement, we observe that all of them were white. How likely is it that 0, 1, or 2 of the three balls are black? How likely is it that the next chosen ball, the $(n + 1)$ th, is white as well?

Problem 2.11. Let A and B be two events in a probability space with $\mathbb{P}(A) > 0$ and $\mathbb{P}(B) > 0$. Under which conditions, do we have

$$\mathbb{P}(A|B) = \mathbb{P}(B|A) ?$$

Problem 2.12. Let A and B certain events in a probability space with $\mathbb{P}(B) > 0$. Do we have

$$\mathbb{P}(A \cup B|B) = \mathbb{P}(A) \quad \text{and/or} \quad \mathbb{P}(A \cap B|B) = \mathbb{P}(A \cap B) ?$$

Give a proof or a counterexample.

Problem 2.13.

(a) Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Suppose we are given events A_1, \dots, A_n with $\mathbb{P}(A_1 \cap \dots \cap A_{n-1}) > 0$. Prove the following *chain rule* for conditional probabilities:

$$\mathbb{P}(A_1 \cap \dots \cap A_n) = \mathbb{P}(A_1)\mathbb{P}(A_2|A_1) \cdots \mathbb{P}(A_n|A_1 \cap A_2 \cap \dots \cap A_{n-1}).$$

Argue why all occurring conditional probabilities are well defined.

- (b) Choose at random three numbers from 1 to 10 without replacement. Find the probability that the first number is even, the second one is odd, and the third one is again even.
- (c) Compare this probability with that of the following event: among three randomly chosen numbers in $\{1, \dots, 10\}$, there are exactly two even and one odd.

Problem 2.14. Three persons, say X , Y , and Z , stand randomly in a row. All orderings are assumed to be equally likely. Event A occurs if Y stands on the right-hand side of X while B occurs in the case that Z is on the right-hand side of X . Hereby, we do not suppose that Y and X or that Z and X stand directly next to each other. Are events A and B independent or dependent?

Problem 2.15. Prove the following generalization of part 1 in Proposition 2.2.7. Let $A \in \mathcal{A}$ be an event with either $\mathbb{P}(A) = 0$ or $\mathbb{P}(A) = 1$. Then for any $B \in \mathcal{A}$, the events A and B are independent.

Problem 2.16. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Given independent events A_1, \dots, A_n in \mathcal{A} , prove that

$$\mathbb{P}\left(\bigcup_{j=1}^n A_j\right) = 1 - \prod_{j=1}^n (1 - \mathbb{P}(A_j)). \quad (2.23)$$

Use $1 - x \leq e^{-x}$, $x \geq 0$, to derive from eq. (2.23) the following:

If independent events⁶ A_1, A_2, \dots satisfy $\sum_{j=1}^{\infty} \mathbb{P}(A_j) = \infty$, then $\mathbb{P}(\bigcup_{j=1}^{\infty} A_j) = 1$.

Problem 2.17. An electric circuit (see Fig. 2.1) contains four switches A, B, C , and D . Each of the switches is independently open or closed (then electricity flows). The switches are

⁶ Compare with Definition 7.1.17: for each $n \in \mathbb{N}$, the events A_1, \dots, A_n are independent.

open with probability $1 - p$ and closed with probability p . Here, $0 \leq p \leq 1$ is given. Find the probability that electricity flows from the left to the right.

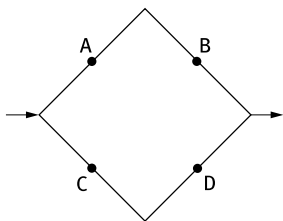


Figure 2.1: An electric circuit with four switches.

Problem 2.18. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Suppose A and B are disjoint events with $\mathbb{P}(A) > 0$ and $\mathbb{P}(B) > 0$. Is it possible that A and B are independent?

Problem 2.19. Let A , B , and C be three independent events.

1. Show that $A \cap B$ and C are independent as well.
2. Even more, show that the independence of A , B , and C implies that of the events $A \cup B$ and C .

Problem 2.20.

1. Suppose that A and C , as well as B and C , are independent. Furthermore, assume $A \cap B = \emptyset$. Show that $A \cup B$ and C are independent as well.
2. Give an example showing that the preceding assertion becomes false without the assumption $A \cap B = \emptyset$.

Hint: To construct such an example, because of Problem 2.19, the events A , B , and C cannot be chosen to be independent. Therefore, the sets defined in Example 2.2.9 are natural candidates for such an example.

Problem 2.21. Suppose $\mathbb{P}(A|B) = \mathbb{P}(A|B^c)$ for some events A and B with $0 < \mathbb{P}(B) < 1$. Does this imply that A and B are independent?

Problem 2.22. Is it possible that an event A is independent of itself? If yes, which events A have this property? Similarly, which A are independent of A^c ?

Problem 2.23. Let A , B , and C be three *independent* events with

$$\mathbb{P}(A) = \mathbb{P}(B) = \mathbb{P}(C) = \frac{1}{3}.$$

Evaluate

$$\mathbb{P}((A \cap B) \cup (A \cap C)).$$

3 Random variables and their distribution

3.1 Transformation of random values

Assume the probability space $(\Omega, \mathcal{A}, \mathbb{P})$ describes a certain random experiment, for example, rolling a die or tossing a coin. If the experiment is executed, a random result $\omega \in \Omega$ shows up. In a second step, we transform this observed result via a mapping $X : \Omega \rightarrow \mathbb{R}$. In this way we obtain a (random) real number $X(\omega)$. Let us point out that X is a fixed, nonrandom function from Ω into \mathbb{R} ; the randomness of $X(\omega)$ stems from the input $\omega \in \Omega$.

Example 3.1.1. Toss a fair coin, labeled on one side with “0” and on the other side with “1,” exactly n times. The appropriate probability space is $(\Omega, \mathcal{P}(\Omega), \mathbb{P})$, where $\Omega = \{0, 1\}^n$ and \mathbb{P} is the uniform distribution on Ω . The result of the experiment is a vector $\omega = (\omega_1, \dots, \omega_n)$ with $\omega_j = 0$ or $\omega_j = 1$. Let X from Ω to \mathbb{R} be defined by

$$X(\omega) = X(\omega_1, \dots, \omega_n) = \omega_1 + \dots + \omega_n.$$

Then $X(\omega)$ tells us how often “1” occurred, but we do no longer know in which order this happened. Of course, $X(\omega)$ is random because, if one tosses the coin another n times, it is very likely that X attains a value different from that in the first trial.

Here we state the most important question in this topic: how are the values of X distributed? As we know, in this case X attains a value $k \leq n$ with probability $\binom{n}{k}2^{-n}$.

Example 3.1.2. Roll a fair die twice. The sample space describing this experiment consists of pairs $\omega = (\omega_1, \omega_2)$, where $\omega_1, \omega_2 \in \{1, \dots, 6\}$. Now define the mapping $X : \Omega \rightarrow \mathbb{R}$ by $X(\omega) := \max\{\omega_1, \omega_2\}$. Thus, instead of recording the values of both rolls, we are only interested in the larger one.

Other possible transformations are, for example, $X_1(\omega) := \min\{\omega_1, \omega_2\}$ or also $X_2(\omega_1, \omega_2) := \omega_1 + \omega_2$.

Let $A \in \mathcal{A}$ be an event. Recall that this event A occurs if and only if we observe an $\omega \in A$. Suppose now $X : \Omega \rightarrow \mathbb{R}$ is a given mapping from Ω into \mathbb{R} and let $B \subseteq \mathbb{R}$ be some event. When do we observe an $\omega \in \Omega$ for which we have $X(\omega) \in B$ or, equivalently, when does the event

$$\{X \in B\} := \{\omega \in \Omega : X(\omega) \in B\}$$

occur? To answer this question, let us recall the definition of the preimage of B with respect to X as given in eq. (A.2):

$$X^{-1}(B) := \{\omega \in \Omega : X(\omega) \in B\}.$$

We observe an $\omega \in \Omega$ for which $X(\omega) \in B$ if and only if $\omega \in X^{-1}(B)$. In other words, the event $\{X \in B\}$ occurs if and only if $X^{-1}(B)$ does. Consequently, the probability to observe

an $\omega \in \Omega$ with $X(\omega) \in B$ should be $\mathbb{P}(X^{-1}(B))$. But to this end, we have to know that $X^{-1}(B) \in \mathcal{A}$; otherwise $\mathbb{P}(X^{-1}(B))$ is not defined at all. Thus, a natural condition for X is $X^{-1}(B) \in \mathcal{A}$ for “sufficiently many” subsets $B \subseteq \mathbb{R}$. The precise mathematical condition reads as follows.

Definition 3.1.3. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. A mapping $X : \Omega \rightarrow \mathbb{R}$ is called a (real-valued) **random variable** (sometimes also **random real number**), provided that it satisfies the following condition:

$$B \in \mathcal{B}(\mathbb{R}) \quad \text{always implies} \quad X^{-1}(B) \in \mathcal{A}. \quad (3.1)$$

Verbally, this condition says that for each Borel set $B \subseteq \mathbb{R}$, its preimage $X^{-1}(B)$ has to be an element of the σ -field \mathcal{A} .

Remark 3.1.4. Condition (3.1) is purely technical and will not be important later on. But, in general, it cannot be avoided, at least if $\mathcal{A} \neq \mathcal{P}(\Omega)$. On the contrary, if $\mathcal{A} = \mathcal{P}(\Omega)$, for example, if either Ω is finite or countably infinite, then *every* mapping $X : \Omega \rightarrow \mathbb{R}$ is a random variable. Indeed, in this case the condition $X^{-1}(B) \in \mathcal{A}$ is trivially always satisfied.

Remark 3.1.5. In order to verify that a given mapping $X : \Omega \rightarrow \mathbb{R}$ is a random variable, it is not necessary to show $X^{-1}(B) \in \mathcal{A}$ for all Borel sets $B \subseteq \mathbb{R}$. It suffices to prove this only for some special Borel sets B . More precisely, the following proposition holds.

Proposition 3.1.6. *A function $X : \Omega \rightarrow \mathbb{R}$ is a random variable if and only if, for all $t \in \mathbb{R}$, we have*

$$X^{-1}((-\infty, t]) = \{\omega \in \Omega : X(\omega) \leq t\} \in \mathcal{A}. \quad (3.2)$$

The assertion remains valid when we replace the intervals $(-\infty, t]$ with intervals of the form $(-\infty, t)$, or we may take intervals $[t, \infty)$ and also (t, ∞) .

Proof. Suppose first that X is a random variable. Given $t \in \mathbb{R}$, the interval $(-\infty, t]$ is a Borel set, hence $X^{-1}((-\infty, t]) \in \mathcal{A}$. Thus, each random variable satisfies condition (3.2).

To prove the converse implication, let X be a mapping from Ω to \mathbb{R} satisfying condition (3.2) for each $t \in \mathbb{R}$. Set

$$\mathcal{C} := \{C \in \mathcal{B}(\mathbb{R}) : X^{-1}(C) \in \mathcal{A}\}.$$

In the first step, one proves¹ that \mathcal{C} is a σ -field. Moreover, (3.2) implies $(-\infty, t] \in \mathcal{C}$ for each $t \in \mathbb{R}$. But $\mathcal{B}(\mathbb{R})$ is the smallest σ -field containing all these intervals. Since \mathcal{C} is another σ -field containing the intervals $(-\infty, t]$, it has to be larger² than the smallest one, that is, we have $\mathcal{C} \supseteq \mathcal{B}(\mathbb{R})$. In other words, every Borel set belongs to \mathcal{C} or, equivalently,

¹ Use Proposition A.2.1 to verify this.

² By the construction of \mathcal{C} , it even coincides with $\mathcal{B}(\mathbb{R})$.

for all $B \in \mathcal{B}(\mathbb{R})$ it follows that $X^{-1}(B) \in \mathcal{A}$. Thus, as asserted, X is a random variable. The proof for intervals of the other types goes along the same lines. Here one has to use that these intervals generate the σ -field of Borel sets as well. \square

Summary: Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. A function $X : \Omega \rightarrow \mathbb{R}$ is a random variable provided that it satisfies the following (purely technical) condition: for any Borel set $B \subseteq \mathbb{R}$,

$$X^{-1}(B) = \{\omega \in \Omega : X(\omega) \in B\} \in \mathcal{A}.$$

This is equivalent to the property that for all $t \in \mathbb{R}$ one has

$$\{\omega \in \Omega : X(\omega) \leq t\} \in \mathcal{A}.$$

3.2 Probability distribution of a random variable

Suppose we are given a random variable $X : \Omega \rightarrow \mathbb{R}$. We define now a mapping \mathbb{P}_X from $\mathcal{B}(\mathbb{R})$ to $[0, 1]$ as follows:

$$\mathbb{P}_X(B) := \mathbb{P}(X^{-1}(B)) = \mathbb{P}\{\omega \in \Omega : X(\omega) \in B\}, \quad B \in \mathcal{B}(\mathbb{R}).$$

Observe that \mathbb{P}_X is well defined. Indeed, since X is a random variable, for all Borel sets $B \subseteq \mathbb{R}$ we have $X^{-1}(B) \in \mathcal{A}$, hence $\mathbb{P}(X^{-1}(B))$ makes sense.

To simplify the notation, given $B \in \mathcal{B}(\mathbb{R})$, we will often write

$$\mathbb{P}\{X \in B\} = \mathbb{P}\{\omega \in \Omega : X(\omega) \in B\}.$$

This is generally used and does not lead to any confusion. Having said this, we may now define \mathbb{P}_X also by

$$\mathbb{P}_X(B) = \mathbb{P}\{X \in B\}.$$

A first easy example shows how \mathbb{P}_X is calculated in concrete cases. Other more interesting examples will follow after some necessary preliminary considerations.

Example 3.2.1. Toss a fair coin, labeled on one side by “0” and on the other side by “1”, three times. The sample space is $\Omega = \{0, 1\}^3$ with the uniform distribution \mathbb{P} describing probability measure. Let the random variable X on Ω be defined by

$$X(\omega) := \omega_1 + \omega_2 + \omega_3 \quad \text{whenever } \omega = (\omega_1, \omega_2, \omega_3) \in \Omega.$$

It follows that

$$\mathbb{P}_X(\{0\}) = \mathbb{P}\{X = 0\} = \mathbb{P}(\{(0, 0, 0)\}) = \frac{1}{8},$$

$$\begin{aligned}\mathbb{P}_X(\{1\}) &= \mathbb{P}\{X = 1\} = \mathbb{P}(\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}) = \frac{3}{8}, \\ \mathbb{P}_X(\{2\}) &= \mathbb{P}\{X = 2\} = \mathbb{P}(\{(1, 1, 0), (0, 1, 1), (1, 0, 1)\}) = \frac{3}{8}, \\ \mathbb{P}_X(\{3\}) &= \mathbb{P}\{X = 3\} = \mathbb{P}(\{(1, 1, 1)\}) = \frac{1}{8}.\end{aligned}$$

Of course, these values describe the distribution of X completely. Indeed, whenever $B \subseteq \mathbb{R}$,

$$\mathbb{P}_X(B) = \sum_{\substack{k=0 \\ k \in B}}^3 \mathbb{P}_X(\{k\}).$$

For example, we have

$$\mathbb{P}_X([-1, 1]) = \mathbb{P}_X(\{0\}) + \mathbb{P}_X(\{1\}) = \frac{1}{8} + \frac{3}{8} = \frac{1}{2}.$$

So, with probability $1/2$, we observe an $\omega = (\omega_1, \omega_2, \omega_3)$ for which

$$-1 \leq \omega_1 + \omega_2 + \omega_3 \leq 1.$$

The proof of the next result heavily depends on properties of the preimage proved in Proposition A.2.1.

Proposition 3.2.2. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. For a random variable $X : \Omega \rightarrow \mathbb{R}$, the mapping $\mathbb{P}_X : \mathcal{B}(\mathbb{R}) \rightarrow [0, 1]$ is a probability measure.*

Proof. Using property (1) in Proposition A.2.1, one easily gets

$$\mathbb{P}_X(\emptyset) = \mathbb{P}(X^{-1}(\emptyset)) = \mathbb{P}(\emptyset) = 0,$$

as well as

$$\mathbb{P}_X(\mathbb{R}) = \mathbb{P}(X^{-1}(\mathbb{R})) = \mathbb{P}(\Omega) = 1.$$

Thus it remains to verify the σ -additivity of \mathbb{P}_X . Take any sequence of disjoint Borel sets B_1, B_2, \dots in \mathbb{R} . Then also $X^{-1}(B_1), X^{-1}(B_2), \dots$ are disjoint subsets of Ω . To see this, apply Proposition A.2.1, which, if $i \neq j$, implies

$$X^{-1}(B_i) \cap X^{-1}(B_j) = X^{-1}(B_i \cap B_j) = X^{-1}(\emptyset) = \emptyset.$$

Another application of Proposition A.2.1 and of the σ -additivity of \mathbb{P} finally gives

$$\mathbb{P}_X\left(\bigcup_{j=1}^{\infty} B_j\right) = \mathbb{P}\left(X^{-1}\left(\bigcup_{j=1}^{\infty} B_j\right)\right) = \mathbb{P}\left(\bigcup_{j=1}^{\infty} X^{-1}(B_j)\right)$$

$$= \sum_{j=1}^{\infty} \mathbb{P}(X^{-1}(B_j)) = \sum_{j=1}^{\infty} \mathbb{P}_X(B_j).$$

Hence, \mathbb{P}_X is a probability measure, as asserted. \square

Definition 3.2.3. The probability measure \mathbb{P}_X on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ defined by

$$\mathbb{P}_X(B) := \mathbb{P}(X^{-1}(B)) = \mathbb{P}\{\omega \in \Omega : X(\omega) \in B\} = \mathbb{P}\{X \in B\}, \quad B \in \mathcal{B}(\mathbb{R}),$$

is called the **probability distribution** of X (with respect to \mathbb{P}) or, in short, the **distribution** of X .

Remark 3.2.4. The distribution \mathbb{P}_X is the most important characteristic of a random variable X . In general, it is completely unimportant how a random variable is defined analytically; only its distribution matters. Thus, two random variables with identical distributions may be regarded as equivalent because they describe the same random experiment.

Remark 3.2.4 leads us to the following definition:

Definition 3.2.5. Two random variables X_1 and X_2 are said to be **identically distributed** provided that $\mathbb{P}_{X_1} = \mathbb{P}_{X_2}$. Hereby, it is not necessary that X_1 and X_2 be defined on the same sample space. Only their distributions have to coincide. In the case of identically distributed X_1 and X_2 , one writes $X_1 \stackrel{d}{=} X_2$.

Example 3.2.6. Toss a fair coin, labeled on each side by “0” or “1,” twice. Let X_1 be the value of the first toss and X_2 that of the second. Then

$$\mathbb{P}\{X_1 = 0\} = \mathbb{P}\{X_2 = 0\} = \frac{1}{2} = \mathbb{P}\{X_1 = 1\} = \mathbb{P}\{X_2 = 1\}.$$

Hence, X_1 and X_2 are identically distributed, or $X_1 \stackrel{d}{=} X_2$. Both random variables describe the same experiment, namely a single toss of a fair coin. Now, toss the coin a third time and let X_3 be the result of the third trial. Then we also have $X_1 \stackrel{d}{=} X_3$, but note that X_1 and X_3 are defined on different sample spaces.

Next, we state and prove some general rules for evaluating the probability distribution of a given random variable. Here we have to distinguish between two different types of random variables, namely between discrete and continuous ones. Let us start with the discrete case.

Definition 3.2.7. A random variable X is **discrete** provided there exists an at most countably infinite set $D \subset \mathbb{R}$ such that $X: \Omega \rightarrow D$.

In other words, a random variable is discrete if it attains at most countably infinitely many different values.

Remark 3.2.8. If a random variable X is discrete with values in $D \subset \mathbb{R}$, then, of course,

$$\mathbb{P}_X(D) = \mathbb{P}\{X \in D\} = 1.$$

Consequently, in this case its probability distribution \mathbb{P}_X is a discrete probability measure on \mathbb{R} . In general, the converse is not valid as the next example shows.

Example 3.2.9. We model the experiment of rolling a fair die by the probability space $(\mathbb{R}, \mathcal{P}(\mathbb{R}), \mathbb{P})$, where $\mathbb{P}(\{1\}) = \dots = \mathbb{P}(\{6\}) = 1/6$ and $\mathbb{P}(\{x\}) = 0$ provided that $x \neq 1, \dots, 6$. If $X : \mathbb{R} \rightarrow \mathbb{R}$ is defined by $X(s) = s^2$ then, of course, \mathbb{P}_X is discrete. Indeed, we have $\mathbb{P}_X(D) = 1$, where $D = \{1, 4, 9, 16, 25, 36\}$. On the other hand, X does not attain values in a countably infinite set; its range is $[0, \infty)$.

Remark 3.2.10. If we look at Example 3.2.9 more thoroughly, then it becomes immediately clear that the values of X outside of $\{1, \dots, 6\}$ are completely irrelevant. With a small change of X , it will attain values in D . More precisely, let $\tilde{X}(\omega) = 1$ if $\omega \neq 1, \dots, 6$ and $\tilde{X}(k) = k^2$, $k = 1, \dots, 6$; then $X \stackrel{d}{=} \tilde{X}$ and \tilde{X} has values in $\{1, 4, 9, 16, 25, 36\}$.

This procedure is also possible in general: if \mathbb{P}_X is discrete with $\mathbb{P}_X(D) = 1$ for some countable set D , then we may change X to \tilde{X} such that $X \stackrel{d}{=} \tilde{X}$ and $\tilde{X} : \Omega \rightarrow D$. Indeed, choose some fixed $d_0 \in D$ and set $\tilde{X}(\omega) = X(\omega)$ if $\omega \in X^{-1}(D)$ and $\tilde{X}(\omega) = d_0$ otherwise. Then $\mathbb{P}_X = \mathbb{P}_{\tilde{X}}$ and \tilde{X} has values in D .

Convention 3.1. Without losing generality, we may always assume the following: if a random variable X has a discrete probability distribution, that is, $\mathbb{P}\{X \in D\} = 1$ for some finite or countably infinite set D , then X attains values in D .

The second type of random variables we investigate is that of continuous ones.³

Definition 3.2.11. A random variable X is said to be **continuous** provided that its distribution \mathbb{P}_X is a continuous probability measure. That is, \mathbb{P}_X possesses a density p . This function p is called the **density function** or, in short, the **density** of the random variable X .

Remark 3.2.12. One should not confuse the continuity of a random variable with the continuity of a function as taught in Calculus. The latter is an (analytic) property of a function, while the former is a property of its distribution. Moreover, whether or not a random variable X is continuous depends not only on X , but also on the underlying probability space.

Remark 3.2.13. Another way to express that a random variable is continuous is as follows: there exists a function $p : \mathbb{R} \rightarrow [0, \infty)$ (the density of X) such that

$$\mathbb{P}\{\omega \in \Omega : X(\omega) \leq t\} = \mathbb{P}\{X \leq t\} = \int_{-\infty}^t p(x) dx, \quad t \in \mathbb{R},$$

³ The precise term would be “absolutely continuous”; but for simplicity let us call them “continuous.”

or, equivalently, for all real numbers $a < b$,

$$\mathbb{P}\{\omega \in \Omega : a \leq X(\omega) \leq b\} = \mathbb{P}\{a \leq X \leq b\} = \int_a^b p(x) dx.$$

How do we determine the probability distribution of a given random variable? To answer this question, let us first consider the case of *discrete* random variables.

Thus, let X be discrete with values in $D = \{x_1, x_2, \dots\} \subset \mathbb{R}$. Then, as observed above, it follows that $\mathbb{P}_X(D) = 1$, and, consequently, \mathbb{P}_X is uniquely determined by the numbers

$$p_j := \mathbb{P}_X(\{x_j\}) = \mathbb{P}\{X = x_j\} = \mathbb{P}\{\omega \in \Omega : X(\omega) = x_j\}, \quad j = 1, 2, \dots \quad (3.3)$$

Moreover, for any $B \subseteq \mathbb{R}$ it follows that

$$\mathbb{P}\{\omega \in \Omega : X(\omega) \in B\} = \mathbb{P}_X(B) = \sum_{x_j \in B} p_j.$$

Consequently, in order to determine \mathbb{P}_X for discrete X , it completely suffices to determine the p_j s defined by eq. (3.3). If we know $(p_j)_{j \geq 1}$, then the probability distribution \mathbb{P}_X of X is completely described.

Remark 3.2.14. In the literature, quite often, one finds a slightly different approach for the description of \mathbb{P}_X . Define $p : \mathbb{R} \rightarrow [0, 1]$ by

$$p(x) = \mathbb{P}\{X = x\}, \quad x \in \mathbb{R}. \quad (3.4)$$

This function p is then called the **probability mass function** of X . Note that $p(x) = 0$ whenever $x \notin D$. This function p satisfies $p(x) \geq 0$, $\sum_{x \in \mathbb{R}} p(x) = 1$, and

$$\mathbb{P}\{X \in B\} = \sum_{x \in B} p(x).$$

In this setting, the numbers p_j in eq. (3.3) coincide with $p(x_j)$.

Example 3.2.15. Roll a fair die twice. Let X on $\{1, \dots, 6\}^2$ be defined by

$$X(\omega) = X(\omega_1, \omega_2) := \omega_1 + \omega_2, \quad \omega = (\omega_1, \omega_2).$$

Which distribution does X possess?

Answer: The very first question one has to answer is always about the possible values of X . In our case, X attains values in $D = \{2, \dots, 12\}$, thus it suffices to determine

$$\mathbb{P}_X(\{k\}) = \mathbb{P}\{X = k\} = \mathbb{P}\{(\omega_1, \omega_2) \in \Omega : X(\omega_1, \omega_2) = k\}, \quad k = 2, \dots, 12.$$

One easily gets

$$\begin{aligned}
\mathbb{P}_X(\{2\}) &= \mathbb{P}\{(\omega_1, \omega_2) : \omega_1 + \omega_2 = 2\} = \frac{|\{(1, 1)\}|}{36} = \frac{1}{36}, \\
\mathbb{P}_X(\{3\}) &= \mathbb{P}\{(\omega_1, \omega_2) : \omega_1 + \omega_2 = 3\} = \frac{|\{(1, 2), (2, 1)\}|}{36} = \frac{2}{36}, \\
&\vdots \\
\mathbb{P}_X(\{7\}) &= \mathbb{P}\{(\omega_1, \omega_2) : \omega_1 + \omega_2 = 7\} = \frac{|\{(1, 6), \dots, (6, 1)\}|}{36} = \frac{6}{36}, \\
&\vdots \\
\mathbb{P}_X(\{12\}) &= \mathbb{P}\{(\omega_1, \omega_2) : \omega_1 + \omega_2 = 12\} = \frac{|\{(6, 6)\}|}{36} = \frac{1}{36},
\end{aligned}$$

hence \mathbb{P}_X is completely described. For example, it follows that

$$\mathbb{P}\{X \leq 4\} = \mathbb{P}_X((-\infty, 4]) = \mathbb{P}_X(\{2\}) + \mathbb{P}_X(\{3\}) + \mathbb{P}_X(\{4\}) = \frac{1}{36} + \frac{2}{36} + \frac{3}{36} = \frac{1}{6}.$$

Example 3.2.16. A biased coin is labeled on one side by “0” and on the other side by “1”: for some $p \in [0, 1]$, number “1” shows up with probability p , thus “0” appears with probability $1 - p$. We toss the coin n times. The result is a sequence $\omega = (\omega_1, \dots, \omega_n)$, where $\omega_i \in \{0, 1\}$, hence the describing sample space is

$$\Omega = \{0, 1\}^n = \{\omega = (\omega_1, \dots, \omega_n) : \omega_i \in \{0, 1\}\}.$$

For $i \leq n$, let $X_i : \Omega \rightarrow \mathbb{R}$ be defined by $X_i(\omega) := \omega_i$. That is, $X_i(\omega)$ is the value of the i th trial. What distribution does X_i possess?

Answer: In Example 1.9.13, we determined the probability measure \mathbb{P} on $\mathcal{P}(\Omega)$, which describes the n -fold tossing of a biased coin. This probability measure was given by

$$\mathbb{P}(\{\omega\}) = p^k (1-p)^{n-k}, \quad k = \sum_{j=1}^n \omega_j \quad \text{where } \omega = (\omega_1, \dots, \omega_n). \quad (3.5)$$

The random variable X_i only attains the values “0” and “1.” Thus, in order to determine \mathbb{P}_{X_i} , it suffices to evaluate $\mathbb{P}_{X_i}(\{0\}) = \mathbb{P}\{\omega \in \Omega : \omega_i = 0\}$. Let $\omega \in \Omega$ be a sequence with $\omega_i = 0$. Then it may contain the value “1” at most $n - 1$ times. Given $k \leq n - 1$, there are exactly $\binom{n-1}{k}$ such sequences ω with $\omega_i = 0$ and with k times “1.” Therefore, we obtain

$$\begin{aligned}
\mathbb{P}_{X_i}(\{0\}) &= \mathbb{P}\{\omega \in \Omega : \omega_i = 0\} = \sum_{k=0}^{n-1} \mathbb{P}\{\omega \in \Omega : \omega_i = 0, \omega_1 + \dots + \omega_n = k\} \\
&= \sum_{k=0}^{n-1} \binom{n-1}{k} p^k (1-p)^{n-k} = (1-p) \sum_{k=0}^{n-1} \binom{n-1}{k} p^k (1-p)^{n-1-k} \\
&= (1-p)[p + (1-p)]^{n-1} = 1-p.
\end{aligned}$$

Of course, this also implies $\mathbb{P}_{X_i}(\{1\}) = p$.

Remark 3.2.17. Note that all X_1, \dots, X_n possess the same distribution, that is,

$$X_1 \stackrel{d}{=} \dots \stackrel{d}{=} X_n.$$

Example 3.2.18. Roll a fair die twice and let ω_1 and ω_2 be the results of the first and second toss, respectively. Define the random variable X by

$$X(\omega_1, \omega_2) = |\omega_1 - \omega_2|, \quad \omega_1, \omega_2 \in \{1, \dots, 6\}.$$

In a first step, we observe that X has values in $D = \{0, \dots, 5\}$. Hence, it suffices to determine $P_X(\{k\}) = \mathbb{P}(X = k)$ for $k = 0, \dots, 5$. Doing so, we easily get

$$\mathbb{P}_X(\{0\}) = \frac{1}{6} \quad \text{and} \quad \mathbb{P}_X(\{1\}) = \frac{5}{18}, \quad \dots, \quad \mathbb{P}_X(\{5\}) = \frac{1}{18}.$$

Summary: Let $X : \Omega \rightarrow \mathbb{R}$ be a *discrete* random variable. In order to describe its distribution \mathbb{P}_X , one has to do two things:

- (1) Determine the finite or countably infinite set $D \subset \mathbb{R}$ for which $\mathbb{P}\{X \in D\} = 1$.
- (2) For each $x \in D$, evaluate

$$\mathbb{P}_X(\{x\}) = \mathbb{P}\{X = x\} = \mathbb{P}\{\omega \in \Omega : X(\omega) = x\}.$$

If $B \subseteq \mathbb{R}$, then it follows that

$$\mathbb{P}\{X \in B\} = \sum_{x \in B \cap D} \mathbb{P}_X(\{x\}) = \sum_{x \in B \cap D} \mathbb{P}\{X = x\}.$$

How do we determine the probability distribution of a random variable if it is *continuous*? For each $x \in \mathbb{R}$, $\mathbb{P}\{X = x\} = 0$, hence the values of $\mathbb{P}\{X = x\}$ cannot be used to describe \mathbb{P}_X as they did in the discrete case. Consequently, a different approach is needed, and this approach is based on the use of distribution functions.

Definition 3.2.19. Let X be a random variable, either discrete or continuous. Then its (**cumulative**) **distribution function** $F_X : \mathbb{R} \rightarrow [0, 1]$ is defined by

$$F_X(t) := \mathbb{P}_X((-\infty, t]) = \mathbb{P}\{X \leq t\}, \quad t \in \mathbb{R}. \quad (3.6)$$

Remark 3.2.20. Observe that for discrete and continuous random variables the distribution function equals

$$F_X(t) = \sum_{x_j \leq t} p_j \quad \text{and} \quad F_X(t) = \int_{-\infty}^t p(x) dx,$$

respectively. Here, in the discrete case, the x_j s and p_j s are as in eq. (3.3), while p denotes the density of X in the continuous case.

Furthermore, note that F_X is nothing else than the distribution function of the probability measure \mathbb{P}_X , as it was introduced in Definition 1.7.1. Consequently, it possesses all properties of a “usual” distribution function as stated in Proposition 1.7.13.

Summary: Let X be a random variable on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. The probability measure \mathbb{P}_X defined by

$$\mathbb{P}_X(B) = \mathbb{P}\{\omega \in \Omega : X(\omega) \in B\}, \quad B \in \mathcal{B}(\mathbb{R}),$$

denotes the probability distribution of X and

$$F_X(t) = \mathbb{P}\{\omega \in \Omega : X(\omega) \leq t\}, \quad t \in \mathbb{R},$$

is its (cumulative) distribution function.

Proposition 3.2.21. *Let F_X be defined by eq. (3.6). Then it possesses the following properties:*

- (1) *The function F_X is nondecreasing.*
- (2) *It follows $F_X(-\infty) = 0$ as well as $F_X(\infty) = 1$.*
- (3) *The function F_X is continuous from the right.*

Furthermore, if $t \in \mathbb{R}$, then

$$\mathbb{P}\{X = t\} = F_X(t) - F_X(t - 0).$$

In particular, if X is continuous, then F_X is a continuous function from \mathbb{R} to $[0, 1]$.

Remark 3.2.22. Note that the converse of the last implication does not hold. Indeed, there exist random variables X for which F_X is continuous, but X does not possess a density. Such random variables are said to be **singularly continuous**. These are exactly those random variables for which the probability measure \mathbb{P}_X is singularly continuous in the sense of Remark 1.7.20.

The next result shows that under slightly stronger conditions about F_X a density of X exists.

Proposition 3.2.23. *Let F_X be continuous and continuously differentiable with the exception of at most finitely many points. Then X is continuous with density $p(t) = \frac{d}{dt} F_X(t)$. Hereby the values of p may be chosen arbitrarily at points where the derivative does not exist; for example, set $p(t) = 0$ for those points.*

Proof. The proof follows from the corresponding properties of distribution functions for probability measures. Recall that F_X is the distribution function of \mathbb{P}_X . \square

The previous proposition provides us with a method to determine the density of a given random variable X . Evaluate the distribution function F_X and differentiate it. The obtained derivative is the density function we are looking for.

The next three examples demonstrate how this method applies.

Example 3.2.24. Let \mathbb{P} be the uniform distribution on a sphere K of radius 1. That is, for each Borel set $B \in \mathcal{B}(\mathbb{R}^2)$, we have

$$\mathbb{P}(B) = \frac{\text{vol}_2(B \cap K)}{\text{vol}_2(K)} = \frac{\text{vol}_2(B \cap K)}{\pi}.$$

Define the random variable $X : \mathbb{R}^2 \rightarrow \mathbb{R}$ by $X(x_1, x_2) := x_1$. Of course, $F_X(t) = 0$ whenever $t < -1$ and $F_X(t) = 1$ when $t > 1$. Thus, it suffices to determine $F_X(t)$ if $-1 \leq t \leq 1$. For those $t \in \mathbb{R}$, we obtain

$$F_X(t) = \frac{\text{vol}_2(S_t \cap K)}{\pi}$$

where S_t is the half-space $\{(x_1, x_2) \in \mathbb{R}^2 : x_1 \leq t\}$ (compare Figure 3.1).

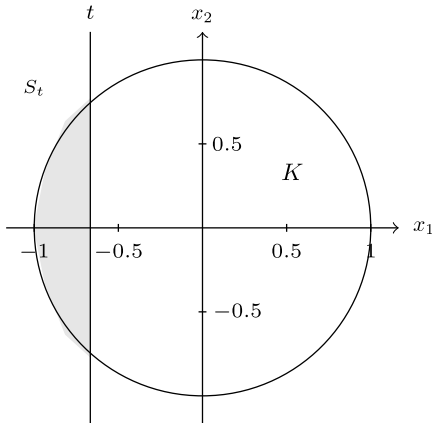


Figure 3.1: The gray shaded set represents the intersection between K and the half-space S_t .

If $|t| \leq 1$, then

$$\text{vol}_2(S_t \cap K) = 2 \int_{-1}^t \sqrt{1-x^2} \, dx,$$

hence,

$$F_X(t) = \frac{2}{\pi} \int_{-1}^t \sqrt{1-x^2} \, dx, \quad |t| \leq 1,$$

and by the fundamental theorem of Calculus, we finally get

$$p(t) = \frac{d}{dt}F_X(t) = \frac{2}{\pi} \sqrt{1-t^2}, \quad |t| \leq 1.$$

Summing up, the random variable X has the density p with

$$p(t) = \begin{cases} \frac{2}{\pi} \sqrt{1-t^2} & \text{if } |t| \leq 1, \\ 0 & \text{if } |t| > 1. \end{cases} \quad (3.7)$$

Example 3.2.25. The probability space is the same as in Example 3.2.24, but this time we define X by

$$X(x_1, x_2) := \sqrt{x_1^2 + x_2^2}, \quad (x_1, x_2) \in \mathbb{R}^2.$$

Of course, it follows $F_X(t) = 0$ if $t < 0$ while $F_X(t) = 1$ if $t > 1$. Take $t \in [0, 1]$. Then

$$F_X(t) = \frac{\text{vol}_2(K(t))}{\text{vol}_2(K(1))} = \frac{t^2\pi}{\pi} = t^2,$$

where $K(t)$ denotes a sphere of radius t . Differentiating F_X with respect to t gives the density

$$p(t) = \begin{cases} 2t & \text{if } 0 \leq t \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Example 3.2.26. Let \mathbb{P} be the uniform distribution on $[0, 1]$ and define the random variable X by $X(s) = \min\{s, 1-s\}$, $s \in \mathbb{R}$. Find the probability distribution of X .

Answer: It is not difficult to see that

$$\mathbb{P}\{X \leq t\} = 0 \quad \text{if } t < 0 \quad \text{and} \quad \mathbb{P}\{X \leq t\} = 1 \quad \text{if } t > 1/2.$$

Thus it remains to evaluate $F_X(t)$ for $0 \leq t \leq 1/2$. Here we obtain

$$\begin{aligned} F_X(t) &= \mathbb{P}\{X \leq t\} = \mathbb{P}\{s \in [0, 1] : 0 \leq s \leq t \text{ or } 1-t \leq s \leq 1\} \\ &= \mathbb{P}\{s \in [0, 1] : 0 \leq s \leq t\} + \mathbb{P}\{s \in [0, 1] : 1-t \leq s \leq 1\} = 2t. \end{aligned}$$

Differentiating gives $F'_X(t) = 2$ if $0 \leq t \leq 1/2$ and $F'_X(t) = 0$ otherwise. Hence \mathbb{P}_X is the uniform distribution on $[0, 1/2]$.

Summary: To determine the density of a *continuous* random variable, proceed as follows:

1. Evaluate (if possible) the distribution function $F_X(t) = \mathbb{P}\{X \leq t\}$.
2. Differentiate F_X . If the derivative $p(t) = F'_X(t)$ is piecewise continuous, then p is the desired density of X .

3.3 Special random variables

We agree upon the following notation: a random variable X is said to be ABC -distributed (or distributed according to ABC) if its probability distribution is a probability measure of type ABC . For example, a random variable is said to be $B_{n,p}$ -distributed (or distributed according to $B_{n,p}$) if $\mathbb{P}_X = B_{n,p}$, that is, if

$$\mathbb{P}\{X = k\} = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, \dots, n.$$

Remark 3.3.1. To shorten the notation, at a few places we also write $X \sim ABC$ whenever \mathbb{P}_X is the probability measure ABC . So, for example, $X \sim B_{n,p}$ tells us that

$$\mathbb{P}\{X = k\} = B_{n,p}(\{k\}), \quad k = 0, \dots, n.$$

In this way, we define the following random variables of special type: X is

- *uniformly distributed* on $\{x_1, \dots, x_N\}$ if

$$\mathbb{P}\{X = x_1\} = \dots = \mathbb{P}\{X = x_N\} = \frac{1}{N},$$

- *Poisson distributed* or Pois_λ -distributed if

$$\mathbb{P}\{X = k\} = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, \dots,$$

- *hypergeometrically distributed* if

$$\mathbb{P}\{X = m\} = \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad m = 0, \dots, n,$$

- G_p -distributed or geometrically distributed if

$$\mathbb{P}\{X = k\} = p(1-p)^{k-1}, \quad k = 1, 2, \dots,$$

- $B_{n,p}^-$ -distributed or negative binomial distributed if

$$\mathbb{P}\{X = k\} = \binom{k-1}{k-n} p^n (1-p)^{k-n} = \binom{k-1}{n-1} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots,$$

or, equivalently, if

$$\mathbb{P}\{X = n+k\} = \binom{-n}{k} p^n (p-1)^k, \quad k = 0, 1, 2, \dots$$

Remark 3.3.2. In view of Convention 3.1, we may suppose that all random variables of the preceding type are discrete. More precisely, we even may assume that X has values in

the (at most countably infinite) set D with $\mathbb{P}_X(D) = 1$. For example, if X is $B_{n,p}$ -distributed, we may suppose that X has values in $\{0, \dots, n\}$.

In quite similar way, we denote specially distributed continuous random variables.

A real-valued random variable X is said to be

- *uniformly distributed* on $[a, \beta]$ if \mathbb{P}_X is the uniform distribution on $[a, \beta]$. That is, if $[a, b] \subseteq [a, \beta]$, then

$$\mathbb{P}\{a \leq X \leq b\} = \frac{b-a}{\beta-a},$$

- *normally distributed* or $\mathcal{N}(\mu, \sigma^2)$ -distributed if

$$\mathbb{P}\{a \leq X \leq b\} = \frac{1}{\sqrt{2\pi}\sigma} \int_a^b e^{-(x-\mu)^2/2\sigma^2} dx,$$

- *standard normally distributed* if it is $\mathcal{N}(0, 1)$ -distributed, that is,

$$\mathbb{P}\{a \leq X \leq b\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx,$$

- *gamma distributed* or $\Gamma_{a,\beta}$ -distributed if for $0 \leq a < b < \infty$,

$$\mathbb{P}\{a \leq X \leq b\} = \frac{1}{a^\beta \Gamma(\beta)} \int_a^b x^{\beta-1} e^{-x/a} dx,$$

- $E_{\lambda,n}$ -distributed or Erlang distributed if it is $\Gamma_{\lambda^{-1},n}$ -distributed, that is, whenever $0 \leq a < b < \infty$, then

$$\mathbb{P}\{a \leq X \leq b\} = \frac{\lambda^n}{(n-1)!} \int_a^b x^{n-1} e^{-\lambda x} dx,$$

or if, equivalently, for any $t \geq 0$,

$$\mathbb{P}\{X \geq t\} = \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t}.$$

- E_λ -distributed or exponentially distributed if for $0 \leq a < b < \infty$,

$$\mathbb{P}\{a \leq X \leq b\} = \lambda \int_a^b e^{-\lambda x} dx = e^{-\lambda a} - e^{-\lambda b},$$

– *arcsine distributed* if, given $0 \leq a < b \leq 1$,

$$\mathbb{P}\{a \leq X \leq b\} = \frac{1}{\pi} \int_a^b \frac{1}{\sqrt{x(1-x)}} dx = \frac{2}{\pi} [\arcsin(\sqrt{b}) - \arcsin(\sqrt{a})].$$

– *Cauchy distributed* if

$$\mathbb{P}\{a \leq X \leq b\} = \frac{1}{\pi} \int_a^b \frac{dx}{1+x^2} = \frac{1}{\pi} [\arctan b - \arctan a].$$

Remark 3.3.3. If a random variable X possesses a special distribution, then all properties of \mathbb{P}_X carry over to X . For example, in this language we may now formulate Poisson's limit theorem (Proposition 1.4.29) as follows.

Let X_n be B_{n,p_n} -distributed and suppose that $np_n \rightarrow \lambda > 0$ as $n \rightarrow \infty$. Then

$$\lim_{n \rightarrow \infty} \mathbb{P}\{X_n = k\} = \mathbb{P}\{X = k\}, \quad k = 0, 1, \dots,$$

where X is Pois_λ -distributed.

Or if X is gamma distributed, then $\mathbb{P}\{X > 0\} = \mathbb{P}_X((0, \infty)) = 1$, and so on.

Remark 3.3.4. A common question is *how does one get a random variable X possessing a certain given distribution*. For example, how do we construct a binomial or a normally distributed random variable? Suppose we want to model the rolling of a die by a random variable X , which is uniformly distributed on $\{1, \dots, 6\}$. The easiest solution is to take $\Omega = \{1, \dots, 6\}$ endowed with the uniform distribution \mathbb{P} and define X by $X(\omega) = \omega$. But this is not the only way to get such a random variable. One may also roll the die n times and choose X as the value of the first (or of the second, etc.) roll. In a similar way, random variables with other probability distributions may be constructed. Further possibilities to model random variables will be investigated in Section 4.4.

Summary: There are two ways to model a random experiment. The classical approach is to construct a *probability space* that describes this experiment. For example, if we toss a fair coin n times and record the number of “heads,” then this may be described by the sample space $\{0, \dots, n\}$ endowed with the probability measure $B_{n,1/2}$. Another way to model a certain random experiment is to choose a *random variable* X so that the probability of the occurrence of an event $B \subseteq \mathbb{R}$ equals $\mathbb{P}\{X \in B\}$. For example, the above experiment of tossing a coin may also be described by a binomial random variable X (with parameters n and $1/2$).

3.4 Random vectors

Suppose we are given n random variables X_1, \dots, X_n defined on a sample space Ω . Our objective is to combine these n variables into a single variable. More precisely, we will investigate the following type of vector-valued mappings.

Definition 3.4.1. Let \vec{X} be a mapping from $\Omega \rightarrow \mathbb{R}^n$ represented as

$$\vec{X}(\omega) = (X_1(\omega), \dots, X_n(\omega)), \quad \omega \in \Omega.$$

Then, \vec{X} is said to be an (n -dimensional) **random vector** or **vector valued random variable**, provided that each of the X_j s is a (real-valued) random variable. The random variables $X_j, 1 \leq j \leq n$, are called the **coordinate mappings** of \vec{X} .

Instead of \vec{X} , we may also write (X_1, \dots, X_n) , that is,

$$(X_1, \dots, X_n)(\omega) = (X_1(\omega), \dots, X_n(\omega)), \quad \omega \in \Omega.$$

A random vector \vec{X} maps Ω into \mathbb{R}^n , that is, we assign to each observed $\omega \in \Omega$ a vector $\vec{X}(\omega)$. The mapping \vec{X} is again fixed and nonrandom. The randomness of $\vec{X}(\omega)$ is caused by the input.

Example 3.4.2. Roll a die two times. Let X_1 be the maximum value, X_2 the minimum, and X_3 the sum of both rolls. The three-dimensional vector $\vec{X} = (X_1, X_2, X_3)$ maps $\Omega = \{1, \dots, 6\}^2$ into \mathbb{R}^3 . For example, the pair $(2, 5)$ is mapped to $(5, 2, 7)$ or the image of $(5, 6)$ is $(6, 5, 11)$.

Example 3.4.3. Suppose there are N people in an auditorium. Enumerate them from 1 to N and choose one person according to the uniform distribution on $\{1, \dots, N\}$. Say we have chosen person k . Let $X_1(k)$ be the height of this person and $X_2(k)$ his or her weight. As a result, we get a random two-dimensional vector (X_1, X_2) mapping k to the vector $(X_1(k), X_2(k))$ in \mathbb{R}^2 .

Example 3.4.4. We place n balls into m urns successively. Hereby, each urn is equally likely. If X_j denotes the number of balls in urn j , then we get an m -dimensional vector $\vec{X} = (X_1, \dots, X_m)$. Observe that the values of \vec{X} lie in the set

$$D = \{(k_1, \dots, k_m) : k_1 + \dots + k_m = n\} \subseteq \mathbb{N}_0^m.$$

Remark 3.4.5. The preceding examples suggest that the values of the coordinate mappings depend on each other. For instance, in Example 3.4.3 larger values of X_1 make also those of X_2 more likely, and vice versa. A basic aim of the following sections is to confirm this guess, that is, we want to find a mathematical formulation that describes whether or not two or more random variables are dependent or independent.

Summary: Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. A mapping $\vec{X} : \Omega \rightarrow \mathbb{R}^n$ is said to be an (n -dimensional) random vector if $\vec{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$ with random variables $X_j : \Omega \rightarrow \mathbb{R}$.

3.5 Joint and marginal distributions

The values of the vector \vec{X} are randomly distributed in \mathbb{R}^n . Consequently, as in the case of random variables, events of the form $\{\vec{X} \in B\}$ occur with certain probabilities. But, in contrast to the case of random variables, the event B is now a subset of \mathbb{R}^n , not of \mathbb{R} as before. More precisely, for events $B \subseteq \mathbb{R}^n$, we are interested in the following quantity:⁴

$$\mathbb{P}\{\omega \in \Omega : \vec{X}(\omega) \in B\} = \mathbb{P}\{\omega \in \Omega : (X_1(\omega), \dots, X_n(\omega)) \in B\}. \quad (3.8)$$

The next definition gives the exact formulation of the problem.

Definition 3.5.1. Let $\vec{X} : \Omega \rightarrow \mathbb{R}^n$ be a random vector with coordinate mappings X_1, \dots, X_n . For each Borel set $B \in \mathcal{B}(\mathbb{R}^n)$, we set

$$\mathbb{P}_{\vec{X}}(B) = \mathbb{P}_{(X_1, \dots, X_n)}(B) = \mathbb{P}\{\vec{X} \in B\}. \quad (3.9)$$

The mapping $\mathbb{P}_{\vec{X}}$ from $\mathcal{B}(\mathbb{R}^n)$ into $[0, 1]$ is said to be the **probability distribution**, or, in short, the **distribution** of \vec{X} . Often, $\mathbb{P}_{\vec{X}} = \mathbb{P}_{(X_1, \dots, X_n)}$ will also be called the **joint distribution** of X_1, \dots, X_n .

In eq. (3.9), we used the shorter expression

$$\mathbb{P}\{\vec{X} \in B\} = \mathbb{P}\{\omega \in \Omega : \vec{X}(\omega) \in B\}.$$

As for random variables, the following is also valid in the case of random vectors.

Proposition 3.5.2. *The mapping $\mathbb{P}_{\vec{X}}$ is a probability measure defined on $\mathcal{B}(\mathbb{R}^n)$.*

Proof. The proof is completely analogous to that of Proposition 3.2.2. Therefore, we decided not to present it here. \square

Let us evaluate $\mathbb{P}_{\vec{X}}(B)$ for special Borel sets $B \subseteq \mathbb{R}^n$. If Q is a box in \mathbb{R}^n as in eq. (1.73), that is, for certain real numbers $a_i < b_i$ we have

$$Q = [a_1, b_1] \times \dots \times [a_n, b_n],$$

then it follows that

$$\mathbb{P}_{\vec{X}}(Q) = \mathbb{P}\{\vec{X} \in Q\} = \mathbb{P}\{\omega \in \Omega : a_1 \leq X_1(\omega) \leq b_1, \dots, a_n \leq X_n(\omega) \leq b_n\}.$$

The later expression may also be written as

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\}.$$

⁴ For random vectors \vec{X} and $B \in \mathcal{B}(\mathbb{R}^n)$, it follows that $\vec{X}^{-1}(B) \in \mathcal{A}$. This can be proved by similar methods as we used in the proof of Proposition 3.1.6. Thus, if $B \in \mathcal{B}(\mathbb{R}^n)$, then eqs. (3.8) and (3.9) are well defined.

Hence, for each box $Q = [a_1, b_1] \times \cdots \times [a_n, b_n]$, we obtain

$$\mathbb{P}_{\vec{X}}(Q) = \mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\}.$$

Thus the quantity $\mathbb{P}_{\vec{X}}(Q)$ is the probability of the occurrence of the following event: X_1 attains a value in $[a_1, b_1]$, and *at the same time* X_2 attains a value in $[a_2, b_2]$, and so on up to X_n attains a value in $[a_n, b_n]$.

Example 3.5.3. Roll a fair die three times. Let X_1, X_2 , and X_3 be the observed values in the first, second, and third roll. If $Q = [1, 2] \times [0, 1] \times [3, 4]$, then

$$\mathbb{P}_{\vec{X}}(Q) = \mathbb{P}\{X_1 \in \{1, 2\}, X_2 = 1, X_3 \in \{3, 4\}\} = \frac{1}{54}.$$

Remark 3.5.4. The previous considerations can easily be generalized to sets $B \subseteq \mathbb{R}^n$ of the form $B = B_1 \times \cdots \times B_n$ with $B_j \in \mathcal{B}(\mathbb{R})$. Then

$$\mathbb{P}_{\vec{X}}(B) = \mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\}. \quad (3.10)$$

Next we introduce the notion of marginal distributions of a random vector.

Definition 3.5.5. Let $\vec{X} = (X_1, \dots, X_n)$ be a random vector. The n probability measures \mathbb{P}_{X_1} to \mathbb{P}_{X_n} are called the **marginal distributions** of \vec{X} .

Observe that each marginal distribution \mathbb{P}_{X_j} is a probability measure on $\mathcal{B}(\mathbb{R})$, while the joint distribution $\mathbb{P}_{(X_1, \dots, X_n)}$ is a probability measure defined on $\mathcal{B}(\mathbb{R}^n)$.

In this context, the following important question arises: does the joint distribution determine the marginal distributions and/or can the joint distribution be derived from the marginal ones?

The next proposition gives the first answer.

Proposition 3.5.6. Let $\vec{X} = (X_1, \dots, X_n)$ be a random vector. If $1 \leq j \leq n$ and $B \in \mathcal{B}(\mathbb{R})$, then

$$\mathbb{P}_{X_j}(B) = \mathbb{P}_{(X_1, \dots, X_n)}\left(\mathbb{R} \times \cdots \times \underbrace{B}_j \times \cdots \times \mathbb{R}\right).$$

In particular, the joint distribution determines the marginal ones.

Proof. The proof is a direct consequence of formula (3.10). Let us apply it to $B_i = \mathbb{R}$ if $i \neq j$ and $B_j = B$. Then, as asserted,

$$\begin{aligned} & \mathbb{P}_{(X_1, \dots, X_n)}\left(\mathbb{R} \times \cdots \times \underbrace{B}_j \times \cdots \times \mathbb{R}\right) \\ &= \mathbb{P}\{X_1 \in \mathbb{R}, \dots, X_j \in B, \dots, X_n \in \mathbb{R}\} = \mathbb{P}\{X_j \in B\} = \mathbb{P}_{X_j}(B). \quad \square \end{aligned}$$

The question whether or not the marginal distributions determine the joint distribution is postponed for a moment. It will be investigated in Example 3.5.8 and, more

thoroughly, in Section 3.6. Before, let us derive some concrete formulas to evaluate the marginal distributions. Here we consider the two cases of discrete and continuous random variables separately.

3.5.1 Marginal distributions: discrete case

To make the results in this subsection easier to understand, we only consider the case of two-dimensional vectors. That is, we investigate two random variables and show how their distributions may be derived from their joint one. We indicate later on how this approach extends to more than two random variables.

In order to avoid confusing notations with many indices, given a two-dimensional random vector, we denote its coordinate mappings by X and Y and not by X_1 and X_2 . This should not lead to mix-ups. Thus, we investigate the random vector (X, Y) with joint distribution $\mathbb{P}_{(X,Y)}$ and marginal distributions \mathbb{P}_X and \mathbb{P}_Y . This random vector maps Ω into \mathbb{R}^2 and acts as follows:

$$(X, Y)(\omega) = (X(\omega), Y(\omega)), \quad \omega \in \Omega.$$

Suppose now that X and Y are discrete. Then, there are finite or countably infinite sets $D = \{x_1, x_2, \dots\}$ and $E = \{y_1, y_2, \dots\}$ such that $X : \Omega \rightarrow D$ as well as $Y : \Omega \rightarrow E$. Consequently, the vector (X, Y) maps Ω into the (at most countably infinite) set $D \times E \subset \mathbb{R}^2$. Observe that

$$D \times E = \{(x_i, y_j) : i, j = 1, 2, \dots\},$$

hence $\mathbb{P}_{(X,Y)}$ is discrete as well and uniquely described by the numbers

$$p_{ij} := \mathbb{P}_{(X,Y)}(\{(x_i, y_j)\}) = \mathbb{P}\{X = x_i, Y = y_j\}, \quad i, j = 1, 2, \dots \quad (3.11)$$

More precisely, given $B \subseteq \mathbb{R}^2$, we have

$$\mathbb{P}_{(X,Y)}(B) = \mathbb{P}\{(X, Y) \in B\} = \sum_{\{(i,j):(x_i,y_j) \in B\}} p_{ij}.$$

We turn now to the description of the marginal distributions \mathbb{P}_X and \mathbb{P}_Y . These are uniquely determined by the numbers

$$q_i := \mathbb{P}_X(\{x_i\}) = \mathbb{P}\{X = x_i\} \quad \text{and} \quad r_j := \mathbb{P}_Y(\{y_j\}) = \mathbb{P}\{Y = y_j\}. \quad (3.12)$$

In other words, if $B, C \subseteq \mathbb{R}$, then it follows that

$$\mathbb{P}_X(B) = \mathbb{P}\{X \in B\} = \sum_{\{i:x_i \in B\}} q_i \quad \text{and} \quad \mathbb{P}_Y(C) = \mathbb{P}\{Y \in C\} = \sum_{\{j:y_j \in C\}} r_j.$$

The next proposition is nothing else than a reformulation of Proposition 3.5.6 in the case of discrete random variables.

Proposition 3.5.7. *Let the probabilities p_{ij} , q_i , and r_j be defined by eqs. (3.11) and (3.12), respectively. Then the q_i s and r_j s may be evaluated by the following equations:*

$$q_i = \sum_{j=1}^{\infty} p_{ij} \quad \text{for } i = 1, 2, \dots \quad \text{and} \quad r_j = \sum_{i=1}^{\infty} p_{ij} \quad \text{for } j = 1, 2, \dots$$

Proof. As already mentioned, Proposition 3.5.7 is a direct consequence of Proposition 3.5.6. But for better understanding, we prefer to give a direct proof.

By virtue of the σ -additivity of \mathbb{P} , it follows that

$$\begin{aligned} q_i &= \mathbb{P}\{X = x_i\} = \mathbb{P}\{X = x_i, Y \in E\} = \mathbb{P}\left\{X = x_i, Y \in \bigcup_{j=1}^{\infty} \{y_j\}\right\} \\ &= \sum_{j=1}^{\infty} \mathbb{P}\{X = x_i, Y \in \{y_j\}\} = \sum_{j=1}^{\infty} \mathbb{P}\{X = x_i, Y = y_j\} = \sum_{j=1}^{\infty} p_{ij}. \end{aligned}$$

This proves the first part. The proof for the r_j s follows exactly along the same line. Here, one uses

$$r_j = \mathbb{P}\{Y = y_j\} = \mathbb{P}\{X \in D, Y = y_j\} = \sum_{i=1}^{\infty} \mathbb{P}\{X = x_i, Y = y_j\} = \sum_{i=1}^{\infty} p_{ij}.$$

This completes the proof. □

The equations in Proposition 3.5.7 may be represented in table form as follows:

$Y \backslash X$	x_1	x_2	x_3	\dots	
y_1	p_{11}	p_{21}	p_{31}	\dots	r_1
y_2	p_{12}	p_{22}	p_{32}	\dots	r_2
y_3	p_{13}	p_{23}	p_{33}	\dots	r_3
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
	q_1	q_2	q_3	\dots	1

The entries in the above matrix are the corresponding probabilities. For example, the entry p_{32} is put into the row marked by x_3 and into the column where one finds y_2 on the left-hand side. This tells us p_{32} is the probability that X attains the value x_3 and, *at the same time*, Y equals y_2 . On the right and in the lower margins,⁵ one finds the corresponding sums of the columns and of the rows, respectively. These numbers describe the marginal distributions (that of X at the bottom and that of Y in the right margin).

⁵ This explains the name “marginal” for the distribution of the coordinate mappings.

Finally, the number “1” at the right lower corner says that both the right column and bottom row have to add up to “1.”

Example 3.5.8. There are four balls in an urn, two labeled with “0” and another two labeled with “1.” Choose two balls without replacing the first. Let X be the value of the first ball and Y that of the second. Direct calculations (use the law of multiplication) lead to

$$\begin{aligned}\mathbb{P}\{X = 0, Y = 0\} &= \frac{1}{6}, & \mathbb{P}\{X = 0, Y = 1\} &= \frac{1}{3}, \\ \mathbb{P}\{X = 1, Y = 0\} &= \frac{1}{3}, & \mathbb{P}\{X = 1, Y = 1\} &= \frac{1}{6}.\end{aligned}$$

In tabular form, this result reads as follows:

$Y \setminus X$	0	1	
0	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{2}$
1	$\frac{1}{3}$	$\frac{1}{6}$	$\frac{1}{2}$
	$\frac{1}{2}$	$\frac{1}{2}$	1

Now suppose that we replace the first ball. This time we denote the values of the first and second ball by X' and Y' , respectively. The corresponding table may now be written as follows:

$Y' \setminus X'$	0	1	
0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$
1	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$
	$\frac{1}{2}$	$\frac{1}{2}$	1

Let us look at Example 3.5.8 more thoroughly. In both cases (nonreplacing and replacing), the marginal distributions coincide, that is, $\mathbb{P}_X = \mathbb{P}_{X'}$ and $\mathbb{P}_Y = \mathbb{P}_{Y'}$. But, on the other hand, the joint distributions are different, that is, we have $\mathbb{P}_{(X,Y)} \neq \mathbb{P}_{(X',Y')}$.

Conclusion. The marginal distributions do *not*, in general, determine the joint distribution. Recall that Proposition 3.5.6 asserts the converse implication: The marginal distributions can be derived from the joint distribution.

Example 3.5.9. Roll a fair die twice. Let X be the minimum value of both rolls and Y the maximum. Then, if $k, l = 1, \dots, 6$, it is easy to see that

$$\mathbb{P}\{X = k, Y = l\} = \begin{cases} 0 & \text{if } k > l, \\ \frac{1}{36} & \text{if } k = l, \\ \frac{1}{18} & \text{if } k < l. \end{cases}$$

Hence, the joint distribution in table form looks as follows:

$Y \setminus X$	1	2	3	4	5	6	
1	$\frac{1}{36}$	0	0	0	0	0	$\frac{1}{36}$
2	$\frac{1}{18}$	$\frac{1}{36}$	0	0	0	0	$\frac{3}{36}$
3	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{36}$	0	0	0	$\frac{5}{36}$
4	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{36}$	0	0	$\frac{7}{36}$
5	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{36}$	0	$\frac{9}{36}$
6	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{36}$	$\frac{11}{36}$
	$\frac{11}{36}$	$\frac{9}{36}$	$\frac{7}{36}$	$\frac{5}{36}$	$\frac{3}{36}$	$\frac{1}{36}$	1

If, for example, $B = \{(4, 5), (5, 4), (6, 5), (5, 6)\}$, then the values in the table imply $\mathbb{P}_{(X,Y)}(B) = 1/9$. In the same way, one gets $\mathbb{P}\{2 \leq X \leq 4\} = (9 + 7 + 5)/36 = 7/12$.

To finish, we shortly go into the case of more than two discrete random variables. Thus, let X_1, \dots, X_n be random variables with $X_j : \Omega \rightarrow D_j$, where the sets D_j are either finite or countably infinite. The set D defined by

$$D = D_1 \times \dots \times D_n = \{(x_1, \dots, x_n), x_j \in D_j\}$$

is at most countably infinite and $\vec{X} : \Omega \rightarrow D$. Consequently, $\mathbb{P}_{\vec{X}}$ is uniquely described by the probabilities

$$p_{x_1, \dots, x_n} = \mathbb{P}\{X_1 = x_1, \dots, X_n = x_n\}, \quad x_j \in D_j.$$

Proposition 3.5.10. For $1 \leq j \leq n$ and $x \in D_j$,

$$\mathbb{P}\{X_j = x\} = \sum_{x_1 \in D_1} \dots \sum_{x_{j-1} \in D_{j-1}} \sum_{x_{j+1} \in D_{j+1}} \dots \sum_{x_n \in D_n} p_{x_1, \dots, x_{j-1}, x, x_{j+1}, \dots, x_n}.$$

Proof. The proof is as that of Proposition 3.5.7. Therefore, we omit it. □

Next, we want to state an important example that shows how Proposition 3.5.10 applies. To do so we need the following definition.

Definition 3.5.11. Let n and m be integers with $m \geq 2$ and let p_1, \dots, p_m be certain success probabilities satisfying $p_j \geq 0$ and $p_1 + \dots + p_m = 1$. An m -dimensional random vector $\vec{X} = (X_1, \dots, X_m)$ has **multinomial distribution** with parameters n and p_1, \dots, p_m if, whenever $k_1 + \dots + k_m = n$,

$$\mathbb{P}\{X_1 = k_1, \dots, X_m = k_m\} = \binom{n}{k_1, \dots, k_m} p_1^{k_1} \dots p_m^{k_m}.$$

Equivalently, a random vector \vec{X} has multinomial distribution if and only if its probability distribution $\mathbb{P}_{\vec{X}}$ is a multinomial distribution as introduced in Definition 1.4.16.

Remark 3.5.12. The m -dimensional random vector \vec{X} in Example 3.4.4 is multinomial distributed with parameters n and $p_j = 1/m$. That means

$$\mathbb{P}\{X_1 = k_1, \dots, X_m = k_m\} = \binom{n}{k_1, \dots, k_m} \left(\frac{1}{m}\right)^n, \quad k_1 + \dots + k_m = n.$$

Example 3.5.13. Let $\vec{X} : \Omega \rightarrow \mathbb{R}^n$ be a multinomial random vector with parameters n and p_1, \dots, p_m . What are the marginal distributions of \vec{X} ?

Answer: To simplify the calculations, we only determine the probability distribution of X_m . The other cases follow in the same way. First note that in the notation of Proposition 3.5.10,

$$p_{k_1, \dots, k_m} = \begin{cases} \binom{n}{k_1, \dots, k_m} p_1^{k_1} \dots p_m^{k_m} & \text{if } k_1 + \dots + k_m = n, \\ 0 & \text{if } k_1 + \dots + k_m \neq n. \end{cases}$$

Consequently, Proposition 3.5.10 leads to

$$\begin{aligned} \mathbb{P}\{X_m = k\} &= \sum_{k_1=0}^n \dots \sum_{k_{m-1}=0}^n p_{k_1, \dots, k_{m-1}, k} \\ &= \sum_{k_1 + \dots + k_{m-1} = n-k} \frac{n!}{k_1! \dots k_{m-1}! k!} p_1^{k_1} \dots p_{m-1}^{k_{m-1}} p_m^k \\ &= \frac{n!}{k! (n-k)!} p_m^k \sum_{k_1 + \dots + k_{m-1} = n-k} \frac{(n-k)!}{k_1! \dots k_{m-1}!} p_1^{k_1} \dots p_{m-1}^{k_{m-1}} \\ &= \binom{n}{k} p_m^k \sum_{k_1 + \dots + k_{m-1} = n-k} \binom{n-k}{k_1, \dots, k_{m-1}} p_1^{k_1} \dots p_{m-1}^{k_{m-1}} \\ &= \binom{n}{k} p_m^k (p_1 + \dots + p_{m-1})^{n-k} = \binom{n}{k} p_m^k (1 - p_m)^{n-k}. \end{aligned}$$

Hereby, in the last step, we used the multinomial theorem (Proposition A.3.20) with $m-1$ summands, with power $n-k$ and entries p_1, \dots, p_{m-1} . Thus X_m is binomial distributed with parameters n and p_m . In the same way, one gets that each X_j is B_{n, p_j} -distributed.

Remark 3.5.14. The previous result can also be seen more directly without using Proposition 3.5.10. Assume we place n particles into m boxes, where p_j is the probability to put a single particle into box j . Fix some $j \leq m$ and let success occur if a particle is placed into box j . Then X_j equals the number of successes, hence it is B_{n, p_j} -distributed. Note that failure occurs if the particle is not placed into box j , and the probability for this is given by $1 - p_j = \sum_{i \neq j}^n p_i$.

Summary: Let $(x_i)_{i \geq 1}$ and $(y_j)_{j \geq 1}$ be the values of the (discrete) random variables X and Y , respectively. If the joint distribution of (X, Y) is given by

$$p_{ij} := \mathbb{P}_{(X, Y)}(\{(x_i, y_j)\}) = \mathbb{P}\{X = x_i, Y = y_j\}, \quad i, j = 1, 2, \dots,$$

then the marginal distributions are described by

$$q_i = \mathbb{P}_X(\{x_i\}) = \mathbb{P}\{X = x_i\} = \sum_{j=1}^{\infty} p_{ij} \quad \text{and} \quad r_j = \mathbb{P}_Y(\{y_j\}) = \mathbb{P}\{Y = y_j\} = \sum_{i=1}^{\infty} p_{ij}.$$

3.5.2 Marginal distributions: continuous case

Let us turn now to the continuous case. Analogous to Definition 3.2.7, a random vector is said to be continuous whenever it possesses a density.⁶ More precisely, we suppose that a random vector shares the following property.

Definition 3.5.15. A random vector $\vec{X} = (X_1, \dots, X_n)$ is said to be **continuous** if there is a function $p : \mathbb{R}^n \rightarrow \mathbb{R}$ such that, for all numbers $a_j < b_j, 1 \leq j \leq n$,

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\} = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} p(x_1, \dots, x_n) dx_n \cdots dx_1.$$

An equivalent formulation is: for all real numbers t_1, \dots, t_n , one has

$$\mathbb{P}\{X_1 \leq t_1, \dots, X_n \leq t_n\} = \int_{-\infty}^{t_1} \cdots \int_{-\infty}^{t_n} p(x_1, \dots, x_n) dx_n \cdots dx_1.$$

The function p is called the **density function** of \vec{X} or the **joint density** of X_1, \dots, X_n .

Remark 3.5.16. Observe that a random vector \vec{X} is continuous if and only if its probability distribution $\mathbb{P}_{\vec{X}}$ is such, that is, the joint distribution of X_1, \dots, X_n is a continuous probability measure on $\mathcal{B}(\mathbb{R}^n)$ in the sense of Definition 1.8.5. Moreover, its density function coincides with the density of $\mathbb{P}_{\vec{X}}$.

In the case of continuous random variables, the marginal distributions are evaluated by the following rule.

Proposition 3.5.17. *If a random vector $\vec{X} = (X_1, \dots, X_n)$ has density $p : \mathbb{R}^n \rightarrow \mathbb{R}$, then for each $j \leq n$ the random variable X_j is continuous with density*

$$p_j(x_j) = \underbrace{\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty}}_{n-1 \text{ integrals}} p(\dots, x_{j-1}, x_j, x_{j+1}, \dots) dx_n \cdots dx_{j+1} dx_{j-1} \cdots dx_1. \quad (3.13)$$

⁶ The following is true: for continuous random variables, the generated vector possesses a density. The proof is far outside the scope of this book. Furthermore, we do not need this assertion because we assume \vec{X} to be continuous, not the X_j s.

Proof. Fix an integer $j \leq n$. An application of Proposition 3.5.6 implies

$$\begin{aligned} \mathbb{P}_{X_j}([a, b]) &= \mathbb{P}_{\vec{X}}\left(\mathbb{R} \times \cdots \times \underbrace{[a, b]}_j \times \cdots \times \mathbb{R}\right) \\ &= \int_{-\infty}^{\infty} \cdots \int_{\underbrace{a}_j}^b \cdots \int_{-\infty}^{\infty} p(x_1, \dots, x_n) dx_n \cdots dx_1 \\ &= \int_a^b \left[\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} p(\dots, x_{j-1}, x_j, x_{j+1}, \dots) dx_n \cdots dx_{j+1} dx_{j-1} \cdots dx_1 \right] dx_j \\ &= \int_a^b p_j(x_j) dx_j \end{aligned}$$

with p_j defined by eq. (3.13). The interchange of the integrals was justified by Fubini's theorem (Proposition A.5.5); note that p is a density, hence it is nonnegative. Since the preceding equation holds for all real numbers $a < b$, the function p_j has to be the density of \mathbb{P}_{X_j} . This completes the proof. \square

In the case $n = 2$, formula (3.13) asserts the following. Let $p : \mathbb{R}^2 \rightarrow \mathbb{R}$ be the joint density of the 2-dimensional vector (X_1, X_2) . Then the functions p_1 and p_2 defined by

$$p_1(x_1) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_2 \quad \text{and} \quad p_2(x_2) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_1$$

are densities of X_1 and X_2 , respectively.

Remark 3.5.18. Another way to formulate Proposition 3.5.17 is as follows: if the function $p : \mathbb{R}^n \rightarrow \mathbb{R}$ is a joint density of $\vec{X} = (X_1, \dots, X_n)$, then p_1, \dots, p_n defined in eq. (3.13) are densities of the random variables X_1, \dots, X_n , respectively.

Example 3.5.19. Choose by random a point $x = (x_1, x_2, x_3)$ in the unit ball of \mathbb{R}^3 . How are the coordinates x_1, x_2 , and x_3 distributed?

Answer: Let $\vec{X} = (X_1, X_2, X_3)$ be uniformly distributed on the unit ball

$$K = \{(x_1, x_2, x_3) : x_1^2 + x_2^2 + x_3^2 \leq 1\}.$$

Then the joint density is given by⁷

⁷ Recall that $\text{vol}_3(K) = \frac{4}{3}\pi$.

$$p(x) = \begin{cases} \frac{3}{4\pi} & \text{if } x \in K, \\ 0 & \text{if } x \notin K. \end{cases}$$

An application of Proposition 3.5.17 leads to $p_1(x_1) = 0$ whenever $|x_1| > 1$ and, if $|x_1| \leq 1$, then it follows that

$$p_1(x_1) = \frac{3}{4\pi} \iint_{x_2^2 + x_3^2 \leq 1 - x_1^2} dx_2 dx_3 = \frac{3}{4\pi} (1 - x_1^2)\pi = \frac{3}{4} (1 - x_1^2).$$

Hence, X_1 has the density (compare Figure 3.2)

$$p_1(s) = \begin{cases} \frac{3}{4}(1 - s^2) & \text{if } -1 \leq s \leq 1, \\ 0 & \text{otherwise.} \end{cases} \tag{3.14}$$

Of course, by symmetry, X_2 and X_3 possess exactly the same distribution densities.

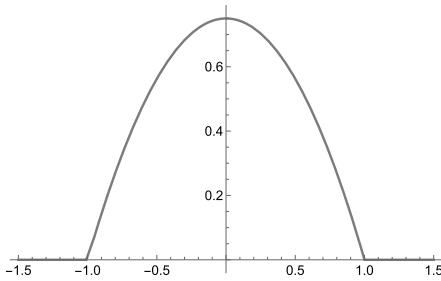


Figure 3.2: The density p_1 defined by eq. (3.14).

Example 3.5.20. Suppose the two-dimensional random vector (X_1, X_2) has the density p defined by⁸

$$p(x_1, x_2) := \begin{cases} 8x_1x_2 & \text{if } 0 \leq x_1 \leq x_2 \leq 1, \\ 0 & \text{otherwise.} \end{cases} \tag{3.15}$$

Then, the density p_1 of X_1 is given by

$$p_1(x_1) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_2 = 8x_1 \int_{x_1}^1 x_2 dx_2 = 4(x_1 - x_1^3), \quad 0 \leq x_1 \leq 1,$$

and $p_1(x_1) = 0$ if $x_1 \notin [0, 1]$.

⁸ Check that p is indeed a probability density.

In the case of p_2 , the density of X_2 , it follows that

$$p_2(x_2) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_1 = 8x_2 \int_0^{x_2} x_1 dx_1 = 4x_2^3, \quad 0 \leq x_2 \leq 1,$$

and $p_2(x_2) = 0$ if $x_2 \notin [0, 1]$. See Figure 3.3 for the joint density of the random vector (X_1, X_2) and its marginal distributions.

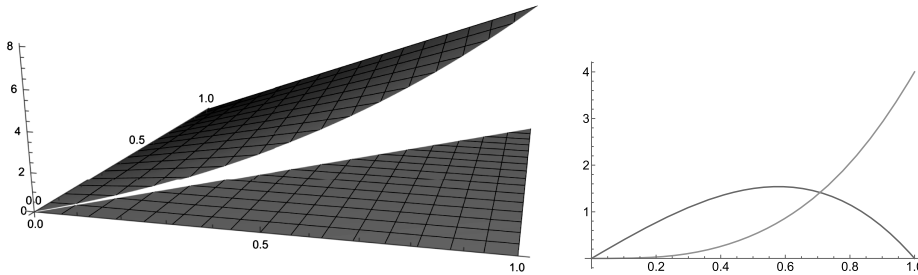


Figure 3.3: The two-dimensional density p defined by eq. (3.15) with marginal densities p_1 and p_2 .

Remark 3.5.21. For $p_1 : \mathbb{R} \rightarrow \mathbb{R}$ and $p_2 : \mathbb{R} \rightarrow \mathbb{R}$ as in Example 3.5.20, define $\tilde{p} : \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$\tilde{p}(x_1, x_2) = p_1(x_1) \cdot p_2(x_2), \quad (x_1, x_2) \in \mathbb{R}^2.$$

In view of Proposition 1.8.8, the function \tilde{p} is a (two-dimensional) density and, moreover, as can be seen easily, its marginal distributions are p_1 and p_2 as well. But note that $p \neq \tilde{p}$. Thus, this is another example showing that the marginal distributions of a random vector do not determine its joint distribution.

Summary: Let $\vec{X} = (X_1, \dots, X_n)$ be a random (n -dimensional) vector. The probability distribution $\mathbb{P}_{\vec{X}}$ on $\mathcal{B}(\mathbb{R}^n)$ is said to be the joint distribution of X_1, \dots, X_n while the probability measures \mathbb{P}_{X_j} , $1 \leq j \leq n$, denote the marginal distributions of \vec{X} . The joint distribution determines the n marginal distributions. In general, the converse implication does not hold.

3.6 Independence of random variables

The central question considered in this section is as follows: when are n given random variables independent? Surely everybody has an intuitive idea about the independence or dependence of random values. But how do we express this property by a mathematical formula? Let us try to approach a solution of this problem with an example.

Example 3.6.1. Roll a fair die twice and define the two random variables X_1 and X_2 as the results of the first and second roll, respectively. These random variables are intuitively independent of each other. But what property of these random variables does this express? Take two subsets $B_1, B_2 \subseteq \{1, \dots, 6\}$ and look at their preimages $A_1 = X_1^{-1}(B_1)$ and $A_2 = X_2^{-1}(B_2)$. Then A_1 occurs if the first result belongs to B_1 while the same is true for A_2 whenever the second result belongs to B_2 . For example, A_1 might indicate that the first result is an even number while A_2 could occur if the second result equals “4.” The basic observation is that no matter how B_1 and B_2 were chosen, the occurrence of their preimages A_1 and A_2 only depends on the first or second roll, respectively. Therefore, they should be independent (as events) in the sense of Definition 2.2.2, that is, the following equation should hold:

$$\begin{aligned} \mathbb{P}\{X_1 \in B_1, X_2 \in B_2\} &= \mathbb{P}(X_1^{-1}(B_1) \cap X_2^{-1}(B_2)) = \mathbb{P}(A_1 \cap A_2) \\ &= \mathbb{P}(A_1) \cdot \mathbb{P}(A_2) = \mathbb{P}(X_1^{-1}(B_1)) \cdot \mathbb{P}(X_2^{-1}(B_2)) = \mathbb{P}\{X_1 \in B_1\} \cdot \mathbb{P}\{X_2 \in B_2\}. \end{aligned}$$

This observation leads us to the following definition of independence.

Definition 3.6.2. Let X_1, \dots, X_n be n random variables mapping Ω into \mathbb{R} . These variables are said to be (stochastically) **independent** if, for all Borel sets $B_j \subseteq \mathbb{R}$,

$$\mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} = \mathbb{P}\{X_1 \in B_1\} \cdots \mathbb{P}\{X_n \in B_n\}. \quad (3.16)$$

Remark 3.6.3. By virtue of Remark 3.5.4, eq. (3.16) may also be written as

$$\mathbb{P}_{(X_1, \dots, X_n)}(B_1 \times \cdots \times B_n) = \mathbb{P}_{X_1}(B_1) \cdots \mathbb{P}_{X_n}(B_n), \quad B_j \in \mathcal{B}(\mathbb{R}).$$

Before proceeding further, we shortly recall Corollary 1.9.9.

Corollary 3.6.4. Given n probability measures $\mathbb{P}_1, \dots, \mathbb{P}_n$ defined on $\mathcal{B}(\mathbb{R})$, there exists a unique probability measure \mathbb{P} on $\mathcal{B}(\mathbb{R}^n)$, the product measure, which is denoted by $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$, such that for all Borel sets $B_j \subseteq \mathbb{R}$,

$$\mathbb{P}(B_1 \times \cdots \times B_n) = \mathbb{P}_1(B_1) \cdots \mathbb{P}_n(B_n). \quad (3.17)$$

Now, we are prepared to state the characterization of independent random variables by properties of their distributions.

Proposition 3.6.5. The random variables X_1, \dots, X_n are independent if and only if their joint distribution coincides with the product probability of the marginal distributions. That is, if and only if

$$\mathbb{P}_{(X_1, \dots, X_n)} = \mathbb{P}_{X_1} \otimes \cdots \otimes \mathbb{P}_{X_n}.$$

Proof. In view of Corollary 3.6.4, the product probability \mathbb{P} of $\mathbb{P}_{X_1}, \dots, \mathbb{P}_{X_n}$ is the unique probability measure on $\mathcal{B}(\mathbb{R}^n)$ satisfying

$$\mathbb{P}(B_1 \times \cdots \times B_n) = \mathbb{P}_{X_1}(B_1) \cdots \mathbb{P}_{X_n}(B_n), \quad B_j \in \mathcal{B}(\mathbb{R}).$$

On the other hand, by Remark 3.6.3, the X_j s are independent if and only if

$$\mathbb{P}_{(X_1, \dots, X_n)}(B_1 \times \cdots \times B_n) = \mathbb{P}_{X_1}(B_1) \cdots \mathbb{P}_{X_n}(B_n), \quad B_j \in \mathcal{B}(\mathbb{R}). \quad (3.18)$$

Consequently, eq. (3.18) holds for all Borel sets B_j if and only if $\mathbb{P}_{(X_1, \dots, X_n)}$ is the product probability $\mathbb{P}_{X_1} \otimes \cdots \otimes \mathbb{P}_{X_n}$. This completes the proof. \square

Corollary 3.6.6. *If X_1, \dots, X_n are independent, the joint distribution $\mathbb{P}_{(X_1, \dots, X_n)}$ is uniquely determined by its marginal distributions $\mathbb{P}_{X_1}, \dots, \mathbb{P}_{X_n}$.*

Proof. Proposition 3.6.5 asserts that $\mathbb{P}_{(X_1, \dots, X_n)} = \mathbb{P}_{X_1} \otimes \cdots \otimes \mathbb{P}_{X_n}$. Hence, the joint distribution is uniquely described by the marginal ones. \square

Another application of Proposition 3.6.5 deals with the existence of independent random variables possessing given distributions.

Proposition 3.6.7. *Let $\mathbb{P}_1, \dots, \mathbb{P}_n$ be given probability measures on the real line, discrete or continuous. Then there is a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and there are independent random variables $X_j : \Omega \rightarrow \mathbb{R}$ such that $\mathbb{P}_{X_j} = \mathbb{P}_j$, $1 \leq j \leq n$. In other words, for all Borel sets B_1, \dots, B_n we have*

$$\mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} = \mathbb{P}\{X_1 \in B_1\} \cdots \mathbb{P}\{X_n \in B_n\} = \mathbb{P}_1(B_1) \cdots \mathbb{P}_n(B_n).$$

Proof. Choose $\Omega = \mathbb{R}^n$ and endow it with the probability distribution (product measure) $\mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n$. Next we define random variables $X_j : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$X_j(\omega) = X_j(\omega_1, \dots, \omega_n) = \omega_j, \quad 1 \leq j \leq n.$$

Thus, $\vec{X} = (X_1, \dots, X_n)$ is the identity, hence $\mathbb{P}_{\vec{X}} = \mathbb{P}$. Consequently, Proposition 1.9.10 implies $\mathbb{P}_{X_j} = \mathbb{P}_j$, $1 \leq j \leq n$, and, moreover, by the choice of \mathbb{P} it follows that

$$\mathbb{P}_{(X_1, \dots, X_n)} = \mathbb{P}_{\vec{X}} = \mathbb{P} = \mathbb{P}_1 \otimes \cdots \otimes \mathbb{P}_n = \mathbb{P}_{X_1} \otimes \cdots \otimes \mathbb{P}_{X_n}.$$

So, in view of Proposition 3.6.5 the random variables X_1 up to X_n are independent, as asserted. \square

Example 3.6.8. Suppose we want to construct n independent Pois_λ -distributed random variables X_1 up to X_n . To do so, choose $\Omega = \mathbb{N}_0^n$ endowed with the product measure \mathbb{P} of n different Pois_λ measures. That is,

$$\mathbb{P}(\{\vec{k}\}) = \frac{\lambda^{k_1 + \cdots + k_n}}{k_1! \cdots k_n!} e^{-n\lambda}, \quad \vec{k} = (k_1, \dots, k_n) \in \mathbb{N}_0^n.$$

Then the random variables $X_j : \Omega \rightarrow \mathbb{N}_0$ with

$$X_j(\vec{k}) = k_j, \quad \vec{k} = (k_1, \dots, k_n) \in \mathbb{N}_0,$$

are the independent n random variables distributed according to Pois_λ .

The next proposition clarifies the relation between the properties “independence of events” and “independence of random variables.” At a first glance, the assertion looks trivial or self-evident, but it is not at all. The reason is that the definition of independence for more than two events, as given in Definition 2.2.12, is more complicated than in the case of two events.

Proposition 3.6.9. *The random variables X_1, \dots, X_n are independent if and only if for all Borel sets B_1, \dots, B_n in \mathbb{R} the events*

$$X_1^{-1}(B_1), \dots, X_n^{-1}(B_n)$$

are stochastically independent in $(\Omega, \mathcal{A}, \mathbb{P})$.

Proof. When are $X_1^{-1}(B_1), \dots, X_n^{-1}(B_n)$ independent? According to Definition 2.2.12, this holds if for all subsets $I \subseteq \{1, \dots, n\}$,

$$\mathbb{P}\left(\bigcap_{i \in I} X_i^{-1}(B_i)\right) = \prod_{i \in I} \mathbb{P}(X_i^{-1}(B_i)). \quad (3.19)$$

On the other hand, by Definition 3.6.2, X_1, \dots, X_n are independent if

$$\begin{aligned} \mathbb{P}\left(\bigcap_{i=1}^n X_i^{-1}(B_i)\right) &= \mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} \\ &= \prod_{i=1}^n \mathbb{P}\{X_i \in B_i\} = \prod_{i=1}^n \mathbb{P}(X_i^{-1}(B_i)). \end{aligned} \quad (3.20)$$

Of course, eq. (3.19) implies eq. (3.20); use eq. (3.19) with $I = \{1, \dots, n\}$. But it is far from clear why, conversely, eq. (3.20) should imply eq. (3.19). As we saw in Example 2.2.10, for fixed sets B_j this is even false. The key observation is that eq. (3.19) has to be valid for *all* Borel sets B_j . This allows us to choose the Borel sets in an appropriate way.

Thus let us assume the validity of eq. (3.20) for all Borel sets in \mathbb{R} . Given $B_j \in \mathcal{B}(\mathbb{R})$ and a subset I of $\{1, \dots, n\}$, we introduce “new” B'_1, \dots, B'_n as follows: $B'_i = B_i$ if $i \in I$ and $B'_i = \mathbb{R}$ if $i \notin I$. This choice of the B'_i implies $X_i^{-1}(B'_i) = \Omega$ whenever $i \notin I$. An application of (3.20) to B'_1, \dots, B'_n leads to (recall $X_i^{-1}(B'_i) = \Omega$ if $i \notin I$)

$$\mathbb{P}\left(\bigcap_{i \in I} X_i^{-1}(B_i)\right) = \mathbb{P}\left(\bigcap_{i=1}^n X_i^{-1}(B'_i)\right) = \prod_{i=1}^n \mathbb{P}(X_i^{-1}(B'_i)) = \prod_{i \in I} \mathbb{P}(X_i^{-1}(B_i)).$$

This proves eq. (3.19) for any subset I of $\{1, \dots, n\}$. So, $X_1^{-1}(B_1), \dots, X_n^{-1}(B_n)$ are independent as asserted. \square

Remark 3.6.10. To verify the independence of X_1, \dots, X_n , it is not necessary to check eq. (3.16) for all Borel sets B_j . It suffices if this is valid for real intervals $[a_j, b_j]$. In other words, X_1, \dots, X_n are independent if and only if, for all $a_j < b_j$,

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\} = \mathbb{P}\{a_1 \leq X_1 \leq b_1\} \cdots \mathbb{P}\{a_n \leq X_n \leq b_n\}.$$

Furthermore, it also suffices to choose the Borel sets as intervals $(-\infty, t_j]$ for $t_j \in \mathbb{R}$, i. e., X_1, \dots, X_n are independent if and only if, for all $t_j \in \mathbb{R}$,

$$\mathbb{P}\{X_1 \leq t_1, \dots, X_n \leq t_n\} = \mathbb{P}\{X_1 \leq t_1\} \cdots \mathbb{P}\{X_n \leq t_n\}.$$

3.6.1 Independence of discrete random variables

As in Section 3.5.1, we restrict ourselves to the case of two random variables. The extension to more than two variables is straightforward and will be shortly considered at the end of this section. We use the same notation as in Section 3.5.1. That is, the two random variables are denoted by X and Y , and they map Ω into $D = \{x_1, x_2, \dots\}$ and $E = \{y_1, y_2, \dots\}$, respectively. The joint distribution of (X, Y) , as well as the marginal distributions, that is, the distributions of X and Y , are described as in eqs. (3.11) and (3.12) by

$$p_{ij} = \mathbb{P}\{X = x_i, Y = y_j\}, \quad q_i = \mathbb{P}\{X = x_i\}, \quad \text{and} \quad r_j = \mathbb{P}\{Y = y_j\}.$$

With these notations the following result is valid.

Proposition 3.6.11. *For the independence of two random variables X and Y , it is necessary and sufficient that*

$$p_{ij} = q_i \cdot r_j, \quad 1 \leq i, j < \infty.$$

Proof. The assertion is an immediate consequence of Propositions 1.9.11 and 3.6.5. But, because of the importance of the result, we give an alternative proof avoiding the direct use of product probabilities; only the techniques are similar.

Let us first show that the condition is necessary. Therefore, choose indices i and j , and put $B_1 := \{x_i\}$ and $B_2 := \{y_j\}$. Then $\{X \in B_1\}$ occurs if and only if $X = x_i$, and, in the same way, the occurrence of $\{Y \in B_2\}$ is equivalent to $Y = y_j$. Since X and Y are assumed to be independent, as claimed,

$$\begin{aligned} p_{ij} &= \mathbb{P}\{X = x_i, Y = y_j\} = \mathbb{P}\{X \in B_1, Y \in B_2\} = \mathbb{P}\{X \in B_1\} \cdot \mathbb{P}\{Y \in B_2\} \\ &= \mathbb{P}\{X = x_i\} \cdot \mathbb{P}\{Y = y_j\} = q_i \cdot r_j. \end{aligned}$$

To prove the converse implication, assume we have $p_{ij} = q_i \cdot r_j$ for all pairs (i, j) of integers. Let B_1 and B_2 be two arbitrary subsets of \mathbb{R} . Then

$$\begin{aligned}
 \mathbb{P}\{X \in B_1, Y \in B_2\} &= \mathbb{P}_{(X,Y)}(B_1 \times B_2) = \sum_{\{(i,j):(x_i,y_j) \in B_1 \times B_2\}} p_{ij} \\
 &= \sum_{\{(i,j):x_i \in B_1, y_j \in B_2\}} q_i \cdot r_j = \sum_{\{i:x_i \in B_1\}} \sum_{\{j:y_j \in B_2\}} q_i \cdot r_j \\
 &= \left(\sum_{\{i:x_i \in B_1\}} q_i \right) \cdot \left(\sum_{\{j:y_j \in B_2\}} r_j \right) = \mathbb{P}_X(B_1) \cdot \mathbb{P}_Y(B_2) \\
 &= \mathbb{P}\{X \in B_1\} \cdot \mathbb{P}\{Y \in B_2\}.
 \end{aligned}$$

Since B_1 and B_2 were arbitrary, the random variables X and Y are independent. This completes the proof. \square

Remark 3.6.12. The previous proposition implies again that for (discrete) independent random variables the joint distribution is determined by the marginal ones. Indeed, in order to know the p_{ij} s, it suffices to know the q_i s and r_j s.

Let us represent the assertion of Proposition 3.6.11 graphically. It tells us that the random variables X and Y are independent if and only if the table describing their joint distribution may be represented as follows:

$Y \backslash X$	x_1	x_2	x_3	\dots	
y_1	$q_1 r_1$	$q_2 r_1$	$q_3 r_1$	\dots	r_1
y_2	$q_1 r_2$	$q_2 r_2$	$q_3 r_2$	\dots	r_2
y_3	$q_1 r_3$	$q_2 r_3$	$q_3 r_3$	\dots	r_3
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
	q_1	q_2	q_3	\dots	1

Example 3.6.13. Proposition 3.6.11 lets us conclude that X and Y in Example 3.5.8 (with-out replacing) are dependent while X' and Y' (with replacement) are independent. Furthermore, by the same argument, the random variables X and Y in Example 3.5.9 (minimum and maximum value when rolling a die twice) are dependent as well.

Example 3.6.14. Let X and Y be two independent Pois_λ -distributed random variables. Then the joint distribution of the vector (X, Y) is determined by

$$\mathbb{P}\{X = k, Y = \ell\} = \frac{\lambda^{k+\ell}}{k! \ell!} e^{-2\lambda}, \quad (k, \ell) \in \mathbb{N}_0 \times \mathbb{N}_0.$$

For example, applying this for $\mathbb{P}_{(X,Y)}(B)$ with $B = \{(k, \ell) : k = \ell\}$ leads to

$$\mathbb{P}\{X = Y\} = \sum_{k=0}^{\infty} \mathbb{P}\{X = k, Y = k\} = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(k!)^2} e^{-2\lambda}.$$

Example 3.6.15. Suppose X and Y are two independent geometrically distributed random variables, with parameters p and q , respectively. Evaluate $\mathbb{P}\{X \leq Y\}$.

Solution: By the independence of X and Y ,

$$\begin{aligned}
 \mathbb{P}\{X \leq Y\} &= \sum_{k=1}^{\infty} \mathbb{P}\{X = k, Y \geq k\} = \sum_{k=1}^{\infty} \mathbb{P}\{X = k\} \cdot \mathbb{P}\{Y \geq k\} \\
 &= \sum_{k=1}^{\infty} p(1-p)^{k-1} \sum_{\ell=k}^{\infty} q(1-q)^{\ell-1} \\
 &= pq \left(\sum_{k=0}^{\infty} (1-p)^k \right) \left(\sum_{\ell=k+1}^{\infty} (1-q)^{\ell-1} \right) \\
 &= pq \left(\sum_{k=0}^{\infty} (1-p)^k \right) \left(\sum_{\ell=k}^{\infty} (1-q)^{\ell} \right) \\
 &= pq \left(\sum_{k=0}^{\infty} (1-p)^k (1-q)^k \right) \left(\sum_{\ell=0}^{\infty} (1-q)^{\ell} \right) \\
 &= \frac{p}{1 - (1-p)(1-q)} = \frac{p}{p+q-pq}.
 \end{aligned}$$

In Example 1.9.12, we investigated the case $p = q$ from a different point of view. The results obtained there let us conclude that

$$\mathbb{P}\{X \leq Y\} = \mathbb{P}\{X < Y\} + \mathbb{P}\{X = Y\} = \frac{p}{2-p} + \frac{1-p}{2-p} = \frac{1}{2-p},$$

which coincides with what we got above if $p = q$.

Example of application: Player A rolls a die and, simultaneously, player B tosses two fair coins labeled with “0” and “1.” Find the probability that player A observes the number “6” for the first time strictly before player B gets a “1” at both coins.

Answer: Let $\{Y = k\}$ be the event that player A observes his first “6” in trial k . Similarly, $\{X = k\}$ occurs if player B has his first two “1” in trial k . Then we ask for the probability $\mathbb{P}\{Y < X\}$. Note that X is geometrically distributed with parameter $p = 1/4$, while the success probability for Y is $q = 1/6$. Hence, by the above calculations,

$$\mathbb{P}\{Y < X\} = 1 - \mathbb{P}\{X \leq Y\} = 1 - \frac{1/4}{1/4 + 1/6 - 1/24} = \frac{1}{3}.$$

The next objective is to investigate in which cases two quite special random variables are independent. To this end, we need the following notation.

Definition 3.6.16. Let Ω be a set and $A \subseteq \Omega$. Then the **indicator function** $\mathbb{1}_A : \Omega \rightarrow \mathbb{R}$ of A is defined by

$$\mathbb{1}_A(\omega) := \begin{cases} 1 & \text{if } \omega \in A, \\ 0 & \text{if } \omega \notin A. \end{cases} \quad (3.21)$$

Let us state some basic properties of indicator functions.

Proposition 3.6.17. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space.*

- (1) *The indicator function of a set $A \subseteq \Omega$ is a random variable if and only if $A \in \mathcal{A}$.*
- (2) *If $A \in \mathcal{A}$, then $\mathbb{1}_A$ is $B_{1,p}$ -distributed (binomial) where $p = \mathbb{P}(A)$.*
- (3) *If $A, B \in \mathcal{A}$, then the random variables $\mathbb{1}_A$ and $\mathbb{1}_B$ are independent if and only if the events A and B are independent.*

Proof. Given $t \in \mathbb{R}$, the event $\{\omega \in \Omega : \mathbb{1}_A(\omega) \leq t\}$ is either empty, A^c , or Ω whenever $t < 0$, $0 \leq t < 1$, or $t \geq 1$, respectively. Consequently, the set $\{\omega \in \Omega : \mathbb{1}_A(\omega) \leq t\}$ is in \mathcal{A} for all $t \in \mathbb{R}$ if and only if $A^c \in \mathcal{A}$. But this happens if and only if $A \in \mathcal{A}$, which proves the first assertion.

To prove the second claim, we first observe that $\mathbb{1}_A$ attains only the values “0” and “1.” Since

$$\mathbb{P}\{\mathbb{1}_A = 1\} = \mathbb{P}\{\omega \in \Omega : \mathbb{1}_A(\omega) = 1\} = \mathbb{P}(A) = p,$$

it is $B_{1,p}$ -distributed with $p = \mathbb{P}(A)$ as claimed.

Let us turn to the last assertion. Given $A, B \in \mathcal{A}$, their joint distribution in table form is⁹

$\mathbb{1}_B \setminus \mathbb{1}_A$	0	1	
0	$\mathbb{P}(A^c \cap B^c)$	$\mathbb{P}(A \cap B^c)$	$\mathbb{P}(B^c)$
1	$\mathbb{P}(A^c \cap B)$	$\mathbb{P}(A \cap B)$	$\mathbb{P}(B)$
	$\mathbb{P}(A^c)$	$\mathbb{P}(A)$	

Consequently, by Proposition 3.6.11, the random variables $\mathbb{1}_A$ and $\mathbb{1}_B$ are independent if and only if the following equations are valid:

$$\begin{aligned} \mathbb{P}(A^c \cap B^c) &= \mathbb{P}(A^c) \cdot \mathbb{P}(B^c), & \mathbb{P}(A^c \cap B) &= \mathbb{P}(A^c) \cdot \mathbb{P}(B), \\ \mathbb{P}(A \cap B^c) &= \mathbb{P}(A) \cdot \mathbb{P}(B^c), & \mathbb{P}(A \cap B) &= \mathbb{P}(A) \cdot \mathbb{P}(B). \end{aligned}$$

Because of Proposition 2.2.7, these four equations are satisfied if and only if the events A and B are independent. This proves the third assertion. □

Finally, let us shortly discuss the independence of more than two discrete random variables. Hereby we use the same notation as in Proposition 3.5.10, that is, the random variables X_1, \dots, X_n satisfy $X_j : \Omega \rightarrow D_j$, where D_j is either finite or countably infinite. Then the following generalization of Proposition 3.6.11 is valid. Its proof is almost identical to that for two variables. Therefore, we omit it.

⁹ Use, for example, that $\mathbb{1}_A(\omega) = 0$ and $\mathbb{1}_B(\omega) = 0$ if and only if $\omega \in A^c \cap B^c$.

Proposition 3.6.18. *The random variables X_1, \dots, X_n are independent if and only if for all $x_j \in D_j$,*

$$\mathbb{P}\{X_1 = x_1, \dots, X_n = x_n\} = \mathbb{P}\{X_1 = x_1\} \cdots \mathbb{P}\{X_n = x_n\}.$$

Example 3.6.19. Let us consider the problem of tossing a biased coin n times. The sample space is $\Omega = \{0, 1\}^n$, and the describing probability measure \mathbb{P} is as in eq. (3.5). The random variables X_j are defined as results of toss j . Then $X_j : \Omega \rightarrow D_j$, where $D_j = \{0, 1\}$. If we choose arbitrary $x_j \in D_j$, then either $x_j = 0$ or $x_j = 1$. Let k be the number of those x_j , which equals 1, that is, $k = x_1 + \dots + x_n$. Formula (3.5) implies

$$\mathbb{P}\{X_1 = x_1, \dots, X_n = x_n\} = \mathbb{P}\{(x_1, \dots, x_n)\} = p^k (1-p)^{n-k}.$$

On the other hand, as shown in Example 3.2.16, the probability distribution of each X_j satisfies

$$\mathbb{P}\{X_j = 0\} = 1-p \quad \text{and} \quad \mathbb{P}\{X_j = 1\} = p.$$

Since exactly k of the x_j s are “1” and $n-k$ are “0,” this implies

$$\mathbb{P}\{X_1 = x_1\} \cdots \mathbb{P}\{X_n = x_n\} = p^k (1-p)^{n-k}.$$

Summing up, for all $x_j \in D_j$,

$$\mathbb{P}\{X_1 = x_1, \dots, X_n = x_n\} = p^k (1-p)^{n-k} = \mathbb{P}\{X_1 = x_1\} \cdots \mathbb{P}\{X_n = x_n\},$$

that is, X_1, \dots, X_n are independent.

Summary: Two discrete random variables X and Y with values $(x_i)_{i \geq 1}$ and $(y_j)_{j \geq 1}$, respectively, are independent if and only if for all $1 \leq i, j < \infty$,

$$\mathbb{P}\{X = x_i, Y = y_j\} = \mathbb{P}\{X = x_i\} \cdot \mathbb{P}\{Y = y_j\}.$$

3.6.2 Independence of continuous random variables

We will consider the question of when *continuous* random variables are independent. Thus, let X_1, \dots, X_n be continuous random variables with distribution densities p_1, \dots, p_n , that is, for $1 \leq j \leq n$ and real numbers $a < b$,

$$\mathbb{P}_{X_j}([a, b]) = \mathbb{P}\{a \leq X_j \leq b\} = \int_a^b p_j(x) dx.$$

With this notation, the independence of the X_j s may be characterized as follows.

Proposition 3.6.20. For random variables X_1, \dots, X_n with densities p_1, \dots, p_n , we define a function $p : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$p(x_1, \dots, x_n) := p_1(x_1) \cdots p_n(x_n), \quad (x_1, \dots, x_n) \in \mathbb{R}^n. \quad (3.22)$$

Then the X_j s are independent if and only if p defined by eq. (3.22) is a distribution density of the random vector $\vec{X} = (X_1, \dots, X_n)$.

Proof. As in the discrete case, the result follows directly from Propositions 1.9.16 and 3.6.5. Without using product probabilities, we may argue as follows.

First, we observe that p defined by eq. (3.22) is a distribution density of \vec{X} if and only if for all $a_j < b_j$,

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\} = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} p_1(x_1) \cdots p_n(x_n) dx_n \cdots dx_1. \quad (3.23)$$

The right-hand side of eq. (3.23) coincides with

$$\left(\int_{a_1}^{b_1} p_1(x_1) dx_1 \right) \cdots \left(\int_{a_n}^{b_n} p_n(x_n) dx_n \right) = \mathbb{P}\{a_1 \leq X_1 \leq b_1\} \cdots \mathbb{P}\{a_n \leq X_n \leq b_n\}.$$

From this we derive that eq. (3.23) is valid for all $a_j < b_j$ if and only if

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\} = \mathbb{P}\{a_1 \leq X_1 \leq b_1\} \cdots \mathbb{P}\{a_n \leq X_n \leq b_n\}.$$

By Remark 3.6.10, this is equivalent to the independence of the X_j s, completing the proof. \square

Example 3.6.21. Throw a dart to a target, which is a circle of radius 1. The center of the circle is the point $(0, 0)$ and $(x_1, x_2) \in K$ denotes the point where the dart hits the target. We assume that the point hit is uniformly distributed on K . The question is whether or not the coordinates x_1 and x_2 of the point hit are dependent or independent of each other.

Answer: Let \mathbb{P} be the uniform distribution on K , and define two random variables X_1 and X_2 by $X_1(x_1, x_2) = x_1$ and $X_2(x_1, x_2) = x_2$. In this notation, the above question is whether the random variables X_1 and X_2 are independent. The density p_1 of X_1 was found in eq. (3.7). By symmetry, p_2 , the density of X_2 , coincides with p_1 , that is, we have

$$p_1(x_1) = \begin{cases} \frac{2}{\pi} \sqrt{1-x_1^2} & \text{if } |x_1| \leq 1, \\ 0 & \text{if } |x_1| > 1, \end{cases} \quad \text{and} \quad p_2(x_2) = \begin{cases} \frac{2}{\pi} \sqrt{1-x_2^2} & \text{if } |x_2| \leq 1, \\ 0 & \text{if } |x_2| > 1. \end{cases}$$

But $p_1(x_1) \cdot p_2(x_2)$ cannot be a distribution density of $\mathbb{P}_{(X_1, X_2)}$. Indeed, the vector $\vec{X} = (X_1, X_2)$ is uniformly distributed on K , thus its (correct) density is p with

$$p(x_1, x_2) = \begin{cases} \frac{1}{\pi} & \text{if } x_1^2 + x_2^2 \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Thus, we conclude that X_1 and X_2 are dependent, hence also the coordinates x_1 and x_2 of the point hit.

Example 3.6.22. We suppose now that the dart does not hit a circle but some rectangle set $R := [\alpha_1, \beta_1] \times [\alpha_2, \beta_2]$. Again we assume that the point $(x_1, x_2) \in R$ is uniformly distributed on R . The posed question is the same as in Example 3.6.21, namely whether x_1 and x_2 are independent of each other.

Answer: Define X_1 and X_2 as in the previous example. By assumption, the vector $\vec{X} = (X_1, X_2)$ is uniformly distributed on R , hence its distribution density p is given by

$$p(x_1, x_2) = \begin{cases} \frac{1}{\text{vol}_2(R)} & \text{if } (x_1, x_2) \in R, \\ 0 & \text{if } (x_1, x_2) \notin R. \end{cases}$$

For the density p_1 of X_1 , we get

$$p_1(x_1) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_2 = \frac{\beta_2 - \alpha_2}{\text{vol}_2(R)} = \frac{1}{\beta_1 - \alpha_1}$$

provided that $\alpha_1 \leq x_1 \leq \beta_1$. Otherwise, we have $p_1(x_1) = 0$. This tells us that X_1 is uniformly distributed on $[\alpha_1, \beta_1]$. In the same way, we obtain for $x_2 \in [\alpha_2, \beta_2]$ that

$$p_2(x_2) = \frac{1}{\beta_2 - \alpha_2}$$

and $p_2(x_2) = 0$ otherwise. Hence, X_2 is also uniformly distributed, but this time on $[\alpha_2, \beta_2]$. From the equations for p_1 and p_2 it follows that for the joint density p one has

$$p(x_1, x_2) = p_1(x_1) \cdot p_2(x_2), \quad (x_1, x_2) \in \mathbb{R}^2.$$

Consequently, by Proposition 3.6.20, the random variables X_1 and X_2 are independent, and so are the coordinates x_1 and x_2 of the point hit.

Example 3.6.23. Let us look at Example 3.6.22 from the opposite side. Now we assume that the coordinates are uniformly distributed, not the vector. Thus let U_1, \dots, U_n be independent random variables with U_j uniformly distributed on the interval $[\alpha_j, \beta_j]$, $1 \leq j \leq n$. Then the random vector $\vec{U} = (U_1, \dots, U_n)$ is (multivariate) uniformly distributed on the box $K = [\alpha_1, \beta_1] \times \dots \times [\alpha_n, \beta_n]$. This is an immediate consequence of Example 1.9.17 combined with Proposition 3.6.5. A direct proof of this fact, without using product measures, is as follows.

The density of U_j is $p_j = \frac{1}{\beta_j - \alpha_j} \mathbb{1}_{[\alpha_j, \beta_j]}$, hence by Proposition 3.6.20 the joint density p of \vec{U} is given by

$$p(x) = p_1(x_1) \cdots p_n(x_n) = \prod_{j=1}^n (\beta_j - \alpha_j)^{-1} = \frac{1}{\text{vol}_n(K)}, \quad x = (x_1, \dots, x_n) \in K,$$

and $p(x) = 0$ if $x \notin K$. Therefore,

$$\mathbb{P}\{\vec{U} \in B\} = \int_B p(x) dx = \frac{\text{vol}_n(K \cap B)}{\text{vol}_n(K)}, \quad B \in \mathcal{B}(\mathbb{R}^n),$$

and \vec{U} is uniformly distributed on K as asserted.

Example 3.6.24. Let X_1, \dots, X_n be independent standard normally distributed. Which joint density does the vector $\vec{X} = (X_1, \dots, X_n)$ possess?

Answer: The density p_j of the X_j , for $j = 1, \dots, n$, is

$$p_j(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad x \in \mathbb{R}.$$

Consequently, by the independence of the X_j s, the joint density p equals

$$\begin{aligned} p(x) &= p_1(x_1) \cdots p_n(x_n) = \frac{1}{(2\pi)^{n/2}} e^{-(x_1^2 + \dots + x_n^2)/2} \\ &= \frac{1}{(2\pi)^{n/2}} e^{-|x|^2/2}, \quad x = (x_1, \dots, x_n). \end{aligned}$$

This tells us that $\mathbb{P}_{\vec{X}} = \mathcal{N}(0, 1)^{\otimes n}$ (cf. Definition 1.9.21) or, equivalently, \vec{X} is n -dimensional standard normally distributed.

Example 3.6.25. If X_1, \dots, X_n are independent E_λ -distributed, then

$$p_j(t) = \begin{cases} 0 & \text{if } t < 0, \\ \lambda e^{-\lambda t} & \text{if } t \geq 0, \end{cases}$$

hence, the random vector $\vec{X} = (X_1, \dots, X_n)$ has the joint density

$$p(t) = \lambda^n e^{-\lambda(t_1 + \dots + t_n)}, \quad t = (t_1, \dots, t_n), \quad t_j \geq 0,$$

and $p(t) = 0$ if one of the t_j s is negative.

Example 3.6.26. Suppose two random variables X_1 and X_2 are independent and E_λ (exponentially) and uniformly on $[0, 1]$ distributed, respectively. Which distribution does the vector $\vec{X} = (X_1, X_2)$ possess?

The density p of \vec{X} equals

$$p(t, s) = \begin{cases} \lambda e^{-\lambda t} \cdot \mathbb{1}_{[0,1]}(s) & \text{if } t \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

For example, if $B = \{(t, s) : 0 \leq t \leq s \leq 1\}$, then

$$\mathbb{P}\{\vec{X} \in B\} = \lambda \int_0^1 \int_0^s e^{-\lambda t} dt ds = \int_0^1 (1 - e^{-\lambda s}) ds = 1 + \frac{1}{\lambda}(e^{-\lambda} - 1).$$

When does this event B occur? Say the lifetime of a component is exponentially distributed with parameter $\lambda > 0$. Then the event B occurs if the component stops working before a randomly chosen time $s \in [0, 1]$. This number s is taken uniformly distributed on $[0, 1]$ and, moreover, independent of the lifetime of the component. For example, first choose the number $s \in [0, 1]$, then check whether or not the component becomes defective before time $s \in [0, 1]$.

Summary: The random variables X_1, \dots, X_n are independent if for all Borel sets B_1, \dots, B_n in \mathbb{R} it follows that

$$\mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} = \mathbb{P}\{X_1 \in B_1\} \cdots \mathbb{P}\{X_n \in B_n\}.$$

In other words, the joint distribution $\mathbb{P}_{(X_1, \dots, X_n)}$ has to coincide with the product measure of the marginal distributions \mathbb{P}_{X_j} , $1 \leq j \leq n$. Another equivalent condition for the independence is: for all real numbers t_1, \dots, t_n , it follows that

$$\mathbb{P}\{X_1 \leq t_1, \dots, X_n \leq t_n\} = \mathbb{P}\{X_1 \leq t_1\} \cdots \mathbb{P}\{X_n \leq t_n\}.$$

If X_1, \dots, X_n are continuous random variables with densities p_1, \dots, p_n , then they are independent if and only if

$$p(x_1, \dots, x_n) = p_1(x_1) \cdots p_n(x_n), \quad (x_1, \dots, x_n) \in \mathbb{R}^n,$$

is a density of the random n -dimensional vector $\vec{X} = (X_1, \dots, X_n)$.

3.7 Order statistics*

This section is devoted to a quite practical problem. Suppose we execute independently of each other the same random experiment n times. Say the results are the real numbers x_1, \dots, x_n . For example, one may think of n different measurements of the same item, and x_1, \dots, x_n are the observed values. After getting the x_j s, we reorder them by their size. These “new” numbers are denoted by x_1^*, \dots, x_n^* and satisfy $x_1^* \leq \dots \leq x_n^*$. In other words, the numbers are the same as before but now in nondecreasing order.

A slightly more precise way to introduce the ordered sample x_1^*, \dots, x_n^* is as follows: There exists at least one permutation π of order n (maybe more than one if some of the x_j s are equal) for which $x_{\pi(1)} \leq x_{\pi(2)} \leq \dots \leq x_{\pi(n)}$. Setting

$$x_1^* = x_{\pi(1)}, \quad x_2^* = x_{\pi(2)}, \quad \text{up to} \quad x_n^* = x_{\pi(n)},$$

the new sample x_1^*, \dots, x_n^* consists of the same values as x_1, \dots, x_n , but now these values are ordered by their size.

For example, if $x_1 = 8.5$, $x_2 = 7.1$, and $x_3 = 7.9$, then we obtain $x_1^* = 7.1$, $x_2^* = 7.9$, and $x_3^* = 8.5$.

In order to apply this ordering procedure to random observations, one needs an algorithm for constructing the x_j^* s. This is easy if $x_i \neq x_j$ for $i \neq j$. Then

$$x_1^* = \min\{x_i : 1 \leq i \leq n\} \quad \text{while} \quad x_k^* = \min\{x_i : x_i > x_{k-1}^*\} \quad \text{if } k \geq 2.$$

For general values, that is, not necessarily different ones, the construction is slightly more complicated. The basic observation is that the ordered values possess the following property: given $t \in \mathbb{R}$, for all $1 \leq k \leq n$,

$$x_k^* \leq t \quad \Leftrightarrow \quad \text{Number of } x_i \leq t \text{ is at least } k \quad \Leftrightarrow \quad |\{i \leq n : x_i \leq t\}| \geq k. \quad (3.24)$$

This implies that

$$x_k^* = \inf\{t \in \mathbb{R} : |\{i \leq n : x_i \leq t\}| \geq k\} = \inf\left\{t \in \mathbb{R} : \sum_{i=1}^n \mathbb{1}_{(-\infty, t]}(x_i) \geq k\right\}. \quad (3.25)$$

For better understanding, here an easy example.

Example 3.7.1. Assume our sample is

$$x_1 = 3, \quad x_2 = 1, \quad x_3 = 3, \quad \text{and} \quad x_4 = 2.$$

Then we get

$$|\{i \leq 4 : x_i \leq t\}| = \begin{cases} 0 & \text{if } -\infty < t < 1, \\ 1 & \text{if } 1 \leq t < 2, \\ 2 & \text{if } 2 \leq t < 3, \\ 4 & \text{if } 3 \leq t < \infty, \end{cases}$$

which implies

$$\inf\{t \in \mathbb{R} : |\{i \leq 4 : x_i \leq t\}| \geq k\} = \begin{cases} 1 & \text{if } k = 1, \\ 2 & \text{if } k = 2, \\ 3 & \text{if } k = 3, \\ 3 & \text{if } k = 4. \end{cases}$$

So we finally arrive at $x_1^* = 1$, $x_2^* = 2$, $x_3^* = 3$, and $x_4^* = 3$.

The basic question treated in this section is now as follows. Suppose the values x_1, \dots, x_n were obtained by n random experiments, independently and according to a

known distribution. How are then the ordered x_k^* s distributed? For example, one rolls a fair die 10 times and observes x_1, \dots, x_{10} in $\{1, \dots, 6\}$. Then, for instance, one may ask how likely it is that the third smallest value x_3^* equals 5. Or which probability does the event $\{x_4^* = 3\}$ possess?

The precise mathematical formulation of this problem is as follows: let X_1, \dots, X_n be n independent identically distributed random variables defined on a sample space Ω . Define random variables X_1^*, \dots, X_n^* as its values ordered by their size. One possible way to define these new random variables is by, for example, using eq. (3.25). That is, given $k \leq n$ and $\omega \in \Omega$,

$$\begin{aligned} X_k^*(\omega) &= \inf\{t \in \mathbb{R} : |\{i \leq n : X_i(\omega) \leq t\}| \geq k\} \\ &= \inf\left\{t \in \mathbb{R} : \sum_{i=1}^n \mathbb{1}_{(-\infty, t]}(X_i(\omega)) \geq k\right\}. \end{aligned}$$

In particular, we get

$$X_1^* = \min\{X_1, \dots, X_n\} \quad \text{and} \quad X_n^* = \max\{X_1, \dots, X_n\}.$$

Remark 3.7.2. Those to whom the definition of the X_k^* s looks too complicated may use the following construction: for each fixed $\omega \in \Omega$, one chooses a permutation π of order n , depending on ω , for which

$$X_{\pi(1)}(\omega) \leq X_{\pi(2)}(\omega) \leq \dots \leq X_{\pi(n)}(\omega).$$

Then one gets the X_k^* s by setting

$$X_1^*(\omega) = X_{\pi(1)}(\omega), \quad X_2^*(\omega) = X_{\pi(2)}(\omega), \quad \text{up to} \quad X_n^*(\omega) = X_{\pi(n)}(\omega).$$

In this connection, it is important that there are at most finitely many (not more than $n!$) different permutations π , depending on ω , for ordering the values $X_1(\omega)$ up to $X_n(\omega)$ by their size. Thus, our sample space splits into at most $n!$ subsets where in each of them one special permutation orders all $X_j(\omega)$ in nondecreasing order.

Definition 3.7.3. The ordered random variables X_1^*, \dots, X_n^* are called the **order statistics** of X_1, \dots, X_n . Similarly, if x_1, \dots, x_n are observed random real numbers, the ordered x_k^* s are the **order statistics** of the x_j s.

Remark 3.7.4. By construction, the random variables X_k^* satisfy $X_1^* \leq \dots \leq X_n^*$. But note that they are no longer independent nor are identically distributed.

Remark 3.7.5. Order statistics play an important role in Mathematical Statistics. For example, suppose at time $t = 0$ we switch on n light bulbs of the same type. Let us record the times $0 < t_1^* < t_2^* < \dots < t_n^*$, where some of the n bulbs burn out. Then these times

are nothing else than the order statistics of the lifetimes t_1, \dots, t_n of the first, second, and so on, light bulb.

Before we state and prove the main result of this section, let us recall that the X_j s are assumed to be identically distributed. Consequently, all of them possess the same distribution function F . That is, for all $j \leq n$, we have

$$F(t) = \mathbb{P}\{X_j \leq t\}, \quad t \in \mathbb{R},$$

Proposition 3.7.6. *Let X_1, \dots, X_n be independent identically distributed random variables with distribution function F . Then for each $k \leq n$, we have*

$$\mathbb{P}\{X_k^* \leq t\} = \sum_{i=k}^n \binom{n}{i} F(t)^i (1-F(t))^{n-i}, \quad t \in \mathbb{R}. \quad (3.26)$$

Proof. Fix $t \in \mathbb{R}$. When does the event $\{X_k^* \leq t\}$ occur? To answer this, for $i \leq n$ introduce disjoint sets A_i as follows: the event A_i occurs if and only if exactly i of the X_j s attain a value in $(-\infty, t]$. More precisely,

$$A_i = \{\omega \in \Omega : |\{j \leq n : X_j(\omega) \leq t\}| = i\}.$$

Using eq. (3.24), the event $\{X_k^* \leq t\}$ occurs if and only if at least k of the X_j s attain a value in $(-\infty, t]$. Thus, by the definition of the A_i s, the event $\{X_k^* \leq t\}$ coincides with $\bigcup_{i=k}^n A_i$. Consequently, since the A_i s are disjoint, it follows that

$$\mathbb{P}\{X_k^* \leq t\} = \sum_{i=k}^n \mathbb{P}(A_i). \quad (3.27)$$

Let $Y_j = \mathbb{1}_{(-\infty, t]}(X_j)$. Then $Y_j = 1$ if and only if $X_j \leq t$ while $Y_j = 0$ otherwise. Hence, the Y_j s are binomial distributed with parameters 1 and p , where

$$p = \mathbb{P}\{Y_j = 1\} = \mathbb{P}\{X_j \leq t\} = F(t).$$

Since the X_j s are independent, so are the Y_j s and their sum¹⁰ $Y_1 + \dots + Y_n$ is binomial distributed with parameters n and $p = F(t)$. Note that the event A_i occurs if and only if $Y_1 + \dots + Y_n = i$, which implies

$$\mathbb{P}(A_i) = \mathbb{P}\{Y_1 + \dots + Y_n = i\} = \binom{n}{i} p^i (1-p)^{n-i} = \binom{n}{i} F(t)^i (1-F(t))^{n-i}. \quad (3.28)$$

Plugging eq. (3.28) into eq. (3.27) proves eq. (3.26). □

¹⁰ Here we already use a result, which will be proved later on in Proposition 4.6.1.

Example 3.7.7. Let us choose independently and according to the uniform distribution n numbers x_1, \dots, x_n out of $\{1, \dots, N\}$. Here, the same number may be chosen more than once. Given integers $m \leq N$ and $k \leq n$, find the probability that the k th largest number x_k^* equals m .

Answer: The distribution function F of the uniform distribution on $\{1, \dots, N\}$ satisfies

$$F(m) = \frac{m}{N}, \quad m = 1, \dots, N.$$

Thus Proposition 3.7.6 implies

$$\mathbb{P}\{x_k^* \leq m\} = \sum_{i=k}^n \binom{n}{i} \left(\frac{m}{N}\right)^i \left(1 - \frac{m}{N}\right)^{n-i}. \quad (3.29)$$

In particular,

$$\mathbb{P}\{x_n^* = 1\} = \mathbb{P}\{x_n^* \leq 1\} = \left(\frac{1}{N}\right)^n \quad \text{and} \quad \mathbb{P}\{x_k^* \leq N\} = 1, \quad k = 1, \dots, n.$$

Because of $\{x_k^* = m\} = \{x_k^* \leq m\} \setminus \{x_k^* \leq m-1\}$, we obtain

$$\begin{aligned} \mathbb{P}\{x_k^* = m\} &= \mathbb{P}\{x_k^* \leq m\} - \mathbb{P}\{x_k^* \leq m-1\} \\ &= \sum_{i=k}^n \binom{n}{i} \left[\left(\frac{m}{N}\right)^i \left(1 - \frac{m}{N}\right)^{n-i} - \left(\frac{m-1}{N}\right)^i \left(1 - \frac{m-1}{N}\right)^{n-i} \right]. \end{aligned} \quad (3.30)$$

Here the right-hand expression vanishes in the case $m = 1$. For instance, we get

$$\mathbb{P}\{x_1^* = 1\} = \sum_{i=1}^n \binom{n}{i} \left(\frac{1}{N}\right)^i \left(1 - \frac{1}{N}\right)^{n-i} = 1 - \left(1 - \frac{1}{N}\right)^n.$$

Another, more direct, approach for this result is as follows: One has $x_1^* > 1$ if and only if all x_j s satisfy $x_j \geq 2$. And among all possible N^n ways to choose the x_j s, there are $(N-1)^n$ possible ones to choose the x_j s in $\{2, \dots, N\}$.

Example 3.7.8. Roll a die four times and order the results in nondecreasing order as $x_1^* \leq x_2^* \leq x_3^* \leq x_4^*$. What is the probability that x_3^* equals m for some $1 \leq m \leq 6$?

Answer: Let us apply formula (3.29) with $N = 6$, $k = 3$, and $n = 4$. For $m \in \{1, \dots, 6\}$, this implies

$$\mathbb{P}\{x_3^* \leq m\} = \sum_{i=3}^4 \binom{4}{i} \left(\frac{m}{6}\right)^i \left(\frac{6-m}{6}\right)^{4-i},$$

hence

$$\mathbb{P}\{x_3^* = m\} = \sum_{i=3}^4 \binom{4}{i} \left[\left(\frac{m}{6}\right)^i \left(\frac{6-m}{6}\right)^{4-i} - \left(\frac{m-1}{6}\right)^i \left(\frac{6-m+1}{6}\right)^{4-i} \right].$$

The probabilities are

m	$\mathbb{P}\{x_3^* \leq m\}$	and	m	$\mathbb{P}\{x_3^* = m\}$
1	0.0162037		1	0.0162037
2	0.1111111		2	0.0949074
3	0.3125		3	0.201389
4	0.592593		4	0.280093
5	0.868056		5	0.275463
6	1		6	0.131944

Thus, the most likely value of x_3^* is the number “4.”

Remark 3.7.9. If we choose as in a lottery 6 numbers out of 49, then Proposition 3.7.6 and/or Example 3.7.7 do not apply. Why? Let x_1, \dots, x_6 be the numbers chosen first, second, and so on. Then they are identically distributed on $\{1, \dots, 49\}$, but they are *not* independent. For example, the probability that the second choice is number “2” given the first number was “1” is not $1/49$. This would be the case if the chosen number were replaced after each choice. We refer to Problem 3.5 for the distribution of the order statistics in the case of lottery numbers.

Let us now turn to the case of *continuous* random variables. That is, we assume that the random variables X_j possess a distribution density p satisfying

$$\mathbb{P}\{X_j \leq t\} = \int_{-\infty}^t p(x) dx, \quad t \in \mathbb{R}.$$

Again we remark that the preceding formula holds for all $j \leq n$. Indeed, the X_j s are identically distributed, hence they all have the same density. A natural question arises: what distribution density does X_k^* possess?

Proposition 3.7.10. *Suppose p is the common density of the X_j s. Let $X_1^* \leq \dots \leq X_n^*$ be the order statistics of the X_j . Then the distribution density p_k of X_k^* is given by*

$$p_k(t) = \frac{n!}{(k-1)!(n-k)!} p(t) F(t)^{k-1} (1-F(t))^{n-k}.$$

Proof. It holds that

$$\begin{aligned} p_k(t) &= \frac{d}{dt} \mathbb{P}\{X_k^* \leq t\} = \frac{d}{dt} \sum_{i=k}^n \binom{n}{i} F(t)^i (1-F(t))^{n-i} \\ &= \sum_{i=k}^n i \binom{n}{i} p(t) F(t)^{i-1} (1-F(t))^{n-i} - \sum_{i=k}^n (n-i) \binom{n}{i} p(t) F(t)^i (1-F(t))^{n-i-1}. \end{aligned} \quad (3.31)$$

In fact, the index i in the second sum of eq. (3.31) runs only from k to $n-1$. Hence, shifting it by 1, this sum becomes

$$\sum_{i=k+1}^n (n-i+1) \binom{n}{i-1} p(t) F(t)^{i-1} (1-F(t))^{n-i}.$$

Because of

$$i \binom{n}{i} = \frac{n!}{(i-1)!(n-i)!} = (n-i+1) \binom{n}{i-1},$$

both sums in eq. (3.31) cancel out for $i = k+1, \dots, n$. Thus, we obtain

$$\begin{aligned} p_k(t) &= k \binom{n}{k} p(t) F(t)^{k-1} (1-F(t))^{n-k} \\ &= \frac{n!}{(k-1)!(n-k)!} p(t) F(t)^{k-1} (1-F(t))^{n-k}, \end{aligned}$$

as asserted. \square

Example 3.7.11. Let us choose independently and according to the uniform distribution on $[0, 1]$ numbers x_1, \dots, x_n . After reordering them, we get $0 \leq x_1^* \leq \dots \leq x_n^* \leq 1$. Which distribution does x_k^* possess?

Answer: The density p of the uniform distribution on $[0, 1]$ is $\mathbb{1}_{[0,1]}$. Furthermore, its distribution function F is given by $F(t) = t$ for $0 \leq t \leq 1$. Thus, by Proposition 3.7.10, the density p_k coincides with

$$p_k(t) = \frac{n!}{(k-1)!(n-k)!} t^{k-1} (1-t)^{n-k}, \quad 0 \leq t \leq 1.$$

As already mentioned in Example 1.6.33, this is nothing else than the density of a beta distribution with parameters k and $n-k+1$. Hence, for all $k = 1, \dots, n$ and all $0 \leq a < b \leq 1$, it follows that

$$\mathbb{P}\{a \leq x_k^* \leq b\} = \mathcal{B}_{k, n-k+1}([a, b]) = \frac{n!}{(k-1)!(n-k)!} \int_a^b x^{k-1} (1-x)^{n-k} dx.$$

Example 3.7.12. Let us investigate here the example that was already mentioned in Remark 3.7.5. At time $t = 0$, we switch on n electric bulbs of the same type. The times $0 < t_1^* \leq \dots \leq t_n^*$ are those where we observe that some of the bulbs burn out. If we assume that the lifetime of each bulb is exponentially distributed, what can we say about the distribution of the t_k^* s?

Answer: Let X_1, \dots, X_n be the lifetimes of the n light bulbs. By assumption, they are independent and exponentially distributed with some parameter $\lambda > 0$. Then the distribution of t_k^* is that of X_k^* . Furthermore, we have $p(t) = \lambda e^{-\lambda t}$ and $F(t) = 1 - e^{-\lambda t}$ for $t \geq 0$.

By Proposition 3.7.10, the distribution density p_k of X_k^* equals

$$p_k(t) = \lambda \frac{n!}{(k-1)!(n-k)!} (1 - e^{-\lambda t})^{k-1} e^{-\lambda t(n-k+1)}, \quad t \geq 0.$$

For example, for t_1^* , the time when we observe the first burnout of any of the bulbs, this implies

$$p_1(t) = \lambda n e^{-\lambda t n}, \quad t \geq 0,$$

that is, t_1^* is $E_{\lambda n}$ -distributed.

Another case of interest is the behavior of t_n^* which is the time of the last outage of one of the n bulbs. Its density is

$$p_n(t) = \lambda n (1 - e^{-\lambda t})^{n-1} e^{-\lambda t}, \quad t > 0.$$

If $F_n(t) = (1 - e^{-\lambda t})^n$, $t > 0$, then $F_n'(t) = p_n(t)$. Moreover, $F_n(0) = 0$ and $F_n(\infty) = 1$, which implies that F_n is a distribution function and p_n is its density. That is,

$$\mathbb{P}\{t_n^* \leq t\} = (1 - e^{-\lambda t})^n, \quad t > 0.$$

Note that there is a more direct way to determine the distribution of t_n^* : use

$$\mathbb{P}\{t_n^* \leq t\} = \mathbb{P}\{X_1 \leq t, \dots, X_n \leq t\}$$

and the independence of X_1, \dots, X_n .

Summary: Let X_1, \dots, X_n be independent identically distributed random variables with distribution function $F: \mathbb{R} \rightarrow \mathbb{R}$. If $X_1^* \leq \dots \leq X_n^*$ is the sequence of the ordered values, then

$$\mathbb{P}\{X_k^* \leq t\} = \sum_{i=k}^n \binom{n}{i} F(t)^i (1 - F(t))^{n-i}, \quad t \in \mathbb{R}.$$

Moreover, if the X_j s are continuous with density p , then for $1 \leq k \leq n$, the density p_k of X_k^* equals

$$p_k(t) = \frac{n!}{(k-1)!(n-k)!} p(t) F(t)^{k-1} (1 - F(t))^{n-k}, \quad t \in \mathbb{R}.$$

3.8 Problems

Problem 3.1. For $n \geq 1$, let S_n be the set of permutations of order n . Suppose \mathbb{P} is the uniform distribution on S_n . That is, all permutations are equally likely. Define now a mapping X from S_n to $\{1, \dots, n\}$ by $X(\pi) = \pi(n)$, $\pi \in S_n$. Determine the probability distribution \mathbb{P}_X of X . What happens if one defines X by $X(\pi) = \pi(k)$ for some fixed $k \leq n$?

Problem 3.2. Roll a fair die twice. Let $\omega = (\omega_1, \omega_2)$ be the observed values. Define two random variables X_1 and X_2 by

$$X_1(\omega) = \max\{\omega_1, \omega_2\} \quad \text{and} \quad X_2(\omega) = \omega_1 + \omega_2.$$

Find the joint distribution of the random vector (X_1, X_2) , as well as its marginal ones. Argue why X_1 and X_2 are not independent.

Problem 3.3. The joint distribution of a random vector $\vec{X} = (X_1, X_2)$ is described by

$X_2 \setminus X_1$	0	1
0	$\frac{1}{10}$	$\frac{2}{5}$
1	$\frac{2}{5}$	$\frac{1}{10}$

Define another vector $\vec{Y} = (Y_1, Y_2)$ by $Y_1 := \min\{X_1, X_2\}$ and $Y_2 := \max\{X_1, X_2\}$. Find the probability distribution of $\vec{Y} = (Y_1, Y_2)$. Are Y_1 and Y_2 independent?

Problem 3.4. Let $\vec{X} = (X_1, X_2)$ be uniformly distributed on the square in \mathbb{R}^2 with corner points $(0, 1)$, $(1, 0)$, $(0, -1)$, and $(-1, 0)$. Find the marginal distributions of \vec{X} .

Problem 3.5. In a lottery, six numbers are chosen out of $\{1, \dots, 49\}$. As usual in lotteries, chosen numbers are not replaced. Let X_1, \dots, X_6 be the chosen numbers as they appeared. That is, X_1 is the number chosen first while X_6 is the number, which appeared last.

1. Determine the joint distribution of the vector $\vec{X} = (X_1, \dots, X_6)$, as well its marginal distributions.
2. Argue why X_1, \dots, X_6 are *not* independent.
3. Reordering the six chosen numbers leads to the order statistics $X_1^* < \dots < X_6^*$. Find the joint distribution of the vector (X_1^*, \dots, X_6^*) , as well as its marginal distributions.

Problem 3.6. A random variable X is geometrically distributed. Given natural numbers k and n , show that

$$\mathbb{P}\{X = k + n \mid X > n\} = \mathbb{P}\{X = k\}.$$

Why is this property called “lack of memory property”?

Problem 3.7. A random variable is exponentially distributed. Prove

$$\mathbb{P}\{X > s + t \mid X > s\} = \mathbb{P}\{X > t\}$$

for all $t, s \geq 0$. Why is this called the “nonaging property” of exponentially distributed lifetimes?

Problem 3.8. Two random variables X and Y are independent and geometrically distributed with parameters p and q for some $0 < p, q < 1$. Evaluate $\mathbb{P}\{X \leq Y \leq 2X\}$.

Problem 3.9. Suppose two independent random variables X and Y satisfy

$$\mathbb{P}\{X = k\} = \mathbb{P}\{Y = k\} = \frac{1}{2^k}, \quad k = 1, 2, \dots$$

Find the probabilities $\mathbb{P}\{X \leq Y\}$ and $\mathbb{P}\{X = Y\}$.

Problem 3.10. Choose two numbers b and c independently, the number b according to the uniform distribution on $[-1, 1]$ and c according to the uniform distribution on $[0, 1]$. Find the probability that the equation

$$x^2 + bx + c = 0$$

does *not* possess a real solution.

Problem 3.11. Use Problem 1.38 to prove the following: If X is a random variable, then the number of points $t \in \mathbb{R}$ with $\mathbb{P}\{X = t\} > 0$ is at most countably infinite.

Problem 3.12. Suppose a fair coin is labeled with “0” and “1.” Toss the coin n times. Let X be the maximum observed value and Y the sum of the n values. Determine the joint distribution of (X, Y) . Argue that X and Y are *not* independent.

Problem 3.13. Suppose a random vector (X, Y) has the joint density function p defined by

$$p(u, v) := \begin{cases} c \cdot uv & \text{if } u, v \geq 0, u + v \leq 1, \\ 0 & \text{if otherwise.} \end{cases}$$

- Find the value of the constant c so that p becomes a density function.
- Determine the density functions of X and Y .
- Evaluate $\mathbb{P}\{X + Y \leq 1/2\}$.
- Are X and Y independent?

Problem 3.14. Gambler A has a biased coin with “heads” having probability p for some $0 < p < 1$, and gambler B ’s coin is biased with “heads” having probability q for some $0 < q < 1$. Gamblers A and B toss their coins simultaneously. Whoever gets “heads” first wins. If both gamblers observe “heads” at the same time, then the game ends in a draw. Evaluate the probability that A wins and the probability that the game ends in a draw.

Problem 3.15. Randomly choose two integers x_1 and x_2 from 1 to 10. Let X be the *minimum* of x_1 and x_2 . Determine the distribution and the probability mass functions of X in the two following cases:

- The number chosen first is replaced.
- The first number is not replaced.

Evaluate in both cases $\mathbb{P}\{2 \leq X \leq 3\}$ and $\mathbb{P}\{X \geq 8\}$.

Problem 3.16. There are four balls labeled with “0” and three balls are labeled with “2” in an urn. Choose three balls without replacement. Let X be the sum of the values on the three chosen balls. Find the distribution of X .

Problem 3.17. As in Example 3.7.7, choose independently n numbers uniformly distributed in $\{1, \dots, N\}$. How likely is it that the largest of the chosen n numbers equals N ? Answer this question by applying formula (3.30). Give also a direct argument for the obtained result.

Problem 3.18. Roll a fair die 5 times. How likely is it that the fourth largest number equals m for some $m = 1, \dots, 6$?

4 Operations on random variables

4.1 Mappings of random variables

This section is devoted to the following problem: let $X : \Omega \rightarrow \mathbb{R}$ be a random variable and let $f : \mathbb{R} \rightarrow \mathbb{R}$ be some function. Set $Y := f(X)$, that is, for all $\omega \in \Omega$ we have $Y(\omega) = f(X(\omega))$. Suppose the distribution of X is known. Then the following task arises:



Determine the distribution of $Y = f(X)$ for a given function $f : \mathbb{R} \rightarrow \mathbb{R}$.

For example, if $f(t) = t^2$, and we know the distribution of X , then we ask for the probability distribution of X^2 . Is it possible to compute this by easy methods?

At the moment it is not clear at all whether $Y = f(X)$ is a random variable. Only if this is valid, the probability distribution \mathbb{P}_Y is well defined. For arbitrary functions f , this need not to be true, they have to satisfy the following additional property.

Definition 4.1.1. A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is called **measurable** if for $B \in \mathcal{B}(\mathbb{R})$ the preimage $f^{-1}(B)$ is a Borel set as well.

Remark 4.1.2. As all previous assumption about σ -fields, random variables, Borel sets, and so on, also this is a purely technical condition for f , which will not play an important role later on. All functions of interest, for example, piecewise continuous, monotone, pointwise limits of continuous functions, and so on, are measurable.

The measurability of f cannot be avoided because it is needed to prove the following result.

Proposition 4.1.3. Let $X : \Omega \rightarrow \mathbb{R}$ be a random variable. If $f : \mathbb{R} \rightarrow \mathbb{R}$ is measurable, then $Y = f(X)$ is a random variable as well.

Proof. Take a Borel set $B \in \mathcal{B}(\mathbb{R})$. Then

$$Y^{-1}(B) = X^{-1}(f^{-1}(B)) = X^{-1}(B'),$$

with $B' := f^{-1}(B)$. We assumed f to be measurable, which implies $B' \in \mathcal{B}(\mathbb{R})$, and hence, since X is a random variable, we conclude $Y^{-1}(B) = X^{-1}(B') \in \mathcal{A}$. The Borel set B was arbitrary, thus, as asserted, Y is a random variable. \square

So let $Y = f(X)$ for some measurable function $f : \mathbb{R} \rightarrow \mathbb{R}$ and some random variable X . Unfortunately, there does not exist a general method for the description of \mathbb{P}_Y in terms of \mathbb{P}_X . Only for some special functions f , for example, for linear functions or for f being strictly monotone and differentiable, there exist general rules for the computation of \mathbb{P}_Y . Nevertheless, quite often we are able to determine \mathbb{P}_Y directly. Mostly the following two approaches turn out to be helpful.

If X is *discrete* with values in $D := \{x_1, x_2, \dots\}$, then $Y = f(X)$ maps the sample space Ω into $f(D) = \{f(x_1), f(x_2), \dots\}$. Problems arise if f is not one-to-one. In this case one has to combine those x_j s that are mapped onto the same element in $f(D)$. For example, if X is uniformly distributed on $D = \{-2, -1, 0, 1, 2\}$ and $f(x) = x^2$, then $Y = X^2$ has values in $f(D) = \{0, 1, 4\}$. Combining -1 and 1 , as well as -2 and 2 , leads to

$$\begin{aligned}\mathbb{P}\{Y = 0\} &= \mathbb{P}\{X = 0\} = \frac{1}{5}, & \mathbb{P}\{Y = 1\} &= \mathbb{P}\{X = -1\} + \mathbb{P}\{X = 1\} = \frac{2}{5}, \\ \mathbb{P}\{Y = 4\} &= \mathbb{P}\{X = -2\} + \mathbb{P}\{X = 2\} = \frac{2}{5}.\end{aligned}$$

The case of one-to-one functions f is easier to handle because then

$$\mathbb{P}\{Y = f(x_j)\} = \mathbb{P}\{X = x_j\}, \quad j = 1, 2, \dots,$$

and the distribution of Y can be directly computed from that of X .

For *continuous* X , one tries to determine the distribution function F_Y of Y . Recall that this was defined as

$$F_Y(t) = \mathbb{P}\{Y \leq t\} = \mathbb{P}\{f(X) \leq t\}.$$

If we are able to compute F_Y , then we are almost done because then we get the distribution density q of Y as the derivative of F_Y .

For instance, if the continuous function f is increasing, one gets F_Y easily by

$$F_Y(t) = \mathbb{P}\{X \leq f^{-1}(t)\} = F_X(f^{-1}(t))$$

with the inverse function f^{-1} (cf. Problem 4.16).

The following examples demonstrate how we compute the distribution of $f(X)$ in some special cases.

Example 4.1.4. Assume the random variable X is $\mathcal{N}(0, 1)$ -distributed. Which distribution does $Y := X^2$ possess?

Answer: Of course, $F_Y(t) = \mathbb{P}\{Y \leq t\} = 0$ when $t \leq 0$. Consequently, it suffices to determine $F_Y(t)$ for $t > 0$. Then

$$\begin{aligned}F_Y(t) &= \mathbb{P}\{X^2 \leq t\} = \mathbb{P}\{-\sqrt{t} \leq X \leq \sqrt{t}\} = \frac{1}{\sqrt{2\pi}} \int_{-\sqrt{t}}^{\sqrt{t}} e^{-s^2/2} ds \\ &= \frac{2}{\sqrt{2\pi}} \int_0^{\sqrt{t}} e^{-s^2/2} ds = h(\sqrt{t}),\end{aligned}$$

where

$$h(u) := \frac{\sqrt{2}}{\sqrt{\pi}} \int_0^u e^{-s^2/2} ds, \quad u \geq 0.$$

Differentiating F_Y with respect to t , the chain rule and the fundamental theorem of Calculus lead to

$$\begin{aligned} q(t) &= F'_Y(t) = \frac{d}{dt}(\sqrt{t}) h'(\sqrt{t}) = \frac{t^{-1/2}}{2} \cdot \frac{\sqrt{2}}{\sqrt{\pi}} e^{-t/2} \\ &= \frac{1}{2^{1/2}\Gamma(1/2)} t^{\frac{1}{2}-1} e^{-t/2}, \quad t > 0. \end{aligned}$$

Hereby, in the last step, we used $\Gamma(1/2) = \sqrt{\pi}$. Consequently, Y possesses the density function

$$q(t) = \begin{cases} 0 & \text{if } t \leq 0, \\ \frac{1}{2^{1/2}\Gamma(1/2)} t^{-1/2} e^{-t/2} & \text{if } t > 0. \end{cases}$$

But this is the density of a $\Gamma_{2, \frac{1}{2}}$ -distribution. Therefore, we obtained the following result, which we, because of its importance, state as a separate proposition.

Proposition 4.1.5. *If X is $\mathcal{N}(0, 1)$ -distributed, then X^2 is $\Gamma_{2, \frac{1}{2}}$ -distributed or, equivalently, distributed according to χ_1^2 .*

Example 4.1.6. Let U be uniformly distributed on $[0, 1]$. Which distribution does the random variable $Y = 1/U$ possess?

Answer: Again we determine F_Y . From $\mathbb{P}\{X \in (0, 1]\} = 1$, we derive $\mathbb{P}\{Y \geq 1\} = 1$, thus, $F_Y(t) = 0$ if $t < 1$. Therefore, we only have to regard numbers $t \geq 1$. Here we have

$$F_Y(t) = \mathbb{P}\left\{\frac{1}{U} \leq t\right\} = \mathbb{P}\left\{U \geq \frac{1}{t}\right\} = 1 - \frac{1}{t}.$$

Hence, the density function q of Y is given by

$$q(t) = F'_Y(t) = \begin{cases} 0 & \text{if } t < 1, \\ \frac{1}{t^2} & \text{if } t \geq 1. \end{cases}$$

Example 4.1.7 (Random walk on \mathbb{Z}). A particle is located at the point 0 of \mathbb{Z} . In the first step, it moves either to -1 or to $+1$. In the second step, it jumps, independently of the first move, again to the left or to the right. Thus, after two steps it is located either at -2 , 0 , or 2 . Hereby we assume that p is the probability for jumps to the right, hence $1 - p$ for jumps to the left. This procedure is repeated arbitrarily often. Let S_n be the position of the particle after n jumps or, equivalently, after n steps.¹ The (random) se-

¹ The value of S_n can also be viewed as the loss or win after n games, where p is the probability to win one dollar in a single game, while $1 - p$ is the probability to lose one dollar.

quence $(S_n)_{n \geq 0}$ is called a (next-neighbor) **random walk** on \mathbb{Z} , where by the construction $\mathbb{P}\{S_0 = 0\} = 1$. See Figure 4.1 for a path of a random walk. Note that it is random, thus, very likely it will look completely different in another experiment.

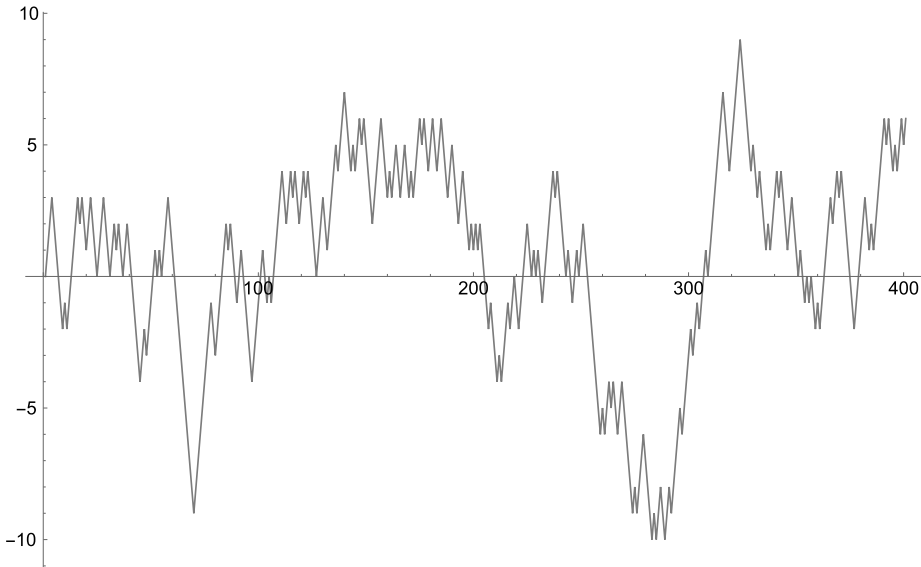


Figure 4.1: The (joined) points (n, S_n) , $n = 0, \dots, 400$, of a symmetric random walk S_n .

After n steps, the possible positions of the particle are in

$$D_n = \{-n, -n + 2, \dots, n - 2, n\}.$$

Note that an integer k belongs to D_n if and only if $|k| \leq n$ and, furthermore, $n + k$ (or equivalently, $n - k$) is even.

Thus, S_n is a random variable with values in D_n . Which distribution does S_n possess? To answer this question define

$$Y_n := \frac{1}{2} (S_n + n).$$

The random variable Y_n attains values in $\{0, 1, \dots, n\}$ and, moreover, $Y_n = m$ if the position of the particle after n steps is $2m - n$, that is, if it jumped m times to the right and $n - m$ times to the left. To see this, take $m = 0$, hence $S_n = -n$, which can only be achieved if all jumps were to the left. If $m = 1$, then $S_n = -n + 2$, that is, there were $n - 1$ jumps to the left and 1 to the right. The same argument applies for all $m \leq n$.

This observation tells us that Y_n is $B_{n,p}$ -distributed, that is,

$$\mathbb{P}\{Y_n = m\} = \binom{n}{m} p^m (1-p)^{n-m}, \quad m = 0, \dots, n.$$

Since $Y_n = \frac{1}{2}(S_n + n)$, if $k \in D_n$, then it follows that

$$\mathbb{P}\{S_n = k\} = \mathbb{P}\left\{Y_n = \frac{1}{2}(k + n)\right\} = \binom{n}{\frac{n+k}{2}} p^{(n+k)/2} (1-p)^{(n-k)/2}. \quad (4.1)$$

For even n , we have $0 \in D_n$, thus one may ask for the probability of $S_n = 0$, that is, for the probability that the particle returns to its starting point after n steps. Applying eq. (4.1) with $k = 0$ gives for even n that

$$\mathbb{P}\{S_n = 0\} = \binom{n}{\frac{n}{2}} p^{n/2} (1-p)^{n/2}.$$

Hence, if $p = 1/2$ (in this case the walk is said to be symmetric), then for even n we obtain

$$\mathbb{P}\{S_n = 0\} = \binom{n}{\frac{n}{2}} 2^{-n} = \frac{n!}{((n/2)!)^2} 2^{-n}.$$

An application of Stirling's formula (1.54) implies

$$\lim_{\substack{n \rightarrow \infty \\ n \text{ even}}} n^{1/2} \mathbb{P}\{S_n = 0\} = \lim_{\substack{n \rightarrow \infty \\ n \text{ even}}} n^{1/2} \frac{\sqrt{2\pi n} (n/e)^n}{[\sqrt{\pi n} (n/2e)^{n/2}]^2} 2^{-n} = \sqrt{\frac{2}{\pi}},$$

that is, if $n \rightarrow \infty$, then for even n it follows that $\mathbb{P}\{S_n = 0\} \sim \sqrt{\frac{2}{\pi}} n^{-1/2}$. Another way to formulate this is

$$\mathbb{P}\{S_{2n} = 0\} \sim \frac{1}{\sqrt{\pi n}}.$$

Example 4.1.8. Suppose X is $B_{n,p}^-$ -distributed, that is,

$$\mathbb{P}\{X = k\} = \binom{k-1}{k-n} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots$$

Let $Y = X - n$. Which probability distribution does Y possess?

Answer: An easy transformation (see formula (1.35)) leads to

$$\mathbb{P}\{Y = k\} = \mathbb{P}\{X = k + n\} = \binom{n+k-1}{k} p^n (1-p)^k = \binom{-n}{k} p^n (p-1)^k \quad (4.2)$$

for all $k = 0, 1, \dots$

Additional question: Which random experiment does Y describe?

Answer: We perform a series of random trials where each time we may obtain either failure or success. Hereby, the success probability is $p \in (0, 1)$. Then the event $\{Y = k\}$ occurs if and only if we observe the n th success in trial $k + n$.

We conclude this section with the investigation of the following problem. Suppose X_1, \dots, X_n are independent random variables. Given n measurable functions f_1, \dots, f_n from \mathbb{R} to \mathbb{R} , we define "new" random variables Y_1, \dots, Y_n by

$$Y_i := f_i(X_i), \quad 1 \leq i \leq n.$$

It is intuitively clear that then Y_1, \dots, Y_n are also independent; the values of Y_i only depend on those of X_i , thus the independence should be preserved. For example, if X_1 and X_2 are independent, then this should also be valid for X_1^2 and $2X_2$.

The next result shows that this is indeed true.

Proposition 4.1.9. *Let X_1, \dots, X_n be independent random variables and let $(f_i)_{i=1}^n$ be measurable functions from \mathbb{R} to \mathbb{R} . Then $f_1(X_1), \dots, f_n(X_n)$ are independent as well.*

Proof. Choose arbitrary Borel sets B_1, \dots, B_n in \mathbb{R} and set $A_i := f_i^{-1}(B_i)$, $1 \leq i \leq n$. With this notation, an $\omega \in \Omega$ satisfies $f_i(X_i(\omega)) \in B_i$ if and only if $X_i(\omega) \in A_i$. Hence, an application of the independence of X_i (use eq. (3.16) with the X_i 's and the A_i 's) leads to

$$\begin{aligned} \mathbb{P}\{f_1(X_1) \in B_1, \dots, f_n(X_n) \in B_n\} &= \mathbb{P}\{X_1 \in A_1, \dots, X_n \in A_n\} \\ &= \mathbb{P}\{X_1 \in A_1\} \cdots \mathbb{P}\{X_n \in A_n\} \\ &= \mathbb{P}\{f_1(X_1) \in B_1\} \cdots \mathbb{P}\{f_n(X_n) \in B_n\}. \end{aligned}$$

The B_i s were chosen arbitrarily, thus the random variables $f_1(X_1), \dots, f_n(X_n)$ are independent as well. \square

Remark 4.1.10. Without proof we still mention that the independence of random variables is preserved whenever they are put together into disjoint groups. For example, if X_1, \dots, X_n are independent, then so are $f(X_1, \dots, X_k)$ and $g(X_{k+1}, \dots, X_n)$ for suitable functions f and g . Assume we roll a die five times and let X_1, \dots, X_5 be the results. Then these random variables are independent, but so are the two random variables $\max\{X_1, X_2\}$ and $X_3 + X_4 + X_5$, or the three $X_1, \max\{X_2, X_3\}$, and $\min\{X_4, X_5\}$.

4.2 Linear transformations

Let a and b real numbers with $a \neq 0$. Given a random variable $Y = aX + b$, that is, Y arises from X by a linear transformation. We ask now for the probability distribution of Y .

Proposition 4.2.1. *Define $Y = aX + b$ with $a, b \in \mathbb{R}$ and $a \neq 0$.*

(a) *Depending on whether $a > 0$ or $a < 0$,*

$$F_Y(t) = F_X\left(\frac{t-b}{a}\right) \quad \text{or} \quad F_Y(t) = 1 - \mathbb{P}\left\{X < \frac{t-b}{a}\right\}.$$

If $a < 0$ and F_X is continuous at $\frac{t-b}{a}$, then

$$F_Y(t) = 1 - F_X\left(\frac{t-b}{a}\right).$$

(b) *Let X be a continuous random variable with density p . Then Y is also continuous with density q given by*

$$q(t) = \frac{1}{|a|} p\left(\frac{t-b}{a}\right), \quad t \in \mathbb{R}. \quad (4.3)$$

Proof. Let us first treat the case $a > 0$. Then we get

$$F_Y(t) = \mathbb{P}\{aX + b \leq t\} = \mathbb{P}\left\{X \leq \frac{t-b}{a}\right\} = F_X\left(\frac{t-b}{a}\right),$$

as asserted.

In the case $a < 0$, we conclude as follows:

$$F_Y(t) = \mathbb{P}\{aX + b \leq t\} = \mathbb{P}\left\{X \geq \frac{t-b}{a}\right\} = 1 - \mathbb{P}\left\{X < \frac{t-b}{a}\right\}.$$

If F_X is continuous at $\frac{t-b}{a}$, then $\mathbb{P}\{X = \frac{t-b}{a}\} = 0$, hence

$$1 - \mathbb{P}\left\{X < \frac{t-b}{a}\right\} = 1 - \mathbb{P}\left\{X \leq \frac{t-b}{a}\right\} = 1 - F_X\left(\frac{t-b}{a}\right),$$

completing the proof of part (a).

Suppose now that p is a density function of X , that is,

$$F_X(t) = \mathbb{P}\{X \leq t\} = \int_{-\infty}^t p(x) dx, \quad t \in \mathbb{R}.$$

If $a > 0$, by part (a) and after the change of variables $x = \frac{y-b}{a}$, we get

$$F_Y(t) = F_X\left(\frac{t-b}{a}\right) = \int_{-\infty}^{\frac{t-b}{a}} p(x) dx = \int_{-\infty}^t \frac{1}{a} p\left(\frac{y-b}{a}\right) dy = \int_{-\infty}^t q(y) dy.$$

Thus, q is a density of Y .

If $a < 0$, the same change of variables² leads to

$$\begin{aligned} F_Y(t) &= 1 - F_X\left(\frac{t-b}{a}\right) = \int_{\frac{t-b}{a}}^{\infty} p(x) dx = - \int_{-\infty}^t \frac{1}{a} p\left(\frac{y-b}{a}\right) dy \\ &= \int_{-\infty}^t \frac{1}{-a} p\left(\frac{y-b}{a}\right) dy = \int_{-\infty}^t \frac{1}{|a|} p\left(\frac{y-b}{a}\right) dy = \int_{-\infty}^t q(y) dy. \end{aligned}$$

This being true for all $t \in \mathbb{R}$ completes the proof. □

² Observe that now $a < 0$, hence the order of integration changes and a minus sign appears.

Example 4.2.2. Let X be $\mathcal{N}(0, 1)$ -distributed. Given $a \neq 0$ and $\mu \in \mathbb{R}$, we ask for the distribution of $Y = aX + \mu$.

Answer: The random variable X is known to be continuous with density

$$p(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}.$$

We apply eq. (4.3) with $b = \mu$ to deduce that the density q of Y equals

$$q(t) = \frac{1}{|a|} p\left(\frac{t - \mu}{|a|}\right) = \frac{1}{\sqrt{2\pi} |a|} e^{-(t - \mu)^2/2a^2}.$$

That is, the random variable Y is $\mathcal{N}(\mu, |a|^2)$ -distributed. In particular, if $\sigma > 0$ and $\mu \in \mathbb{R}$, then $\sigma X + \mu$ is distributed according to $\mathcal{N}(\mu, \sigma^2)$.

Additional question: Suppose Y is $\mathcal{N}(\mu, \sigma^2)$ -distributed. Which probability distribution does $X := \frac{Y - \mu}{\sigma}$ possess?

Answer: Formula (4.3) immediately implies that X has a standard normal distribution.

Because of the importance of the previous observation, we formulate it as proposition.

Proposition 4.2.3. *Suppose $\mu \in \mathbb{R}$ and $\sigma > 0$. Then the following are equivalent:*

X is $\mathcal{N}(0, 1)$ -distributed $\iff \sigma X + \mu$ is distributed according to $\mathcal{N}(\mu, \sigma^2)$.



Corollary 4.2.4. *Let Φ be the Gaussian Φ -function introduced in eq. (1.70). For each interval $[a, b]$,*

$$\mathcal{N}(\mu, \sigma^2)([a, b]) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right).$$

Proof. This is a direct consequence of Proposition 4.2.3. Indeed, if X has a standard normal distribution, then

$$\begin{aligned} \mathcal{N}(\mu, \sigma^2)([a, b]) &= \mathbb{P}\{a \leq \sigma X + \mu \leq b\} = \mathbb{P}\left\{\frac{a - \mu}{\sigma} \leq X \leq \frac{b - \mu}{\sigma}\right\} \\ &= \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right), \end{aligned}$$

as asserted. □

Let X be an $\mathcal{N}(\mu, \sigma^2)$ -distributed random variable. The next result shows that X with high probability (more than 99.7%) attains values in $[\mu - 3\sigma, \mu + 3\sigma]$. Therefore, in most cases, one may assume that X maps into $[\mu - 3\sigma, \mu + 3\sigma]$. This observation is usually called “**Three Sigma Rule**”.

Corollary 4.2.5 (Three Sigma Rule). *If X is distributed according to $\mathcal{N}(\mu, \sigma^2)$, then*

$$\mathbb{P}\{|X - \mu| \leq 2\sigma\} \geq 0.954 \quad \text{and} \quad \mathbb{P}\{|X - \mu| \leq 3\sigma\} \geq 0.997.$$

Proof. By virtue of Corollary 4.2.4, for each $c > 0$,

$$\mathbb{P}\{|X - \mu| \leq c\sigma\} = \Phi(c) - \Phi(-c),$$

hence the desired estimates follow from

$$\Phi(2) - \Phi(-2) = 2\Phi(2) - 1 > 0.9545 \quad \text{and} \quad \Phi(3) - \Phi(-3) = 2\Phi(3) - 1 > 0.9973. \quad \square$$

Example 4.2.6. Let U be uniformly distributed on $[0, 1]$. What is the probability distribution of $aU + b$ if $a \neq 0$ and $b \in \mathbb{R}$?

Answer: The distribution density p of U is given by $p(t) = 1$ if $0 \leq t \leq 1$ and $p(t) = 0$ otherwise. Therefore, the density q of $aU + b$ equals

$$q(t) = \begin{cases} \frac{1}{|a|} & \text{if } 0 \leq \frac{t-b}{a} \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Assume first $a > 0$. Then $q(t) = 1/a$ if and only if $b \leq t \leq a+b$ and $q(t) = 0$ otherwise. Consequently, $aU + b$ is uniformly distributed on $[b, a+b]$.

If, in contrast, $a < 0$, then $q(t) = 1/|a|$ if and only if $a+b \leq t \leq b$ and $q(t) = 0$ otherwise. Hence, now $aU + b$ is uniformly distributed on $[a+b, b]$.

It is easy to see that the reversed implications are also true. That is, we have



U uniformly distributed on $[0, 1] \Leftrightarrow aU + b$ uniformly distributed on $\begin{cases} [b, a+b] & \text{if } a > 0, \\ [a+b, b] & \text{if } a < 0. \end{cases}$

Corollary 4.2.7. *A random variable U is uniformly distributed on $[0, 1]$ if and only if $1 - U$ is such. In other words, for a uniformly distributed U on $[0, 1]$ it follows that $U \stackrel{d}{=} 1 - U$.*

Example 4.2.8. Suppose a random variable X is $\Gamma_{\alpha, \beta}$ -distributed for some $\alpha, \beta > 0$ and let $a > 0$. Which distribution does aX possess?

Answer: The distribution density p of X satisfies $p(t) = 0$ if $t \leq 0$ and, if $t > 0$, then

$$p(t) = \frac{1}{\alpha^\beta \Gamma(\beta)} t^{\beta-1} e^{-t/\alpha}.$$

An application of eq. (4.3) implies that the density q of aX is given by $q(t) = 0$ if $t \leq 0$ and, if $t > 0$, then

$$q(t) = \frac{1}{a} p\left(\frac{t}{a}\right) = \frac{1}{a \alpha^\beta \Gamma(\beta)} \left(\frac{t}{a}\right)^{\beta-1} e^{-t/aa} = \frac{1}{(aa)^\beta \Gamma(\beta)} t^{\beta-1} e^{-t/aa}.$$

Thus, aX is $\Gamma_{aa,\beta}$ -distributed.

In particular, we have that X is $\Gamma_{1,\beta}$ -distributed if and only if for some (each) $a > 0$ it follows that aX is $\Gamma_{\alpha,\beta}$ -distributed.

In the case of the exponential distribution $E_\lambda = \Gamma_{1/\lambda,1}$, the previous result implies the following: if $a > 0$, then a random variable X is E_λ -distributed if and only if aX possesses an $E_{\lambda/a}$ distribution.

4.3 Coin tossing versus uniform distribution

4.3.1 Binary fractions

We start this section with the following statement: each real number $x \in [0, 1)$ may be represented as binary fraction $x = 0.x_1x_2\dots$, where $x_k \in \{0, 1\}$. This is a shortened way to express that

$$x = \sum_{k=1}^{\infty} \frac{x_k}{2^k}.$$

The representation of x as binary fraction is in general not unique. For example,

$$\frac{1}{2} = 0.1000\dots, \quad \text{but also} \quad \frac{1}{2} = 0.0111\dots$$

Check this by computing the infinite sums in both cases.

It is not difficult to prove that two different representations admit exactly those $x \in [0, 1)$ which may be written as $x = k/2^n$ for some $n \in \mathbb{N}$ and some $k = 1, 3, 5, \dots, 2^n - 1$. Those numbers are usually called **dyadic rational numbers**.

To make the binary representation unique, we declare the following:

Convention 4.1. *If a number $x \in [0, 1)$ admits the representations*

$$x = 0.x_1\dots x_{n-1}1000\dots \quad \text{and} \quad x = 0.x_1\dots x_{n-1}0111\dots,$$

then we always choose the former one. In other words, there do not exist numbers $x \in [0, 1)$ whose binary representation consists only of 1s from a certain point onward.

How do we get the binary fraction for a given $x \in [0, 1)$?

The procedure is not difficult. First, one checks whether $x < \frac{1}{2}$ or $x \geq \frac{1}{2}$. In the former case, one takes $x_1 = 0$ and in the latter $x_1 = 1$.

With this choice, it follows that $0 \leq x - \frac{x_1}{2} < \frac{1}{2}$. In the next step, one asks whether $x - \frac{x_1}{2} < \frac{1}{4}$ or $x - \frac{x_1}{2} \geq \frac{1}{4}$. Depending on this, one chooses either $x_2 = 0$ or $x_2 = 1$. This

choice implies $0 \leq x - \frac{x_1}{2} - \frac{x_2}{2^2} < \frac{1}{4}$, and if this difference belongs either to $[0, \frac{1}{8})$ or $[\frac{1}{8}, \frac{1}{4})$, then $x_3 = 0$ or $x_3 = 1$, respectively. Proceeding further in this way leads to the binary fraction representing x .

After this heuristic method, we now present a mathematically more exact way. To this end, for each $n \geq 1$, we divide the interval $[0, 1)$ into 2^n intervals of length 2^{-n} .

We start with $n = 1$ and divide $[0, 1)$ into two intervals,

$$I_0 := \left[0, \frac{1}{2} \right) \quad \text{and} \quad I_1 := \left[\frac{1}{2}, 1 \right).$$

In the second step, we divide each of the two intervals I_0 and I_1 further into two parts of equal length. In this way, we obtain the four intervals

$$I_{00} := \left[0, \frac{1}{4} \right), \quad I_{01} := \left[\frac{1}{4}, \frac{1}{2} \right), \quad I_{10} := \left[\frac{1}{2}, \frac{3}{4} \right), \quad \text{and} \quad I_{11} := \left[\frac{3}{4}, 1 \right).$$

Observe that the left end point of $I_{a_1 a_2}$ equals $a_1/2 + a_2/4$, that is,

$$I_{a_1 a_2} = \left[\sum_{j=1}^2 \frac{a_j}{2^j}, \sum_{j=1}^2 \frac{a_j}{2^j} + \frac{1}{2^2} \right), \quad a_1, a_2 \in \{0, 1\}.$$

It is clear now how to proceed. Given $n \geq 1$ and numbers $a_1, \dots, a_n \in \{0, 1\}$, set

$$I_{a_1 \dots a_n} = \left[\sum_{j=1}^n \frac{a_j}{2^j}, \sum_{j=1}^n \frac{a_j}{2^j} + \frac{1}{2^n} \right). \tag{4.4}$$

In this way, we obtain 2^n disjoint intervals of length 2^{-n} where the left corner points are $0.a_1 a_2 \dots a_n$ (see Figure 4.2 for the first dyadic intervals).

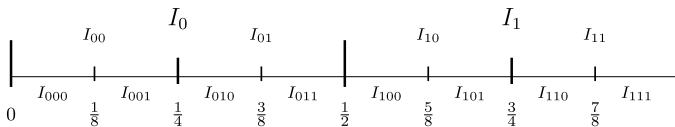


Figure 4.2: The first dyadic intervals of $[0, 1)$.

The following lemma makes the above heuristic method more precise.

Lemma 4.3.1. *For all $a_1, \dots, a_n \in \{0, 1\}$, the intervals in (4.4) are characterized by*

$$I_{a_1 \dots a_n} = \{x \in [0, 1) : x = 0.a_1 a_2 \dots a_n \dots \}.$$

Verbally, a number in $[0, 1)$ belongs to $I_{a_1 \dots a_n}$ if and only if its first n digits in the binary fraction are a_1, \dots, a_n .

Proof. Assume first $x \in I_{a_1 \dots a_n}$. If $a := 0.a_1 \dots a_n$ denotes the left end point of $I_{a_1 \dots a_n}$, by definition $a \leq x < a + 1/2^n$ or, equivalently, $0 \leq x - a < 1/2^n$. Therefore, the binary fraction of $x - a$ is of the form $0.00 \dots 0b_{n+1} \dots$ with certain numbers $b_{n+1}, b_{n+2}, \dots \in \{0, 1\}$. This yields

$$x = a + (x - a) = 0.a_1 \dots a_n b_{n+1} \dots$$

Thus, as asserted, the first n digits in the representation of x are a_1, \dots, a_n .

Conversely, if x can be written as $x = 0.x_1 x_2 \dots$ with $x_1 = a_1, \dots, x_n = a_n$, then $a \leq x$ where, as above, a denotes the left end point of $I_{a_1 \dots a_n}$. Moreover, by Convention 4.1, at least one of the x_k s, $k > n$, has to be zero. Consequently,

$$x - a = \sum_{k=n+1}^{\infty} \frac{x_k}{2^k} < \sum_{k=n+1}^{\infty} \frac{1}{2^k} = \frac{1}{2^n},$$

that is, we have $a \leq x < a + \frac{1}{2^n}$ or, equivalently, $x \in I_{a_1 \dots a_n}$ as asserted. \square

A direct consequence of Lemma 4.3.1 is as follows.

Corollary 4.3.2. *For each $n \geq 1$, the 2^n sets $I_{a_1 \dots a_n}$ form a disjoint partition of $[0, 1)$, that is,*

$$\bigcup_{a_1, \dots, a_n \in \{0, 1\}} I_{a_1 \dots a_n} = [0, 1) \quad \text{and} \quad I_{a_1 \dots a_n} \cap I_{a'_1 \dots a'_n} = \emptyset$$

provided that $(a_1, \dots, a_n) \neq (a'_1, \dots, a'_n)$. Furthermore,

$$\{x \in [0, 1) : x_k = 0\} = \bigcup_{a_1, \dots, a_{k-1} \in \{0, 1\}} I_{a_1 \dots a_{k-1} 0}.$$

4.3.2 Binary fractions of random numbers

We saw above each number $x \in [0, 1)$ admits a representation $x = 0.x_1 x_2 \dots$ with certain $x_k \in \{0, 1\}$. What does happen if we choose a number x randomly, say according to the uniform distribution on $[0, 1]$? Then the x_k s in the binary fraction are also random, with values in $\{0, 1\}$. How are they distributed?

The mathematical formulation of this question is as follows: let $U : \Omega \rightarrow \mathbb{R}$ be a random variable uniformly distributed on $[0, 1]$. If $\omega \in \Omega$, write³

$$U(\omega) = 0.X_1(\omega)X_2(\omega) \dots = \sum_{k=1}^{\infty} \frac{X_k(\omega)}{2^k}. \quad (4.5)$$

In this way, we obtain infinitely many random variables $X_k : \Omega \rightarrow \{0, 1\}$.

³ Note that $\mathbb{P}\{U \in [0, 1)\} = 1$. Thus, without losing generality, we may assume $U(\omega) \in [0, 1)$.

Which distribution do these random variables possess? Answer gives the next proposition.

Proposition 4.3.3. *If $k \in \mathbb{N}$, then*

$$\mathbb{P}\{X_k = 0\} = \mathbb{P}\{X_k = 1\} = \frac{1}{2}. \quad (4.6)$$

Furthermore, given $n \geq 1$, the random variables X_1, \dots, X_n are independent.

Proof. By assumption, \mathbb{P}_U is the uniform distribution on $[0, 1]$. Thus, the finite additivity of \mathbb{P}_U , Corollary 4.3.2 and eq. (1.46) imply

$$\begin{aligned} \mathbb{P}\{X_k = 0\} &= \mathbb{P}_U\left(\bigcup_{a_1, \dots, a_{k-1} \in \{0,1\}} I_{a_1 \dots a_{k-1} 0}\right) \\ &= \sum_{a_1, \dots, a_{k-1} \in \{0,1\}} \mathbb{P}_U(I_{a_1 \dots a_{k-1} 0}) = \sum_{a_1, \dots, a_{k-1} \in \{0,1\}} \frac{1}{2^k} = \frac{2^{k-1}}{2^k} = \frac{1}{2}. \end{aligned}$$

Since X_k attains only two different values, $\mathbb{P}\{X_k = 1\} = 1/2$ as well, proving the first part.

We want to verify that for all $n \geq 1$ the random variables X_1, \dots, X_n are independent. Equivalently, according to Proposition 3.6.18, the following has to be proven: if $a_1, \dots, a_n \in \{0, 1\}$, then

$$\mathbb{P}\{X_1 = a_1, \dots, X_n = a_n\} = \mathbb{P}\{X_1 = a_1\} \cdots \mathbb{P}\{X_n = a_n\}. \quad (4.7)$$

By eq. (4.6), the right-hand side of eq. (4.7) equals

$$\mathbb{P}\{X_1 = a_1\} \cdots \mathbb{P}\{X_n = a_n\} = \underbrace{\frac{1}{2} \cdots \frac{1}{2}}_n = \frac{1}{2^n}.$$

To compute the left-hand side of eq. (4.7), note that Lemma 4.3.1 implies that we have $X_1 = a_1$ up to $X_n = a_n$ if and only if U attains a value in $I_{a_1 \dots a_n}$. The intervals $I_{a_1 \dots a_n}$ are of length 2^{-n} , hence by eq. (1.46) (recall that \mathbb{P}_U is the uniform distribution on $[0, 1]$),

$$\mathbb{P}\{X_1 = a_1, \dots, X_n = a_n\} = \mathbb{P}\{U \in I_{a_1 \dots a_n}\} = \mathbb{P}_U(I_{a_1 \dots a_n}) = \frac{1}{2^n}.$$

Thus, for all $a_1, \dots, a_n \in \{0, 1\}$, eq. (4.7) is valid, and, as asserted, the random variables X_1, \dots, X_n are independent. \square

To formulate the previous result in a different way, let us introduce the following notation.

Definition 4.3.4. An infinite sequence X_1, X_2, \dots of random variables is said to be **independent** provided that any finite collection of the X_k s is independent.

Remark 4.3.5. Since any subcollection of independent random variables is independent as well, the independence of X_1, X_2, \dots is equivalent to the following. For all $n \geq 1$, the random variables X_1, \dots, X_n are independent, that is, for all $n \geq 1$ and all Borel sets B_1, \dots, B_n , it follows that

$$\mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} = \mathbb{P}\{X_1 \in B_1\} \cdots \mathbb{P}\{X_n \in B_n\}.$$

Remark 4.3.6. In view of Definition 4.3.4, the basic observation in Example 3.6.19 may now be formulated in the following way. If we toss a (maybe biased) coin, labeled with “0” and “1,” infinitely often and let X_1, X_2, \dots be the results of the single tosses, then this infinite sequence of random variables is independent with $\mathbb{P}\{X_k = 0\} = 1 - p$ and $\mathbb{P}\{X_k = 1\} = p$. In particular, for a fair coin, the X_k s possess the following properties:

1. If $k \in \mathbb{N}$, then $\mathbb{P}\{X_k = 0\} = \mathbb{P}\{X_k = 1\} = 1/2$.
2. X_1, X_2, \dots is an infinite sequence of independent random variables.

This observation leads us to the following definition.

Definition 4.3.7. An infinite sequence X_1, X_2, \dots of independent random variables with values in $\{0, 1\}$ satisfying

$$\mathbb{P}\{X_k = 0\} = \mathbb{P}\{X_k = 1\} = 1/2, \quad k = 1, 2, \dots,$$

is said to be a **model for tossing a fair coin infinitely often**.

Consequently, Proposition 4.3.3 asserts that the random variables X_1, X_2, \dots defined by eq. (4.5) serve as a model for tossing a fair coin infinitely often.

Example 4.3.8. Suppose a randomly chosen number in $[0, 1]$ is $x = 0.1657432763$. Its binary expansion is $0.001010100110111000100110101111100111101100010 \dots$. Thus, translating this into a result of tossing a coin, the first observations are $0, 0, 1, 0, 1, \dots$. Of course, since x is only approximately chosen uniformly, also only the first finitely many digits of the binary expansion simulate the tossing of a fair coin.

Summary: Let U be some random variable distributed according to the uniform distribution on $[0, 1]$. Represent the values of U as binary fraction $0.X_1X_2 \dots$ where the X_j s attain values in $\{0, 1\}$. Then the random variables X_j s are independent with $\mathbb{P}\{X_j = 0\} = \mathbb{P}\{X_j = 1\} = \frac{1}{2}$. Thus, in order to generate an infinite independent sequence of equiprobable zeroes and ones, choose a number uniformly distributed on $[0, 1]$ and expand it as a binary fraction.

4.3.3 Random numbers generated by coin tossing

We saw in Proposition 4.3.3 that choosing a random number in $[0, 1]$ leads to a model for tossing a fair coin infinitely often. Our aim is now to investigate the opposite question. That is, we are given an infinite random sequence of zeroes and ones and want to

construct a uniformly distributed number in $[0, 1]$. The precise mathematical question is as follows: suppose we are given an infinite sequence $(X_k)_{k \geq 1}$ of independent random variables with

$$\mathbb{P}\{X_k = 0\} = \mathbb{P}\{X_k = 1\} = 1/2, \quad k = 1, 2, \dots \quad (4.8)$$

Is it possible to construct from these X_k s a uniformly distributed U ? The next proposition answers this question to the affirmative.

Proposition 4.3.9. *Let X_1, X_2, \dots be an arbitrary sequence of independent random variables satisfying eq. (4.8). If U is defined by*

$$U(\omega) := \sum_{k=1}^{\infty} \frac{X_k(\omega)}{2^k}, \quad \omega \in \Omega,$$

then this random variable is uniformly distributed on $[0, 1]$.

Proof. In order to prove that U is uniformly distributed on $[0, 1]$, we have to show that, if $t \in [0, 1]$, then

$$\mathbb{P}\{U \leq t\} = t. \quad (4.9)$$

We start the proof of eq. (4.9) with the following observation: suppose the binary fraction of some $t \in [0, 1]$ is $0.t_1t_2\dots$ for certain $t_i \in \{0, 1\}$. If $s = 0.s_1s_2\dots$, then $s < t$ if and only if there is an $n \in \mathbb{N}$ so that the following is satisfied.⁴

$$s_1 = t_1, \dots, s_{n-1} = t_{n-1}, \quad s_n = 0, \quad \text{and} \quad t_n = 1.$$

Fix $t \in [0, 1]$ for a moment and set

$$A_n(t) := \{s \in [0, 1] : s_1 = t_1, \dots, s_{n-1} = t_{n-1}, s_n < t_n\}.$$

Of course, $A_n(t) \cap A_m(t) = \emptyset$ whenever $n \neq m$ and, moreover, $A_n(t) \neq \emptyset$ if and only if $t_n = 1$. Furthermore, by the previous remark

$$[0, t) = \bigcup_{n=1}^{\infty} A_n(t) = \bigcup_{\{n:t_n=1\}} A_n(t).$$

Finally, if $A_n(t) \neq \emptyset$, that is, if $t_n = 1$, then

$$\begin{aligned} \mathbb{P}\{U \in A_n(t)\} &= \mathbb{P}\{X_1 = t_1, \dots, X_{n-1} = t_{n-1}, X_n = 0\} \\ &= \mathbb{P}\{X_1 = t_1\} \cdots \mathbb{P}\{X_{n-1} = t_{n-1}\} \cdot \mathbb{P}\{X_n = 0\} = \frac{1}{2^n}. \end{aligned}$$

⁴ In the case $n = 1$, this says $s_1 = 0$ and $t_1 = 1$.

In the last step, we used both properties of the X_k s, that is, they are independent and satisfy $\mathbb{P}\{X_k = 0\} = \mathbb{P}\{X = 1\} = 1/2$.

Summing up, we get

$$\begin{aligned}\mathbb{P}\{U < t\} &= \mathbb{P}\left\{U \in \bigcup_{\{n:t_n=1\}} A_n(t)\right\} = \sum_{\{n:t_n=1\}} \mathbb{P}\{U \in A_n(t)\} \\ &= \sum_{\{n:t_n=1\}} \frac{1}{2^n} = \sum_{n=1}^{\infty} \frac{t_n}{2^n} = t.\end{aligned}$$

This “almost” proves eq. (4.9). It remains to show that $\mathbb{P}\{U < t\} = \mathbb{P}\{U \leq t\}$ or, equivalently, $\mathbb{P}\{U = t\} = 0$. To verify this, we use the continuity of \mathbb{P} from above. Then

$$\begin{aligned}\mathbb{P}\{U = t\} &= \mathbb{P}\{X_1 = t_1, X_2 = t_2, \dots\} \\ &= \lim_{n \rightarrow \infty} \mathbb{P}\{X_1 = t_1, \dots, X_n = t_n\} = \lim_{n \rightarrow \infty} \frac{1}{2^n} = 0.\end{aligned}$$

Consequently, eq. (4.9) holds for all $t \in [0, 1]$ and, as asserted, U is uniformly distributed on $[0, 1]$. \square

Remark 4.3.10. Another possibility to write U is as binary fraction

$$U(\omega) = 0.X_1(\omega)X_2(\omega)\dots, \quad \omega \in \Omega.$$

Consequently, in order to construct a random number u in $[0, 1]$ one may proceed as follows: toss a fair coin with “0” and “1” infinitely often and take the obtained sequence as binary fraction of u . The u obtained in this way is uniformly distributed on $[0, 1]$.

Of course, in practice one tosses a coin not infinitely often. One stops the procedure after N trials for some “large” N . In this way, one gets a number u , which is “almost” uniformly distributed on $[0, 1]$.

Example 4.3.11. Suppose tossing a fair coin led to the sequence 0, 1, 1, 0, 0, 1, 0, 1. Then the generated number $u \in [0, 1]$ equals $u = 0.39453125$. In the same way, the sequence 1, 1, 1, 0, 1, 1, 0, 1 when tossing leads to 0.92578125 as a randomly chosen number in $[0, 1]$.

Then how does one construct n independent numbers u_1, \dots, u_n , all uniformly distributed on $[0, 1]$? The answer is quite obvious. Take n coins and toss them. As functions of independent observations the generated u_1, \dots, u_n are independent as well and, by construction, each of these numbers is uniformly distributed on $[0, 1]$. Another way is to toss the same coin n times “infinitely often,” thus getting n infinite sequences of zeroes and ones.

Summary: If X_1, X_2, \dots is an arbitrary sequence of independent random variables such that $\mathbb{P}\{X_j = 0\} = \mathbb{P}\{X_j = 1\} = \frac{1}{2}$, then $U = \sum_{k=1}^{\infty} \frac{X_k}{2^k} = 0.X_1X_2\dots$ is uniformly distributed on $[0, 1]$. Thus, in order to obtain a number u uniformly distributed on $[0, 1]$, toss a fair coin “infinitely” often and define u as binary fraction according to the observed sequence of zeroes and ones.

4.4 Simulation of random variables

Proposition 4.3.9 provides us with a technique to simulate a uniformly distributed random variable U by tossing a fair coin. The aim of this section is to find a suitable function $f : [0, 1] \rightarrow \mathbb{R}$, so that the transformed random variable $X = f(U)$ possesses a given probability distribution.

Remark 4.4.1. Typical questions of this kind are as follows:

- Find a function f so that $X = f(U)$ is standard normally distributed.
- Does there exist a function $g : [0, 1] \rightarrow \mathbb{R}$ for which $g(U)$ is $B_{n,p}$ -distributed?

Suppose for a moment we already found such functions f and g . According to Remark 4.3.10, we construct independent numbers u_1, \dots, u_n , uniformly distributed on $[0, 1]$, and set $x_i = f(u_i)$ and $y_i = g(u_i)$. In this way, we get either n standard normally distributed numbers x_1, \dots, x_n or n binomial distributed numbers y_1, \dots, y_n . Moreover, by Proposition 4.1.9, these numbers are independent. In this way, we may simulate independent random numbers possessing a given probability distribution.

We start with simulating *discrete* random variables. Thus suppose we are given real numbers x_1, x_2, \dots and $p_k \geq 0$ with $\sum_{k=1}^{\infty} p_k = 1$, and we look for a random variable $X = f(U)$ such that

$$\mathbb{P}\{X = x_k\} = p_k, \quad k = 1, 2, \dots$$

One possible way to find such a function f is as follows: divide $[0, 1]$ into disjoint intervals I_1, I_2, \dots of length $|I_k| = p_k$ where $k = 1, 2, \dots$. Since $\sum_{k=1}^{\infty} p_k = 1$, such intervals exist. For example, take $I_1 = [0, p_1)$ and

$$I_k = \left[\sum_{i=1}^{k-1} p_i, \sum_{i=1}^k p_i \right), \quad k = 2, 3, \dots$$

With these intervals I_k , we define $f : [0, 1] \rightarrow \mathbb{R}$ by

$$f(x) := x_k \quad \text{if } x \in I_k, \quad (4.10)$$

or, equivalently,

$$f(x) = \sum_{k=1}^{\infty} x_k \mathbb{1}_{I_k}(x). \quad (4.11)$$

Then the following is true.

Proposition 4.4.2. *Let U be uniformly distributed on $[0, 1]$, and set $X = f(U)$ with f defined by eq. (4.10) or eq. (4.11). Then*

$$\mathbb{P}\{X = x_k\} = p_k, \quad k = 1, 2, \dots$$

Proof. Using that U is uniformly distributed on $[0, 1]$, this is an easy consequence of eq. (1.46) in view of

$$\mathbb{P}\{X = x_k\} = \mathbb{P}\{f(U) = x_k\} = \mathbb{P}\{U \in I_k\} = |I_k| = p_k. \quad \square$$

Remark 4.4.3. Note that the concrete shape of the intervals⁵ I_k is not important at all. They only have to satisfy $|I_k| = p_k$ for all $k = 1, 2, \dots$. Moreover, these intervals need not necessarily be disjoint; a “small” overlap does not influence the assertion. Indeed, it suffices that $\mathbb{P}\{U \in I_k \cap I_l\} = 0$ whenever $k \neq l$. For example, if always $|I_k \cap I_l| < \infty$, $k \neq l$, then the construction works as well. In particular, we may choose also $I_1 = [0, p_1]$ and $I_k = [\sum_{i=1}^{k-1} p_i, \sum_{i=1}^k p_i]$ if $k \geq 2$.

Example 4.4.4. We want to simulate a random variable X , which is uniformly distributed on $\{x_1, \dots, x_N\}$. How to proceed?

Answer: Divide the interval $[0, 1]$ into N intervals I_1, \dots, I_N of length $\frac{1}{N}$. For example, choose $I_k := [\frac{k-1}{N}, \frac{k}{N}]$, $k = 1, \dots, N$. If $f = \sum_{k=1}^N x_k \mathbb{1}_{I_k}$, then $X = f(U)$ is uniformly distributed on $\{x_1, \dots, x_N\}$.

Example 4.4.5. Suppose we want to simulate a number $k \in \mathbb{N}_0$, which is Pois_λ -distributed. Set

$$I_k := \left[\sum_{j=0}^{k-1} \frac{\lambda^j}{j!} e^{-\lambda}, \sum_{j=0}^k \frac{\lambda^j}{j!} e^{-\lambda} \right), \quad k = 0, 1, \dots,$$

where the left-hand sum is supposed to be zero if $k = 0$. Choose randomly a number $u \in [0, 1]$ and take the k with $u \in I_k$. Then k is the number we are interested in.

Example 4.4.6. Finally, let us simulate $B_{n,p}$ -distributed random numbers. One possible way is to divide $[0, 1]$ into $n + 1$ disjoint intervals as follows:

$$I_k = \left[\sum_{j=0}^{k-1} \binom{n}{j} p^j (1-p)^{n-j}, \sum_{j=0}^k \binom{n}{j} p^j (1-p)^{n-j} \right), \quad k = 0, \dots, n.$$

If $k = 0$, the left-hand sum is taken as zero, that is, $I_0 = [0, (1-p)^n)$. Next choose a random number u uniformly distributed on $[0, 1]$. For example, use the technique presented in Remark 4.3.10. If this $u \in I_k$, then $k \in \{0, \dots, n\}$ is the desired $B_{n,p}$ -distributed integer.

For instance, if $p = 1/2$ and $n = 4$, the five intervals are

$$I_0 = \left[0, \frac{1}{16} \right), \quad I_1 = \left[\frac{1}{16}, \frac{5}{16} \right), \quad I_2 = \left[\frac{5}{16}, \frac{11}{16} \right), \quad I_3 = \left[\frac{11}{16}, \frac{15}{16} \right), \quad I_4 = \left[\frac{15}{16}, 1 \right).$$

⁵ They do not even need to be intervals.

After fixing the intervals, let us sufficiently often toss a fair coin labeled with “0” and “1.” Say we observed the sequence 0, 1, 1, 0, 1, 0. Since the dyadic number 0.011010 equals $u = 0.40625$ in decimal representation, one notes that $u \in I_2$. Thus, the randomly chosen number is $k = 2$. It is distributed according to $B_{4,0.5}$. Which number $k \in \{0, 1, 2, 3, 4\}$ do we select if we toss 1, 1, 0, 1, 1, 1?

Our next aim is to simulate *continuous* random variables. More precisely, suppose we are given a probability density p . Then we look for a function $f : [0, 1] \rightarrow \mathbb{R}$ such that p is the density of $X = f(U)$, that is, we have to have

$$\mathbb{P}\{f(U) \leq t\} = \mathbb{P}\{X \leq t\} = \int_{-\infty}^t p(x) \, dx, \quad t \in \mathbb{R}. \quad (4.12)$$

To this end, set

$$F(t) = \int_{-\infty}^t p(x) \, dx, \quad t \in \mathbb{R}. \quad (4.13)$$

Thus, F is the distribution function of the random variable X , which we are going to construct.

Suppose first that F is one-to-one on a finite or infinite interval (a, b) , so that $F(x) = 0$ if $x < a$, and $F(x) = 1$ if $x > b$. For example, this is valid if $p(t) > 0$ for all $t \in (a, b)$. Since F is continuous, the inverse function F^{-1} exists and maps $(0, 1)$ to (a, b) .

Proposition 4.4.7. *Let p be a probability density and define F by eq. (4.12). Suppose F satisfies the above condition. If $X = F^{-1}(U)$, then*

$$\mathbb{P}\{X \leq t\} = \int_{-\infty}^t p(x) \, dx, \quad t \in \mathbb{R},$$

that is, p is a density of X .

Proof. First note that the assumptions about F imply that it is increasing on (a, b) . Hence, if $t \in \mathbb{R}$, then

$$\mathbb{P}\{X \leq t\} = \mathbb{P}\{F^{-1}(U) \leq t\} = \mathbb{P}\{U \leq F(t)\} = F(t) = \int_{-\infty}^t p(x) \, dx, \quad t \in \mathbb{R}.$$

Here we used $0 \leq F(t) \leq 1$ and $\mathbb{P}\{U \leq s\} = s$ whenever $0 \leq s \leq 1$. This completes the proof. \square

But what do we do if F does not satisfy the above assumption? For example, this happens if $p(x) = 0$ on an interval $I = (\alpha, \beta)$ and $p(x) > 0$ on some left- and right-hand intervals⁶ of I . In this case F^{-1} does not exist, and we have to modify the construction.⁷

Definition 4.4.8. Let F be defined by eq. (4.13). Then we set

$$F^{-}(s) = \inf\{t \in \mathbb{R} : F(t) = s\}, \quad 0 \leq s < 1.$$

The function F^{-} , mapping $[0, 1)$ to $[-\infty, \infty)$, is called the **pseudoinverse** of F .

Remark 4.4.9. If $0 < s < 1$, then $F^{-}(s) \in \mathbb{R}$ while $F^{-}(0) = -\infty$. Moreover, if F is increasing on some interval I , then $F^{-}(s) = F^{-1}(s)$ for $s \in I$.

Lemma 4.4.10. The pseudoinverse function F^{-} possesses the following properties:

1. If $s \in (0, 1)$ and $t \in \mathbb{R}$, then

$$F(F^{-}(s)) = s \quad \text{and} \quad F^{-}(F(t)) \leq t.$$

2. Given $t \in (0, 1)$, we have

$$F^{-}(s) \leq t \iff s \leq F(t). \quad (4.14)$$

Proof. The equality $F(F^{-}(s)) = s$ is a direct consequence of the continuity of F . Indeed, if there are $t_n \searrow F^{-}(s)$ with $F(t_n) = s$, then

$$s = \lim_{n \rightarrow \infty} F(t_n) = F(F^{-}(s)).$$

The second part of the first assertion follows by the definition of F^{-} .

Now let us come to the proof of property (4.14). If $F^{-}(s) \leq t$, then the monotonicity of F and $F(F^{-}(s)) = s$ lead to $s = F(F^{-}(s)) \leq F(t)$.

Conversely, if $s \leq F(t)$, then $F^{-}(s) \leq F^{-}(F(t)) \leq t$ by the first part, thus, property (4.14) is proved. \square

Now choose a uniformly distributed U and set $X = F^{-}(U)$. Since $\mathbb{P}\{U = 0\} = 0$, we may assume that X attains values in \mathbb{R} .

Proposition 4.4.11. Let p be a probability density, that is, we have $p(x) \geq 0$ and $\int_{-\infty}^{\infty} p(x) dx = 1$. Define F by eq. (4.13) and let F^{-} be its pseudoinverse. Take U to be uniformly distributed on $[0, 1]$ and set $X = F^{-}(U)$. Then p is a distribution density of the random variable X .

⁶ Take, for instance, p with $p(x) = \frac{1}{2}$ if $x \in [0, 1]$ and if $x \in [1, 2]$, and $p(x) = 0$ otherwise.

⁷ All subsequent distribution functions F possess an inverse function on a suitable interval (a, b) . Thus, Proposition 4.4.7 applies in almost all cases of interest. Therefore, if the statements about pseudoinverse functions look too complicated, you may skip them.

Proof. Using property (4.14), it follows that

$$\begin{aligned} F_X(t) &= \mathbb{P}\{X \leq t\} = \mathbb{P}\{\omega \in \Omega : F^-(U(\omega)) \leq t\} \\ &= \mathbb{P}\{\omega \in \Omega : U(\omega) \leq F(t)\} = F(t), \end{aligned}$$

which completes the proof. \square

Remark 4.4.12. Since $F^- = F^{-1}$ whenever the inverse function exists, Proposition 4.4.7 is a special case of Proposition 4.4.11.

Example 4.4.13. Let us simulate an $\mathcal{N}(0, 1)$ -distributed random variable, that is, we are looking for a function $f : (0, 1) \rightarrow \mathbb{R}$ such that for a uniformly distributed U ,

$$\mathbb{P}\{f(U) \leq t\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx, \quad t \in \mathbb{R}.$$

The distribution function

$$\Phi(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx$$

is one-to-one from $\mathbb{R} \rightarrow (0, 1)$, hence Proposition 4.4.7 applies, and $\Phi^{-1}(U)$ is a standard normal random variable.

How does one get an $\mathcal{N}(\mu, \sigma^2)$ -distributed random variable? If X is standard normal, by Proposition 4.2.3 the transformed variable $\sigma X + \mu$ is $\mathcal{N}(\mu, \sigma^2)$ -distributed. Consequently, $\sigma \Phi^{-1}(U) + \mu$ possesses the desired distribution.

How do we find n independent $\mathcal{N}(\mu, \sigma^2)$ -distributed numbers x_1, \dots, x_n ? To achieve this, choose u_1, \dots, u_n in $[0, 1]$ according to the construction presented in Remark 4.3.10 and set $x_i = \sigma \Phi^{-1}(u_i) + \mu$, $1 \leq i \leq n$.

Example 4.4.14. Our next aim is to simulate an E_λ -distributed (exponentially distributed) random variable. Here

$$F(t) = \begin{cases} 0 & \text{if } t \leq 0, \\ 1 - e^{-\lambda t} & \text{if } t > 0, \end{cases}$$

which satisfies the assumptions of Proposition 4.4.7 on the interval $(0, \infty)$. Its inverse F^{-1} maps $(0, 1)$ to $(0, \infty)$ and equals

$$F^{-1}(s) = -\frac{\ln(1-s)}{\lambda}, \quad 0 < s < 1.$$

Therefore, if U is uniformly distributed on $[0, 1]$, then $X = -\frac{\ln(1-U)}{\lambda}$ is E_λ -distributed. This is true for any uniformly distributed random variable U . By Corollary 4.2.7, the random variable $1 - U$ has the same distribution as U , hence, setting

$$Y = -\frac{\ln(1 - (1 - U))}{\lambda} = -\frac{\ln(U)}{\lambda},$$

the random variable Y is E_λ -distributed as well.

Example 4.4.15. Let us simulate a random variable with Cauchy distribution (see Definition 1.6.37). The distribution function F is given by

$$F(t) = \frac{1}{\pi} \int_{-\infty}^t \frac{1}{1+x^2} dx = \frac{1}{\pi} \arctan(t) + \frac{1}{2}, \quad t \in \mathbb{R},$$

hence $X := \tan(\pi U - \frac{\pi}{2})$ possesses a Cauchy distribution.

Example 4.4.16. Finally, let us give an example where Proposition 4.4.11 applies and Proposition 4.4.7 does not. Suppose we want to simulate a random variable X with distribution function F defined by

$$F(t) = \begin{cases} 0 & \text{if } t < 0, \\ \frac{t}{2} & \text{if } 0 \leq t < 1, \\ \frac{1}{2} & \text{if } 1 \leq t < 2, \\ \frac{1}{2} + \frac{t-2}{2} & \text{if } 2 \leq t < 3, \\ 1 & \text{if } t \geq 3. \end{cases} \quad (4.15)$$

Direct computations imply (see Figure 4.3) that

$$F^-(s) = \begin{cases} 2s & \text{if } 0 < s < \frac{1}{2}, \\ 1 & \text{if } s = \frac{1}{2}, \\ 2s + 1 & \text{if } \frac{1}{2} < s \leq 1. \end{cases}$$

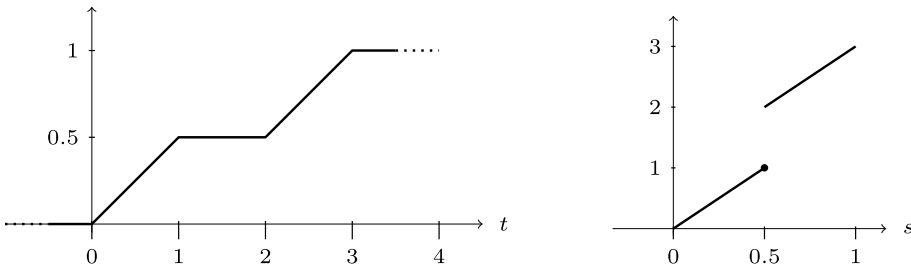


Figure 4.3: The functions F and F^- in Example 4.4.16.

Hence, if X is defined by

$$X = F^-(U) = 2U \mathbb{1}_{(0, \frac{1}{2})}(U) + (2U + 1) \mathbb{1}_{(\frac{1}{2}, 1)}(U),$$

then $\mathbb{P}\{X \leq t\} = F(t)$ with F defined by eq. (4.15). In other words, X is acting as follows. Choose by random a number $u \in [0, 1]$. If $u \leq \frac{1}{2}$, then $X(u) = 2u$, while for $u > \frac{1}{2}$ we take $X(u) = 2u + 1$.

Summary: Let \mathbb{P} be a probability measure on the Borel sets of \mathbb{R} . In order to simulate n independent numbers x_1, \dots, x_n distributed according to \mathbb{P} , one proceeds in the following way. In the first step, one constructs n independent numbers u_1, \dots, u_n uniformly distributed on $[0, 1]$. To this end, one tosses a fair coin sufficiently often as explained in Remark 4.3.10.

The second step is then as follows: If \mathbb{P} is discrete, one defines the numbers x_1, \dots, x_n by $x_j = f(u_j)$ with the function f constructed in (4.11). The probability that one (or each) x_j belongs to a set B equals $\mathbb{P}(B)$. In case of continuous \mathbb{P} , define x_j as $F^{-1}(u_j)$ provided the distribution function F of \mathbb{P} is invertible. Otherwise replace F^{-1} by the pseudoinverse F^- introduced in Definition 4.4.8. As before, the x_j s are independent and distributed according to the given \mathbb{P} .

4.5 Addition of random variables

Suppose we are given two random variables X and Y , both mapping from Ω into \mathbb{R} . As usual, their sum $X + Y$ is defined by

$$(X + Y)(\omega) := X(\omega) + Y(\omega), \quad \omega \in \Omega.$$

The main question we investigate in this section is as follows: suppose we know the probability distributions of X and Y . Is there a way to compute the distribution of $X + Y$? For example, if we roll a die twice, X is the result of the first roll, Y that of the second, then we know \mathbb{P}_X and \mathbb{P}_Y . But how do we get \mathbb{P}_{X+Y} ?

Before we treat this question, we have to be sure that $X + Y$ is also a random variable. This is not obvious at all. Otherwise, the probability distribution of $X + Y$ is not defined and our question does not make sense.

Proposition 4.5.1. *If X and Y are random variables, then so is $X + Y$.*

Proof. We start the proof with the following observation. For two real numbers a and b , one has $a < b$ if and only if there is a rational number $q \in \mathbb{Q}$ such that $a < q$ and $b > q$. Therefore, given $t \in \mathbb{R}$, it follows that

$$\begin{aligned} \{\omega \in \Omega : X(\omega) + Y(\omega) < t\} &= \{\omega \in \Omega : X(\omega) < t - Y(\omega)\} \\ &= \bigcup_{q \in \mathbb{Q}} [\{\omega : X(\omega) < q\} \cap \{\omega : q < t - Y(\omega)\}]. \end{aligned} \quad (4.16)$$

By assumption, X and Y are random variables. Hence, for each $q \in \mathbb{Q}$,

$$A_q := \{\omega : X(\omega) < q\} \in \mathcal{A} \quad \text{and} \quad B_q := \{\omega : Y(\omega) < t - q\} \in \mathcal{A},$$

which, by the properties of σ -fields, implies $C_q := A_q \cap B_q \in \mathcal{A}$. With this notation, we may write eq. (4.16) as

$$\{\omega \in \Omega : X(\omega) + Y(\omega) < t\} = \bigcup_{q \in \mathbb{Q}} C_q.$$

The σ -field \mathcal{A} is closed under taking countable unions, thus, since \mathbb{Q} is countably infinite and $C_q \in \mathcal{A}$, it follows that $\bigcup_{q \in \mathbb{Q}} C_q \in \mathcal{A}$. Therefore, we have proven that, if $t \in \mathbb{R}$, then

$$\{\omega \in \Omega : X(\omega) + Y(\omega) < t\} \in \mathcal{A}.$$

Proposition 3.1.6 lets us conclude that, as asserted, $X + Y$ is a random variable. \square

Remark 4.5.2. In view of Proposition 4.5.1, the following *question* makes sense: does there exist a general approach to evaluate \mathbb{P}_{X+Y} by virtue of \mathbb{P}_X and of \mathbb{P}_Y ?

Answer: Such a general way does not exist. The deeper reason behind this is that, in order to get \mathbb{P}_{X+Y} , one has to know the joint distribution of (X, Y) . And as we saw in Section 3.5, in general, the knowledge of \mathbb{P}_X and \mathbb{P}_Y does not suffice to determine their joint distribution, hence generally we also do not know \mathbb{P}_{X+Y} .

The next example emphasizes the previous remark.

Example 4.5.3. Let X, Y, X' , and Y' be as in Example 3.5.8, that is, we choose two balls out of an urn where two balls are labeled by “0” and two by “1,” once without replacing the first ball and once with replacing. The joint distributions of (X, Y) and (X', Y') are

$$\begin{aligned} \mathbb{P}\{X = 0, Y = 0\} &= \frac{1}{6}, & \mathbb{P}\{X = 0, Y = 1\} &= \frac{1}{3}, \\ \mathbb{P}\{X = 1, Y = 0\} &= \frac{1}{3}, & \mathbb{P}\{X = 1, Y = 1\} &= \frac{1}{6}, \\ \mathbb{P}\{X' = 0, Y' = 0\} &= \frac{1}{4}, & \mathbb{P}\{X' = 0, Y' = 1\} &= \frac{1}{4}, \\ \mathbb{P}\{X' = 1, Y' = 0\} &= \frac{1}{4}, & \mathbb{P}\{X' = 1, Y' = 1\} &= \frac{1}{4}. \end{aligned}$$

Then $\mathbb{P}_X = \mathbb{P}_{X'}$ and $\mathbb{P}_Y = \mathbb{P}_{Y'}$, but

$$\mathbb{P}\{X + Y = 0\} = \frac{1}{6}, \quad \mathbb{P}\{X + Y = 1\} = \frac{2}{3}, \quad \text{and} \quad \mathbb{P}\{X + Y = 2\} = \frac{1}{6},$$

while

$$\mathbb{P}\{X' + Y' = 0\} = \frac{1}{4}, \quad \mathbb{P}\{X' + Y' = 1\} = \frac{1}{2}, \quad \text{and} \quad \mathbb{P}\{X' + Y' = 2\} = \frac{1}{4}.$$

Consequently, if we do not replace the chosen ball, the probability that the sum of both choices equals “1” is $2/3$ while it is $1/2$ if we replace the first ball. And this happens although in both cases the numbers “0” and “1” occur every time with probability $1/2$.

On the other hand, as we saw in Proposition 3.6.5, the joint distribution is uniquely determined by the marginal ones, provided the random variables are independent. Therefore, for *independent* random variables X and Y , the distribution of $X + Y$ is determined by those of X and Y . The question remains, how \mathbb{P}_{X+Y} can be computed.

4.5.1 Sums of discrete random variables

We first consider an important special case, namely that X and Y attain values in \mathbb{Z} . Here we have

Proposition 4.5.4 (Convolution formula for \mathbb{Z} -valued random variables). *Let X and Y be two independent random variables with values in \mathbb{Z} . If $k \in \mathbb{Z}$, then*



$$\mathbb{P}\{X + Y = k\} = \sum_{i=-\infty}^{\infty} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = k - i\}.$$

Proof. Fix $k \in \mathbb{Z}$ and define $B_k \subseteq \mathbb{Z} \times \mathbb{Z}$ by

$$B_k := \{(i, j) \in \mathbb{Z} \times \mathbb{Z} : i + j = k\}.$$

Then we get

$$\mathbb{P}\{X + Y = k\} = \mathbb{P}\{(X, Y) \in B_k\} = \mathbb{P}_{(X, Y)}(B_k) \quad (4.17)$$

with joint distribution $\mathbb{P}_{(X, Y)}$. Proposition 3.6.11 asserts that for independent X and Y and $B \subseteq \mathbb{Z} \times \mathbb{Z}$,

$$\mathbb{P}_{(X, Y)}(B) = \sum_{(i, j) \in B} \mathbb{P}_X(\{i\}) \cdot \mathbb{P}_Y(\{j\}) = \sum_{(i, j) \in B} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = j\}.$$

We apply this formula with $B = B_k$ and, from eq. (4.17), obtain

$$\begin{aligned} \mathbb{P}\{X + Y = k\} &= \sum_{(i, j) \in B_k} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = j\} \\ &= \sum_{\{(i, j) : i + j = k\}} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = j\} \\ &= \sum_{i=-\infty}^{\infty} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = k - i\}, \end{aligned}$$

as asserted. □

Example 4.5.5. Two independent random variables X and Y are distributed according to $\mathbb{P}\{X = j\} = \mathbb{P}\{Y = j\} = 1/2^j, j = 1, 2, \dots$. Determine the probability distribution of $X - Y$.

Solution: First note that $\mathbb{P}\{X = j\} = \mathbb{P}\{Y = j\} = 0$ for $j \leq 0$. Hence, given $k \in \mathbb{Z}$, an application of Proposition 4.5.4 to X and $-Y$ yields

$$\mathbb{P}\{X - Y = k\} = \sum_{i=-\infty}^{\infty} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{-Y = k - i\} = \sum_{i=1}^{\infty} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = i - k\}.$$

If $k \geq 0$, then $\mathbb{P}\{Y = i - k\} = 0$ for $i \leq k$, thus

$$\begin{aligned} \mathbb{P}\{X - Y = k\} &= \sum_{i=k+1}^{\infty} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = i - k\} = \sum_{i=k+1}^{\infty} \frac{1}{2^i} \cdot \frac{1}{2^{i-k}} \\ &= 2^k \sum_{i=k+1}^{\infty} \frac{1}{2^{2i}} = 2^k \cdot 2^{-2k-2} \cdot \sum_{i=0}^{\infty} \frac{1}{2^{2i}} = 2^{-k-2} \cdot \frac{4}{3} = \frac{2^{-k}}{3}. \end{aligned}$$

For $k < 0$, it follows that

$$\mathbb{P}\{X - Y = k\} = \sum_{i=1}^{\infty} \frac{1}{2^i} \cdot \frac{1}{2^{i-k}} = 2^k \sum_{i=1}^{\infty} \frac{1}{2^{2i}} = 2^k \sum_{i=1}^{\infty} \frac{1}{4^i} = \frac{2^k}{3}.$$

We combine both cases and obtain

$$\mathbb{P}\{X - Y = k\} = \frac{2^{-|k|}}{3}, \quad k \in \mathbb{Z}.$$

Which random experiment does $X - Y$ describe? Suppose player A and B both toss a fair coin. Let X be the number of necessary trials for A to observe the first “heads.” Similarly, Y describes how often B has to toss his coin to get the first “heads.” Thus, the value of $X - Y$ tells us how many trials later (or earlier if $X - Y$ is negative) player A got his first “heads” compared to B .

For example, if B got his first “heads” one trial earlier than A , then $X - Y = 1$. The probability that this occurs equals $1/6$.

One special case of Proposition 4.5.4 is of particular interest.

Proposition 4.5.6 (Convolution formula for \mathbb{N}_0 -valued random variables). *Let X and Y be two independent random variables with values in \mathbb{N}_0 . If $k \in \mathbb{N}_0$, then it follows that*

$$\mathbb{P}\{X + Y = k\} = \sum_{i=0}^k \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = k - i\}.$$



Proof. Regard X and Y as \mathbb{Z} -valued random variables with $\mathbb{P}\{X = i\} = \mathbb{P}\{Y = i\} = 0$ for all $i = -1, -2, \dots$ If $k \in \mathbb{N}_0$, then Proposition 4.5.4 lets us conclude that

$$\mathbb{P}\{X + Y = k\} = \sum_{i=-\infty}^{\infty} \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = k - i\} = \sum_{i=0}^k \mathbb{P}\{X = i\} \cdot \mathbb{P}\{Y = k - i\}.$$

Here we used $\mathbb{P}\{X = i\} = 0$ for $i < 0$ and $\mathbb{P}\{Y = k - i\} = 0$ if $i > k$. For $k < 0$, it follows that $\mathbb{P}\{X + Y = k\} = 0$ because in this case $\mathbb{P}\{Y = k - i\} = 0$ for all $i \geq 0$. This completes the proof. \square

Example 4.5.7. Let X and Y be two independent random variables, both uniformly distributed on $\{1, 2, \dots, N\}$. Which probability distribution does $X + Y$ possess?

Answer: Of course, $X + Y$ attains only values in the set $\{2, 3, \dots, 2N\}$. Hence, $\mathbb{P}\{X + Y = k\}$ is only of interest for $2 \leq k \leq 2N$. Here we get

$$\mathbb{P}\{X + Y = k\} = \frac{|I_k|}{N^2}, \quad (4.18)$$

where I_k is defined by

$$I_k := \{i \in \{1, \dots, N\} : 1 \leq k - i \leq N\} = \{i \in \{1, \dots, N\} : k - N \leq i \leq k - 1\}.$$

To verify eq. (4.18), use that for $i \notin I_k$ either $\mathbb{P}\{X = i\} = 0$ or $\mathbb{P}\{Y = k - i\} = 0$. It is not difficult to prove that

$$|I_k| = \begin{cases} k - 1 & \text{if } 2 \leq k \leq N + 1, \\ 2N - k + 1 & \text{if } N + 1 < k \leq 2N, \end{cases}$$

which leads to

$$\mathbb{P}\{X + Y = k\} = \begin{cases} \frac{k-1}{N^2} & \text{if } 2 \leq k \leq N + 1, \\ \frac{2N-k+1}{N^2} & \text{if } N + 1 < k \leq 2N, \\ 0 & \text{otherwise} \end{cases}$$

If $N = 6$, then $X + Y$ may be viewed as the sum of two rolls of a die. Here the above formula leads to the values of $\mathbb{P}\{X + Y = k\}$, $k = 2, \dots, 12$, which we, by a direct approach, already computed in Example 3.2.15. For another example with $N = 8$ see Figure 4.4.

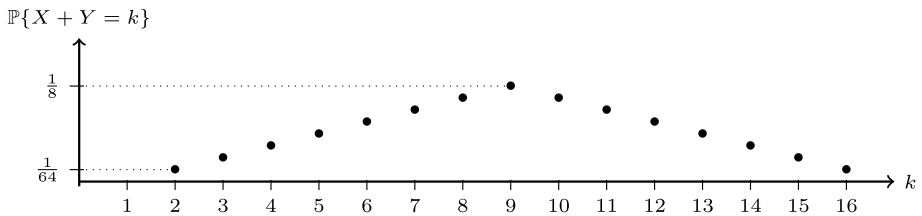


Figure 4.4: The sum of independent X and Y , both uniformly distributed on $\{1, \dots, 8\}$.

Finally, let us shortly discuss the case of two arbitrary independent discrete random variables. Assume that X and Y have values in at most countable infinite sets D and E ,

respectively. Then $X + Y$ maps into

$$D + E := \{x + y : x \in D, y \in E\}.$$

Note that $D + E$ is also at most countably infinite.

Under these assumptions, the following is valid.

Proposition 4.5.8. *Suppose X and Y are two independent discrete random variables with values in the (at most) countably infinite sets D and E , respectively. For $z \in D + E$, it follows that*

$$\mathbb{P}\{X + Y = z\} = \sum_{\{(x,y) \in D \times E : x+y=z\}} \mathbb{P}\{X = x\} \cdot \mathbb{P}\{Y = y\}.$$

Proof. For a fixed $z \in D + E$ define $B_z \subseteq D \times E$ by $B_z := \{(x, y) : x + y = z\}$. Using this notation, we get

$$\mathbb{P}\{X + Y = z\} = \mathbb{P}\{(X, Y) \in B_z\} = \mathbb{P}_{(X,Y)}(B_z),$$

where again $\mathbb{P}_{(X,Y)}$ denotes the joint distribution of X and Y . Now we may proceed as in the proof of Proposition 4.5.4. The independence of X and Y implies

$$\mathbb{P}_{(X,Y)}(B_z) = \sum_{\{(x,y) \in D \times E : x+y=z\}} \mathbb{P}\{X = x\} \cdot \mathbb{P}\{Y = y\},$$

proving the proposition. □

Remark 4.5.9. If $D = E = \mathbb{Z}$, then Proposition 4.5.8 implies Proposition 4.5.4, while for $D = E = \mathbb{N}_0$ we rediscover Proposition 4.5.6.

Example 4.5.10. Suppose X is uniformly distributed on $D = \{1, 2, 3, 4\}$ while Y is uniformly distributed on $E = \{5, 6, 7, 8\}$. Their sum attains its values in $\{6, \dots, 12\}$. If X and Y are independent, then, for example,

$$\begin{aligned} \mathbb{P}\{X + Y = 7\} &= \sum_{\substack{x \in D, y \in E \\ x+y=7}} \mathbb{P}\{X = x\} \mathbb{P}\{Y = y\} \\ &= \mathbb{P}\{X = 1\} \cdot \mathbb{P}\{Y = 6\} + \mathbb{P}\{X = 2\} \cdot \mathbb{P}\{Y = 5\} \\ &= \frac{1}{4} \cdot \frac{1}{4} + \frac{1}{4} \cdot \frac{1}{4} = \frac{1}{8}. \end{aligned}$$

4.5.2 Sums of continuous random variables

In this section we investigate the following question: let X and Y be two continuous random variables with density functions p and q . Is $X + Y$ continuous as well, and if this is so, how do we compute its density?

To answer this question, we need a special type of composing two functions.

Definition 4.5.11. Let f and g be two Riemann integrable functions from \mathbb{R} to \mathbb{R} . Their **convolution** $f * g$ is defined by

$$(f * g)(x) := \int_{-\infty}^{\infty} f(x-y)g(y) dy, \quad x \in \mathbb{R}. \quad (4.19)$$

Remark 4.5.12. The convolution is a commutative operation, that is,

$$f * g = g * f.$$

This follows by the change of variables $u = x - y$ in eq. (4.19), thus

$$(f * g)(x) = \int_{-\infty}^{\infty} f(x-y)g(y) dy = \int_{-\infty}^{\infty} f(u)g(x-u) du = (g * f)(x), \quad x \in \mathbb{R}.$$

Remark 4.5.13. For general functions f and g , the integral in eq. (4.19) does not always exist for all $x \in \mathbb{R}$. The investigation of this question requires facts and notations⁸ from Measure Theory; therefore, we will not treat it here. We only state a special case, which suffices for our later purposes. Moreover, for concrete functions f and g , it is mostly easy to check for which $x \in \mathbb{R}$ the value $(f * g)(x)$ exists.

Proposition 4.5.14. Let p and q be two probability densities and suppose that at least one of them is bounded. Then $(p * q)(x)$ exists for all $x \in \mathbb{R}$.

Proof. Say p is bounded, that is, there is a constant $c \geq 0$ such that $0 \leq p(z) \leq c$ for all $z \in \mathbb{R}$. Since $q(y) \geq 0$, if $x \in \mathbb{R}$, then

$$0 \leq \int_{-\infty}^{\infty} p(x-y)q(y) dy \leq c \int_{-\infty}^{\infty} q(y) dy = c < \infty.$$

This proves that $(p * q)(x)$ exists for all $x \in \mathbb{R}$.

Since $p * q = q * p$, the same argument applies if q is bounded. \square

The next result provides us with a formula for the evaluation of the density function of $X + Y$ for independent continuous X and Y .

Proposition 4.5.15 (Convolution formula for continuous random variables). Let X and Y be two independent random variables with distribution densities p and q . Then $X + Y$ is continuous as well, and its density r may be computed by

⁸ For example, “exists almost everywhere.”

$$r(x) = (p \star q)(x) = \int_{-\infty}^{\infty} p(y) q(x-y) dy.$$



Proof. We have to show that $r = p \star q$ satisfies

$$\mathbb{P}\{X + Y \leq t\} = \int_{-\infty}^t r(x) dx, \quad t \in \mathbb{R}. \quad (4.20)$$

Fix $t \in \mathbb{R}$ for a moment and define $B_t \subseteq \mathbb{R}^2$ by

$$B_t := \{(u, y) \in \mathbb{R}^2 : u + y \leq t\}.$$

Then we get

$$\mathbb{P}\{X + Y \leq t\} = \mathbb{P}\{(X, Y) \in B_t\} = \mathbb{P}_{(X, Y)}(B_t). \quad (4.21)$$

To compute the right-hand side of eq. (4.21), we use Proposition 3.6.20. It asserts that the joint distribution $\mathbb{P}_{(X, Y)}$ of independent X and Y has density $(u, y) \mapsto p(u)q(y)$, that is, if $B \subseteq \mathbb{R}^2$, then

$$\mathbb{P}_{(X, Y)}(B) = \iint_B p(u)q(y) dy du.$$

Choosing $B = B_t$ in the last formula, eq. (4.21) may now be written as

$$\mathbb{P}\{X + Y \leq t\} = \iint_{B_t} p(u) q(y) dy du = \int_{-\infty}^{\infty} \left[\int_{-\infty}^{t-y} p(u) du \right] q(y) dy. \quad (4.22)$$

Next we change the variables in the inner integral as follows:⁹ $u = x - y$, hence $du = dx$. Then the right-hand integrals in eq. (4.22) coincide with

$$\begin{aligned} \int_{-\infty}^{\infty} \left[\int_{-\infty}^t p(x-y) dx \right] q(y) dy &= \int_{-\infty}^t \left[\int_{-\infty}^{\infty} p(x-y) q(y) dy \right] dx \\ &= \int_{-\infty}^t (p \star q)(x) dx. \end{aligned}$$

Hereby we used that p and q are nonnegative, so that we may interchange the integrals by virtue of Proposition A.5.5. Thus, eq. (4.20) is satisfied, which completes the proof. \square

⁹ Note that in the inner integral y is a constant.

4.6 Sums of certain random variables

Let us start with the investigation of the sum of independent *binomial distributed* random variables. Here the following is valid.

Proposition 4.6.1. *Let X and Y be two independent random variables, accordingly $B_{n,p}$ - and $B_{m,p}$ -distributed for some $n, m \geq 1$, and some $p \in [0, 1]$. Then $X + Y$ is $B_{n+m,p}$ -distributed.*

Proof. By Proposition 4.5.6, we get that for $0 \leq k \leq m + n$,

$$\begin{aligned} \mathbb{P}\{X + Y = k\} &= \sum_{j=0}^k \left[\binom{n}{j} p^j (1-p)^{n-j} \right] \cdot \left[\binom{m}{k-j} p^{k-j} (1-p)^{m-(k-j)} \right] \\ &= p^k (1-p)^{n+m-k} \sum_{j=0}^k \binom{n}{j} \binom{m}{k-j}. \end{aligned}$$

To evaluate the sum, we apply Vandermonde's identity (Proposition A.3.9), which asserts

$$\sum_{j=0}^k \binom{n}{j} \binom{m}{k-j} = \binom{n+m}{k}.$$

This leads to

$$\mathbb{P}\{X + Y = k\} = \binom{n+m}{k} p^k (1-p)^{m+n-k},$$

and $X + Y$ is $B_{n+m,p}$ -distributed. \square

Interpretation: In the first experiment, we toss a biased coin n times and in the second m times. We combine these two experiments into one and toss the coin now $n + m$ times. Then we observe “heads” exactly k times during the $n + m$ trials if there is some $j \leq k$ so that we had j “heads” among the first n trials and $k - j$ among the second m . Finally, we have to sum the probabilities of all these events over $j \leq k$.

Corollary 4.6.2. *Let X_1, \dots, X_n be independent $B_{1,p}$ -distributed, that is,*

$$\mathbb{P}\{X_j = 0\} = 1 - p \quad \text{and} \quad \mathbb{P}\{X_j = 1\} = p, \quad j = 1, \dots, n.$$

Then their sum $X_1 + \dots + X_n$ is $B_{n,p}$ -distributed.

Proof. Apply Proposition 4.6.1 successively, first to X_1 and X_2 , then to $X_1 + X_2$ and X_3 , and so on. \square

Remark 4.6.3. Observe that

$$X_1 + \dots + X_n = |\{j \leq n : X_j = 1\}|.$$

Corollary 4.6.2 justifies the interpretation of the binomial distribution given in Section 1.4.3. Indeed, the event $\{X_j = 1\}$ occurs if in trial j we observe success. Thus, the sum $X_1 + \cdots + X_n$ describes the number of successes in n independent trials. Hereby, the success probability is $\mathbb{P}\{X_j = 1\} = p$.

In the literature, the following notion is common.

Definition 4.6.4. A finite or infinite sequence X_1, X_2, \dots of independent $B_{1,p}$ -distributed random variables is called a **Bernoulli trial** or **Bernoulli process**, or also **binomial trial**, with success probability $p \in [0, 1]$.

With these notations, Corollary 4.6.2 may now be formulated as follows:

Let X_1, X_2, \dots , be a Bernoulli trial with success probability p . Then for $n \geq 1$,

$$\mathbb{P}\{X_1 + \cdots + X_n = k\} = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, \dots, n.$$

Let X and Y be two independent *Poisson distributed* random variables. Which distribution does $X + Y$ possess? The next result answers this question.

Proposition 4.6.5. *Let X and Y be independent Pois_λ - and Pois_μ -distributed for some $\lambda > 0$ and $\mu > 0$, respectively. Then $X + Y$ is $\text{Pois}_{\lambda+\mu}$ -distributed.*

Proof. Proposition 4.5.6 and the binomial theorem (see Proposition A.3.8) imply

$$\begin{aligned} \mathbb{P}\{X + Y = k\} &= \sum_{j=0}^k \left[\frac{\lambda^j}{j!} e^{-\lambda} \right] \left[\frac{\mu^{k-j}}{(k-j)!} e^{-\mu} \right] \\ &= \frac{e^{-(\lambda+\mu)}}{k!} \sum_{j=0}^k \frac{k!}{j! (k-j)!} \lambda^j \mu^{k-j} \\ &= \frac{e^{-(\lambda+\mu)}}{k!} \sum_{j=0}^k \binom{k}{j} \lambda^j \mu^{k-j} = \frac{(\lambda + \mu)^k}{k!} e^{-(\lambda+\mu)}. \end{aligned}$$

Consequently, as asserted, $X + Y$ is $\text{Pois}_{\lambda+\mu}$ -distributed. □

Interpretation: The numbers of phone calls arriving per day at some call centers A and B are Poisson distributed with parameters¹⁰ λ and μ . Suppose that these two centers have different customers, that is, we assume that the number of calls in A and B is independent of each other. Proposition 4.6.5 asserts that the number of calls arriving per day either in A or in B is again Poisson distributed, but now with parameter $\lambda + \mu$.

Example 4.6.6. This example deals with the distribution of raisins in a set of dough. More precisely, suppose we have N pounds of dough and therein are n uniformly dis-

¹⁰ Later on, in Proposition 5.1.16, we will see that λ and μ are the mean values of arriving calls per day.

tributed raisins. Choose at random a one-pound piece of dough. Find the probability that there are $k \geq 0$ raisins in the chosen piece.

Approach 1: Since the raisins are uniformly distributed in the dough, the probability that a single raisin is in the chosen piece equals $1/N$. Hence, if X is the number of raisins in that piece, it is $B_{n,p}$ -distributed with $p = 1/N$. Assuming that N is big, the random variable X is approximately Pois_λ -distributed with $\lambda = n/N$, that is,

$$\mathbb{P}\{X = k\} = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, \dots$$

Note that $\lambda = n/N$ coincides with the average number of raisins per pound dough.

Approach 2: Assume that we took in the previous model $N \rightarrow \infty$, that is, we have an “infinite” amount of dough and “infinitely” many raisins. Which distribution does X , the number of raisins in a one-pound piece, now possess?

First, we have to determine what it means that the amount of dough is “infinite” and that the raisins are uniformly distributed¹¹ therein. This is expressed by the following conditions:

- (a) The mass of dough is unbelievably huge, hence, whenever we choose two different pieces, the numbers of raisins in the pieces are independent of each other.
- (b) The fact that the raisins are uniformly distributed is expressed by the following condition: suppose the number of raisins in a one-pound piece is $n \geq 0$. If this piece is split into two pieces, say K_1 and K_2 of weight α and $1 - \alpha$ pounds, then the probability that a single raisin (out of n) is in K_1 equals α , and the probability that it is in K_2 is $1 - \alpha$.

Fix $0 < \alpha < 1$ and choose in the first step a piece K_1 of α pounds and in the second one another piece K_2 of weight $1 - \alpha$. Let X_1 and X_2 be the numbers of raisins in each of the two pieces. By condition (a), the random variables X_1 and X_2 are independent. If $X = X_1 + X_2$, then X is the number of raisins in a randomly chosen one-pound piece. Suppose now $X = n$, that is, there are n raisins in the one-pound piece. Then by condition (b), the probability for k raisins in K_1 is described by the binomial distribution $B_{n,\alpha}$. Recall that the success probability for a single raisin is α , thus, $X_1 = k$ means, we have k successes. This may be formulated as follows: for all $0 \leq k \leq n$,

$$\mathbb{P}\{X_1 = k | X = n\} = B_{n,\alpha}(\{k\}) = \binom{n}{k} \alpha^k (1 - \alpha)^{n-k}. \quad (4.23)$$

Rewriting eq. (4.23) leads to

$$\mathbb{P}\{X_1 = k, X_2 = n - k\} = \mathbb{P}\{X_1 = k, X = n\}$$

¹¹ Note that the multivariate uniform distribution only makes sense (cf. Definition 1.8.13) if the underlying set has a finite volume.

$$= \mathbb{P}\{X_1 = k | X = n\} \cdot \mathbb{P}\{X = n\} = \mathbb{P}\{X = n\} \cdot \binom{n}{k} \alpha^k (1 - \alpha)^{n-k}. \quad (4.24)$$

Observe that in contrast to eq. (4.23), eq. (4.24) remains valid also if $\mathbb{P}\{X = n\} = 0$. Indeed, if $\mathbb{P}\{X = n\} = 0$, by Proposition 4.5.6, the event $\{X_1 = k, X_2 = n - k\}$ has zero probability as well.

The independence of X_1 and X_2 and eq. (4.24) imply that, if $n = 0, 1, 2, \dots$ and $k = 0, \dots, n$, then

$$\mathbb{P}\{X_1 = k\} \cdot \mathbb{P}\{X_2 = n - k\} = \mathbb{P}\{X = n\} \cdot \binom{n}{k} \alpha^k (1 - \alpha)^{n-k}.$$

Setting $k = n$, we get

$$\mathbb{P}\{X_1 = n\} \cdot \mathbb{P}\{X_2 = 0\} = \mathbb{P}\{X = n\} \cdot \alpha^n, \quad (4.25)$$

while for $n \geq 1$ and $k = n - 1$ we obtain

$$\mathbb{P}\{X_1 = n - 1\} \cdot \mathbb{P}\{X_2 = 1\} = \mathbb{P}\{X = n\} \cdot n \cdot \alpha^{n-1} (1 - \alpha). \quad (4.26)$$

In particular, from eq. (4.25) $\mathbb{P}\{X_2 = 0\} > 0$ follows. If this probability were zero, then this would imply $\mathbb{P}\{X = n\} = 0$ for all $n \in \mathbb{N}_0$, which is impossible in view of $\mathbb{P}\{X \in \mathbb{N}_0\} = 1$.

In the next step, we solve eqs. (4.25) and (4.26) with respect to $\mathbb{P}\{X = n\}$ and make them equal. Doing so, for $n \geq 1$ we get

$$\begin{aligned} \mathbb{P}\{X_1 = n\} &= \frac{\alpha}{n} (1 - \alpha)^{-1} \cdot \frac{\mathbb{P}\{X_2 = 1\}}{\mathbb{P}\{X_2 = 0\}} \cdot \mathbb{P}\{X_1 = n - 1\} \\ &= \frac{\alpha \lambda}{n} \cdot \mathbb{P}\{X_1 = n - 1\}, \end{aligned} \quad (4.27)$$

where $\lambda \geq 0$ is defined by

$$\lambda := (1 - \alpha)^{-1} \cdot \frac{\mathbb{P}\{X_2 = 1\}}{\mathbb{P}\{X_2 = 0\}}. \quad (4.28)$$

Do we have $\lambda > 0$? If $\lambda = 0$, then $\mathbb{P}\{X_2 = 1\} = 0$ and, by eq. (4.26), $\mathbb{P}\{X = n\} = 0$ for $n \geq 1$. Consequently, $\mathbb{P}\{X = 0\} = 1$, which says that there are no raisins in the dough. We exclude this trivial case, thus it follows that $\lambda > 0$.

Finally, successive application of eq. (4.27) implies for $n \in \mathbb{N}_0$ ¹² that

$$\mathbb{P}\{X_1 = n\} = \frac{(\alpha \lambda)^n}{n!} \cdot \mathbb{P}\{X_1 = 0\}, \quad (4.29)$$

¹² If $n = 0$, the equation holds trivially.

leading to

$$1 = \sum_{n=0}^{\infty} \mathbb{P}\{X_1 = n\} = \mathbb{P}\{X_1 = 0\} \cdot \sum_{n=0}^{\infty} \frac{(a\lambda)^n}{n!} = \mathbb{P}\{X_1 = 0\} e^{a\lambda},$$

that is, we have $\mathbb{P}\{X_1 = 0\} = e^{-a\lambda}$. Plugging this into eq. (4.29) gives

$$\mathbb{P}\{X_1 = n\} = \frac{(a\lambda)^n}{n!} e^{-a\lambda},$$

and so X_1 is Poisson distributed with parameter $a\lambda$.

Let us interchange now the roles of X_1 and X_2 , hence also of a and $1 - a$. An application of the first step to X_2 tells us that it is Poisson distributed, but now with parameter $(1 - a)\lambda'$, where in view of eq. (4.28) λ' is given by¹³

$$\lambda' = a^{-1} \cdot \frac{\mathbb{P}\{X_1 = 1\}}{\mathbb{P}\{X_1 = 0\}} = a^{-1} \frac{a\lambda e^{-a\lambda}}{e^{-a\lambda}} = \lambda.$$

Thus, X_2 is $\text{Pois}_{(1-a)\lambda}$ -distributed.

Since X_1 and X_2 are independent, Proposition 4.6.5 applies, hence $X = X_1 + X_2$ is Pois_{λ} -distributed or, equivalently,

$$\mathbb{P}\{\text{There are } k \text{ raisins in a one-pound piece}\} = \frac{\lambda^k}{k!} e^{-\lambda}.$$

Remark 4.6.7. Which role does the parameter $\lambda > 0$ play in this model? As already mentioned, Proposition 5.1.16 will tell us that λ is the average number of raisins per pound dough. Thus, if $\rho > 0$ and we ask for the number of raisins in a piece of ρ pounds, then this number is $\text{Pois}_{\rho\lambda}$ -distributed,¹⁴ that is,

$$\mathbb{P}\{k \text{ raisins in } \rho \text{ pounds of dough}\} = \frac{(\rho\lambda)^k}{k!} e^{-\rho\lambda}.$$

Assume that dough contains on average 20 raisins per pound. Let X be the number of raisins in a piece of bread baked from five pounds of dough. Then X is Pois_{100} -distributed and

$$\begin{aligned} \mathbb{P}(\{95 \leq X \leq 105\}) &= 0.4176, & \mathbb{P}(\{90 \leq X \leq 110\}) &= 0.7065, \\ \mathbb{P}(\{85 \leq X \leq 115\}) &= 0.8793, & \mathbb{P}(\{80 \leq X \leq 120\}) &= 0.9599, \\ \mathbb{P}(\{75 \leq X \leq 125\}) &= 0.9892, & \mathbb{P}(\{70 \leq X \leq 130\}) &= 0.9976. \end{aligned}$$

¹³ Observe that we have to replace X_2 by X_1 and $1 - a$ by a .

¹⁴ Because on average there are $\rho\lambda$ raisins in a piece of ρ pounds.

Additional question: Suppose we buy two loaves of bread baked from ρ pounds dough each. What is the probability that one of these two loaves contains more than twice as many raisins as the other?

Answer: Let X be the number of raisins in the first loaf, and Y the number of raisins in the second. By assumption, X and Y are independent, and both are $\text{Pois}_{\rho\lambda}$ -distributed, where as before $\lambda > 0$ is the average number of raisins per pound dough. The probability we are interested in is (use Proposition 1.2.4 as well as that $(X, Y) \stackrel{d}{=} (Y, X)$)

$$\begin{aligned} \mathbb{P}\{X > 2Y \text{ or } Y > 2X\} &= \mathbb{P}\{X > 2Y\} + \mathbb{P}\{Y > 2X\} = 2 \mathbb{P}\{X > 2Y\} \\ &= 2 \sum_{k=0}^{\infty} \mathbb{P}\{Y = k, X > 2k\} = 2 \sum_{k=0}^{\infty} \mathbb{P}\{Y = k\} \cdot \mathbb{P}\{X > 2k\} \\ &= 2 \sum_{k=0}^{\infty} \mathbb{P}\{Y = k\} \cdot \sum_{j=2k+1}^{\infty} \mathbb{P}\{X = j\} = 2 e^{-2\rho\lambda} \sum_{k=0}^{\infty} \frac{(\rho\lambda)^k}{k!} \sum_{j=2k+1}^{\infty} \frac{(\rho\lambda)^j}{j!}. \end{aligned}$$

If the average number of raisins per pound is $\lambda = 20$, and if the loaves are baked from $\rho = 5$ pounds dough, then this probability is approximately

$$\mathbb{P}\{X > 2Y \text{ or } Y > 2X\} \approx 3.17061 \times 10^{-6}.$$

If $\rho = 1$, that is, the loaves are made from one-pound dough each, then

$$\mathbb{P}\{X > 2Y \text{ or } Y > 2X\} \approx 0.0430079.$$

Next we investigate the distribution of the sum of two independent *negative binomial distributed* random variables. Recall that X is $B_{n,p}^-$ -distributed if

$$\mathbb{P}\{X = k\} = \binom{k-1}{k-n} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots$$

Proposition 4.6.8. *Let X and Y be independent and respectively $B_{n,p}^-$ - and $B_{m,p}^-$ -distributed for some $m, n \geq 1$. Then $X + Y$ is $B_{n+m,p}^-$ -distributed.*

Proof. We derive from Example 4.1.8 that, if $k \in \mathbb{N}_0$, then

$$\mathbb{P}\{X - n = k\} = \binom{-n}{k} p^n (p-1)^k \quad \text{and} \quad \mathbb{P}\{Y - m = k\} = \binom{-m}{k} p^m (p-1)^k.$$

An application of Proposition 4.5.6 to $X - n$ and $Y - m$ implies

$$\begin{aligned} \mathbb{P}\{X + Y - (n + m) = k\} &= \sum_{j=0}^k \left[\binom{-n}{j} p^n (p-1)^j \right] \left[\binom{-m}{k-j} p^m (p-1)^{k-j} \right] \\ &= p^{n+m} (p-1)^k \sum_{j=0}^k \binom{-n}{j} \binom{-m}{k-j}. \end{aligned}$$

To compute the last sum, we use Proposition A.5.3, which asserts that

$$\sum_{j=0}^k \binom{-n}{j} \binom{-m}{k-j} = \binom{-n-m}{k}.$$

Consequently, for each $k \in \mathbb{N}_0$,

$$\mathbb{P}\{X + Y - (n + m) = k\} = \binom{-n-m}{k} p^{n+m} (p-1)^k.$$

Another application of eq. (4.2) (this time with $n + m$) leads to

$$\mathbb{P}\{X + Y = k\} = \binom{k-1}{k-(n+m)} p^{n+m} (1-p)^{k-(n+m)}, \quad k = n+m, n+m+1, \dots,$$

completing the proof. \square

Corollary 4.6.9. *Let X_1, \dots, X_n be independent G_p -distributed (geometrically distributed) random variables. Then their sum $X_1 + \dots + X_n$ is $B_{n,p}^-$ -distributed.*

Proof. Use $G_p = B_{1,p}^-$ and apply Proposition 4.6.8 n times. \square

Interpretation: The following two experiments are completely equivalent: one is to play the same game until one observes success for the n th time. The other experiment is, after each success to start a new game, until one observes success in the n th (and last) game. Here we assume that all n games are executed independently and possess the same success probability.

Let U and V be two independent random variables, both *uniformly distributed* on $[0, 1]$. Which distribution density does $U + V$ possess?

Proposition 4.6.10. *The sum of two independent random variables U and V , uniformly distributed on $[0, 1]$, has the density r defined by*

$$r(x) = \begin{cases} x & \text{if } 0 \leq x < 1, \\ 2-x & \text{if } 1 \leq x \leq 2, \\ 0 & \text{otherwise.} \end{cases} \quad (4.30)$$

Proof. The distribution densities p and q of U and V are given by $p(x) = q(x) = 1$ if $0 \leq x \leq 1$ and $p(x) = q(x) = 0$ otherwise. Proposition 4.5.15 asserts that $U + V$ has density $r = p * q$ computed by

$$r(x) = \int_{-\infty}^{\infty} p(x-y) q(y) dy = \int_0^1 p(x-y) dy.$$

But $p(x-y) = 1$ if and only if $0 \leq x-y \leq 1$ or, equivalently, if and only if $x-1 \leq y \leq x$. Taking into account the restriction $0 \leq y \leq 1$, it follows that $p(x-y)q(y) = 1$ if and only

if $y \in [\max\{x - 1, 0\}, \min\{x, 1\}]$. In particular, $r(x) = 0$ for $x \notin [0, 2]$, and if $0 \leq x \leq 2$, then

$$r(x) = \min\{x, 1\} - \max\{x - 1, 0\}.$$

It is not difficult to see (treat the cases $0 \leq x \leq 1$ and $1 \leq x \leq 2$ separately) that $r(x)$ may be written as stated in eq. (4.30). This completes the proof. \square

Application. Suppose we choose independently and according to the uniform distribution two numbers u_1 and u_2 in $[0, 1]$. Then the probability that $a \leq u_1 + u_2 \leq b$ equals $\int_a^b r(x) dx$ with r given by eq. (4.30). For example (see Figure 4.5),

$$\mathbb{P}\left\{\frac{1}{2} \leq u_1 + u_2 \leq \frac{3}{2}\right\} = \int_{1/2}^1 x dx + \int_1^{3/2} (2 - x) dx = \frac{3}{4}.$$

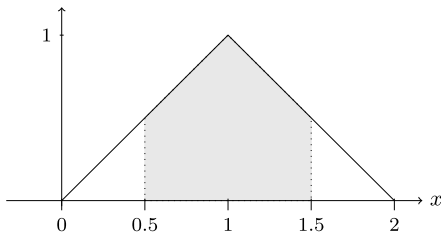


Figure 4.5: The area of the gray-shaded set equals $\mathbb{P}\{1/2 \leq U + V \leq 3/2\}$. Here the random variables U and V are independent and both uniformly distributed on $[0, 1]$.

We investigate now the sum of two *gamma distributed* random variables. Recall that the density of a $\Gamma_{\alpha,\beta}$ -distributed random variable is given by

$$p_{\alpha,\beta}(x) = \frac{1}{\alpha^\beta \Gamma(\beta)} x^{\beta-1} e^{-x/\alpha}$$

if $x > 0$, while $p_{\alpha,\beta}(x) = 0$ otherwise.

Proposition 4.6.11. *Let X_1 and X_2 be two independent random variables distributed according to Γ_{α,β_1} and Γ_{α,β_2} , respectively. Then $X_1 + X_2$ is $\Gamma_{\alpha,\beta_1+\beta_2}$ -distributed.*

Proof. If r denotes the density of $X_1 + X_2$, Proposition 4.5.15 implies

$$r(x) = (p_{\alpha,\beta_1} * p_{\alpha,\beta_2})(x) = \int_{-\infty}^{\infty} p_{\alpha,\beta_1}(x-y)p_{\alpha,\beta_2}(y) dy, \quad x \in \mathbb{R}, \quad (4.31)$$

and we have to show that $r = p_{\alpha,\beta_1+\beta_2}$.

It is easy to see that $r(x) = 0$ if $x \leq 0$, hence it suffices to evaluate eq. (4.31) for $x > 0$. Since $p_{\alpha, \beta_2}(x - y) = 0$ if $y > x$,

$$\begin{aligned} r(x) &= \frac{1}{\alpha^{\beta_1 + \beta_2} \Gamma(\beta_1) \Gamma(\beta_2)} \int_0^x y^{\beta_1 - 1} (x - y)^{\beta_2 - 1} e^{-y/\alpha} e^{-(x-y)/\alpha} dy \\ &= \frac{1}{\alpha^{\beta_1 + \beta_2} \Gamma(\beta_1) \Gamma(\beta_2)} x^{\beta_1 + \beta_2 - 2} e^{-x/\alpha} \int_0^x \left(\frac{y}{x}\right)^{\beta_1 - 1} \left(1 - \frac{y}{x}\right)^{\beta_2 - 1} dy. \end{aligned}$$

Changing the variable as $u := y/x$, hence $dy = x du$, we obtain

$$\begin{aligned} r(x) &= \frac{1}{\alpha^{\beta_1 + \beta_2} \Gamma(\beta_1) \Gamma(\beta_2)} x^{\beta_1 + \beta_2 - 1} e^{-x/\alpha} \int_0^1 u^{\beta_1 - 1} (1 - u)^{\beta_2 - 1} du \\ &= \frac{B(\beta_1, \beta_2)}{\alpha^{\beta_1 + \beta_2} \Gamma(\beta_1) \Gamma(\beta_2)} \cdot x^{\beta_1 + \beta_2 - 1} e^{-x/\alpha}, \end{aligned} \quad (4.32)$$

where B denotes the beta function defined by eq. (1.61). Equation (1.62) yields

$$\frac{B(\beta_1, \beta_2)}{\Gamma(\beta_1) \Gamma(\beta_2)} = \frac{1}{\Gamma(\beta_1 + \beta_2)}, \quad (4.33)$$

hence, if $x > 0$, then, from eqs. (4.32) and (4.33), it follows that $r(x) = p_{\alpha, \beta_1 + \beta_2}(x)$. This completes the proof. \square

Recall that the *Erlang distribution* is defined as $E_{\lambda, n} = \Gamma_{\lambda^{-1}, n}$. Thus, Proposition 4.6.11 implies the following corollary.

Corollary 4.6.12. *Let X and Y be independent and distributed according to $E_{\lambda, n}$ and $E_{\lambda, m}$, respectively. Then their sum $X + Y$ is $E_{\lambda, n+m}$ -distributed.*

Another corollary of Proposition 4.5.15 (or of Corollary 4.6.12) describes the sum of independent *exponentially distributed* random variables.

Corollary 4.6.13. *Let X_1, \dots, X_n be independent E_λ -distributed random variables. Then their sum $X_1 + \dots + X_n$ is Erlang distributed with parameters λ and n .*

Proof. Recall that $E_\lambda = E_{\lambda, 1}$. By Corollary 4.6.12, the sum $X_1 + X_2$ is distributed according to $E_{\lambda, 2}$. Proceeding in this way, every time applying Corollary 4.6.12 leads to the desired result. \square

Example 4.6.14. The lifetime of light bulbs is assumed to be E_λ -distributed for a certain $\lambda > 0$. At time zero, we switch on the first bulb. At the moment it burns out, we replace it by the second one of the same type. If the second burns out, we replace it by the third, and so on. Let S_n be the moment when the n th light bulb burns out. Which distribution does S_n possess?

Answer: Let X_1, X_2, \dots be the lifetimes of the first, second, and so on, light bulb. Then $S_n = X_1 + \dots + X_n$. Since the light bulbs are assumed to be of the same type, the random variables X_j are all E_λ -distributed. Furthermore, the different lifetimes do not influence each other, thus, we may assume that the X_j s are independent. Now Corollary 4.6.13 lets us conclude that S_n is Erlang distributed with parameters λ and n , hence, if $t > 0$, by Proposition 1.6.26, we get

$$\mathbb{P}\{S_n \leq t\} = \frac{\lambda^n}{(n-1)!} \int_0^t x^{n-1} e^{-\lambda x} dx = 1 - \sum_{j=0}^{n-1} \frac{(\lambda t)^j}{j!} e^{-\lambda t}. \quad (4.34)$$

Example 4.6.15. We continue the preceding example, but ask now a different question. How often do we have to change light bulbs before some given time $T > 0$?

Answer: Let Y be the number of changes necessary until time T . Then for $n \geq 0$ the event $\{Y = n\}$ occurs if and only if $S_n \leq T$, but $S_{n+1} > T$. Hereby, we use the notation of Example 4.6.14. In other words,

$$\mathbb{P}\{Y = n\} = \mathbb{P}\{S_n \leq T, S_{n+1} > T\} = \mathbb{P}(\{S_n \leq T\} \setminus \{S_{n+1} \leq T\}), \quad n = 0, 1, \dots$$

Since $\{S_{n+1} \leq T\} \subseteq \{S_n \leq T\}$, from eq. (4.34) it follows that

$$\begin{aligned} \mathbb{P}\{Y = n\} &= \mathbb{P}\{S_n \leq T\} - \mathbb{P}\{S_{n+1} \leq T\} \\ &= \left[1 - \sum_{j=0}^{n-1} \frac{(\lambda T)^j}{j!} e^{-\lambda T} \right] - \left[1 - \sum_{j=0}^n \frac{(\lambda T)^j}{j!} e^{-\lambda T} \right] \\ &= \frac{(\lambda T)^n}{n!} e^{-\lambda T} = \text{Pois}_{\lambda T}(\{n\}). \end{aligned}$$

Summing up, the number of necessary replacements of burned out light bulbs until time $T > 0$ is Poisson distributed with parameter λT where $1/\lambda > 0$ is the average lifetime of a single bulb (compare with Example 5.1.30).

Let us still mention an important equivalent random “experiment”: customers arrive at the post office randomly. We assume that the times between their arrivals are independent and E_λ -distributed. Then S_n is the time when the n th customer arrives. Hence, under these assumptions, the number of arriving customers until a certain time $T > 0$ is Poisson distributed with parameter λT .

We investigate now the sum of two independent *chi-squared distributed* random variables. Recall Definition 1.6.27: A random variable X is χ_n^2 -distributed if it is $\Gamma_{2, \frac{n}{2}}$ -distributed. Hence, Proposition 4.6.11 implies the following result.

Proposition 4.6.16. *Suppose that X is χ_n^2 -distributed and that Y is χ_m^2 -distributed for some $n, m \geq 1$. If X and Y are independent, then $X + Y$ is χ_{n+m}^2 -distributed.*

Proof. Because of Proposition 4.6.11, the sum $X + Y$ is $\Gamma_{2, \frac{n}{2} + \frac{m}{2}} = \chi_{n+m}^2$ -distributed. This proves the assertion. \square

Proposition 4.6.16 has the following important consequence.

Proposition 4.6.17. *Let X_1, \dots, X_n be a sequence of independent $\mathcal{N}(0, 1)$ -distributed random variables. Then $X_1^2 + \dots + X_n^2$ is χ_n^2 -distributed.*

Proof. Proposition 4.1.5 asserts that the random variables X_j^2 are χ_1^2 -distributed. Furthermore, because of Proposition 4.1.9, they are also independent. Thus successive application of Proposition 4.6.16 proves the assertion. \square

Our next and final aim in this section is to investigate the distribution of the sum of two independent *normally distributed* random variables. Here the following important result is valid.

Proposition 4.6.18. *Let X_1 and X_2 be two independent random variables distributed according to $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$. Then $X_1 + X_2$ is $\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$ -distributed.*

Proof. In the first step, we treat a special case, namely $\mu_1 = \mu_2 = 0$ and $\sigma_1 = 1$. To simplify the notation, set $\lambda = \sigma_2$. Thus we have to prove the following: if X_1 and X_2 are $\mathcal{N}(0, 1)$ - and $\mathcal{N}(0, \lambda^2)$ -distributed, then $X_1 + X_2$ is $\mathcal{N}(0, 1 + \lambda^2)$ -distributed.

Let $p_{0,1}$ and p_{0,λ^2} be the corresponding densities introduced in eq. (1.49). Then we have to prove that

$$p_{0,1} * p_{0,\lambda^2} = p_{0,1+\lambda^2}. \quad (4.35)$$

To verify this start with

$$\begin{aligned} (p_{0,1} * p_{0,\lambda^2})(x) &= \frac{1}{2\pi\lambda} \int_{-\infty}^{\infty} e^{-(x-y)^2/2} e^{-y^2/2\lambda^2} dy \\ &= \frac{1}{2\pi\lambda} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(x^2 - 2xy + (1+\lambda^{-2})y^2)} dy. \end{aligned} \quad (4.36)$$

We use

$$\begin{aligned} x^2 - 2xy + (1 + \lambda^{-2})y^2 &= ((1 + \lambda^{-2})^{1/2}y - (1 + \lambda^{-2})^{-1/2}x)^2 - x^2 \left(\frac{1}{1 + \lambda^{-2}} - 1 \right) \\ &= ((1 + \lambda^{-2})^{1/2}y - (1 + \lambda^{-2})^{-1/2}x)^2 + \frac{x^2}{1 + \lambda^2} \\ &= \left(ay - \frac{x}{a} \right)^2 + \frac{x^2}{1 + \lambda^2} \end{aligned}$$

with $a := (1 + \lambda^{-2})^{1/2}$. Plugging this transformation into eq. (4.36) leads to

$$(p_{0,1} * p_{0,\lambda^2})(x) = \frac{e^{-x^2/2(1+\lambda^2)}}{2\pi\lambda} \int_{-\infty}^{\infty} e^{-(ay - \frac{x}{a})^2/2} dy. \quad (4.37)$$

Next change the variables by $u := ay - x/a$, thus, $dy = du/a$, and observe that $a\lambda = (1 + \lambda^2)^{1/2}$. Then the right-hand side of eq. (4.37) transforms to

$$(p_{0,1} * p_{0,\lambda^2})(x) = \frac{e^{-x^2/2(1+\lambda^2)}}{2\pi(1+\lambda^2)^{1/2}} \int_{-\infty}^{\infty} e^{-u^2/2} du = p_{0,1+\lambda^2}(x).$$

Hereby, we used Proposition 1.6.7 asserting $\int_{-\infty}^{\infty} e^{-u^2/2} du = \sqrt{2\pi}$. This proves the validity of eq. (4.35).

In the second step, we treat the general case, that is, X_1 is $\mathcal{N}(\mu_1, \sigma_1^2)$ - and X_2 is $\mathcal{N}(\mu_2, \sigma_2^2)$ -distributed. Set

$$Y_1 := \frac{X_1 - \mu_1}{\sigma_1} \quad \text{and} \quad Y_2 := \frac{X_2 - \mu_2}{\sigma_2}.$$

By Proposition 4.2.3, the random variables Y_1 and Y_2 are standard normal and, moreover, because of Proposition 4.1.9, also independent. Thus, the sum $X_1 + X_2$ may be represented as

$$X_1 + X_2 = \mu_1 + \mu_2 + \sigma_1 Y_1 + \sigma_2 Y_2 = \mu_1 + \mu_2 + \sigma_1 Z,$$

where $Z = Y_1 + \lambda Y_2$ with $\lambda = \sigma_2/\sigma_1$.

An application of the first step shows that Z is $\mathcal{N}(0, 1 + \lambda^2)$ -distributed. Hence, Proposition 4.2.3 implies the existence of a standard normally distributed Z_0 such that $Z = (1 + \lambda^2)^{1/2} Z_0$. Summing up, $X_1 + X_2$ may now be written as

$$X_1 + X_2 = \mu_1 + \mu_2 + \sigma_1 (1 + \lambda^2)^{1/2} Z_0 = \mu_1 + \mu_2 + (\sigma_1^2 + \sigma_2^2)^{1/2} Z_0,$$

and another application of Proposition 4.2.3 lets us conclude that, as asserted, the sum $X_1 + X_2$ is $\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$ -distributed. \square

Summary: Let X and Y be independent random variables. As we agreed upon in Remark 3.3.1, by “ \sim ” we mean that X , Y , and $X + Y$ possess the stated distribution. Then the following are valid:

1. $X \sim B_{m,p}$ and $Y \sim B_{n,p} \Rightarrow X + Y \sim B_{m+n,p}$.
2. $X \sim \text{Pois}_\lambda$ and $Y \sim \text{Pois}_\mu \Rightarrow X + Y \sim \text{Pois}_{\lambda+\mu}$.
3. $X \sim B_{m,p}^-$ and $Y \sim B_{n,p}^- \Rightarrow X + Y \sim B_{m+n,p}^-$.
4. $X \sim \Gamma_{\alpha,\beta_1}$ and $Y \sim \Gamma_{\alpha,\beta_2} \Rightarrow X + Y \sim \Gamma_{\alpha,\beta_1+\beta_2}$.
5. $X \sim E_{\lambda,m}$ and $Y \sim E_{\lambda,n} \Rightarrow X + Y \sim E_{\lambda,m+n}$.
6. $X \sim \chi_m^2$ and $Y \sim \chi_n^2 \Rightarrow X + Y \sim \chi_{m+n}^2$.
7. $X \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $Y \sim \mathcal{N}(\mu_2, \sigma_2^2) \Rightarrow X + Y \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$.
8. X_1, \dots, X_n independent and $\mathcal{N}(0, 1)$ -distributed $\Rightarrow X_1^2 + \dots + X_n^2 \sim \chi_n^2$.

4.7 Products and quotients of random variables

Let X and Y be two random variables mapping a sample space Ω into \mathbb{R} . Then their product $X \cdot Y$ and their quotient X/Y (assume $Y(\omega) \neq 0$ for $\omega \in \Omega$) are defined by

$$(X \cdot Y)(\omega) := X(\omega) \cdot Y(\omega) \quad \text{and} \quad \left(\frac{X}{Y}\right)(\omega) := \frac{X(\omega)}{Y(\omega)}, \quad \omega \in \Omega.$$

The aim of this section is to investigate the distribution of such products and quotients. We restrict ourselves to continuous X and Y because, later on, we will only deal with products and quotients of those random variables. Furthermore, we omit the proof of the fact that products and fractions are random variables as well. The proofs of these permanent properties are not complicated and follow the ideas used in the proof of Proposition 4.5.1. Thus, our interest are products $X \cdot Y$ and quotients X/Y for independent X and Y , where, to simplify the computations, we suppose $\mathbb{P}\{Y > 0\} = 1$.

We start with the investigation of *products* of continuous random variables. Thus, let X and Y be two random variables with distribution densities p and q . Since we assumed $\mathbb{P}\{Y > 0\} = 1$, we may choose the density q such that $q(x) = 0$ if $x \leq 0$.

Proposition 4.7.1. *Let X and Y be two independent random variables possessing the stated properties. Then $X \cdot Y$ is continuous as well, and its density r may be calculated by*

$$r(x) = \int_0^{\infty} p\left(\frac{x}{y}\right) \frac{q(y)}{y} dy, \quad x \in \mathbb{R}. \quad (4.38)$$

Proof. For $t \in \mathbb{R}$, we evaluate $\mathbb{P}\{X \cdot Y \leq t\}$. To this end, fix $t \in \mathbb{R}$ and set

$$A_t := \{(u, y) \in \mathbb{R} \times (0, \infty) : u \cdot y \leq t\}. \quad (4.39)$$

As in the proof of Proposition 4.5.15, it follows that

$$\mathbb{P}\{X \cdot Y \leq t\} = \mathbb{P}_{(X, Y)}(A_t) = \int_0^{\infty} \left[\int_{-\infty}^{t/y} p(u) du \right] q(y) dy. \quad (4.40)$$

In the inner integral of eq. (4.40), we change the variables by $x = uy$, hence we get $dx = y du$. Notice that in the inner integral y is a constant. After this change of variables, the right-hand integral in eq. (4.40) becomes¹⁵

$$\int_0^{\infty} \left[\int_{-\infty}^t p\left(\frac{x}{y}\right) dx \right] \frac{q(y)}{y} dy = \int_{-\infty}^t \left[\int_0^{\infty} p\left(\frac{x}{y}\right) \frac{q(y)}{y} dy \right] dx = \int_{-\infty}^t r(x) dx.$$

This being valid for all $t \in \mathbb{R}$, the function r is a density of $X \cdot Y$. □

¹⁵ The interchange of the integrals is justified by Proposition A.5.5. Note that p and q are nonnegative.

Example 4.7.2. Let U and V be two independent random variables uniformly distributed on $[0, 1]$. Which probability distribution does $U \cdot V$ possess?

Answer: We have $p(y) = q(y) = 1$ if $0 \leq y \leq 1$, and $p(y) = q(y) = 0$ otherwise. Furthermore, $0 \leq U \cdot V \leq 1$, hence its density r satisfies $r(x) = 0$ if $x \notin [0, 1]$. For $x \in [0, 1]$, we apply formula (4.38) and obtain (see Figure 4.6)

$$r(x) = \int_0^{\infty} p\left(\frac{x}{y}\right) \frac{q(y)}{y} dy = \int_x^1 \frac{1}{y} dy = -\ln(x) = \ln\left(\frac{1}{x}\right), \quad 0 < x \leq 1. \quad (4.41)$$

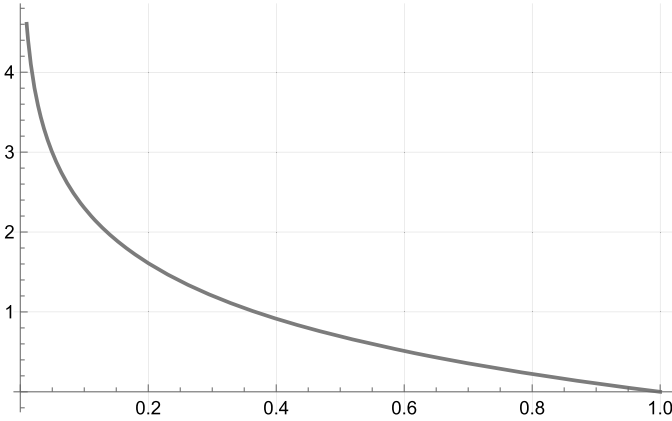


Figure 4.6: The density r of $U \cdot V$ given by eq. (4.41).

Consequently, if $0 < a < b \leq 1$, then

$$\mathbb{P}\{a \leq U \cdot V \leq b\} = - \int_a^b \ln(x) dx = -[x \ln x - x]_a^b = a \ln(a) - b \ln(b) + b - a.$$

In particular, it follows that

$$\text{vol}_2(B_t) = \mathbb{P}\{U \cdot V \leq t\} = t - t \ln t, \quad 0 < t \leq 1. \quad (4.42)$$

Here B_t is defined as in Fig. 4.7. That is,

$$B_t = \{(u, y) \in [0, 1]^2 : uy \leq t\}, \quad 0 \leq t \leq 1.$$

In other words, $B_t = A_t \cap [0, 1]^2$ with A_t defined by eq. (4.39). Furthermore, the random vector (U, V) is uniformly distributed on $[0, 1]^2$, so we get

$$\mathbb{P}\{(U, V) \in A_t\} = \mathbb{P}\{(U, V) \in B_t\} = \text{vol}_2(B_t).$$

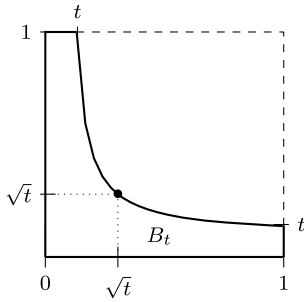


Figure 4.7: The set $B_t = \{(u, y) \in [0, 1]^2 : uy \leq t\}$ for a given $0 < t < 1$.

Our next objective are *quotients* of random variables X and Y . We denote their densities by p and q , thereby assuming $q(x) = 0$ if $x \leq 0$. Then we get

Proposition 4.7.3. *Let X and Y be independent with $\mathbb{P}\{Y > 0\} = 1$. Then their quotient X/Y has the density r given by*

$$r(x) = \int_0^{\infty} y p(xy) q(y) dy, \quad x \in \mathbb{R}.$$

Proof. The proof of Proposition 4.7.3 is similar to that of Proposition 4.7.1. Therefore, we present only the main steps. Setting now

$$A_t := \{(u, y) \in \mathbb{R} \times (0, \infty) : u \leq ty\},$$

we obtain

$$\mathbb{P}\{(X/Y) \leq t\} = \mathbb{P}_{(X,Y)}(A_t) = \int_0^{\infty} \left[\int_{-\infty}^{ty} p(u) du \right] q(y) dy. \quad (4.43)$$

We change the variables in the inner integral of eq. (4.43) by putting $x = u/y$. After that, we interchange the integrals and arrive at

$$\mathbb{P}\{(X/Y) \leq t\} = \int_{-\infty}^t r(x) dx$$

for all $t \in \mathbb{R}$. This proves that r is the density of X/Y . □

Example 4.7.4. Let U and V be as in Example 4.7.2. We investigate now their quotient U/V . By Proposition 4.7.3, its density r can be computed by

$$r(x) = \int_0^{\infty} y p(xy) q(y) dy = \int_0^1 y p(xy) dy = \int_0^1 y dy = \frac{1}{2}$$

in the case $0 \leq x \leq 1$. If $1 \leq x < \infty$, then $p(xy) = 0$ if $y > 1/x$, and it follows that

$$r(x) = \int_0^{1/x} y \, dy = \frac{1}{2x^2}$$

for those x . Combining both cases, the density r of U/V may be written as (see Figure 4.8)

$$r(x) = \begin{cases} \frac{1}{2} & \text{if } 0 \leq x \leq 1, \\ \frac{1}{2x^2} & \text{if } 1 < x < \infty, \\ 0 & \text{otherwise.} \end{cases} \quad (4.44)$$

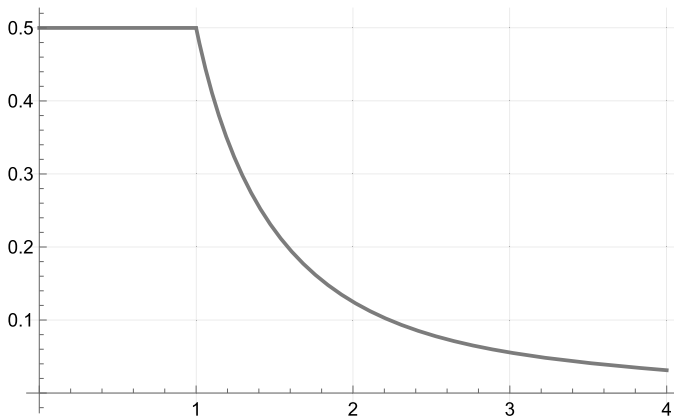


Figure 4.8: The density r defined by eq. (4.44).

Question: Does there exist an easy geometric explanation for $r(x) = \frac{1}{2}$ in the case $0 \leq x \leq 1$?

Answer: If $t > 0$, then

$$F_{U/V}(t) = \mathbb{P}\{U/V \leq t\} = \mathbb{P}\{U \leq tV\} = \mathbb{P}_{(U,V)}(B_t),$$

where now

$$B_t := \{(u, v) \in [0, 1]^2 : 0 \leq u \leq vt\}. \quad (4.45)$$

If $0 < t \leq 1$, then B_t is a triangle in $[0, 1]^2$ with area $\text{vol}_2(B_t) = \frac{t}{2}$. The independence of U and V implies (cf. Example 3.6.23) that $\mathbb{P}_{(U,V)}$ is the uniform distribution on $[0, 1]^2$, hence

$$F_{U/V}(t) = \mathbb{P}_{(U,V)}(B_t) = \text{vol}_2(B_t) = \frac{t}{2}, \quad 0 < t \leq 1,$$

leading to $r(t) = F'_{U|V}(t) = \frac{1}{2}$ for those t . See Figure 4.9 for a geometric explanation of this fact.

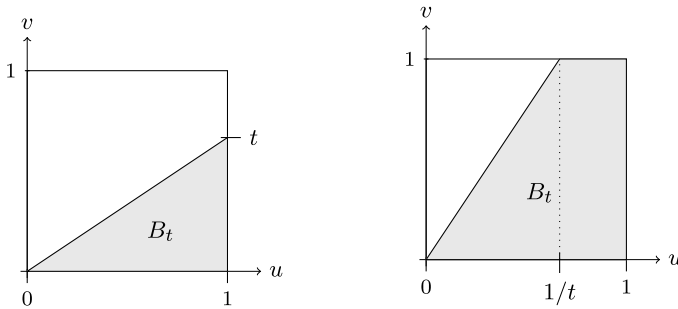


Figure 4.9: The set B_t , defined in eq. (4.45). Here $0 \leq t \leq 1$ in the left-hand figure, and $1 < t < \infty$ in the right-hand one. The area of B_t is either $\frac{t}{2}$ or $1 - \frac{1}{2t}$, respectively.

4.7.1 Student's t -distribution

Let us use Proposition 4.7.3 to compute the density of a distribution which plays a crucial role in Mathematical Statistics.

Proposition 4.7.5. *Let X be $\mathcal{N}(0,1)$ -distributed and Y be independent of X and χ_n^2 -distributed for some $n \geq 1$. Define the random variable Z as*

$$Z := \frac{X}{\sqrt{Y/n}}.$$

Then Z possesses the density r given by (see Figure 4.10 for the graphs of these functions in the cases $n = 1$, $n = 2$, and $n = 8$)

$$r(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2}, \quad x \in \mathbb{R}. \tag{4.46}$$

Proof. In the first step, we determine the density of \sqrt{Y} with Y distributed according to χ_n^2 . If $t > 0$, then

$$F_{\sqrt{Y}}(t) = \mathbb{P}\{\sqrt{Y} \leq t\} = \mathbb{P}\{Y \leq t^2\} = \frac{1}{2^{n/2} \Gamma(\frac{n}{2})} \int_0^{t^2} x^{n/2-1} e^{-x/2} dx.$$

Thus, if $t > 0$, then the density q of \sqrt{Y} equals

$$q(t) = \frac{d}{dt} F_{\sqrt{Y}}(t) = \frac{1}{2^{n/2} \Gamma(\frac{n}{2})} (2t) t^{n-2} e^{-t^2/2} = \frac{1}{2^{n/2-1} \Gamma(\frac{n}{2})} t^{n-1} e^{-t^2/2}. \tag{4.47}$$

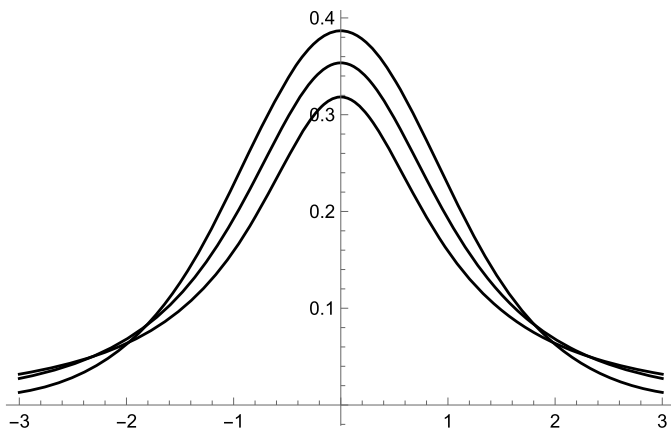


Figure 4.10: From bottom to top, these are the densities of t_1 , t_2 , and t_8 distributions.

Of course, we have $q(t) = 0$ if $t \leq 0$.

In the second step, we determine the density \tilde{r} of $\tilde{Z} = Z/\sqrt{n} = X/\sqrt{Y}$. An application of Proposition 4.7.3 for $p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ and q given in eq. (4.47) leads to

$$\begin{aligned} \tilde{r}(x) &= \int_0^{\infty} y \left[\frac{1}{\sqrt{2\pi}} e^{-(xy)^2/2} \right] \left[\frac{1}{2^{n/2-1} \Gamma(\frac{n}{2})} y^{n-1} e^{-y^2/2} \right] dy \\ &= \frac{1}{\sqrt{\pi} 2^{n/2-1/2} \Gamma(\frac{n}{2})} \int_0^{\infty} y^n e^{-(1+x^2)y^2/2} dy. \end{aligned} \quad (4.48)$$

Change the variables in the last integral by setting $v = \frac{y^2}{2}(1+x^2)$. Then $y = \frac{\sqrt{2v}}{(1+x^2)^{1/2}}$ and, consequently, $dy = \frac{1}{\sqrt{2}} v^{-1/2} (1+x^2)^{-1/2} dv$. Inserting this into eq. (4.48) shows that

$$\begin{aligned} \tilde{r}(x) &= \frac{1}{\sqrt{\pi} 2^{n/2} \Gamma(\frac{n}{2})} \int_0^{\infty} \frac{2^{n/2} v^{n/2-1/2} e^{-v}}{(1+x^2)^{n/2+1/2}} dv \\ &= \frac{\Gamma(\frac{n+1}{2})}{\sqrt{\pi} \Gamma(\frac{n}{2})} (1+x^2)^{-n/2-1/2}. \end{aligned} \quad (4.49)$$

In the third step, we finally obtain the density r of Z . Since $Z = \sqrt{n}\tilde{Z}$, formula (4.3) applies with $b = 0$ and $a = \sqrt{n}$. Thus, by eq. (4.49) for \tilde{r} , as asserted,

$$r(x) = \frac{1}{\sqrt{n}} \tilde{r}\left(\frac{x}{\sqrt{n}}\right) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2}. \quad \square$$

Definition 4.7.6. The probability measure on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ with density r , given by eq. (4.46), is called the **t_n -distribution** or **Student's t -distribution** with n degrees of freedom. A random variable Z is said to be t_n -distributed (or t -distributed with n degrees of freedom), provided its probability distribution is a t_n -distribution, that is, for $a < b$,

$$\mathbb{P}\{a \leq Z \leq b\} = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \int_a^b \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} dx.$$

Remark 4.7.7. The t_1 -distribution coincides with the Cauchy distribution introduced in Section 1.6.8. Observe that $\Gamma(1/2) = \sqrt{\pi}$ and $\Gamma(1) = 1$.

In view of Definition 4.7.6, we may now formulate Proposition 4.7.5 as follows.

Proposition 4.7.8. *If X and Y are independent and $\mathcal{N}(0, 1)$ and χ_n^2 distributed, then $\frac{X}{\sqrt{Y/n}}$ is t_n -distributed.*

Proposition 4.6.17 leads still to another version of Proposition 4.7.5.

Proposition 4.7.9. *If X, X_1, \dots, X_n are independent $\mathcal{N}(0, 1)$ -distributed, then*

$$\frac{X}{\sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}}$$

is t_n -distributed.

Corollary 4.7.10. *If X and Y are independent and $\mathcal{N}(0, 1)$ -distributed, then $X/|Y|$ possesses a Cauchy distribution.*

Proof. An application of Proposition 4.7.9 with $n = 1$ and $X_1 = Y$ implies that $X/|Y|$ is t_1 -distributed. We saw in Remark 4.7.7 the t_1 and the Cauchy distributions coincide, thus, $X/|Y|$ is also Cauchy distributed. \square

4.7.2 F-distribution

We present now another important class of probability measures or probability distributions playing a central role in Mathematical Statistics.

Proposition 4.7.11. *For two natural numbers m and n , let X and Y be independent and χ_m^2 - and χ_n^2 -distributed. Then $Z := \frac{X/m}{Y/n}$ has the distribution density r defined by*

$$r(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ m^{m/2} n^{n/2} \cdot \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \cdot \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}} & \text{if } x > 0. \end{cases} \quad (4.50)$$

Proof. We first evaluate the density \tilde{r} of $\tilde{Z} = X/Y$. To this end, we apply Proposition 4.7.3 with functions p and q given by

$$p(x) = \frac{1}{2^{m/2} \Gamma(m/2)} x^{m/2-1} e^{-x/2} \quad \text{and} \quad q(y) = \frac{1}{2^{n/2} \Gamma(n/2)} y^{n/2-1} e^{-y/2}$$

whenever $x, y > 0$. Then we get

$$\begin{aligned} \tilde{r}(x) &= \frac{1}{2^{(m+n)/2} \Gamma(m/2) \Gamma(n/2)} \int_0^\infty y (xy)^{m/2-1} y^{n/2-1} e^{-xy/2} e^{-y/2} dy \\ &= \frac{x^{m/2-1}}{2^{(m+n)/2} \Gamma(m/2) \Gamma(n/2)} \int_0^\infty y^{(m+n)/2-1} e^{-y(1+x)/2} dy. \end{aligned} \quad (4.51)$$

We replace in eq. (4.51) the variable y by $u = y(1+x)/2$, thus, $dy = \frac{2}{1+x} du$. Inserting this into eq. (4.51), the last expression transforms to

$$\begin{aligned} \tilde{r}(x) &= \frac{x^{m/2-1}}{\Gamma(m/2) \Gamma(n/2)} (1+x)^{-(m+n)/2} \int_0^\infty u^{(m+n)/2-1} e^{-u} du \\ &= \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2}) \Gamma(\frac{n}{2})} \cdot \frac{x^{m/2-1}}{(1+x)^{(m+n)/2}}. \end{aligned}$$

Because of $Z = \frac{n}{m} \cdot \tilde{Z}$, we obtain the density r of Z by Proposition 1.7.21. Indeed, then

$$r(x) = \frac{m}{n} \tilde{r}\left(\frac{mx}{n}\right) = m^{m/2} n^{n/2} \cdot \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2}) \Gamma(\frac{n}{2})} \cdot \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}},$$

as asserted. □

Remark 4.7.12. Using relation (1.62) between the beta and gamma functions, the density r of Z may also be written as

$$r(x) = \frac{m^{m/2} n^{n/2}}{B(\frac{m}{2}, \frac{n}{2})} \cdot \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}}, \quad x > 0.$$

For the behavior of $r(x)$ as $x \rightarrow 0$, we refer to Problem 4.17. For certain parameters m and n the graphs of these densities can be found in Figure 4.11.

Definition 4.7.13. The probability measure on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ with density r defined by eq. (4.50) is called the **Fisher–Snedecor distribution** or **F-distribution** (with m and n degrees of freedom).

A random variable Z is F-distributed (with m and n degrees of freedom), provided its probability distribution is an F-distribution. Equivalently, if $0 \leq a < b$, then

$$\mathbb{P}\{a \leq Z \leq b\} = m^{m/2} n^{n/2} \cdot \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2}) \Gamma(\frac{n}{2})} \int_a^b \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}} dx.$$

The random variable Z is also said to be **F_{m,n}-distributed**.

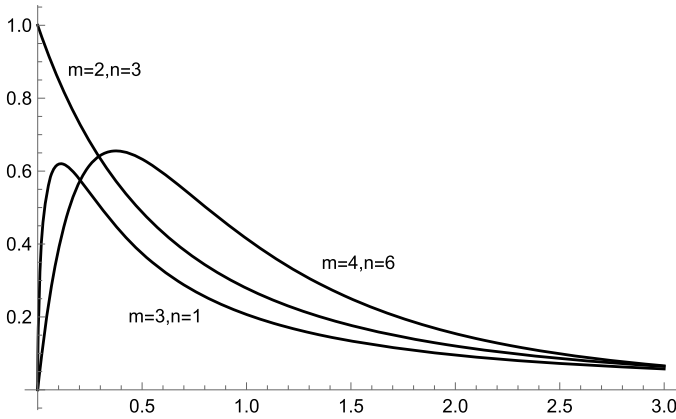


Figure 4.11: Densities of $F_{2,3}$, $F_{3,1}$, and $F_{4,6}$ distributed random variables.

With this notation, Proposition 4.7.11 may now be formulated as follows:

Proposition 4.7.14. *If two independent random variables X and Y are χ_m^2 and χ_n^2 distributed, then $\frac{X/m}{Y/n}$ is $F_{m,n}$ -distributed.*

Finally, Proposition 4.6.17 implies the following version of the previous result.

Proposition 4.7.15. *Let $X_1, \dots, X_m, Y_1, \dots, Y_n$ be independent $\mathcal{N}(0, 1)$ -distributed. Then*

$$Z = \frac{(X_1^2 + \dots + X_m^2)/m}{(Y_1^2 + \dots + Y_n^2)/n} = \frac{\frac{1}{m} \sum_{i=1}^m X_i^2}{\frac{1}{n} \sum_{j=1}^n Y_j^2}$$

is $F_{m,n}$ -distributed.

Corollary 4.7.16. *If a random variable Z is $F_{m,n}$ -distributed, then $1/Z$ possesses an $F_{n,m}$ distribution.*

Proof. This is an immediate consequence of Proposition 4.7.11. □

Summary: A random variable X is t_n -distributed provided that for all $a < b$ we have

$$\mathbb{P}\{a \leq X \leq b\} = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \int_a^b \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} dx.$$

If X, X_1, \dots, X_n are independent $\mathcal{N}(0, 1)$ -distributed, then

$$\frac{X}{\sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}} \text{ is } t_n\text{-distributed.}$$

A random variable X is $F_{m,n}$ -distributed if, whenever $0 \leq a < b < \infty$, one has

$$\mathbb{P}\{a \leq X \leq b\} = m^{m/2} n^{n/2} \cdot \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_a^b \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}} dx.$$

If $X_1, \dots, X_m, Y_1, \dots, Y_n$ are independent $\mathcal{N}(0, 1)$ -distributed, then

$$\frac{\frac{1}{m} \sum_{i=1}^m X_i^2}{\frac{1}{n} \sum_{j=1}^n Y_j^2} \text{ is } F_{m,n}\text{-distributed.}$$

4.8 Problems

Problem 4.1. Let U be uniformly distributed on $[0, 1]$. Which distributions do the following random variables possess:

$$\min\{U, 1 - U\}, \quad \max\{U, 1 - U\}, \quad |2U - 1|, \quad \text{and} \quad \left|U - \frac{1}{3}\right|?$$

Problem 4.2 (Generating functions). Let X be a random variable with values in \mathbb{N}_0 . For $k \in \mathbb{N}_0$, let $p_k = \mathbb{P}\{X = k\}$. Then its **generating function** φ_X is defined by

$$\varphi_X(t) = \sum_{k=0}^{\infty} p_k t^k.$$

1. Show that $\varphi_X(t)$ exists if $|t| \leq 1$.
2. Let X and Y be two independent random variables with values in \mathbb{N}_0 . Prove that then

$$\varphi_{X+Y} = \varphi_X \cdot \varphi_Y.$$

3. Compute φ_X in each of the following cases:
 - (a) X is uniformly distributed on $\{1, \dots, N\}$ for some $N \geq 1$.
 - (b) X is $B_{n,p}$ -distributed for some $n \geq 1$ and $p \in [0, 1]$.
 - (c) X is Pois_λ -distributed for some $\lambda > 0$.
 - (d) X is G_p -distributed for a certain $0 < p < 1$.
 - (e) X is $B_{n,p}^-$ -distributed.

Problem 4.3. Roll two dice simultaneously. Let X be the result of the first die and Y that of the second. Is it possible to falsify these two dice in such a way that $X + Y$ is uniformly distributed on $\{2, \dots, 12\}$? It is not assumed that both dice are falsified in the same way.

Hint: One possible way to answer this question is as follows: investigate the generating functions of X and Y and compare their product with the generating function of the uniform distribution on $\{2, \dots, 12\}$.

Problem 4.4. Let X_1, \dots, X_n be a sequence of independent identically distributed random variables with common distribution function F and distribution density p , that is,

$$\mathbb{P}\{X_j \leq t\} = F(t) = \int_{-\infty}^t p(x) dx, \quad j = 1, \dots, n.$$

Define random variables X_* and X^* by

$$X_* := \min\{X_1, \dots, X_n\} \quad \text{and} \quad X^* := \max\{X_1, \dots, X_n\}.$$

1. Determine the distribution functions and densities of X_* and X^* directly, that is, without using general results about order statistics.
2. Describe the distribution of the random variables X_* and X^* in the case that the X_j s are exponentially distributed with parameter $\lambda > 0$.
3. Suppose now the X_j s are uniformly distributed on $[0, 1]$. Describe the distribution of X_* and X^* in this case.

Problem 4.5. Let F be the distribution function of the uniform distribution on $\{1, \dots, N\}$ for some $N \in \mathbb{N}$. Determine its pseudoinverse function F^- . Do the same if F is either the distribution function of a binomial or a Poisson distribution. Given U uniformly distributed on $[0, 1]$, how is $F^-(U)$ distributed in each of these cases?

Problem 4.6. Find a function f from $(0, 1)$ to \mathbb{R} such that

$$\mathbb{P}\{f(U) = k\} = \frac{1}{2^k}, \quad k = 1, 2, \dots$$

for U uniformly distributed on $[0, 1]$.

Problem 4.7. Let U be uniform distributed on $[0, 1]$. Find functions f and g such that $X = f(U)$ and $Y = g(U)$ have the distribution densities p and q with

$$p(x) := \begin{cases} 0 & \text{if } x \notin (0, 1], \\ \frac{x^{-1/2}}{2} & \text{if } x \in (0, 1] \end{cases} \quad \text{and} \quad q(x) := \begin{cases} 0 & \text{if } |x| > 1, \\ x + 1 & \text{if } -1 \leq x \leq 0, \\ 1 - x & \text{if } 0 < x \leq 1. \end{cases}$$

Problem 4.8. Let X and Y be independent random variables with

$$\mathbb{P}\{X = k\} = \mathbb{P}\{Y = k\} = \frac{1}{2^k}, \quad k = 1, 2, \dots$$

How is $X + Y$ distributed?

Problem 4.9. The number of customers visiting a shop per day is Poisson distributed with parameter $\lambda > 0$. The probability that a single customer buys something equals p for a given $0 \leq p \leq 1$. Let X be the number of customers per day buying some goods. Determine the distribution of X .

Remark: We assume that the decision whether or not a single customer buys something is independent of the number of daily visitors.

A different way to formulate the above question is as follows: let X_0, X_1, \dots be independent random variables with $\mathbb{P}\{X_0 = 0\} = 1$,

$$\mathbb{P}\{X_j = 1\} = p, \quad \text{and} \quad \mathbb{P}\{X_j = 0\} = 1 - p, \quad j = 1, 2, \dots,$$

for a certain $p \in [0, 1]$. Furthermore, let Y be a Poisson-distributed random variable with parameter $\lambda > 0$, independent of the X_j s. Determine the distribution of

$$X := \sum_{j=0}^Y X_j.$$

Hint: Use the “infinite” version of the law of total probability as stated in Problem 2.5.

Problem 4.10. Suppose X and Y are independent and exponentially distributed with parameter $\lambda > 0$. Find the distribution densities of $X - Y$ and X/Y .

Problem 4.11. Two random variables U and V are independent and uniformly distributed on $[0, 1]$. Given $n \in \mathbb{N}$, find the distribution density of $U + nV$.

Problem 4.12. Let X and Y be independent random variable distributed according to Pois_λ and Pois_μ , respectively. Given $n \in \mathbb{N}_0$ and some $k \in \{0, \dots, n\}$, prove

$$\mathbb{P}\{X = k \mid X + Y = n\} = \binom{n}{k} \left(\frac{\lambda}{\lambda + \mu} \right)^k \left(\frac{\mu}{\lambda + \mu} \right)^{n-k} = B_{n,p}(\{k\})$$

with $p = \frac{\lambda}{\lambda + \mu}$.

Reformulation of the preceding problem: An owner of two stores, say store A and store B , observes that the number of customers in each of these stores is independent and Pois_λ and Pois_μ distributed. One day he was told that there were n customers in both stores together. What is the probability that k of the n customers were in store A , hence $n - k$ in store B ?

Problem 4.13. Let X and Y be independent standard normal variables. Show that X/Y is Cauchy distributed.

Hint: Use Corollary 4.7.10 and the fact that the vectors (X, Y) , $(-X, Y)$, $(X, -Y)$, and $(-X, -Y)$ are identically distributed. Note that the probability distribution of each of these two-dimensional vectors is the (two-dimensional) standard normal distribution.

Problem 4.14. Let X and Y be independent G_p -distributed. Find the probability distribution of $X - Y$.

Hint: Compare with Example 4.5.5. There we evaluated the distribution of $X - Y$ in the case $p = \frac{1}{2}$.

Problem 4.15. Let U and V be independent random variables, both uniformly distributed on $[-1, 1]$. Determine the density of $U \cdot V$.

Hint: Use the technique presented in Example 4.7.2 and its geometric explanation in Figure 4.7. To this end, treat the cases $UV > 0$ and $UV \leq 0$ separately.

Problem 4.16. Suppose X is a random variable with values in $(a, b) \subseteq \mathbb{R}$ and with density p . Let f from $(a, b) \rightarrow \mathbb{R}$ be (strictly) monotone and differentiable. Give a formula for q , the density of $f(X)$.

Use your result to evaluate the densities of e^X and e^{-X} where X is distributed according to $\mathcal{N}(0, 1)$.

Problem 4.17. Let r be the density of an $F_{m,n}$ -distributed random variable given by eq. (4.50). Argue why $r(x) \rightarrow \infty$ as $x \rightarrow 0$ if $m = 1$ and $n \geq 1$. Why do we have $r(0) = 1$ if $m = 2$ and $n \geq 1$? Evaluate $r(0)$ if $n = 1, 2, \dots$ and $m \geq 3$.

5 Expected value, variance, and covariance

5.1 Expected value

5.1.1 Expected value of discrete random variables

What is an expected value (also called mean value or expectation) of a random variable? How is it defined? Which property of the random variable does it describe and how it can be computed? Does every random variable possess an expected value? To answer these questions, let us start with an example.

Example 5.1.1. Suppose N students attend a certain exam. The number of possible points is 100. Given $j = 0, 1, \dots, 100$, let n_j be the number of students who achieved j points. Now choose randomly, according to the uniform distribution (a single student is chosen with probability $1/N$), one student. Name him or her ω , and define $X(\omega)$ as the number of points that the chosen student achieved. Then X is a random variable with values in $D = \{0, 1, \dots, 100\}$. How is X distributed? Since X has values in D , its distribution is described by the probabilities

$$p_j = \mathbb{P}\{X = j\} = \frac{n_j}{N}, \quad j = 0, 1, \dots, 100. \quad (5.1)$$

As expected value of X , we take the average number A of points in this exam. How is A evaluated? The easiest way to do this is

$$A = \frac{1}{N} \sum_{j=0}^{100} j \cdot n_j = \sum_{j=0}^{100} j \cdot \frac{n_j}{N} = \sum_{j=0}^{100} j \cdot p_j,$$

where the p_j s are defined by eq. (5.1). If we write $\mathbb{E}X$ for the expected value (or mean value) of X , and if we assume that this value coincides with A , then the preceding equation says

$$\mathbb{E}X = \sum_{j=0}^{100} j p_j = \sum_{j=0}^{100} j \mathbb{P}\{X = j\} = \sum_{j=0}^{100} x_j \mathbb{P}\{X = x_j\},$$

where the $x_j = j$ with $j = 0, \dots, 100$ denote the possible values of X .

In view of this example, the following definition for the expected value of a discrete random variable X looks feasible. Suppose X has values in $D = \{x_1, x_2, \dots\}$, and let $p_j = \mathbb{P}\{X = x_j\}$, $j = 1, 2, \dots$. Then the expected value $\mathbb{E}X$ of X is given by

$$\mathbb{E}X = \sum_{j=1}^{\infty} x_j p_j = \sum_{j=1}^{\infty} x_j \mathbb{P}\{X = x_j\}. \quad (5.2)$$

Unfortunately, the sum in eq. (5.2) does not always exist. In order to overcome this difficulty, let us recall some basic facts about infinite series of real numbers.

A sequence $(a_j)_{j \geq 1}$ of real numbers is called **summable**, provided its sequence of partial sums $(s_n)_{n \geq 1}$ with $s_n = a_1 + \cdots + a_n$ converges in \mathbb{R} . Then one defines

$$\sum_{j=1}^{\infty} a_j = \lim_{n \rightarrow \infty} s_n = \lim_{n \rightarrow \infty} \sum_{j=1}^n a_j.$$

If the sequence of partial sums diverges, nevertheless, in some cases we may assign to the infinite series a limit. If either $\lim_{n \rightarrow \infty} s_n = -\infty$ or $\lim_{n \rightarrow \infty} s_n = \infty$, then we write $\sum_{j=1}^{\infty} a_j = -\infty$ or $\sum_{j=1}^{\infty} a_j = \infty$, respectively. In particular, if $a_j \geq 0$ for $j \in \mathbb{N}$, then the sequence of partial sums is nondecreasing, which implies that only two different cases may occur: Either $\sum_{j=1}^{\infty} a_j < \infty$ (in this case the sequence is summable) or $\sum_{j=1}^{\infty} a_j = \infty$.

Let $(a_j)_{j \geq 1}$ be an arbitrary sequence of real numbers. If $\sum_{j=1}^{\infty} |a_j| < \infty$, then it is called **absolutely summable**. Note that each absolutely summable sequence is summable. This is a direct consequence of Cauchy's convergence criterion. The converse implication is wrong, as can be seen by considering $((-1)^n/n)_{n \geq 1}$.

Now we are prepared to define the expected value of a nonnegative random variable.

Definition 5.1.2. Let X be a discrete random variable with values in $\{x_1, x_2, \dots\}$ for some $x_j \geq 0$. Equivalently, the random variable X is discrete with $X \geq 0$. Then the **expected value** of X is defined by

$$\mathbb{E}X := \sum_{j=1}^{\infty} x_j \mathbb{P}\{X = x_j\}. \quad (5.3)$$

Remark 5.1.3. Since $x_j \mathbb{P}\{X = x_j\} \geq 0$ for nonnegative X , for those random variables the sum in eq. (5.3) is always well defined, but may be infinite. That is, each nonnegative discrete random variable X possesses an expected value $\mathbb{E}X \in [0, \infty]$.

Let us now turn to the case of arbitrary (not necessarily nonnegative) random variables. The next example shows which problems may arise.

Example 5.1.4. We consider the probability measure introduced in Example 1.3.6 and choose a random variable X with values in \mathbb{Z} distributed according to the probability measure in this example. In other words,

$$\mathbb{P}\{X = k\} = \frac{3}{\pi^2} \frac{1}{k^2}, \quad k \in \mathbb{Z} \setminus \{0\}.$$

If we try to evaluate the expected value of X by formula (5.2), then this leads to the undetermined expression

$$\begin{aligned} \frac{3}{\pi^2} \sum_{\substack{k=-\infty \\ k \neq 0}}^{\infty} \frac{k}{k^2} &= \frac{3}{\pi^2} \lim_{\substack{n \rightarrow \infty \\ m \rightarrow \infty}} \sum_{k=-m}^n \frac{1}{k} = \frac{3}{\pi^2} \left[\lim_{n \rightarrow \infty} \sum_{k=1}^n \frac{1}{k} + \lim_{m \rightarrow \infty} \sum_{k=-m}^{-1} \frac{1}{k} \right] \\ &= \frac{3}{\pi^2} \left[\lim_{n \rightarrow \infty} \sum_{k=1}^n \frac{1}{k} - \lim_{m \rightarrow \infty} \sum_{k=1}^m \frac{1}{k} \right] = \infty - \infty. \end{aligned}$$

To exclude phenomenons as in Example 5.1.4, we suppose that a random variable has to meet the following condition.

Definition 5.1.5. Let X be discrete with values in $\{x_1, x_2, \dots\} \subset \mathbb{R}$. Then the **expected value of X exists**, provided that

$$\mathbb{E}|X| = \sum_{j=1}^{\infty} |x_j| \mathbb{P}\{X = x_j\} < \infty. \quad (5.4)$$

We mentioned above that an absolutely summable sequence is summable. Hence, under assumption (5.4), the sum in the subsequent definition is a well-defined real number.

Definition 5.1.6. Let X be a discrete random variable satisfying $\mathbb{E}|X| < \infty$. Then its **expected value** is defined as

$$\mathbb{E}X = \sum_{j=1}^{\infty} x_j \mathbb{P}\{X = x_j\}. \quad (5.5)$$

As before, the numbers x_1, x_2, \dots in formula (5.5) are the possible values of X . For example, if X attains values in $\{1, 4, 9\}$, then we may choose $x_1 = 1, x_2 = 4$, and $x_3 = 9$. But we could also take $x_1 = 4, x_2 = 9$, and $x_3 = 1$, and so on. Every time we get the same expected value.

Example 5.1.7. We start with an easy example that demonstrates how to compute the expected value in concrete cases. If the distribution of a random variable X is defined as $\mathbb{P}\{X = -1\} = 1/6, \mathbb{P}\{X = 0\} = 1/8, \mathbb{P}\{X = 1\} = 3/8$, and $\mathbb{P}\{X = 2\} = 1/3$, then its expected value equals

$$\begin{aligned} \mathbb{E}X &= (-1) \cdot \mathbb{P}\{X = -1\} + 0 \cdot \mathbb{P}\{X = 0\} + 1 \cdot \mathbb{P}\{X = 1\} + 2 \cdot \mathbb{P}\{X = 2\} \\ &= -\frac{1}{6} + \frac{3}{8} + \frac{2}{3} = \frac{7}{8}. \end{aligned}$$

Example 5.1.8. The next example shows that $\mathbb{E}X = \infty$ may occur even for quite natural random variables. Thus, let us come back to the model presented in Example 1.4.47. There we developed a strategy how to always win one dollar in a series of games. The basic idea was, after losing a game, next time one doubles the amount in the pool. As in Example 1.4.47, let $X(k)$ be the amount of money needed when winning for the first time in game k . We obtained

$$\mathbb{P}\{X = 2^k - 1\} = p(1-p)^{k-1}, \quad k = 1, 2, \dots$$

Recall that $0 < p < 1$ is the probability to win a single game. We ask for the expected value of money needed to apply this strategy. It follows that

$$\mathbb{E}X = \sum_{k=1}^{\infty} (2^k - 1) \mathbb{P}\{X = 2^k - 1\} = p \sum_{k=1}^{\infty} (2^k - 1)(1-p)^{k-1}. \quad (5.6)$$

If the game is fair, that is, if $p = 1/2$, then this leads to

$$\mathbb{E}X = \sum_{k=1}^{\infty} \frac{2^k - 1}{2^k} = \infty,$$

because of $(2^k - 1)/2^k \rightarrow 1$ as $k \rightarrow \infty$. This yields $\mathbb{E}X = \infty$ for all¹ $p \leq 1/2$.

Let us sum up: if $p \leq 1/2$ (which is the case in all provided games), the obtained result tells us that the average amount of money needed, to use this strategy, is arbitrarily large. The owners of gambling casinos know this strategy as well. Therefore, they limit the possible amount of money in the pool. For example, if the largest possible stakes is N dollars, then the strategy breaks down as soon as one loses n games for some n with $2^n > N$. And, as our calculations show, on average this always happens.

Remark 5.1.9. If $p > 1/2$, then the average amount of money needed is finite, and it can be calculated by

$$\begin{aligned} \mathbb{E}X &= p \sum_{k=1}^{\infty} (2^k - 1)(1-p)^{k-1} = 2p \sum_{k=0}^{\infty} (2 - 2p)^k - p \sum_{k=0}^{\infty} (1-p)^k \\ &= \frac{2p}{1 - (2 - 2p)} - \frac{p}{1 - (1-p)} = \frac{2p}{2p - 1} - 1 = \frac{1}{2p - 1}. \end{aligned}$$

5.1.2 Expected value of certain discrete random variables

The aim of this section is to compute the expected value of the most interesting discrete random variables. We start with uniformly distributed ones.

Proposition 5.1.10. *Let X be uniformly distributed on the set $\{x_1, \dots, x_N\}$ of real numbers. Then it follows that*

$$\mathbb{E}X = \frac{1}{N} \sum_{j=1}^N x_j. \quad (5.7)$$

That is, $\mathbb{E}X$ is the arithmetic mean of the x_j s.

Proof. This is an immediate consequence of $\mathbb{P}\{X = x_j\} = 1/N$, implying

$$\mathbb{E}X = \sum_{j=1}^N x_j \cdot \mathbb{P}\{X = x_j\} = \sum_{j=1}^N x_j \cdot \frac{1}{N}. \quad \square$$

Remark 5.1.11. For general discrete random variables X with values x_1, x_2, \dots , their expected value $\mathbb{E}X$ may be regarded as a weighted (the weights are the p_j s) mean of the x_j s.

¹ If $p \leq 1/2$ then $1 - p \geq 1/2$, hence the sum in eq. (5.6) becomes bigger and, therefore, it also diverges.

Example 5.1.12. Let X be uniformly distributed on $\{1, \dots, 6\}$. Then X is a model for rolling a fair die. Its expected value is, as is well known,

$$\mathbb{E}X = \frac{1 + \dots + 6}{6} = \frac{21}{6} = \frac{7}{2}.$$

Next we determine the expected value of a binomial distributed random variable.

Proposition 5.1.13. *Let X be binomial distributed with parameters n and p . Then we get*

$$\mathbb{E}X = np. \quad (5.8)$$

Proof. The possible values of X are $0, \dots, n$. Thus, it follows that

$$\begin{aligned} \mathbb{E}X &= \sum_{k=0}^n k \cdot \mathbb{P}\{X = k\} = \sum_{k=1}^n k \cdot \binom{n}{k} p^k (1-p)^{n-k} \\ &= \sum_{k=1}^n \frac{n!}{(k-1)!(n-k)!} p^k (1-p)^{n-k} \\ &= np \sum_{k=1}^n \frac{(n-1)!}{(k-1)!(n-k)!} p^{k-1} (1-p)^{n-k}. \end{aligned}$$

Shifting the index from $k-1$ to k in the last sum implies

$$\begin{aligned} \mathbb{E}X &= np \sum_{k=0}^{n-1} \frac{(n-1)!}{k!(n-1-k)!} p^k (1-p)^{n-1-k} \\ &= np \sum_{k=0}^{n-1} \binom{n-1}{k} p^k (1-p)^{n-1-k} = np [p + (1-p)]^{n-1} = np. \end{aligned}$$

This completes the proof. \square

Remark 5.1.14. The previous result tells us the following: if we perform n independent trials of an experiment with success probability $0 \leq p \leq 1$, then on average we will observe np successes.

Example 5.1.15. One kilogram of a radioactive material consists of N atoms. The atoms decay independently of each other and, moreover, the lifetime of each of the atoms is exponentially distributed with some parameter $\lambda > 0$. We ask for the time $T_0 > 0$, at which, on average, half of the atoms are decayed; T_0 is usually called radioactive half-life.

Answer: If $T > 0$, then the probability that a single atom decays before time T is given by

$$p(T) = E_\lambda([0, T]) = 1 - e^{-\lambda T}.$$

Since the atoms decay independently, the number of atoms decaying before time T is $B_{N,p(T)}$ -distributed. Therefore, by Proposition 5.1.13, the expected value of decayed atoms equals $N \cdot p(T) = N(1 - e^{-\lambda T})$. Hence, T_0 has to satisfy

$$N(1 - e^{-\lambda T_0}) = \frac{N}{2},$$

leading to $T_0 = \ln 2/\lambda$. Conversely, if we know T_0 and want to determine λ , then $\lambda = \ln 2/T_0$. Consequently, the probability that a single atom decays before time $T > 0$ can also be described by

$$E_\lambda([0, T]) = 1 - e^{-T \ln 2/T_0} = 1 - 2^{-T/T_0}.$$

Next, we determine the expected value of *Poisson distributed* random variables.

Proposition 5.1.16. *For some $\lambda > 0$, let X be distributed according to Pois_λ . Then it follows that $\mathbb{E}X = \lambda$.*

Proof. The possible values of X are $0, 1, \dots$. Hence, the expected value is given by

$$\mathbb{E}X = \sum_{k=0}^{\infty} k \cdot \mathbb{P}\{X = k\} = \sum_{k=1}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda} = \lambda \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} e^{-\lambda},$$

which transforms by a shift of the index to

$$\lambda \left[\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \right] e^{-\lambda} = \lambda [e^\lambda] e^{-\lambda} = \lambda,$$

proving the assertion. □

Interpretation: Proposition 5.1.16 explains the role of the parameter λ in the definition of the Poisson distribution. Whenever certain numbers are Poisson distributed, then $\lambda > 0$ is the average of the observed values. For example, if the number of accidents per week is known to be Pois_λ -distributed, then the parameter λ is determined by the average number of accidents per week in the past. Or, as we already mentioned in Example 4.6.6, the number of raisins in a piece of ρ pounds of dough is $\text{Pois}_{\lambda\rho}$ -distributed, where λ is the proportion of raisins per pound dough, hence $\lambda\rho$ is the average number of raisins per ρ pounds.

Example 5.1.17. Let us once more take a look at Example 4.6. There we considered light bulbs with E_λ -distributed lifetime. Every time a bulb burned out, we replaced it by a new one of the same type. It turned out that the number of necessary replacements until time $T > 0$ was $\text{Pois}_{\lambda T}$ -distributed. Consequently, by Proposition 5.1.16, on average, until time T we have to change the light bulbs λT times.

Finally, we compute the expected value of a negative binomial distributed random variable. According to Definition 1.4.49, a random variable X is $B_{n,p}^-$ -distributed if

$$\mathbb{P}\{X = k\} = \binom{k-1}{k-n} p^n (1-p)^{k-n}, \quad k = n, n+1, \dots$$

or, equivalently, if

$$\mathbb{P}\{X = k+n\} = \binom{-n}{k} p^n (p-1)^k, \quad k = 0, 1, \dots \quad (5.9)$$

Proposition 5.1.18. *Suppose X is $B_{n,p}^-$ -distributed for some $n \geq 1$ and $p \in (0, 1)$. Then*

$$\mathbb{E}X = \frac{n}{p}.$$

Proof. Using eq. (5.9), the expected value of X is computed as

$$\begin{aligned} \mathbb{E}X &= \sum_{k=n}^{\infty} k \mathbb{P}\{X = k\} = \sum_{k=0}^{\infty} (k+n) \mathbb{P}\{X = k+n\} \\ &= p^n \sum_{k=1}^{\infty} k \binom{-n}{k} (p-1)^k + n p^n \sum_{k=0}^{\infty} \binom{-n}{k} (p-1)^k. \end{aligned} \quad (5.10)$$

To evaluate the two sums in eq. (5.10), we use Proposition A.5.2, which asserts

$$\frac{1}{(1+x)^n} = \sum_{k=0}^{\infty} \binom{-n}{k} x^k \quad (5.11)$$

for $|x| < 1$. Applying this with $x = p-1$ (recall $0 < p < 1$) yields

$$n p^n \sum_{k=0}^{\infty} \binom{-n}{k} (p-1)^k = n p^n \frac{1}{(1+(p-1))^n} = n. \quad (5.12)$$

Next we differentiate eq. (5.11) with respect to x and obtain

$$\frac{-n}{(1+x)^{n+1}} = \sum_{k=1}^{\infty} k \binom{-n}{k} x^{k-1},$$

which, multiplying both sides by x , gives

$$\frac{-nx}{(1+x)^{n+1}} = \sum_{k=1}^{\infty} k \binom{-n}{k} x^k. \quad (5.13)$$

Letting $x = p-1$ in eq. (5.13), the first sum in eq. (5.10) becomes

$$p^n \sum_{k=1}^{\infty} k \binom{-n}{k} (p-1)^k = p^n \frac{-n(p-1)}{(1+(p-1))^{n+1}} = \frac{n(1-p)}{p}. \quad (5.14)$$

Finally, we combine eqs. (5.10), (5.12), and (5.14) to obtain

$$\mathbb{E}X = \frac{n(1-p)}{p} + n = \frac{n}{p},$$

as claimed. \square

Remark 5.1.19. Proposition 5.1.18 asserts that, on average, the n th success occurs in trial n/p . For example, rolling a die, on average, the first appearance of number “6” will be in trial 6, the second one in trial 12, and so on.

Corollary 5.1.20. *If X is geometrically distributed with parameter p , then*

$$\mathbb{E}X = \frac{1}{p}. \quad (5.15)$$

Proof. Recall that $G_p = B_{1,p}^-$, hence X is $B_{1,p}^-$ -distributed, and $\mathbb{E}X = \frac{1}{p}$ by Proposition 5.1.18. \square

Because of its beauty, let us give another, more direct proof of Corollary 5.1.20.

Suppose X is G_p -distributed. Then we write

$$\begin{aligned} \mathbb{E}X &= p \sum_{k=1}^{\infty} k(1-p)^{k-1} = p \sum_{k=0}^{\infty} (k+1)(1-p)^k \\ &= (1-p) \sum_{k=0}^{\infty} k p (1-p)^{k-1} + p \sum_{k=0}^{\infty} (1-p)^k = (1-p) \mathbb{E}X + 1. \end{aligned}$$

Solving this equation with respect to $\mathbb{E}X$ proves eq. (5.15), as asserted. Observe that this alternative proof is based upon the knowledge of $\mathbb{E}X < \infty$. Otherwise, we could not solve the equation for $\mathbb{E}X$. But, because of $0 < p < 1$, this fact is an easy consequence of

$$\mathbb{E}X = p \sum_{k=1}^{\infty} k(1-p)^{k-1} < \infty.$$

Summary: Let X be a discrete random variable with values x_1, x_2, \dots . Then

$$\mathbb{E}X = \sum_{j=1}^{\infty} x_j \mathbb{P}\{X = x_j\} \quad \text{provided that} \quad \mathbb{E}|X| = \sum_{j=1}^{\infty} |x_j| \mathbb{P}\{X = x_j\} < \infty.$$

5.1.3 Expected value of continuous random variables

Let X be a continuous random variable with distribution density p , that is, if $t \in \mathbb{R}$, then

$$\mathbb{P}\{X \leq t\} = \int_{-\infty}^t p(x) dx.$$

How to define $\mathbb{E}X$ in this case?

To answer this question, let us present formula (5.3) in an equivalent way. Suppose X maps Ω into a set $D \subset \mathbb{R}$, which is either finite or countably infinite. Let $p : \mathbb{R} \rightarrow [0, 1]$ be the probability mass function of X introduced in eq. (3.4). Then the expected value of X may also be written as

$$\mathbb{E}X = \sum_{x \in \mathbb{R}} x p(x).$$

In this form, the preceding formula suggests that in the continuous case the sum should be replaced by an integral. This can be made more precise by approximating continuous random variables by discrete ones. But this is only a heuristic explanation; for a precise approach, deeper convergence theorems for random variables are needed. Therefore, we do not give more details here, we simply replace sums by integrals.

Doing so, for continuous random variables the following approach for the definition of $\mathbb{E}X$ might be taken. If $p : \mathbb{R} \rightarrow [0, \infty)$ is the distribution density of X , set

$$\mathbb{E}X := \int_{-\infty}^{\infty} x p(x) dx. \quad (5.16)$$

However, here we have a similar problem as in the discrete case, namely that the integral in eq. (5.16) need not exist. Therefore, let us give a short digression about the integrability of real functions.

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a function such that for all $a < b$ the integral $\int_a^b f(x) dx$ is a well-defined real number. Then

$$\int_{-\infty}^{\infty} f(x) dx := \lim_{\substack{a \rightarrow -\infty \\ b \rightarrow \infty}} \int_a^b f(x) dx, \quad (5.17)$$

provided *both* limits on the right-hand side of eq. (5.17) exist. In this case we call f **integrable** (in the Riemann sense) on \mathbb{R} .

Remark 5.1.21. Let us point out that the numbers a and b in eq. (5.17) tend independently to $-\infty$ and ∞ , respectively. For instance, as the example $f(x) = x$ shows, it does not suffice that $\lim_{b \rightarrow \infty} \int_{-b}^b f(x) dx$ exists. Another way to express the existence of the integral in eq. (5.17) is as follows: the two limits

$$\lim_{a \rightarrow -\infty} \int_a^0 f(x) dx \quad \text{and} \quad \lim_{b \rightarrow \infty} \int_0^b f(x) dx \quad (5.18)$$

have to exist, and if this is so, the integral $\int_{-\infty}^{\infty} f(x) dx$ is defined as the sum of the two limits in eq. (5.18).

If $f(x) \geq 0$, $x \in \mathbb{R}$, then the limit $\lim_{\substack{a \rightarrow -\infty \\ b \rightarrow \infty}} \int_a^b f(x) dx$ always exists in a generalized sense, that is, it may be finite (then f is integrable) or infinite, then this is expressed by $\int_{-\infty}^{\infty} f(x) dx = \infty$.

If $\int_{-\infty}^{\infty} |f(x)| dx < \infty$, then f is said to be **absolutely integrable**, and as in the case of infinite series, absolutely integrable function are integrable. Note that $x \mapsto \sin x/x$ is integrable, but not absolutely integrable.

After this preparation, we come back to the definition of the expected value for continuous random variables.

Definition 5.1.22. Let X be a random variable with distribution density p . If $p(x) = 0$ for $x < 0$, or, equivalently, $\mathbb{P}\{X \geq 0\} = 1$, then the **expected value** of X is defined by

$$\mathbb{E}X := \int_0^{\infty} x p(x) dx. \quad (5.19)$$

Observe that under these conditions on p or X , we have $x p(x) \geq 0$. Therefore, the integral in eq. (5.19) is always well defined, but might be infinite. In this case we write $\mathbb{E}X = \infty$.

Let us turn now to the case of \mathbb{R} -valued random variables. The following example shows that the integral in eq. (5.16) may not exist, hence, in general, without an additional assumption the expected value cannot be defined by eq. (5.16).

Example 5.1.23. A random variable X is supposed to possess the density (check that this is indeed a density function)

$$p(x) = \begin{cases} 0 & \text{if } -1 < x < 1, \\ \frac{1}{2x^2} & \text{if } |x| \geq 1. \end{cases}$$

If we try to evaluate $\mathbb{E}X$ by virtue of eq. (5.16), then, because of

$$\begin{aligned} \int_{-\infty}^{\infty} x p(x) dx &= \lim_{\substack{a \rightarrow -\infty \\ b \rightarrow \infty}} \int_{-a}^b x p(x) dx = \frac{1}{2} \left[\lim_{b \rightarrow \infty} \int_1^b \frac{dx}{x} + \lim_{a \rightarrow \infty} \int_{-a}^{-1} \frac{dx}{x} \right] \\ &= \frac{1}{2} \left[\lim_{b \rightarrow \infty} \int_1^b \frac{dx}{x} - \lim_{a \rightarrow \infty} \int_1^a \frac{dx}{x} \right] = \infty - \infty, \end{aligned}$$

we observe an undetermined expression. Thus, there is no meaningful way to introduce an expected value for X .

We enforce the existence of the integral by the following condition.

Definition 5.1.24. Let X be a (real-valued) random variable with distribution density p . We say the **expected value of X exists**, provided p satisfies the following integrability condition:

$$\mathbb{E}|X| := \int_{-\infty}^{\infty} |x| p(x) dx < \infty. \quad (5.20)$$

Remark 5.1.25. At this point, it is not clear that the right-hand integral in (5.20) is indeed the expected value of $|X|$. This will follow later on by Proposition 5.1.38. Nevertheless, we use this notation before giving a proof.

Condition (5.20) says nothing but that $f(x) := xp(x)$ is absolutely integrable. Hence, as mentioned above, f is integrable, and the integral in the following definition is well defined.

Definition 5.1.26. Suppose condition (5.20) is satisfied. Then the **expected value** of X is defined by

$$\mathbb{E}X := \int_{-\infty}^{\infty} x p(x) dx.$$

Summary: Let X be a continuous random variable with density $p : \mathbb{R} \rightarrow \mathbb{R}$. Then

$$\mathbb{E}X = \int_{-\infty}^{\infty} x p(x) dx \quad \text{provided that} \quad \mathbb{E}|X| = \int_{-\infty}^{\infty} |x| p(x) dx < \infty.$$

5.1.4 Expected value of certain continuous random variables

We start with computing the expected value of a *uniformly distributed* (continuous) random variable.

Proposition 5.1.27. Let X be uniformly distributed on the finite interval $I = [a, \beta]$. Then

$$\mathbb{E}X = \frac{a + \beta}{2},$$

that is, the expected value is the midpoint of the interval I .

Proof. The distribution density of X is the function p defined as $p(x) = (\beta - a)^{-1}$ if $x \in I$, and $p(x) = 0$ if $x \notin I$. Of course, X possesses an expected value,² which can be evaluated by

² The product $|x|p(x)$ is bounded and nonzero only on a finite interval.

$$\mathbb{E}X = \int_{-\infty}^{\infty} xp(x) dx = \int_{\alpha}^{\beta} \frac{x}{\beta - \alpha} dx = \frac{1}{\beta - \alpha} \left[\frac{x^2}{2} \right]_{\alpha}^{\beta} = \frac{1}{2} \cdot \frac{\beta^2 - \alpha^2}{\beta - \alpha} = \frac{\alpha + \beta}{2}.$$

This proves the proposition. \square

Next we determine the expected value of a *gamma distributed* random variable.

Proposition 5.1.28. *Suppose X is $\Gamma_{\alpha, \beta}$ -distributed with $\alpha, \beta > 0$. Then its expected value is*

$$\mathbb{E}X = \alpha\beta.$$

Proof. Because of $\mathbb{P}\{X \geq 0\} = 1$, its expected value is well defined and computed by

$$\begin{aligned} \mathbb{E}X &= \int_0^{\infty} xp(x) dx = \frac{1}{\alpha^{\beta} \Gamma(\beta)} \int_0^{\infty} x \cdot x^{\beta-1} e^{-x/\alpha} dx \\ &= \frac{1}{\alpha^{\beta} \Gamma(\beta)} \int_0^{\infty} x^{\beta} e^{-x/\alpha} dx. \end{aligned} \quad (5.21)$$

The change of variables $u := x/\alpha$ transforms eq. (5.21) into

$$\mathbb{E}X = \frac{\alpha^{\beta+1}}{\alpha^{\beta} \Gamma(\beta)} \int_0^{\infty} u^{\beta} e^{-u} du = \frac{\alpha^{\beta+1}}{\alpha^{\beta} \Gamma(\beta)} \cdot \Gamma(\beta + 1) = \alpha\beta,$$

where we used eq. (1.50) in the last step. This completes the proof. \square

Corollary 5.1.29. *Let X be E_{λ} -distributed for a certain $\lambda > 0$. Then*

$$\mathbb{E}X = \frac{1}{\lambda}.$$

Proof. Note that $E_{\lambda} = \Gamma_{\lambda^{-1}, 1}$. \square

Example 5.1.30. The lifetime of a special type of light bulbs is exponentially distributed. Suppose the average lifetime constitutes 100 units of time. This implies $\lambda = 1/100$, hence, if X describes the lifetime, then

$$\mathbb{P}\{X \leq t\} = 1 - e^{-t/100}, \quad t \geq 0.$$

For example, the probability that the light bulb burns longer than 200 time units equals

$$\mathbb{P}\{X \geq 200\} = e^{-200/100} = e^{-2} = 0.135335 \dots$$

Remark 5.1.31. If we evaluate in the previous example

$$\mathbb{P}\{X \geq \mathbb{E}X\} = \mathbb{P}\{X \geq 100\} = e^{-1},$$

then we see that in general $\mathbb{P}\{X \geq \mathbb{E}X\} \neq 1/2$. Thus, in this case, the expected value is different from the **median** of X defined as a real number M satisfying $\mathbb{P}\{X \geq M\} \geq 1/2$ and $\mathbb{P}\{X \leq M\} \geq 1/2$. In particular, if F_X satisfies the condition of Proposition 4.4.7, then the median is uniquely determined by $M = F_X^{-1}(1/2)$, i. e., by $\mathbb{P}\{X \leq M\} = 1/2$. It is easy to see that the above phenomenon appears for all exponentially distributed random variables. Indeed, if X is E_λ -distributed, then $M = \ln 2/\lambda$ while, as we saw, $\mathbb{E}X = 1/\lambda$, compare Figure 5.1.

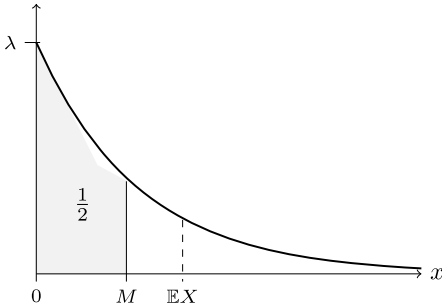


Figure 5.1: The expected value $\mathbb{E}X = 1/\lambda$ and the median $M = \ln 2/\lambda$ of an E_λ -distributed random variable X .

Corollary 5.1.32. *If X is χ_n^2 -distributed, then*

$$\mathbb{E}X = n.$$

Proof. Since $\chi_n^2 = \Gamma_{2,n/2}$, by Proposition 5.1.28 it follows that $\mathbb{E}X = 2 \cdot n/2 = n$. \square

Which expected value does a beta distributed random variable possess? The next proposition answers this question.

Proposition 5.1.33. *Let X be $\mathcal{B}_{\alpha,\beta}$ -distributed for certain $\alpha, \beta > 0$. Then*

$$\mathbb{E}X = \frac{\alpha}{\alpha + \beta}.$$

Proof. Using eq. (1.63), from eq. (5.19) we obtain, as asserted,

$$\begin{aligned} \mathbb{E}X &= \frac{1}{B(\alpha, \beta)} \int_0^1 x \cdot x^{\alpha-1} (1-x)^{\beta-1} dx \\ &= \frac{1}{B(\alpha, \beta)} \int_0^1 x^\alpha (1-x)^{\beta-1} dx = \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} = \frac{\alpha}{\alpha + \beta}. \end{aligned} \quad \square$$

Example 5.1.34. Suppose we choose independently n numbers x_1, \dots, x_n uniformly distributed on $[0, 1]$ and order them by their size. Then we get the order statistics

$0 \leq x_1^* \leq \dots \leq x_n^* \leq 1$. According to Example 3.7.11, if $1 \leq k \leq n$, then the number x_k^* is $\mathcal{B}_{k, n-k+1}$ -distributed. Thus Proposition 5.1.33 implies that the average value of x_k^* , that is, of the k th largest number, equals

$$\frac{k}{k + (n - k + 1)} = \frac{k}{n + 1}.$$

In particular, the expected value of the smallest number is $\frac{1}{n+1}$ while that of the largest one is $\frac{n}{n+1}$.

Does a Cauchy distributed random variable possess an expected value? Here we obtain the following.

Proposition 5.1.35. *If X is Cauchy distributed, then $\mathbb{E}X$ does not exist.*

Proof. First, observe that we may not use Definition 5.1.22. The distribution density of X is given by $p(x) = \frac{1}{\pi} \cdot \frac{1}{1+x^2}$, hence, it does not satisfy $p(x) = 0$ for $x < 0$. Consequently, we have to check whether condition (5.20) is satisfied. Here we get

$$\mathbb{E}|X| = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{|x|}{1+x^2} dx = \frac{2}{\pi} \int_0^{\infty} \frac{x}{1+x^2} dx = \frac{1}{\pi} [\ln(1+x^2)]_0^{\infty} = \infty.$$

Thus, $\mathbb{E}|X| = \infty$, that is, X does not possess an expected value. \square

Finally, we determine the expected value of normally distributed random variables.

Proposition 5.1.36. *If X is $\mathcal{N}(\mu, \sigma^2)$ -distributed, then*

$$\mathbb{E}X = \mu.$$

Proof. First, we check whether the expected value exists. The density of X is given by eq. (1.49), hence

$$\begin{aligned} \int_{-\infty}^{\infty} |x| p_{\mu, \sigma}(x) dx &= \frac{1}{\sqrt{2\pi} \sigma} \int_{-\infty}^{\infty} |x| e^{-(x-\mu)^2/2\sigma^2} dx = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} |\sqrt{2}\sigma u + \mu| e^{-u^2} du \\ &\leq \sigma \frac{2\sqrt{2}}{\sqrt{\pi}} \int_0^{\infty} u e^{-u^2} du + |\mu| \frac{2}{\sqrt{\pi}} \int_0^{\infty} e^{-u^2} du < \infty, \end{aligned}$$

where we used the well-known fact³ that for all $k \in \mathbb{N}_0$,

$$\int_0^{\infty} u^k e^{-u^2} du < \infty.$$

³ See either [Spi08] or use that for all $k \geq 1$ one has $\sup_{u>0} u^k e^{-u} < \infty$.

The expected value $\mathbb{E}X$ is now evaluated in a similar way as

$$\begin{aligned}\mathbb{E}X &= \int_{-\infty}^{\infty} x p_{\mu, \sigma}(x) dx = \frac{1}{\sqrt{2\pi} \sigma} \int_{-\infty}^{\infty} x e^{-(x-\mu)^2/2\sigma^2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\sigma v + \mu) e^{-v^2/2} dv \\ &= \sigma \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} v e^{-v^2/2} dv + \mu \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-v^2/2} dv.\end{aligned}\quad (5.22)$$

The function $f(v) := ve^{-v^2/2}$ is odd, that is, $f(-v) = -f(v)$, thus $\int_{-\infty}^{\infty} f(v) dv = 0$, and the first integral in eq. (5.22) vanishes. To compute the second integral, use Proposition 1.6.7 and obtain

$$\mu \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-v^2/2} dv = \mu \frac{1}{\sqrt{2\pi}} \sqrt{2\pi} = \mu.$$

This completes the proof. \square

Remark 5.1.37. Proposition 5.1.36 justifies the notation “expected value” for the parameter μ in the definition of the probability measure $\mathcal{N}(\mu, \sigma^2)$.

Summary: Let X be some random variable. Using the notation introduced in Remark 3.3.1, the following are valid:

1. X uniformly distributed on $\{x_1, \dots, x_N\} \Rightarrow \mathbb{E}X = \frac{x_1 + \dots + x_N}{N}$.
 2. $X \sim B_{n,p} \Rightarrow \mathbb{E}X = np$.
 3. $X \sim \text{Pois}_\lambda \Rightarrow \mathbb{E}X = \lambda$.
 4. $X \sim B_{n,p}^- \Rightarrow \mathbb{E}X = \frac{n}{p}$.
 5. $X \sim G_p \Rightarrow \mathbb{E}X = \frac{1}{p}$.
 6. X uniformly distributed on $[\alpha, \beta] \Rightarrow \mathbb{E}X = \frac{\alpha + \beta}{2}$.
 7. $X \sim \Gamma_{\alpha, \beta} \Rightarrow \mathbb{E}X = \alpha\beta$.
 8. $X \sim E_{\lambda, n} \Rightarrow \mathbb{E}X = \frac{n}{\lambda}$.
 9. $X \sim \chi_n^2 \Rightarrow \mathbb{E}X = n$.
 10. $X \sim \mathcal{B}_{\alpha, \beta} \Rightarrow \mathbb{E}X = \frac{\alpha}{\alpha + \beta}$.
 11. $X \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow \mathbb{E}X = \mu$.
 12. X Cauchy distributed $\Rightarrow \mathbb{E}X$ does not exist.
-

5.1.5 Properties of the expected value

In this section we summarize the main properties of the expected value. They are valid for both discrete and continuous random variables. But, unfortunately, within

the framework of this book it is not possible to prove most of them in full generality. To do so, one needs an integral (Lebesgue integral) $\int_{\Omega} f d\mathbb{P}$ of functions $f : \Omega \rightarrow \mathbb{R}$ for some probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Then $\mathbb{E}X = \int_{\Omega} X d\mathbb{P}$, and all subsequent properties of $X \mapsto \mathbb{E}X$ follow from those of the (Lebesgue) integral. We refer to [LG22] for a thorough presentation of this (wonderful) topic.

Proposition 5.1.38. *The expected value of random variables has the following properties:*

- (1) *The expected value of X only depends on its probability distribution \mathbb{P}_X , not on the way how X is defined. That is, if $X \stackrel{d}{=} Y$ for two random variables X and Y , then $\mathbb{E}X = \mathbb{E}Y$.*
- (2) *If X is constant with probability one, that is, there is some $c \in \mathbb{R}$ with $\mathbb{P}\{X = c\} = 1$, then $\mathbb{E}X = c$.*
- (3) *The expected value is linear: let X and Y be two random variables possessing an expected value and let $a, b \in \mathbb{R}$. Then $\mathbb{E}(aX + bY)$ exists as well and, moreover,*

$$\mathbb{E}(aX + bY) = a \mathbb{E}X + b \mathbb{E}Y.$$

- (4) *Suppose X is a discrete random variable with values x_1, x_2, \dots . Given a function f from \mathbb{R} to \mathbb{R} , the expected value $\mathbb{E}f(X)$ exists if and only if*

$$\sum_{i=1}^{\infty} |f(x_i)| \mathbb{P}\{X = x_i\} < \infty,$$

and, moreover, then

$$\mathbb{E}f(X) = \sum_{i=1}^{\infty} f(x_i) \mathbb{P}\{X = x_i\}. \quad (5.23)$$

- (5) *If X is continuous with density p , then for any measurable function $f : \mathbb{R} \rightarrow \mathbb{R}$ the expected value $\mathbb{E}f(X)$ exists if and only if*

$$\int_{-\infty}^{\infty} |f(x)| p(x) dx < \infty.$$

In this case it follows that

$$\mathbb{E}f(X) = \int_{-\infty}^{\infty} f(x) p(x) dx. \quad (5.24)$$

- (6) *For independent X and Y possessing an expected value, the expected value of their product⁴ XY exists as well and, moreover,*

$$\mathbb{E}[XY] = \mathbb{E}X \cdot \mathbb{E}Y.$$

⁴ Recall that $(XY)(\omega) = X(\omega) \cdot Y(\omega)$ for all $\omega \in \Omega$.

- (7) Write $X \leq Y$ provided that $X(\omega) \leq Y(\omega)$ for all $\omega \in \Omega$. If in this sense $|X| \leq Y$ for some Y with $\mathbb{E}Y < \infty$, then $\mathbb{E}|X| < \infty$ and, hence, $\mathbb{E}X$ exists.
- (8) Suppose $\mathbb{E}X$ and $\mathbb{E}Y$ exist. Then $X \leq Y$ implies $\mathbb{E}X \leq \mathbb{E}Y$. In particular, if $X \geq 0$, then $\mathbb{E}X \geq 0$.

Proof. We only prove properties (1), (2), (4), and (8). Some of the other properties are not difficult to verify in the case of discrete random variables, for example, (3), but because the proofs are incomplete, we do not present them here. We refer to [Bil12, Dur19] or [Kho07] for the proofs of the remaining properties.

We begin with the proof of (1). If X and Y are identically distributed, then either both are discrete or both are continuous. If they are discrete, and $\mathbb{P}_X(D) = 1$ for an at most countably infinite set D , then $X \stackrel{d}{=} Y$ implies $\mathbb{P}_Y(D) = 1$. Moreover, by the same argument $\mathbb{P}_X(\{x\}) = \mathbb{P}_Y(\{x\})$ for any $x \in D$. Hence, in view of Definition 5.1.2, $\mathbb{E}X$ exists if and only if $\mathbb{E}Y$ does. Moreover, if this is valid, then $\mathbb{E}X = \mathbb{E}Y$ by the same argument.

In the continuous case, we argue as follows. Let p be a density of X . Due to $X \stackrel{d}{=} Y$, it follows that

$$\int_{-\infty}^t p(x) \, dx = \mathbb{P}_X((-\infty, t]) = \mathbb{P}_Y((-\infty, t]), \quad t \in \mathbb{R}.$$

Thus, p is also a distribution density of Y and, consequently, in view of Definition 5.1.24, the expected value of X exists if and only if this is the case for Y . Moreover, by Definition 5.1.26, we get $\mathbb{E}X = \mathbb{E}Y$.

Next we show that (2) is valid. Thus, suppose $\mathbb{P}\{X = c\} = 1$ for some $c \in \mathbb{R}$. Then X is discrete with $\mathbb{P}_X(D) = 1$ where $D = \{c\}$, and by Definition 5.1.2 we obtain

$$\mathbb{E}X = c \cdot \mathbb{P}\{X = c\} = c \cdot 1 = c,$$

as asserted.

To prove (4), we assume that X has values in $D = \{x_1, x_2, \dots\}$. Then $Y = f(X)$ maps into $f(D) = \{y_1, y_2, \dots\}$. Given $j \in \mathbb{N}$, let $D_j = \{x_i : f(x_i) = y_j\}$. Thus,

$$\mathbb{P}\{Y = y_j\} = \mathbb{P}\{X \in D_j\} = \sum_{x_i \in D_j} \mathbb{P}\{X = x_i\}.$$

Consequently, since $D_j \cap D_{j'} = \emptyset$ if $j \neq j'$, due to $\bigcup_{j=1}^{\infty} D_j = D$, we get

$$\begin{aligned} \mathbb{E}|Y| &= \sum_{j=1}^{\infty} |y_j| \mathbb{P}\{Y = y_j\} = \sum_{j=1}^{\infty} \sum_{x_i \in D_j} |y_j| \mathbb{P}\{X = x_i\} \\ &= \sum_{j=1}^{\infty} \sum_{x_i \in D_j} |f(x_i)| \mathbb{P}\{X = x_i\} = \sum_{i=1}^{\infty} |f(x_i)| \mathbb{P}\{X = x_i\}. \end{aligned}$$

This proves the first part of (4). The second part follows by exactly the same arguments (replace $|y_j|$ by y_j). Therefore, we omit its proof.

We finally prove (8). To this end, we first show the second part, that is, $\mathbb{E}X \geq 0$ for $X \geq 0$. If X is discrete, then X attains values in D , where D consists only of nonnegative real numbers. Hence, $x_j \mathbb{P}\{X = x_j\} \geq 0$, which implies $\mathbb{E}X \geq 0$. If X is continuous, in view of $X \geq 0$, we may choose its density p such that $p(x) = 0$ if $x < 0$. Then $\mathbb{E}X = \int_0^\infty p(x) dx \geq 0$.

Suppose now $X \leq Y$. Setting $Z = Y - X$, from the first step we get $\mathbb{E}Z \geq 0$. But, property (3) implies $\mathbb{E}Z = \mathbb{E}Y - \mathbb{E}X$, from which we derive $\mathbb{E}X \leq \mathbb{E}Y$, as asserted. Note that by assumption $\mathbb{E}X$ and $\mathbb{E}Y$ are real numbers, so that $\mathbb{E}Y - \mathbb{E}X$ is not an undetermined expression. \square

Remark 5.1.39. Properties (4) and (5) of the previous proposition, applied with the function $f(x) = |x|$, lead to

$$\mathbb{E}|X| = \sum_{j=1}^{\infty} |x_j| \mathbb{P}\{X = x_j\} \quad \text{or} \quad \mathbb{E}|X| = \int_{-\infty}^{\infty} |x| p(x) dx,$$

as we already stated in conditions (5.4) and (5.20), respectively.

Corollary 5.1.40. *If $\mathbb{E}X$ exists, then shifting X by $\mu = \mathbb{E}X$, it becomes centralized. In other words, if $\mu = \mathbb{E}X$, then*

$$\mathbb{E}(X - \mu) = 0.$$

Proof. If $Y = X - \mu$, then properties (2) and (3) of Proposition 5.1.38 imply

$$\mathbb{E}Y = \mathbb{E}(X - \mu) = \mathbb{E}X - \mathbb{E}\mu = \mu - \mu = 0,$$

as asserted. \square

An important consequence of (8) in Proposition 5.1.38 reads as follows.

Corollary 5.1.41. *If $\mathbb{E}X$ exists, then*

$$|\mathbb{E}X| \leq \mathbb{E}|X|.$$

Proof. For all $\omega \in \Omega$, it follows that

$$-|X(\omega)| \leq X(\omega) \leq |X(\omega)|,$$

that is, we have $-|X| \leq X \leq |X|$. We apply now (3) and (8) of Proposition 5.1.38 and conclude that

$$-\mathbb{E}|X| = \mathbb{E}(-|X|) \leq \mathbb{E}X \leq \mathbb{E}|X|. \quad (5.25)$$

Since $|a| \leq c$ for $a, c \in \mathbb{R}$ is equivalent to $-c \leq a \leq c$, the desired estimate is a consequence of inequalities (5.25) with $a = \mathbb{E}X$ and $c = \mathbb{E}|X|$. \square

We now present some examples that show how Proposition 5.1.38 may be used to evaluate certain expected values.

Example 5.1.42. Suppose we roll n fair dice. Let S_n be the sum of the observed values. What is the expected value of S_n ?

Answer: If X_j denotes the value of the j th die, then X_1, \dots, X_n are uniformly distributed on $\{1, \dots, 6\}$ with $\mathbb{E}X_j = 7/2$ and, moreover, $S_n = X_1 + \dots + X_n$. Thus, property (3) lets us conclude that

$$\mathbb{E}S_n = \mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}X_1 + \dots + \mathbb{E}X_n = \frac{7n}{2}.$$

Example 5.1.43. In Example 4.1.7, we investigated the random walk of a particle on \mathbb{Z} . Each time it jumped with probability p either one step to the right or with probability $1-p$ one step to the left. There S_n denoted the position of the particle after n steps. What is the expected position after n steps?

Answer: We proved that $S_n = 2Y_n - n$ with a $B_{n,p}$ -distributed random variable Y_n . Proposition 5.1.13 implies $\mathbb{E}Y_n = np$, hence the linearity of the expected value leads to

$$\mathbb{E}S_n = 2\mathbb{E}Y_n - n = 2np - n = n(2p - 1). \quad (5.26)$$

For $p = 1/2$, we obtain the (not very surprising) result $\mathbb{E}S_n = 0$. But note that eq. (5.26) can also be proved directly by using $\mathbb{E}X_j = (-1)(1-p) + 1 \cdot p = 2p - 1$.

Remark 5.1.44. If we regard S_n as the position of a particle after n jumps, then since $2p - 1 > 0$ if $p > 1/2$ it follows that in this case the particle drifts on average to ∞ . On the contrary, if $p < 1/2$, the particle tends on average to $-\infty$.

On the other hand, if we interpret S_n as the loss or win after n games, we get the following conclusion: whenever one plays a series of games with success probability $p < 1/2$ (for example, roulette), in the long run one will lose on average an arbitrarily big amount of money.

The next example demonstrates how property (4) of Proposition 5.1.38 may be used.

Example 5.1.45. Let X be Pois_λ -distributed. Find $\mathbb{E}X^2$.

Solution: Property (4) of Proposition 5.1.38 implies

$$\mathbb{E}X^2 = \sum_{k=0}^{\infty} k^2 \mathbb{P}\{X = k\} = \sum_{k=1}^{\infty} k^2 \frac{\lambda^k}{k!} e^{-\lambda} = \lambda \sum_{k=1}^{\infty} k \frac{\lambda^{k-1}}{(k-1)!} e^{-\lambda}.$$

We shift the index of summation in the right-hand sum by 1 and get

$$\mathbb{E}X^2 = \lambda \sum_{k=0}^{\infty} (k+1) \frac{\lambda^k}{k!} e^{-\lambda} = \lambda \sum_{k=0}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda} + \lambda \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda}.$$

By Proposition 5.1.16, the first sum coincides with $\lambda \mathbb{E}X = \lambda^2$, while the second gives $\lambda \text{Pois}_\lambda(\mathbb{N}_0) = \lambda \cdot 1 = \lambda$. Adding both values leads to

$$\mathbb{E}X^2 = \lambda^2 + \lambda.$$

The next example rests upon an application of properties (3), (4), and (6) in Proposition 5.1.38.

Example 5.1.46. Compute $\mathbb{E}X^2$ for X being $B_{n,p}$ -distributed.

Solution: Let X_1, \dots, X_n be independent $B_{1,p}$ -distributed random variables. Then Corollary 4.6.2 asserts that $X = X_1 + \dots + X_n$ is $B_{n,p}$ -distributed. Therefore, it suffices to evaluate $\mathbb{E}X^2$ with $X = X_1 + \dots + X_n$. Thus, property (3) of Proposition 5.1.38 implies

$$\mathbb{E}X^2 = \mathbb{E}(X_1 + \dots + X_n)^2 = \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}X_i X_j.$$

If $i \neq j$, then X_i and X_j are independent, hence property (6) applies and yields

$$\mathbb{E}X_i X_j = \mathbb{E}X_i \cdot \mathbb{E}X_j = p \cdot p = p^2.$$

For $i = j$, property (4) gives

$$\mathbb{E}X_j^2 = 0^2 \cdot \mathbb{P}\{X_j = 0\} + 1^2 \cdot \mathbb{P}\{X_j = 1\} = p.$$

Combining both cases leads to

$$\mathbb{E}X^2 = \sum_{i \neq j} \mathbb{E}X_i \cdot \mathbb{E}X_j + \sum_{j=1}^n \mathbb{E}X_j^2 = n(n-1)p^2 + np = n^2 p^2 + np(1-p).$$

Example 5.1.47. Let X be G_p -distributed. Compute $\mathbb{E}X^2$.

Solution: We claim that

$$\mathbb{E}X^2 = \frac{2-p}{p^2}. \quad (5.27)$$

To prove this, let us start with

$$\mathbb{E}X^2 = \sum_{k=1}^{\infty} k^2 \mathbb{P}\{X = k\} = p \sum_{k=1}^{\infty} k^2 (1-p)^{k-1}. \quad (5.28)$$

We evaluate the right-hand sum by the following approach. If $|x| < 1$, then

$$\frac{1}{1-x} = \sum_{k=0}^{\infty} x^k.$$

Differentiating both sides of this equation leads to

$$\frac{1}{(1-x)^2} = \sum_{k=1}^{\infty} k x^{k-1}.$$

Next we multiply this equation by x and arrive at

$$\frac{x}{(1-x)^2} = \sum_{k=1}^{\infty} k x^k.$$

Another differentiation of both functions on $\{x \in \mathbb{R} : |x| < 1\}$ implies

$$\frac{1}{(1-x)^2} + \frac{2x}{(1-x)^3} = \sum_{k=1}^{\infty} k^2 x^{k-1}.$$

If we use the last equation with $x = 1 - p$, then, by eq. (5.28),

$$\mathbb{E}X^2 = p \left[\frac{1}{(1 - (1-p))^2} + \frac{2(1-p)}{(1 - (1-p))^3} \right] = \frac{2-p}{p^2},$$

as we claimed in eq. (5.27).

In the next example we use property (5) of Proposition 5.1.38.

Example 5.1.48. Let U be uniformly distributed on $[0, 1]$. Which expected value does \sqrt{U} possess?

Solution: By property (5), it follows that

$$\mathbb{E}\sqrt{U} = \int_{-\infty}^{\infty} \sqrt{x} \cdot \mathbb{1}_{[0,1]}(x) \, dx = \int_0^1 \sqrt{x} \, dx = \frac{2}{3} [x^{3/2}]_0^1 = \frac{2}{3}.$$

Another approach is as follows. Because of

$$F_{\sqrt{U}}(t) = \mathbb{P}\{\sqrt{U} \leq t\} = \mathbb{P}\{U \leq t^2\} = t^2$$

for $0 \leq t \leq 1$, the density q of \sqrt{U} is given by $q(x) = 2x$, $0 \leq x \leq 1$, and $q(x) = 0$ otherwise. Thus,

$$\mathbb{E}\sqrt{U} = \int_0^1 x \cdot 2x \, dx = 2 \left[\frac{x^3}{3} \right]_0^1 = \frac{2}{3}.$$

Let us present now an interesting example called **Coupon collector's problem**. It was first mentioned in 1708 by A. De Moivre. We formulate it in a present-day language.

Example 5.1.49. A company produces cornflakes. Each pack contains a picture. We assume that there are n different pictures and that they are equally likely. That is, when buying a pack, the probability to get a certain fixed picture is $1/n$. How many packs of cornflakes have to be bought on average before one gets all possible n pictures?

An equivalent formulation of the problem is as follows. In an urn there are n balls numbered from 1 to n . One chooses balls out of the urn with replacement. How many balls have to be chosen on average before one observes all n numbers?

Answer: Assume we already have k different pictures for some $k = 0, 1, \dots, n-1$. Let X_k be the number of necessary purchases to obtain a new picture, that is, to get one which we do not have. Since each pack contains a picture,

$$\mathbb{P}\{X_0 = 1\} = 1.$$

If $k \geq 1$, then there are still $n - k$ pictures that one does not possess. Hence, X_k is geometrically distributed with success probability $p_k = (n - k)/n$. If $S_n = X_0 + \dots + X_{n-1}$, then S_n is the totality of necessary purchases. By Corollary 5.1.20, we obtain

$$\mathbb{E}X_k = \frac{1}{p_k} = \frac{n}{n - k}, \quad k = 0, \dots, n - 1.$$

Note that $\mathbb{E}X_0 = 1$, thus the previous formula also holds in this case. Then the linearity of the expected value implies

$$\begin{aligned} \mathbb{E}S_n &= 1 + \mathbb{E}X_1 + \dots + \mathbb{E}X_{n-1} = 1 + \frac{1}{p_1} + \dots + \frac{1}{p_{n-1}} \\ &= 1 + \frac{n}{n-1} + \frac{n}{n-2} + \dots + \frac{n}{1} = n \sum_{k=1}^n \frac{1}{k}. \end{aligned}$$

Consequently, on average, one needs $n \sum_{k=1}^n \frac{1}{k}$ purchases to obtain a complete collection of all pictures.

For example, if $n = 50$, on average, we have to buy 225 packs; if $n = 100$, then 519; for $n = 200$, on average, there are 1176 purchases necessary; if $n = 300$, then 1885; if $n = 400$, we have to buy 2628 packs; and, finally, if $n = 500$, we need to buy 3397.

Remark 5.1.50. As $n \rightarrow \infty$, the harmonic series $\sum_{k=1}^n \frac{1}{k}$ behaves like $\ln n$. More precisely (cf. [Lag13] or [Spi08], Problem 12, Chapter 22)

$$\lim_{n \rightarrow \infty} \left[\sum_{k=1}^n \frac{1}{k} - \ln n \right] = \gamma, \quad (5.29)$$

where γ denotes Euler's constant, which is approximately 0.57721. Therefore, for large n , the average number of necessary purchases is approximately $n[\ln n + \gamma]$. For example, if $n = 300$, then the approximative value is 1884.29, leading also to 1885 necessary purchases.

Summary: The most important properties of the expected value are (we assume that all expected values are well defined):

1. For all X, Y and $a, b \in \mathbb{R}$, it follows that $\mathbb{E}[aX + bY] = a \mathbb{E}X + b \mathbb{E}Y$.
2. If $f: \mathbb{R} \rightarrow \mathbb{R}$, then $\mathbb{E}[f(X)] = \sum_{k=1}^{\infty} f(x_k) \mathbb{P}\{X = x_k\}$, resp. $\mathbb{E}[f(X)] = \int_{-\infty}^{\infty} f(x) p(x) dx$.
3. For all random variables X , it follows that $|\mathbb{E}X| \leq \mathbb{E}|X|$.
4. If X and Y are independent, then $\mathbb{E}[XY] = \mathbb{E}X \cdot \mathbb{E}Y$.

5.2 Variance

5.2.1 Higher moments of random variables

Definition 5.2.1. Let $n \geq 1$ be some integer. A random variable X possesses an **n th moment**, provided that $\mathbb{E}|X|^n < \infty$. We also say X has a finite **absolute n th moment**. If this is so, then $\mathbb{E}X^n$ exists, and it is called the **n th moment** of X .

Remark 5.2.2. Because of $|X|^n = |X^n|$, the assumption $\mathbb{E}|X|^n < \infty$ implies the existence of the n th moment $\mathbb{E}X^n$.

Note that a random variable X has a first moment if and only if the expected value of X exists, cf. conditions (5.4) and (5.20). Moreover, then the first moment coincides with $\mathbb{E}X$.

Proposition 5.2.3. Let X be either a discrete random variable with values in $\{x_1, x_2, \dots\}$ and with $p_j = \mathbb{P}\{X = x_j\}$, or let X be continuous with density p . If $n \geq 1$, then

$$\mathbb{E}|X|^n = \sum_{j=1}^{\infty} |x_j|^n \cdot p_j \quad \text{or} \quad \mathbb{E}|X|^n = \int_{-\infty}^{\infty} |x|^n p(x) dx. \quad (5.30)$$

Consequently, X possesses a finite absolute n th moment if and only if either the sum or the integral in eq. (5.30) are finite. If this is satisfied, then these moments are given by

$$\mathbb{E}X^n = \sum_{j=1}^{\infty} x_j^n \cdot p_j \quad \text{or} \quad \mathbb{E}X^n = \int_{-\infty}^{\infty} x^n p(x) dx.$$

Proof. Apply properties (4) and (5) in Proposition 5.1.38 with $f(x) = |x|^n$ or with $f(x) = x^n$, respectively. \square

Example 5.2.4. Let U be uniformly distributed on $[0,1]$. Then

$$\mathbb{E}|U|^n = \mathbb{E}U^n = \int_0^1 x^n dx = \frac{1}{n+1}.$$

For the subsequent investigations, we need the following elementary lemma.

Lemma 5.2.5. *If $0 < \alpha < \beta$, then for all $x \geq 0$,*

$$x^\alpha \leq x^\beta + 1.$$

Proof. If $0 \leq x \leq 1$, from $x^\beta \geq 0$ it follows that

$$x^\alpha \leq 1 \leq x^\beta + 1,$$

and the inequality is valid.

If $x > 1$, then $\alpha < \beta$ implies $x^\alpha < x^\beta$, hence also for those x we arrive at

$$x^\alpha < x^\beta \leq x^\beta + 1,$$

which proves the lemma. □

Proposition 5.2.6. *Suppose a random variable X has a finite absolute n th moment. Then X possesses all m th moments with $m < n$.*

Proof. Suppose $\mathbb{E}|X|^n < \infty$ and choose an $m < n$. For a fixed $\omega \in \Omega$, we apply Lemma 5.2.5 with $\alpha = m$, $\beta = n$, and $x = |X(\omega)|$. Doing so, we obtain

$$|X(\omega)|^m \leq |X(\omega)|^n + 1,$$

and this being true for all $\omega \in \Omega$ implies $|X|^m \leq |X|^n + 1$. Hence, property (7) of Proposition 5.1.38 yields

$$\mathbb{E}|X|^m \leq \mathbb{E}(|X|^n + 1) = \mathbb{E}|X|^n + 1 < \infty.$$

Consequently, as asserted, X possesses also an absolute m th moment. □

Remark 5.2.7. There exist much stronger estimates between different absolute moments of X . For example, Lyapunov's inequality, a special case of Jensen's inequality, asserts that for any $0 < \alpha \leq \beta$,

$$[\mathbb{E}|X|^\alpha]^{1/\alpha} \leq [\mathbb{E}|X|^\beta]^{1/\beta}.$$

The case $n = 2$ and $m = 1$ in Proposition 5.2.6 is of special interest. Here we get the following useful result.

Corollary 5.2.8. *If X possesses a finite second moment, then $\mathbb{E}|X| < \infty$, that is, its expected value exists.*

Let us state another important consequence of Proposition 5.2.6.

Corollary 5.2.9. *Suppose X has a finite absolute n th moment. Then for any $b \in \mathbb{R}$, we also have $\mathbb{E}|X + b|^n < \infty$.*

Proof. An application of the binomial theorem (Proposition A.3.8) implies

$$|X + b|^n \leq (|X| + |b|)^n = \sum_{k=0}^n \binom{n}{k} |X|^k |b|^{n-k}.$$

Hence, using properties (3) and (7) of Proposition 5.1.38, we obtain

$$\mathbb{E}|X + b|^n \leq \sum_{k=0}^n \binom{n}{k} |b|^{n-k} \mathbb{E}|X|^k < \infty.$$

Note that Proposition 5.2.6 implies $\mathbb{E}|X|^k < \infty$ for all $k < n$. This ends the proof. \square

Example 5.2.10. Let X be $\Gamma_{\alpha, \beta}$ -distributed with parameters $\alpha, \beta > 0$. Which moments does X possess, and how can they be computed?

Answer: In view of $X \geq 0$, it suffices to investigate $\mathbb{E}X^n$. For all $n \geq 1$, it follows that

$$\begin{aligned} \mathbb{E}X^n &= \frac{1}{\alpha^\beta \Gamma(\beta)} \int_0^\infty x^{n+\beta-1} e^{-x/\alpha} dx = \frac{\alpha^{n+\beta}}{\alpha^\beta \Gamma(\beta)} \int_0^\infty y^{n+\beta-1} e^{-y} dy \\ &= \alpha^n \frac{\Gamma(\beta + n)}{\Gamma(\beta)} = \alpha^n (\beta + n - 1)(\beta + n - 2) \cdots (\beta + 1)\beta. \end{aligned}$$

In particular, X has moments of any order $n \geq 1$.

In the case of an E_λ -distributed random variable X , we have $\alpha = 1/\lambda$ and $\beta = 1$, hence

$$\mathbb{E}X^n = \frac{n!}{\lambda^n}.$$

Example 5.2.11. Suppose a random variable is t_n -distributed. Which moments does X possess?

Answer: We already know that a t_1 -distributed random variable does not possess a first moment. Recall that X is t_1 -distributed if it is Cauchy distributed. And in Proposition 5.1.35 we proved $\mathbb{E}|X| = \infty$ for Cauchy distributed random variables.

But what can be said if $n \geq 2$?

According to Definition 4.6, the random variable X has the density p with

$$p(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2}, \quad x \in \mathbb{R}.$$

If $m \in \mathbb{N}$, then

$$\mathbb{E}|X|^m = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \int_{-\infty}^{\infty} |x|^m \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} dx.$$

Hence, X has an m th moment if and only if the integral

$$\int_{-\infty}^{\infty} |x|^m \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} dx = 2 \int_0^{\infty} x^m \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} dx \quad (5.31)$$

is finite. Note that

$$\lim_{x \rightarrow \infty} x^{n+1} \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} = \lim_{x \rightarrow \infty} \left(x^{-2} + \frac{1}{n}\right)^{-n/2-1/2} = n^{n/2+1/2},$$

thus, there are constants $0 < c_1 < c_2$ (depending on n , but not on x) such that

$$\frac{c_1}{x^{n-m+1}} \leq x^m \left(1 + \frac{x^2}{n}\right)^{-n/2-1/2} \leq \frac{c_2}{x^{n-m+1}} \quad (5.32)$$

for large x , that is, if $x > x_0$ for a suitable $x_0 \in \mathbb{R}$.

Recall that $\int_1^{\infty} x^{-a} dx < \infty$ if and only if $a > 1$. Having this in mind, in view of eq. (5.31) and by the estimates in (5.32), we get $\mathbb{E}|X|^m < \infty$ if and only if $n - m + 1 > 1$, that is, if and only if $m < n$.

Summing up, a t_n -distributed random variable has moments of order $1, \dots, n - 1$, but no moments of order greater than or equal to n .

Finally, let us investigate the moments of normally distributed random variables.

Example 5.2.12. How do we calculate $\mathbb{E}X^n$ for an $\mathcal{N}(0, 1)$ -distributed random variable?

Answer: Well-known properties of the exponential function imply

$$\mathbb{E}|X|^n = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} |x|^n e^{-x^2/2} dx = \frac{2}{\sqrt{2\pi}} \int_0^{\infty} x^n e^{-x^2/2} dx < \infty$$

for all $n \in \mathbb{N}$. Thus, a normally distributed random variable possesses moments of any order. These moments are evaluated by

$$\mathbb{E}X^n = \int_{-\infty}^{\infty} x^n p_{0,1}(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^n e^{-x^2/2} dx.$$

If n is an odd integer, then $x \mapsto x^n e^{-x^2/2}$ is an odd function, hence $\mathbb{E}X^n = 0$ for odd integers n .

Therefore, it suffices to investigate even $n = 2m$ with $m \in \mathbb{N}$. Here we get

$$\mathbb{E}X^{2m} = 2 \cdot \frac{1}{\sqrt{2\pi}} \int_0^{\infty} x^{2m} e^{-x^2/2} dx,$$

which, by the change of variables $y := x^2/2$, thus $x = \sqrt{2y}$ with $dx = \frac{1}{\sqrt{2}} y^{-1/2} dy$, transforms into

$$\mathbb{E}X^{2m} = \frac{1}{\sqrt{\pi}} 2^m \int_0^{\infty} y^{m-1/2} e^{-y} dy = \frac{2^m}{\sqrt{\pi}} \Gamma\left(m + \frac{1}{2}\right).$$

Since $\Gamma(1/2) = \sqrt{\pi}$ and by an application of eq. (1.50), we finally obtain

$$\begin{aligned} \mathbb{E}X^{2m} &= \frac{2^m}{\sqrt{\pi}} \left(m - \frac{1}{2}\right) \Gamma\left(m - \frac{1}{2}\right) = \frac{2^m}{\sqrt{\pi}} \left(m - \frac{1}{2}\right) \left(m - \frac{3}{2}\right) \Gamma\left(m - \frac{3}{2}\right) \\ &= \frac{2^m \Gamma(1/2) \cdot 1/2 \cdot 3/2 \cdots (m-1/2)}{\sqrt{\pi}} \\ &= (2m-1)(2m-3) \cdots 3 \cdot 1 = (2m-1)!! . \end{aligned}$$

Summary: A random variable X possess an n th moment if $\mathbb{E}|X|^n < \infty$. Then its n th moment is defined by

$$\mathbb{E}X^n = \sum_{j=1}^{\infty} x_j^n \mathbb{P}\{X = x_j\} \quad \text{or} \quad \mathbb{E}X^n = \int_{-\infty}^{\infty} x^n p(x) dx ,$$

respectively. Here the x_j s are in the discrete case the possible values of X or p denotes in the continuous case its density.

If X possesses an n th moment, then this also so for all moments of order $m \leq n$. In particular, for each random variable with second moment its expected value exists.

5.2.2 Variance of random variables

Let X be a random variable with finite second moment. As we saw in Corollary 5.2.8, then its expected value $\mu := \mathbb{E}X$ exists. Furthermore, letting $b = -\mu$, by Corollary 5.2.9, we also have $\mathbb{E}|X - \mu|^2 < \infty$. After this preparation, we can introduce the variance of a random variable.

Definition 5.2.13. Let X be a random variable possessing a finite second moment. If $\mu := \mathbb{E}X$, then its **variance** is defined as

$$\mathbb{V}X := \mathbb{E}|X - \mu|^2 = \mathbb{E}|X - \mathbb{E}X|^2 .$$

Interpretation: The expected value μ of a random variable is its main characteristic. It tells us around which value the observations of X have to be expected. But it does not tell us how far away from μ these observations will be on average. Are they concentrated around μ or are they widely dispersed? This behavior is described by the variance. It is defined as the average quadratic distance of X to its mean value. If $\mathbb{V}X$ is small, then we will observe realizations of X quite near to its mean. Otherwise, if $\mathbb{V}X$ is large, then it is very likely to observe values of X far away from its expected value.

How do we evaluate the variance in concrete cases? We answer this question for discrete and continuous random variables separately.

Proposition 5.2.14. *Let X be a random variable with finite second moment and let $\mu \in \mathbb{R}$ be its expected value. Then it follows that*

$$\mathbb{V}X = \sum_{j=1}^{\infty} (x_j - \mu)^2 \cdot p_j \quad \text{and} \quad \mathbb{V}X = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) \, dx \quad (5.33)$$

in the discrete and continuous case, respectively. Hereby, x_1, x_2, \dots are the possible values of X and $p_j = \mathbb{P}\{X = x_j\}$ in the discrete case, while p denotes the density of X in the continuous case.

Proof. The assertion follows directly by an application of properties (4) and (5) of Proposition 5.1.38 to $f(x) = (x - \mu)^2$. \square

Before we present concrete examples, let us state and prove certain properties of the variance, which will simplify the calculations later on.

Proposition 5.2.15. *Assume X and Y are random variables with finite second moment. Then the following are valid:*

(i) *We have*

$$\mathbb{V}X = \mathbb{E}X^2 - (\mathbb{E}X)^2. \quad (5.34)$$

(ii) *If $\mathbb{P}\{X = c\} = 1$ for some $c \in \mathbb{R}$, then⁵ $\mathbb{V}X = 0$.*

(iii) *For $a, b \in \mathbb{R}$ follows that*

$$\mathbb{V}(aX + b) = a^2 \mathbb{V}X.$$

(iv) *In the case of independent X and Y , one has*

$$\mathbb{V}(X + Y) = \mathbb{V}X + \mathbb{V}Y.$$

Proof. Let us begin with the proof of (i). With $\mu = \mathbb{E}X$, we obtain

$$\begin{aligned} \mathbb{V}X &= \mathbb{E}(X - \mu)^2 = \mathbb{E}[X^2 - 2\mu X + \mu^2] = \mathbb{E}X^2 - 2\mu \mathbb{E}X + \mu^2 \\ &= \mathbb{E}X^2 - 2\mu^2 + \mu^2 = \mathbb{E}X^2 - \mu^2. \end{aligned}$$

This proves (i).

To verify (ii), we use property (2) in Proposition 5.1.38. Then $\mu = \mathbb{E}X = c$, hence $\mathbb{P}\{X - \mu = 0\} = 1$. Another application of property (2) leads to

$$\mathbb{V}X = \mathbb{E}(X - \mu)^2 = 0$$

as asserted.

⁵ The converse implication is also true. If $\mathbb{V}X = 0$, then X is constant with probability 1.

Next we prove (iii). If $a, b \in \mathbb{R}$, then $\mathbb{E}(aX + b) = a\mathbb{E}X + b$ by the linearity of the expected value. Consequently,

$$\mathbb{V}(aX + b) = \mathbb{E}[aX + b - (a\mathbb{E}X + b)]^2 = a^2 \mathbb{E}(X - \mathbb{E}X)^2 = a^2 \mathbb{V}X.$$

Thus (iii) is valid.

To prove (iv), observe that, if $\mu := \mathbb{E}X$ and $\nu := \mathbb{E}Y$, then $\mathbb{E}(X + Y) = \mu + \nu$, and hence

$$\begin{aligned} \mathbb{V}(X + Y) &= \mathbb{E}[(X - \mu) + (Y - \nu)]^2 \\ &= \mathbb{E}(X - \mu)^2 + 2\mathbb{E}[(X - \mu)(Y - \nu)] + \mathbb{E}(Y - \nu)^2 \\ &= \mathbb{V}X + 2\mathbb{E}[(X - \mu)(Y - \nu)] + \mathbb{V}Y. \end{aligned} \quad (5.35)$$

By Proposition 4.1.9, the independence of X and Y implies that of $X - \mu$ and $Y - \nu$. Therefore, from property (6) in Proposition 5.1.38, we derive

$$\mathbb{E}[(X - \mu)(Y - \nu)] = \mathbb{E}(X - \mu) \cdot \mathbb{E}(Y - \nu) = (\mathbb{E}X - \mu) \cdot (\mathbb{E}Y - \nu) = 0 \cdot 0 = 0.$$

Plugging this into eq. (5.35) completes the proof of (iv). \square

Summary: Let X be a random variable with finite second moment. If μ denotes the expected value of X , its variance is defined by $\mathbb{V}X = \mathbb{E}(X - \mu)^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$. The variance may be evaluated in the discrete, respectively continuous case as follows:

$$\mathbb{V}X = \sum_{k=1}^{\infty} (x_k - \mu)^2 \mathbb{P}\{X = x_k\} \quad \text{and} \quad \mathbb{V}X = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx.$$

Here the x_k s and p are the possible values of X and its density, respectively. The basic properties are $\mathbb{V}(aX + b) = a^2\mathbb{V}X$ and $\mathbb{V}[X + Y] = \mathbb{V}X + \mathbb{V}Y$ for *independent* X and Y .

5.2.3 Variance of certain random variables

Our first objective is to describe the variance of a random variable uniformly distributed on a finite set.

Proposition 5.2.16. *If X is uniformly distributed on $\{x_1, \dots, x_N\}$, then*

$$\mathbb{V}X = \frac{1}{N} \sum_{j=1}^N (x_j - \mu)^2,$$

where μ is given by $\mu = \frac{1}{N} \sum_{j=1}^N x_j$.

Proof. Because of $p_j = \frac{1}{N}$, $1 \leq j \leq N$, this is a direct consequence of eq. (5.33). Recall that μ was computed in eq. (5.7). \square

Example 5.2.17. Suppose X is uniformly distributed on $\{1, \dots, 6\}$. Then $\mathbb{E}X = 7/2$, and we get

$$\begin{aligned}\mathbb{V}X &= \frac{(1 - \frac{7}{2})^2 + (2 - \frac{7}{2})^2 + (3 - \frac{7}{2})^2 + (4 - \frac{7}{2})^2 + (5 - \frac{7}{2})^2 + (6 - \frac{7}{2})^2}{6} \\ &= \frac{\frac{25}{4} + \frac{9}{4} + \frac{1}{4} + \frac{1}{4} + \frac{9}{4} + \frac{25}{4}}{6} = \frac{35}{12}.\end{aligned}$$

Thus, when rolling a die once, the variance is given by $\frac{35}{12}$.

Now assume that we roll the die n times. Let X_1, \dots, X_n be the results of the single rolls. The X_j s are independent, hence, if $S_n = X_1 + \dots + X_n$ denotes the sum of the n trials, then, by (iv) in Proposition 5.2.15, it follows that

$$\mathbb{V}S_n = \mathbb{V}(X_1 + \dots + X_n) = \mathbb{V}X_1 + \dots + \mathbb{V}X_n = \frac{35n}{12}.$$

The next proposition examines the variance of binomial distributed random variables.

Proposition 5.2.18. *If X is $B_{n,p}$ -distributed, then*

$$\mathbb{V}X = np(1-p). \quad (5.36)$$

Proof. Let X be $B_{n,p}$ -distributed. In Example 5.1.47, we found $\mathbb{E}X^2 = n^2p^2 + np(1-p)$. Moreover, $\mathbb{E}X = np$ by Proposition 5.1.13. Thus, using formula (5.34) we derive

$$\mathbb{V}X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = n^2p^2 + np(1-p) - (np)^2 = np(1-p),$$

as asserted. □

Remark 5.2.19. An alternative proof of Proposition 5.2.18 is as follows: if $n = 1$, then we get

$$\mathbb{V}X = (1-p)(0-p)^2 + p(1-p)^2 = p(1-p)[p + (1-p)] = p(1-p).$$

Thus, if $X = X_1 + \dots + X_n$ with X_j s independent $B_{1,p}$ -distributed, then, on the one hand, X is $B_{n,p}$ -distributed and, on the other hand, we obtain

$$\mathbb{V}X = \mathbb{V}(X_1 + \dots + X_n) = \mathbb{V}X_1 + \dots + \mathbb{V}X_n = n\mathbb{V}X_1 = np(1-p).$$

Corollary 5.2.20. *Binomial distributed random variables have maximal variance (with n fixed) if $p = 1/2$.*

Proof. The function $p \mapsto np(1-p)$ becomes maximal for $p = \frac{1}{2}$. In the extreme cases $p = 0$ and $p = 1$, the variance is zero. □

Corollary 5.2.21. Let $(S_n)_{n \geq 0}$ be a random walk, that is, $S_0 = 0$ and $S_n = X_1 + \cdots + X_n$ if $n \geq 1$ where the X_i s are independent and attaining the values -1 and 1 with probabilities $1 - p$ and p , respectively. Then, if $n \geq 1$, it follows that

$$\mathbb{V}S_n = 4np(1 - p). \quad (5.37)$$

Proof. As shown in Example 4.1.7, the random variables $Y_n = (S_n + n)/2$ are $B_{n,p}$ -distributed. Applying (iii) of Proposition 5.2.15, together with eq. (5.36), implies

$$\mathbb{V}S_n = \mathbb{V}(2Y_n - n) = 4\mathbb{V}Y_n = 4np(1 - p),$$

as asserted. □

Next we determine the variance of Poisson distributed random variables.

Proposition 5.2.22. Let X be Pois_λ -distributed for some $\lambda > 0$. Then

$$\mathbb{V}X = \lambda.$$

Proof. In Example 5.1.45, we computed $\mathbb{E}X^2 = \lambda^2 + \lambda$. Furthermore, by Proposition 5.1.16, we know that $\mathbb{E}X = \lambda$. Thus, by eq. (5.34), we obtain, as asserted,

$$\mathbb{V}X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \lambda^2 + \lambda - \lambda^2 = \lambda. \quad \square$$

Next, we compute the variance of a geometrically distributed random variable.

Proposition 5.2.23. Let X be G_p -distributed for some $0 < p < 1$. Then its variance equals

$$\mathbb{V}X = \frac{1 - p}{p^2}.$$

Proof. In Example 5.1.47, we found $\mathbb{E}X^2 = \frac{2-p}{p^2}$, and by eq. (5.15) we have $\mathbb{E}X = \frac{1}{p}$. Consequently, formula (5.34) implies

$$\mathbb{V}X = \frac{2-p}{p^2} - \left(\frac{1}{p}\right)^2 = \frac{1-p}{p^2},$$

as asserted. □

Corollary 5.2.24. If X is $B_{n,p}^-$ -distributed, then

$$\mathbb{V}X = n \frac{1-p}{p^2}$$

Proof. Let X_1, \dots, X_n be independent G_p -distributed random variables. By Corollary 4.6, their sum $X = X_1 + \cdots + X_n$ is $B_{n,p}^-$ -distributed, hence property (iv) in Proposition 5.2.15

lets us conclude that

$$\mathbb{V}X = \mathbb{V}(X_1 + \cdots + X_n) = \mathbb{V}X_1 + \cdots + \mathbb{V}X_n = n \mathbb{V}X_1 = n \frac{1-p}{p^2}. \quad \square$$

Interpretation: The smaller the p , the bigger the variance of a geometrically or negative binomial distributed random variable (for n fixed). This is not surprising, because the smaller the p , the larger the expected value, and so the values of X may be very far from $1/p$ (success is very unlikely).

We consider now variances of continuous random variables. Let us begin with uniformly distributed ones.

Proposition 5.2.25. *Let X be uniformly distributed on an interval $[a, \beta]$. Then it follows that*

$$\mathbb{V}X = \frac{(\beta - a)^2}{12}.$$

Proof. We know by Proposition 5.1.27 that $\mathbb{E}X = (\alpha + \beta)/2$. In order to apply formula (5.34), we still have to compute the second moment $\mathbb{E}X^2$. Here we obtain

$$\mathbb{E}X^2 = \frac{1}{\beta - \alpha} \int_{\alpha}^{\beta} x^2 dx = \frac{1}{3} \cdot \frac{\beta^3 - \alpha^3}{\beta - \alpha} = \frac{\beta^2 + \alpha\beta + \alpha^2}{3}.$$

Consequently, formula (5.34) lets us conclude that

$$\begin{aligned} \mathbb{V}X &= \mathbb{E}X^2 - (\mathbb{E}X)^2 = \frac{\beta^2 + \alpha\beta + \alpha^2}{3} - \left(\frac{\alpha + \beta}{2}\right)^2 \\ &= \frac{\beta^2 + \alpha\beta + \alpha^2}{3} - \frac{\alpha^2 + 2\alpha\beta + \beta^2}{4} \\ &= \frac{\alpha^2 - 2\alpha\beta + \beta^2}{12} = \frac{(\beta - \alpha)^2}{12}. \end{aligned}$$

This completes the proof. □

In the case of gamma distributed random variables, the following is valid.

Proposition 5.2.26. *If X is $\Gamma_{\alpha, \beta}$ -distributed, then*

$$\mathbb{V}X = \alpha^2 \beta.$$

Proof. Recall that $\mathbb{E}X = \alpha\beta$ by Proposition 5.1.28. Furthermore, in Example 5.2.10 we evaluated $\mathbb{E}X^n$ for a gamma distributed X . Taking $n = 2$ yields

$$\mathbb{E}X^2 = \alpha^2 (\beta + 1) \beta,$$

and, hence, by eq. (5.34),

$$\mathbb{V}X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \alpha^2(\beta + 1)\beta - (\alpha\beta)^2 = \alpha^2\beta,$$

as asserted. \square

Corollary 5.2.27. *If X is E_λ -distributed, then*

$$\mathbb{V}X = \frac{1}{\lambda^2}.$$

Proof. Because of $E_\lambda = \Gamma_{\frac{1}{\lambda}, 1}$, this directly follows from Proposition 5.2.26. \square

Corollary 5.2.28. *For a χ_n^2 -distributed X , it holds that*

$$\mathbb{V}X = 2n.$$

Proof. Let us give two alternative proofs of the assertion. The first uses Proposition 5.2.26 and $\chi_n^2 = \Gamma_{2, \frac{n}{2}}$.

The second proof is longer, but maybe more interesting. Let X_1, \dots, X_n be independent $\mathcal{N}(0, 1)$ -distributed random variables. Proposition 4.6.10 implies that $X_1^2 + \dots + X_n^2$ is χ_n^2 -distributed, thus property (iv) of Proposition 5.2.15 applies and leads to

$$\mathbb{V}X = \mathbb{V}X_1^2 + \dots + \mathbb{V}X_n^2 = n\mathbb{V}X_1^2.$$

In Example 5.2.12, we evaluated the moments of an $\mathcal{N}(0, 1)$ -distributed random variable. In particular, $\mathbb{E}X_1^2 = 1$ and $\mathbb{E}(X_1^2)^2 = \mathbb{E}X_1^4 = 3!! = 3$, hence

$$\mathbb{V}X = n\mathbb{V}X_1^2 = n(\mathbb{E}X_1^4 - (\mathbb{E}X_1^2)^2) = (3 - 1)n = 2n,$$

as claimed. \square

Finally, we determine the variance of a normal random variable.

Proposition 5.2.29. *If X is $\mathcal{N}(\mu, \sigma^2)$ -distributed, then it follows that*

$$\mathbb{V}X = \sigma^2.$$

Proof. Of course, this could be proven by computing the integral

$$\mathbb{V}X = \int_{-\infty}^{\infty} (x - \mu)^2 p_{\mu, \sigma}(x) dx.$$

We prefer a different approach that avoids the calculation of integrals. Because of Proposition 4.2.3, the random variable X may be represented as $X = \sigma X_0 + \mu$ for a standard normal X_0 . Applying (iii) in Proposition 5.2.15 gives

$$\mathbb{V}X = \sigma^2 \mathbb{V}X_0. \quad (5.38)$$

But $\mathbb{E}X_0 = 0$, and by Example 5.2.12 we have $\mathbb{E}X_0^2 = 1$, thus

$$\mathbb{V}X_0 = 1 - 0 = 1.$$

Plugging this into eq. (5.38) proves $\mathbb{V}X = \sigma^2$. \square

Remark 5.2.30. The previous result explains why the parameter σ^2 is called the “variance” of an $\mathcal{N}(\mu, \sigma^2)$ -distributed random variable. Recall that the other parameter μ denotes its expected value. Moreover, it shows that the smaller the $\sigma^2 > 0$, the more are the values of an $\mathcal{N}(\mu, \sigma^2)$ -distributed random variable concentrated around the expected value μ ; see Figure 5.2.

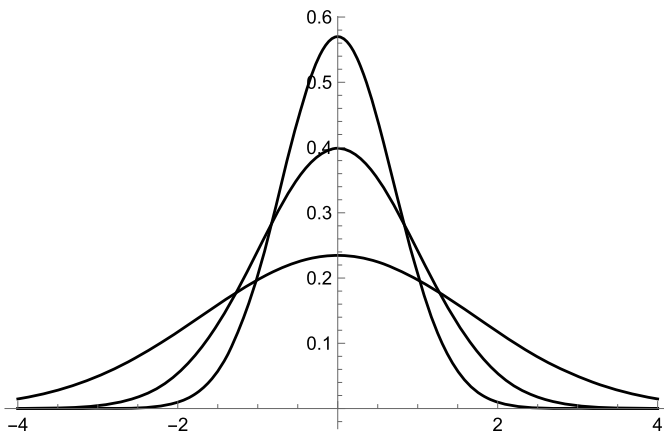


Figure 5.2: Densities of normal distributions with mean value 0 and variances (from bottom to top) $\sigma^2 = 1.7$, $\sigma^2 = 1$, and $\sigma^2 = 0.7$. The larger the σ^2 , the more likely are events away from zero.

Summary: Let X be some random variable. Then we have

1. X uniformly distributed on $\{x_1, \dots, x_N\} \Rightarrow \mathbb{V}X = \frac{(x_1 - \mu)^2 + \dots + (x_N - \mu)^2}{N}$ with $\mu = \mathbb{E}X$.
 2. $X \sim B_{n,p} \Rightarrow \mathbb{V}X = np(1-p)$.
 3. $X \sim \text{Pois}_\lambda \Rightarrow \mathbb{V}X = \lambda$.
 4. $X \sim B_{n,p}^- \Rightarrow \mathbb{V}X = \frac{n(1-p)}{p^2}$.
 5. X uniformly distributed on $[\alpha, \beta] \Rightarrow \mathbb{V}X = \frac{(\beta - \alpha)^2}{12}$.
 6. $X \sim \Gamma_{\alpha,\beta} \Rightarrow \mathbb{V}X = \alpha^2 \beta$.
 7. $X \sim E_{\lambda,n} \Rightarrow \mathbb{V}X = \frac{n}{\lambda^2}$.
 8. $X \sim \chi_n^2 \Rightarrow \mathbb{V}X = 2n$.
 9. $X \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow \mathbb{V}X = \sigma^2$.
-

5.3 Covariance and correlation

5.3.1 Covariance

Suppose we know or conjecture that two given random variables X and Y are dependent. The aim of this section is to introduce a quantity that measures their degree of dependence. Such a quantity should tell us whether the random variables are strongly or only weakly dependent. Furthermore, we want to know what kind of dependence we observe. Do larger values of X trigger larger values of Y or is it the other way round? To illustrate these questions, let us come back to the experiment presented in Example 2.2.5.

Example 5.3.1. In an urn are n balls labeled with “0” and another n balls labeled with “1.” Choose two balls out of the urn *without* replacement. Let X be the number appearing on the first ball and Y that on the second. Then X and Y are dependent (check this), but it is intuitively clear that if n becomes larger, then their dependence diminishes. We ask for a quantity that tells us their degree of dependence. This measure should decrease as n increases and it should tend to zero as $n \rightarrow \infty$.

Moreover, if $X = 1$ occurred, then there remained in the urn more balls with “0” than with “1,” and the probability of the event $Y = 0$ increases. Thus, larger values of X make smaller values of Y more likely.

Before we are able to introduce such a “measure of dependence,” we need some preparation.

Proposition 5.3.2. *If two random variables X and Y possess a finite second moment, then the expected value of their product XY exists.*

Proof. We use the elementary estimate $|ab| \leq \frac{a^2+b^2}{2}$ valid for $a, b \in \mathbb{R}$. Thus, if $\omega \in \Omega$, then

$$|X(\omega)Y(\omega)| \leq \frac{X(\omega)^2}{2} + \frac{Y(\omega)^2}{2},$$

that is, we have

$$|XY| \leq \frac{X^2}{2} + \frac{Y^2}{2}. \quad (5.39)$$

By assumption,

$$\mathbb{E} \left[\frac{X^2}{2} + \frac{Y^2}{2} \right] = \frac{1}{2} [\mathbb{E}X^2 + \mathbb{E}Y^2] < \infty,$$

consequently, because of estimate (5.39), property (7) in Proposition 5.1.38 applies and tells us that $\mathbb{E}|XY| < \infty$. Thus, $\mathbb{E}[XY]$ exists as asserted. \square

How do we compute $\mathbb{E}[XY]$ for given X and Y ? In Section 4.5 we observed that the distribution of $X + Y$ does not only depend on the distributions of X and Y . We have to know their *joint distribution*, that is, the distribution of the vector (X, Y) . And the same is true for products and the expected value of the product.

Example 5.3.3. Let us again investigate the random variables X, Y, X' , and Y' introduced in Example 3.5.8. Recall that they satisfied

$$\begin{aligned}\mathbb{P}\{X = 0, Y = 0\} &= \frac{1}{6}, & \mathbb{P}\{X = 0, Y = 1\} &= \frac{1}{3}, \\ \mathbb{P}\{X = 1, Y = 0\} &= \frac{1}{3}, & \mathbb{P}\{X = 1, Y = 1\} &= \frac{1}{6}, \\ \mathbb{P}\{X' = 0, Y' = 0\} &= \frac{1}{4}, & \mathbb{P}\{X' = 0, Y' = 1\} &= \frac{1}{4}, \\ \mathbb{P}\{X' = 1, Y' = 0\} &= \frac{1}{4}, & \mathbb{P}\{X' = 1, Y' = 1\} &= \frac{1}{4}.\end{aligned}$$

Then $\mathbb{P}_X = \mathbb{P}_{X'}$ and $\mathbb{P}_Y = \mathbb{P}_{Y'}$, but

$$\begin{aligned}\mathbb{E}[XY] &= \frac{1}{6}(0 \cdot 0) + \frac{1}{3}(1 \cdot 0) + \frac{1}{3}(0 \cdot 1) + \frac{1}{6}(1 \cdot 1) = \frac{1}{6} \quad \text{and} \\ \mathbb{E}[X'Y'] &= \frac{1}{4}(0 \cdot 0) + \frac{1}{4}(1 \cdot 0) + \frac{1}{4}(0 \cdot 1) + \frac{1}{4}(1 \cdot 1) = \frac{1}{4}.\end{aligned}$$

This example tells us that we have to know the joint distribution in order to compute $\mathbb{E}[XY]$. The knowledge of the marginal distributions does not suffice.

To evaluate $\mathbb{E}[XY]$, we need the following two-dimensional generalization of formulas (5.23) and (5.24).

Proposition 5.3.4. *Let X and Y be two random variables and let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be some function.*

1. *Suppose X and Y are discrete with values in $\{x_1, x_2, \dots\}$ and in $\{y_1, y_2, \dots\}$. Set $p_{ij} = \mathbb{P}\{X = x_i, Y = y_j\}$. If*

$$\mathbb{E}|f(X, Y)| = \sum_{i,j=1}^{\infty} |f(x_i, y_j)| p_{ij} < \infty, \quad (5.40)$$

then $\mathbb{E}f(X, Y)$ exists and can be computed by

$$\mathbb{E}f(X, Y) = \sum_{i,j=1}^{\infty} f(x_i, y_j) p_{ij}.$$

2. *Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be continuous.⁶ If $p : \mathbb{R}^2 \rightarrow \mathbb{R}$ is the joint density of (X, Y) (as intro-*

⁶ In fact, we need only a measurability in the sense of Definition 4.1.1, but this time for functions f from \mathbb{R}^2 to \mathbb{R} . For our purposes “continuity” of f suffices.

duced in Definition 3.5.15), then

$$\mathbb{E}|f(X, Y)| = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |f(x, y)| p(x, y) dx dy < \infty \quad (5.41)$$

implies the existence of $\mathbb{E}f(X, Y)$, which can be evaluated by

$$\mathbb{E}f(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) p(x, y) dx dy. \quad (5.42)$$

Remark 5.3.5. The previous formulas extend easily to higher dimensions. That is, if $\vec{X} = (X_1, \dots, X_n)$ is an n -dimensional random vector with (joint) distribution density $p : \mathbb{R}^n \rightarrow \mathbb{R}$, then for continuous⁷ $f : \mathbb{R}^n \rightarrow \mathbb{R}$ one has

$$\mathbb{E}f(\vec{X}) = \mathbb{E}f(X_1, \dots, X_n) = \int_{\mathbb{R}^n} f(x_1, \dots, x_n) p(x_1, \dots, x_n) dx_n \cdots dx_1$$

provided the integral exists. The case of discrete X_1, \dots, X_n is treated in a similar way. If \vec{X} maps into the finite or countably infinite set $D \subset \mathbb{R}^n$, then

$$\mathbb{E}f(\vec{X}) = \mathbb{E}f(X_1, \dots, X_n) = \sum_{x \in D} f(x) \mathbb{P}\{\vec{X} = x\}.$$

If we apply Proposition 5.3.4 with $f : (x, y) \mapsto x \cdot y$, then we obtain the following formulas for the evaluation of $\mathbb{E}[XY]$. Hereby, we assume that conditions (5.40) or (5.41) are satisfied.

Corollary 5.3.6. *In the notation of Proposition 5.3.4, the following are valid:*

$$\mathbb{E}[XY] = \sum_{i,j=1}^{\infty} (x_i \cdot y_j) p_{ij} \quad \text{and} \quad \mathbb{E}[XY] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x \cdot y) p(x, y) dx dy$$

in the discrete and continuous case, respectively.

After all these preparations, we are now in a position to introduce the covariance of two random variables.

Definition 5.3.7. Let X and Y be two random variables with finite second moments. Setting $\mu = \mathbb{E}X$ and $\nu = \mathbb{E}Y$, the **covariance** of X and Y is defined as

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu)(Y - \nu)].$$

⁷ Compare with the footnote for $n = 2$ in part 2 of Proposition 5.3.4.

Remark 5.3.8. Apply Corollary 5.2.9 and Proposition 5.3.2 to see that the covariance is well defined for random variables with a finite second moment. Furthermore, in view of Proposition 5.3.4, the covariance may be computed as

$$\text{Cov}(X, Y) = \sum_{i,j=1}^{\infty} (x_i - \mu)(y_j - \nu) p_{ij}$$

in the discrete case (recall that $p_{ij} = \mathbb{P}\{X = x_i, Y = y_j\}$), and as

$$\text{Cov}(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu)(y - \nu) p(x, y) dx dy$$

in the continuous case.

Example 5.3.9. Let us once more consider the random variables X, Y, X' , and Y' in Example 3.5.8 or Example 5.3.3, respectively. Each of the four random variables has the expected value $1/2$. Therefore, we obtain

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{1}{6} \left(0 - \frac{1}{2}\right) \cdot \left(0 - \frac{1}{2}\right) + \frac{1}{3} \left(1 - \frac{1}{2}\right) \cdot \left(0 - \frac{1}{2}\right) \\ &\quad + \frac{1}{3} \left(0 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) + \frac{1}{6} \left(1 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) = -\frac{1}{12}, \end{aligned}$$

while

$$\begin{aligned} \text{Cov}(X', Y') &= \frac{1}{4} \left(0 - \frac{1}{2}\right) \cdot \left(0 - \frac{1}{2}\right) + \frac{1}{4} \left(1 - \frac{1}{2}\right) \cdot \left(0 - \frac{1}{2}\right) \\ &\quad + \frac{1}{4} \left(0 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) + \frac{1}{4} \left(1 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) = 0. \end{aligned}$$

The following proposition summarizes the main properties of the covariance.

Proposition 5.3.10. *Let X and Y be random variables with finite second moments. Then the following are valid:*

- (1) $\text{Cov}(X, Y) = \text{Cov}(Y, X)$.
- (2) $\text{Cov}(X, X) = \mathbb{V}X$.
- (3) *The covariance is bilinear; that is, for X_1, X_2 and real numbers a_1 and a_2 ,*

$$\text{Cov}(a_1 X_1 + a_2 X_2, Y) = a_1 \text{Cov}(X_1, Y) + a_2 \text{Cov}(X_2, Y)$$

and, analogously,

$$\text{Cov}(X, b_1 Y_1 + b_2 Y_2) = b_1 \text{Cov}(X, Y_1) + b_2 \text{Cov}(X, Y_2),$$

for random variables Y_1, Y_2 and real numbers b_1, b_2 .

(4) *The covariance may also be evaluated by*

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - (\mathbb{E}X)(\mathbb{E}Y). \quad (5.43)$$

(5) $\text{Cov}(X, Y) = 0$ for independent X and Y .

Proof. Properties (1) and (2) follow directly from the definition of the covariance.

Let us verify (3). Setting $\mu_1 = \mathbb{E}X_1$ and $\mu_2 = \mathbb{E}X_2$, the linearity of the expected value implies

$$\mathbb{E}(a_1X_1 + a_2X_2) = a_1\mu_1 + a_2\mu_2.$$

Hence, if $\nu = \mathbb{E}Y$, then

$$\begin{aligned} \text{Cov}(a_1X_1 + a_2X_2, Y) &= \mathbb{E}[(a_1(X_1 - \mu_1) + a_2(X_2 - \mu_2))(Y - \nu)] \\ &= a_1\mathbb{E}[(X_1 - \mu_1)(Y - \nu)] + a_2\mathbb{E}[(X_2 - \mu_2)(Y - \nu)] \\ &= a_1\text{Cov}(X_1, Y) + a_2\text{Cov}(X_2, Y). \end{aligned}$$

This proves the first part of (3). The second part can be proven in the same way or one uses $\text{Cov}(X, Y) = \text{Cov}(Y, X)$ and the first part of (3).

Next we prove eq. (5.43). With $\mu = \mathbb{E}X$ and $\nu = \mathbb{E}Y$, from

$$(X - \mu)(Y - \nu) = XY - \mu Y - \nu X + \mu\nu,$$

we get that

$$\begin{aligned} \text{Cov}(X, Y) &= \mathbb{E}[XY - \mu Y - \nu X + \mu\nu] = \mathbb{E}[XY] - \mu\mathbb{E}Y - \nu\mathbb{E}X + \mu\nu \\ &= \mathbb{E}[XY] - \mu\nu. \end{aligned}$$

This proves (4) by the definition of μ and ν .

Finally, we verify (5). If X and Y are independent, then, by Proposition 4.1.9, this is also true for $X - \mu$ and $Y - \nu$. Thus, property (6) of Proposition 5.1.38 applies and leads to

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu)(Y - \nu)] = \mathbb{E}(X - \mu) \mathbb{E}(Y - \nu) = [\mathbb{E}X - \mu] [\mathbb{E}Y - \nu] = 0.$$

Therefore, the proof is completed. \square

Remark 5.3.11. Quite often the computation of $\text{Cov}(X, Y)$ can be simplified by the use of eq. (5.43). For example, consider X and Y in Example 3.5.8. In Example 5.3.3 we found $\mathbb{E}[XY] = 1/6$. Since $\mathbb{E}X = \mathbb{E}Y = 1/2$, from eq. (5.43) we immediately get

$$\text{Cov}(X, Y) = \frac{1}{6} - \frac{1}{4} = -\frac{1}{12}.$$

We obtained the same result in Example 5.3.9 with slightly more effort.

Property (5) in Proposition 5.3.10 is of special interest. It asserts $\text{Cov}(X, Y) = 0$ for independent X and Y . One may ask now whether this characterizes independent random variables. More precisely, are the random variables X and Y independent if and only if $\text{Cov}(X, Y) = 0$?

The answer is negative as the next example shows.

Example 5.3.12. The joint distribution of X and Y is given by the following table:

$Y \setminus X$	-1	0	1	
-1	$\frac{1}{10}$	$\frac{1}{10}$	$\frac{1}{10}$	$\frac{3}{10}$
0	$\frac{1}{10}$	$\frac{2}{10}$	$\frac{1}{10}$	$\frac{2}{5}$
1	$\frac{1}{10}$	$\frac{1}{10}$	$\frac{1}{10}$	$\frac{3}{10}$
	$\frac{3}{10}$	$\frac{2}{5}$	$\frac{3}{10}$	

Of course, $\mathbb{E}X = \mathbb{E}Y = 0$ and, moreover,

$$\mathbb{E}[XY] = \frac{1}{10}((-1)(-1) + (-1)(+1) + (+1)(-1) + (+1)(+1)) = 0,$$

which by eq. (5.43) implies $\text{Cov}(X, Y) = 0$. On the other hand, Proposition 3.6.11 tells us that X and Y are *not* independent. For example,

$$\mathbb{P}\{X = 0, Y = 0\} = \frac{1}{5} \quad \text{while} \quad \mathbb{P}\{X = 0\}\mathbb{P}\{Y = 0\} = \frac{4}{25}.$$

Example 5.3.12 shows that $\text{Cov}(X, Y) = 0$ is, in general, weaker than the independence of X and Y . Therefore, the following definition makes sense.

Definition 5.3.13. Two random variables X and Y satisfying $\text{Cov}(X, Y) = 0$ are said to be **uncorrelated**. Otherwise, if $\text{Cov}(X, Y) \neq 0$, then X and Y are **correlated**.

More generally, a sequence X_1, \dots, X_n of random variables is called (pairwise) **uncorrelated**, if $\text{Cov}(X_i, X_j) = 0$ whenever $i \neq j$.

Using this notation, property (5) in Proposition 5.3.10 may now be formulated in the following way:

! X and Y independent $\begin{matrix} \Rightarrow \\ \nRightarrow \end{matrix}$ X and Y uncorrelated.

Example 5.3.14. Let $A, B \in \mathcal{A}$ be two events in a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and let $\mathbb{1}_A$ and $\mathbb{1}_B$ be their indicator functions as introduced in Definition 3.6.16. How can we compute $\text{Cov}(\mathbb{1}_A, \mathbb{1}_B)$?

Answer: Since $\mathbb{E}\mathbb{1}_A = \mathbb{P}(A)$, we get

$$\begin{aligned}\text{Cov}(\mathbb{1}_A, \mathbb{1}_B) &= \mathbb{E}[\mathbb{1}_A \mathbb{1}_B] - (\mathbb{E}\mathbb{1}_A)(\mathbb{E}\mathbb{1}_B) = \mathbb{E}\mathbb{1}_{A \cap B} - \mathbb{P}(A) \mathbb{P}(B) \\ &= \mathbb{P}(A \cap B) - \mathbb{P}(A) \mathbb{P}(B).\end{aligned}$$

This tells us that $\mathbb{1}_A$ and $\mathbb{1}_B$ are uncorrelated if and only if the events A and B are independent. But as we saw in Proposition 3.6.17, this happens if and only if the random variables $\mathbb{1}_A$ and $\mathbb{1}_B$ are independent. In other words, two indicator functions are independent if and only if they are uncorrelated.

Finally, we consider the covariance of two continuous random variables.

Example 5.3.15. Suppose a random vector (X, Y) is uniformly distributed on the unit ball of \mathbb{R}^2 . Then the joint density of (X, Y) is given by

$$p(x, y) = \begin{cases} \frac{1}{\pi} & \text{if } x^2 + y^2 \leq 1, \\ 0 & \text{if } x^2 + y^2 > 1. \end{cases}$$

We proved in Example 3.5.19 that X and Y possess the distribution densities

$$q(x) = \begin{cases} \frac{2}{\pi} \sqrt{1-x^2} & \text{if } |x| \leq 1, \\ 0 & \text{if } |x| > 1 \end{cases} \quad \text{and} \quad r(y) = \begin{cases} \frac{2}{\pi} \sqrt{1-y^2} & \text{if } |y| \leq 1, \\ 0 & \text{if } |y| > 1. \end{cases}$$

The function $y \mapsto y(1-y^2)^{1/2}$ is odd. Consequently, because we integrate over an interval symmetric around the origin,

$$\mathbb{E}X = \mathbb{E}Y = \frac{2}{\pi} \int_{-1}^1 y(1-y^2)^{1/2} dy = 0.$$

By the same argument, we obtain

$$\mathbb{E}[XY] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x \cdot y) p(x, y) dx dy = \frac{1}{\pi} \int_{-1}^1 y \left[\int_{-\sqrt{1-y^2}}^{\sqrt{1-y^2}} x dx \right] dy = 0,$$

and these two assertions imply $\text{Cov}(X, Y) = 0$. Hence, X and Y are uncorrelated, but as we already observed in Example 3.6.21, they are not independent.

Summary: Let X and Y be two random variables with second moment and expected values μ and ν , respectively. Then

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu)(Y - \nu)] = \mathbb{E}[XY] - (\mathbb{E}X)(\mathbb{E}Y)$$

denotes the covariance of X and Y . It can be evaluated in the discrete case as follows:

$$\text{Cov}(X, Y) = \sum_{i,j=1}^{\infty} (x_i - \mu)(y_j - \nu) \mathbb{P}\{X = x_i, Y = y_j\}.$$

If X and Y are continuous with joint density p , then

$$\text{Cov}(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu)(y - \nu) p(x, y) dx dy.$$

5.3.2 Correlation coefficient

The question arises whether or not the covariance is the quantity that we are looking for; that is, which measures the degree of dependence. The answer is only partially affirmative. Why? Suppose X and Y are dependent. If a is a nonzero real number, then a natural demand is that the degree of dependence between X and Y should be the same as that between aX and Y . But

$$\text{Cov}(aX, Y) = a \text{Cov}(X, Y),$$

thus, if $a \neq 1$, then the measure of dependence would increase or decrease. To overcome this drawback, we normalize the covariance in the following way.

Definition 5.3.16. Let X and Y be random variables with finite second moments. Furthermore, we assume that neither X nor Y are constant with probability 1, that is, we have $\mathbb{V}X > 0$ and $\mathbb{V}Y > 0$. Then the quotient

$$\rho(X, Y) := \frac{\text{Cov}(X, Y)}{(\mathbb{V}X)^{1/2}(\mathbb{V}Y)^{1/2}} \quad (5.44)$$

is called the **correlation coefficient** of X and Y .

To verify a crucial property of the correlation coefficient, we need the following version of the Cauchy–Schwarz inequality.

Proposition 5.3.17 (Cauchy–Schwarz inequality). *For any two random variables X and Y with finite second moments, it follows that*

$$|\mathbb{E}(XY)| \leq (\mathbb{E}X^2)^{1/2} (\mathbb{E}Y^2)^{1/2}. \quad (5.45)$$

Proof. By property (8) of Proposition 5.1.38, we have

$$0 \leq \mathbb{E}(|X| - \lambda|Y|)^2 = \mathbb{E}X^2 - 2\lambda\mathbb{E}|XY| + \lambda^2 \mathbb{E}Y^2 \quad (5.46)$$

for any $\lambda \in \mathbb{R}$. To proceed further, we have to assume⁸ $\mathbb{E}X^2 > 0$ and $\mathbb{E}Y^2 > 0$. The latter assumption allows us to choose λ as

$$\lambda := \frac{(\mathbb{E}X^2)^{1/2}}{(\mathbb{E}Y^2)^{1/2}}.$$

If we apply inequality (5.46) with this λ , then we obtain

$$0 \leq \mathbb{E}X^2 - 2 \frac{(\mathbb{E}X^2)^{1/2}}{(\mathbb{E}Y^2)^{1/2}} \mathbb{E}|XY| + \mathbb{E}X^2 = 2\mathbb{E}X^2 - 2 \frac{(\mathbb{E}X^2)^{1/2}}{(\mathbb{E}Y^2)^{1/2}} \mathbb{E}|XY|,$$

which easily implies (recall that we assumed $\mathbb{E}X^2 > 0$)

$$\mathbb{E}|XY| \leq (\mathbb{E}X^2)^{1/2} (\mathbb{E}Y^2)^{1/2}.$$

To complete the proof, we use Corollary 5.1.41 and get

$$|\mathbb{E}(XY)| \leq \mathbb{E}|XY| \leq (\mathbb{E}X^2)^{1/2} (\mathbb{E}Y^2)^{1/2},$$

as asserted. □

Remark 5.3.18. An analogue inequality as (5.45) for vectors in \mathbb{R}^n is as follows: let $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ be two elements in \mathbb{R}^n . Then one has

$$\left| \sum_{j=1}^n x_j y_j \right| \leq \left(\sum_{j=1}^n x_j^2 \right)^{1/2} \left(\sum_{j=1}^n y_j^2 \right)^{1/2}.$$

In the language of scalar products and Euclidean distance, this says

$$|\langle x, y \rangle| \leq |x| |y|.$$

Corollary 5.3.19. *The correlation coefficient satisfies*

$$-1 \leq \rho(X, Y) \leq 1.$$

Proof. Let as before $\mu = \mathbb{E}X$ and $\nu = \mathbb{E}Y$. Applying inequality (5.45) to $X - \mu$ and $Y - \nu$ leads to

$$\begin{aligned} |\text{Cov}(X, Y)| &= |\mathbb{E}(X - \mu)(Y - \nu)| \leq (\mathbb{E}(X - \mu)^2)^{1/2} (\mathbb{E}(Y - \nu)^2)^{1/2} \\ &= (\mathbb{V}X)^{1/2} (\mathbb{V}Y)^{1/2}, \end{aligned}$$

⁸ The Cauchy–Schwarz inequality remains valid for $\mathbb{E}X^2 = 0$ or $\mathbb{E}Y^2 = 0$. In this case, it follows that $\mathbb{P}\{X = 0\} = 1$ or $\mathbb{P}\{Y = 0\} = 1$, hence $\mathbb{P}\{XY = 0\} = 1$ and $\mathbb{E}[XY] = 0$.

or, equivalently,

$$-(\mathbb{V}X)^{1/2} (\mathbb{V}Y)^{1/2} \leq \text{Cov}(X, Y) \leq (\mathbb{V}X)^{1/2} (\mathbb{V}Y)^{1/2}.$$

By the definition of $\rho(X, Y)$ given in eq. (5.44), this implies $-1 \leq \rho(X, Y) \leq 1$, as asserted. \square

Interpretation: For uncorrelated X and Y , we have $\rho(X, Y) = 0$. In particular, this is valid if X and Y are independent. On the contrary, $\rho(X, Y) \neq 0$ tells us that X and Y are dependent. Thereby, values near to zero correspond to weak dependence, while $\rho(X, Y)$ near 1 or -1 indicate a strong dependence. The strongest possible dependence is when $Y = aX$ for some $a \neq 0$. Then $\rho(X, Y) = 1$ if $a > 0$ while $\rho(X, Y) = -1$ for $a < 0$.

Definition 5.3.20. Two random variables X and Y are said to be **positively correlated** if $\rho(X, Y) > 0$. In the case that $\rho(X, Y) < 0$, they are said to be **negatively correlated**.

Interpretation: Random variables X and Y are positively correlated, provided that larger (or smaller) values of X make larger (or smaller) values of Y more likely. This does *not* mean that a larger X -value always implies a larger Y -value, only that the probability for those larger values increases. And in the same way, if X and Y are negatively correlated, then larger values of X make smaller Y -values more likely.

Let us explain this with two typical examples. Choose by random a person ω in the audience. Let $X(\omega)$ be his or her height and $Y(\omega)$ his or her weight. Then X and Y will surely be positively correlated. But this does not necessarily mean that each taller person has a bigger weight. Another example of negatively correlated random variables could be as follows: X is the average number of cigarettes that a randomly chosen person smokes per day and Y is his lifetime.

Example 5.3.21. Let us come back to Example 5.3.1: in an urn there are n balls labeled with “0” and n labeled with “1.” One chooses two balls without replacement. Then X is the value of the first ball, Y that of the second. How does the correlation coefficient of X and Y depend on n ?

Answer: The joint distribution of X and Y is given by the following table:

$Y \backslash X$	0	1	
0	$\frac{n-1}{4n-2}$	$\frac{n}{4n-2}$	$\frac{1}{2}$
1	$\frac{n}{4n-2}$	$\frac{n-1}{4n-2}$	$\frac{1}{2}$
	$\frac{1}{2}$	$\frac{1}{2}$	

Direct computations show $\mathbb{E}X = \mathbb{E}Y = 1/2$ and $\mathbb{V}X = \mathbb{V}Y = 1/4$. Moreover, it easily follows $\mathbb{E}[XY] = \frac{n-1}{4n-2}$, hence

$$\text{Cov}(X, Y) = \frac{n-1}{4n-2} - \frac{1}{4} = \frac{-1}{8n-4},$$

and the correlation coefficient equals

$$\rho(X, Y) = \frac{-1}{\frac{8n-4}{\sqrt{\frac{1}{4}}\sqrt{\frac{1}{4}}}} = \frac{-1}{2n-1}.$$

If $n \rightarrow \infty$, then $\rho(X, Y)$ is of order $\frac{-1}{2n}$. Hence, if n is large, then the random variables X and Y are “almost” uncorrelated.

Since $\rho(X, Y) < 0$, the two random variables are negatively correlated. Why? This was already explained in Example 5.3.1: an occurrence of $X = 1$ makes $Y = 0$ more likely, while the occurrence of $X = 0$ increases the likelihood of $Y = 1$. In the case $n = 1$, the value of Y is completely determined by that of X , expressed by $\rho(X, Y) = -1$.

Summary: Let X and Y be two random variables with finite second moment. Then

$$\rho(X, Y) := \frac{\text{Cov}(X, Y)}{(\mathbb{V}X)^{1/2}(\mathbb{V}Y)^{1/2}}$$

is said to be their correlation coefficient. The random variables are said to be positively correlated if $\rho(X, Y) > 0$, negatively correlated if $\rho(X, Y) < 0$, and uncorrelated if $\rho(X, Y) = 0$. The basic properties are $-1 \leq \rho(X, Y) \leq 1$ and $\rho(X, Y) = 0$ for independent X and Y . But note that uncorrelated X and Y need not be independent.

5.4 Some paradoxes and examples

The aim of this section is to present a few more comprehensive examples of special interest having a long history.

5.4.1 Boy or girl paradox

In 1959 Martin Gardner⁹ phrased the following two questions:

1. Mr. Jones has two children. The older child is a girl. What is the probability that both children are girls?
2. Mr. Smith has two children. At least one of them is a boy. What is the probability that both children are boys?

The basic assumptions for answering these questions are:

- (a) Each child is either a boy or a girl.

⁹ Martin Gardner, *Problems involving questions of probability and ambiguity*, Scientific American, October 1959.

- (b) Boys and girls occur equally likely.
 (c) The gender of the two children is independent of each other.

Answers: There are 4 different possibilities for the gender of the two children:

$$(G, G), (B, G), (G, B), \text{ and } (B, B).$$

For example, (B, G) means that the older child is a boy while the younger one is a girl. In view of the basic assumptions, all 4 elementary events are equally likely, hence their probability is $1/4$.

Solution of the first question: Here one asks for the probability of $\{(G, G)\}$ under the condition $\{(G, B), (G, G)\}$. So we get

$$\mathbb{P}(\{(G, G)\} | \{(G, B), (G, G)\}) = \frac{\mathbb{P}(\{(G, G)\})}{\mathbb{P}(\{(G, B), (G, G)\})} = \frac{1/4}{1/2} = \frac{1}{2}.$$

First solution of the second question: One chooses by random a family with two children. Let the event A occur if the chosen family has at least one boy. That is,

$$A = \{(B, G), (G, B), (B, B)\},$$

hence $\mathbb{P}(A) = 3/4$. In question 2, we asked for the probability of the occurrence of $\{(B, B)\}$ under the condition A . Then we get

$$\mathbb{P}(\{(B, B)\} | A) = \frac{\mathbb{P}(\{(B, B)\})}{\mathbb{P}(A)} = \frac{1/4}{3/4} = \frac{1}{3}.$$

We obtained this result by restricting the sample space and ruling out families with two girls.

Alternative solution of the second question: We split the problem into two steps. In the *first* step, one chooses at random a family with two children. But we do not have any information about the gender of these children. Any of the four configurations (G, G) , (G, B) , (B, G) , and (B, B) is possible, each occurring with probability $1/4$.

Next, in the *second* step, we choose equiprobable one of the two children of the family, that is, of the family chosen in the first step. The result may be “B” or “G.” Set

$$S = \{\text{Second step leads to “B”}\}.$$

Thus, if the randomly chosen family with two children has a boy and a girl, then the second step will lead equally likely to “B” and to “G.” Consequently,

$$\mathbb{P}(S | \{(B, G)\}) = \mathbb{P}(S | \{(G, B)\}) = \frac{1}{2}.$$

On the other hand, if the chosen family has either two girls or two boys, then it follows that

$$\mathbb{P}(S|\{(G, G)\}) = 0 \quad \text{or} \quad \mathbb{P}(S|\{(B, B)\}) = 1,$$

respectively.

By symmetry (recall that “B” and “G” are equally likely), one should have that $\mathbb{P}(S) = 1/2$. A rigorous proof is based on the law of total probability. It implies

$$\begin{aligned} \mathbb{P}(S) &= \mathbb{P}\{(G, G)\} \cdot \mathbb{P}(S|\{(G, G)\}) + \mathbb{P}\{(B, G)\} \cdot \mathbb{P}(S|\{(B, G)\}) \\ &\quad + \mathbb{P}\{(G, B)\} \cdot \mathbb{P}(S|\{(G, B)\}) + \mathbb{P}\{(B, B)\} \cdot \mathbb{P}(S|\{(B, B)\}) \\ &= \frac{1}{4} \cdot 0 + \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{4} \cdot 1 = \frac{1}{2}. \end{aligned}$$

This lets us conclude that

$$\mathbb{P}(\{(B, B)\}|S) = \mathbb{P}(S|\{(B, B)\}) \cdot \frac{\mathbb{P}(\{(B, B)\})}{\mathbb{P}(S)} = 1 \cdot \frac{1/4}{1/2} = \frac{1}{2}.$$

Consequently, the answer to the second question is: the probability for both children being boys knowing that at least one child is a boy equals $1/2$.

For the sake of completeness, we also state the other probabilities. These are $\mathbb{P}(\{(G, G)\}|S) = 0$ and

$$\mathbb{P}(\{(G, B)\}|S) = \mathbb{P}(\{(B, G)\}|S) = \mathbb{P}(S|\{(B, G)\}) \cdot \frac{\mathbb{P}(\{(B, G)\})}{\mathbb{P}(S)} = \frac{1}{2} \cdot \frac{1/4}{1/2} = \frac{1}{4}.$$

The obtained result may also be phrased as follows: the *a priori* probabilities of getting (G, G) , (G, B) , (B, G) , and (B, B) (the probabilities before choosing randomly a child) are all $1/4$ while the *a posteriori* probabilities (after observing a “B” in the second step) are 0 , $1/4$, $1/4$, and $1/2$, respectively.

Remark 5.4.1. One may wonder how it is possible that the answer to question 2 is at the same time $1/3$ and $1/2$. Maybe there is some error in the calculations. This is not so. The reason lies in the ambiguity of the way to use the information “At least one child is a boy.” In the first approach, we use this information to rule out families with two girls from the very beginning. That is, if we chose at random a family with two girls, we discard it and make another trial.

The alternative approach looks more natural, at least to us. Choose at random any family with two children. Then all four possibilities of the distribution of boys and girls may occur. Having the family fixed, we check whether or not a randomly chosen child (maybe the older, maybe the younger) is a boy. For example, we chose the family of Mr. Smith, and we see him walking with one of his children who is a boy. There is a 50% chance that this is his older child, but, of course, it also could be the younger one.

5.4.2 Randomly chosen entries

Let us extend the boy or girl paradox to a more general setting. It reads as follows: toss a fair coin labeled by “0” and “1” exactly $n \geq 2$ times, without recording the obtained results. After that, one gets the information that one of n the entries equals one. Given some $1 \leq k \leq n$, the question is now how likely is it that the observed sequence contains exactly k times the number “1.”

Note that in the case $n = k = 2$ this is exactly question 2 in Section 5.4.1. To see this link “B” with “1” and “G” with “0.”

But also here, in this generalized setting, the information that *one of the entries equals “1”* can be interpreted in different ways.

One possible interpretation of the problem is to discard from the very beginning sequences without “1.” So the restricted sample space contains $2^n - 1$ elements. Consequently, if A_k is the event to observe k times “1,” then

$$\mathbb{P}(A_k | \{\text{At least one “1”}\}) = \frac{\binom{n}{k}}{2^n - 1} \quad k = 1, \dots, n.$$

Of course, $\mathbb{P}(A_0 | \{\text{At least one “1”}\}) = 0$.

Another interpretation of the problem is as follows: one chooses at random a sequence $x = (x_1, \dots, x_n)$ of “0”s and “1”s. The probability of its occurrence is $1/2^n$. Next one fixes the observed $x = (x_1, \dots, x_n)$ and chooses at random one of its entries (all entries are equally likely), say x_j for some $j \leq n$. Then the event $x_j = 1$ occurs with probability k/n , where k is the number of “1”s in the chosen sequence x . Thus, if

$$S := \{\text{The randomly chosen entry equals “1”}\},$$

it follows that

$$\mathbb{P}(S | \{x\}) = \frac{k}{n} \quad \text{whenever } x \in A_k.$$

Recall that A_k denotes the set of all sequences of “0” and “1” with k times “1.” The law of total probability yields (see the proof of Proposition 5.1.13 for the evaluation of the sum)

$$\begin{aligned} \mathbb{P}(S) &= \sum_{k=1}^n \sum_{x \in A_k} \mathbb{P}(\{x\}) \cdot \mathbb{P}(S | \{x\}) = \sum_{k=1}^n \sum_{x \in A_k} \frac{1}{2^n} \cdot \frac{k}{n} \\ &= \frac{1}{n2^n} \sum_{k=1}^n k \binom{n}{k} = \frac{1}{n2^n} \cdot (n2^{n-1}) = \frac{1}{2}. \end{aligned}$$

This result is not surprising at all and could also be obtained heuristically. Indeed, by symmetry the occurrence of “0” and “1” has to be equally likely.

Consequently, if $x \in A_k$, then

$$\mathbb{P}(\{x\}|S) = \mathbb{P}(S|\{x\}) \cdot \frac{\mathbb{P}(\{x\})}{\mathbb{P}(S)} = \frac{k}{n} \cdot \frac{1}{2^n} \cdot \frac{1}{2^{-1}} = \frac{k}{n 2^{n-1}}.$$

Finally, since $|A_k| = \binom{n}{k}$, the additivity of conditional probabilities leads to

$$\mathbb{P}(A_k|S) = \binom{n}{k} \cdot \frac{k}{n 2^{n-1}} = \binom{n-1}{k-1} \cdot \frac{1}{2^{n-1}}, \quad k = 1, \dots, n.$$

Summing up, we received the following result.

Proposition 5.4.2. *Let $n \geq 1$. Tossing a fair coin n times, then for all $1 \leq k \leq n$ it follows that*

$$\mathbb{P}\{\text{Observe } k \text{ times "1" | Randomly chosen entry is "1"}\} = \binom{n-1}{k-1} \cdot \frac{1}{2^{n-1}}.$$

Example 5.4.3. If we apply Proposition 5.4.2 to the case $n = k = 2$, we rediscover the answer given in Section 5.4.1. Indeed, then

$$\mathbb{P}\{\text{Both children are boys | A randomly chosen child is a boy}\} = \frac{1}{2}.$$

Another case of interest is as follows: one asks for the probability that a family with three children has two boys, provided a randomly chosen child is a boy. Then $n = 3$ and $k = 2$, hence the probability for this event equals $2 \cdot 1/4 = 1/2$. Equivalently, under the given condition, each of the events $\{(G, B, B)\}$, $\{(B, G, B)\}$, and $\{(B, B, G)\}$ occurs with probability $1/6$. On the other hand, if a randomly chosen child of the three ones is a boy, the occurrence of three boys possesses the probability $1/4$. Without that knowledge, the probability of $\{(B, B, B)\}$ equals $1/8$.

Remark 5.4.4. Proposition 5.4.2 tells us that the left-hand probability coincides with that of the occurrence of $k-1$ times “1” when tossing a fair coin $n-1$ times. There is a straightforward explanation of this coincidence. In the presented setting, we first toss the coin and next choose randomly an entry of the observed sequence of zeroes and ones. But we could do it also the other way round: first choosing randomly a number from 1 to n and after that tossing the coin. We leave the details as problem for the interested reader (see Problem 5.25).

5.4.3 Secretary problem

This is a well-known problem in Probability Theory, also called “marriage problem” or the “sultan’s dowry problem.”

A company wants to hire a secretary. There are n applicants interviewed in random order. Immediately after the interview, the administrator decides whether or not the

candidate is rejected or hired. Once rejected, the candidate cannot be recalled. The goal of the company is to get the best applicant. To this end, every interviewed candidate is ranked in a linear order. But note that the administrator has no information about the quality of the unseen applicants.

Now the strategy of the administrator is as follows: Choose a number $0 \leq r \leq n - 1$, interview the first r applicants and reject them all. After that choose the first candidate who is better than all of the r rejected ones. If $r = 0$, then nobody is interviewed and the first applicant is hired.

Questions: How likely is it that this strategy leads to the employment of the overall best candidate? What is the optimal choice of the number r ?

Before proceeding further, let us explain the problem with an easy example. Suppose there are three candidates enumerated by 1, 2, and 3. We assume that applicant 2 is better than 1 and that candidate 3 is better than 2. Then there are $3! = 6$ ways of interviewing the applicants:

$$a = (1, 2, 3), \quad b = (1, 3, 2), \quad c = (2, 1, 3), \quad d = (2, 3, 1), \quad e = (3, 1, 2), \quad f = (3, 2, 1).$$

If $r = 0$, only orderings (e) and (f) lead to the best candidate. So, the chance for hiring the best one is $2/6 = 1/3$.

In the case $r = 1$, the company hires the best one in the cases (b) , (c) , and (d) . Hence, the chance to get the best one equals $1/2$.

Finally, if $r = 2$, only (a) and (c) are successful. Thus, also here the chance to get the best applicant is $1/3$.

Summing up, if $n = 3$, then the optimal strategy is to choose $r = 1$. That is, reject the first applicant and then choose the next one who is better than the first. Of course, it may happen that there is nobody better than the first, which occurs in the cases (e) and (f) . Then nobody is hired and the administrator failed to get the best applicant.

Let us transform the general problem into a mathematical setting. Name the candidates by numbers from $1, \dots, n$ and, without losing generality, let us assume that the applicant n is the best, that applicant $n - 1$ is the second best, and so on. Suppose now the candidates are interviewed in the order $\pi(1), \pi(2), \dots, \pi(n)$ for some permutation $\pi \in S_n$. In this context, the best applicant appears at position k if and only if $\pi(k) = n$ or, equivalently, $\pi^{-1}(n) = k$.

It is assumed that all orderings of the applicants are equally likely, so we have to endow the set S_n of permutations with the uniform distribution \mathbb{P} . That is, for all $A \subseteq S_n$ we have

$$\mathbb{P}(A) = \frac{|A|}{|S_n|} = \frac{|A|}{n!}.$$

In particular, if

$$A_k := \{\pi \in S_n : \pi(k) = n\},$$

it follows that

$$\mathbb{P}(A_k) = \frac{(n-1)!}{n!} = \frac{1}{n}, \quad k = 1, \dots, n. \tag{5.47}$$

Recall that $\pi \in A_k$ if and only if the best applicant is at position k in the queue of candidates.

Let $P(r)$ be the probability to hire the best candidate when choosing the strategy of rejecting the first r candidates.

If $r = 0$, then the best candidate is hired if and only if $\pi(1) = n$ or, equivalently, if and only if $\pi \in A_1$. Hence, we get in this case

$$P(0) = \mathbb{P}(A_1) = \frac{1}{n}.$$

Consider now an arbitrary $1 \leq r \leq n-1$ and suppose that $\pi \in A_k$ for a certain $1 \leq k \leq n$, that is, the permutation satisfies $\pi(n) = k$. When does the administrator choose the best applicant? This happens if and only if $k > r$ and, moreover,

$$\pi(r+1), \dots, \pi(k-1) < \max\{\pi(1), \dots, \pi(r)\}. \tag{5.48}$$

Another way to formulate property (5.48) is as follows:

$$\text{If } \pi(a) = \max\{\pi(1), \dots, \pi(k-1)\}, \text{ then necessarily } 1 \leq a \leq r. \tag{5.49}$$

This may also be expressed as follows: whenever $\pi \in A_k$ for some $k \geq 2$, then the second best candidate among the first k candidates has to be at position a for some $a \leq r$. Recall that the best one occurs at position k . Compare the three possible situations discussed in Figure 5.3.

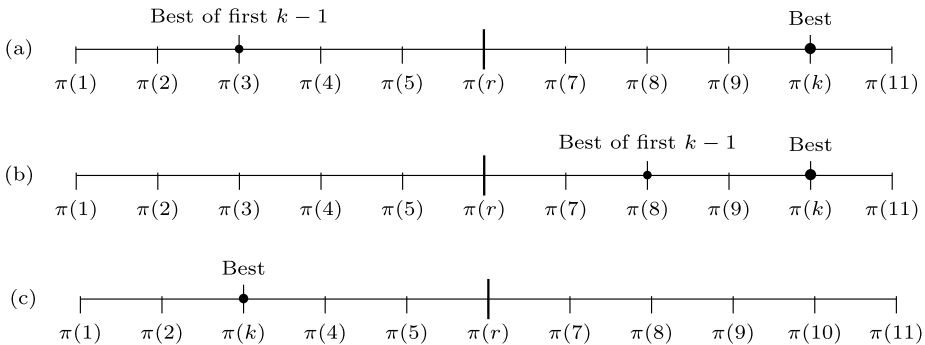


Figure 5.3: In case (a), the company hires the best secretary while it fails to do so in the cases (b) and (c).

To proceed further, given $1 \leq a \leq k-1$, define disjoint subsets S_k^a of S_n by

$$S_k^a := \{\pi \in S_n : \pi(a) = \max\{\pi(1), \dots, \pi(k-1)\}\}.$$

Verbally said, a permutation π belongs to S_k^a if and only if it attains its maximal value in $\{1, \dots, k-1\}$ at the given number a .

We claim now that if $2 \leq k \leq n$, then

$$\mathbb{P}(S_k^a | A_k) = \frac{1}{k-1}, \quad a = 1, \dots, k-1. \tag{5.50}$$

Of course,¹⁰ $\mathbb{P}(S_k^1 | A_k) = \dots = \mathbb{P}(S_k^{k-1} | A_k)$. Moreover, because $k \geq 2$, there always exists a best candidate among positions 1 and $k-1$, thus

$$1 = \mathbb{P}\left(\bigcup_{a=1}^{k-1} S_k^a | A_k\right) = \sum_{a=1}^{k-1} \mathbb{P}(S_k^a | A_k).$$

Clearly, these two properties prove eq. (5.50).

Summing up, in view of assertion (5.49), for a given order $\pi(1), \dots, \pi(n)$ the strategy leads to the best candidate if and only if

$$\exists k > r, \quad \exists 1 \leq a \leq r, \quad \pi \in S_k^a \cap A_k \iff \pi \in \bigcup_{k=r+1}^n \left[A_k \cap \left(\bigcup_{a=1}^r S_k^a \right) \right]. \tag{5.51}$$

Now we are prepared to evaluate $P(r)$, the probability to choose the best applicant when rejecting the first r candidates. Using (5.51), an application of the law of multiplication together with eqs. (5.47) and (5.50) implies (recall that the A_k s and the S_k^a s are disjoint for fixed k)

$$\begin{aligned} P(r) &= \mathbb{P}\left(\bigcup_{k=r+1}^n \left[A_k \cap \left(\bigcup_{a=1}^r S_k^a \right) \right]\right) = \sum_{k=r+1}^n \sum_{a=1}^r \mathbb{P}(A_k \cap S_k^a) \\ &= \sum_{k=r+1}^n \sum_{a=1}^r \mathbb{P}(S_k^a | A_k) \mathbb{P}(A_k) = \frac{r}{n} \sum_{k=r+1}^n \frac{1}{k-1} = \frac{r}{n} \sum_{k=r}^{n-1} \frac{1}{k}. \end{aligned}$$

Conclusion: The optimal choice of the number $1 \leq r < n$ is that for which

$$r \mapsto \frac{r}{n} \sum_{k=r}^{n-1} \frac{1}{k} \tag{5.52}$$

becomes maximal. For example, if $n = 30$, the maximal value is attained at $r = 11$ and one has $P(11) = 0.378651$. But note that also choices of r near to 11 lead to reasonably large probabilities. For example, we have $P(9) = 0.373139$, $P(10) = 0.377562$, $P(12) = 0.376711$, and $P(13) = 0.371992$. One should compare these values with the random choice of a single candidate where the probability to get the best applicant equals $1/30 = 0.0\bar{3}$. See also Figure 5.4.

¹⁰ If $1 \leq a, b \leq k-1$, define a bijection between S_k^a and S_k^b by $\pi \mapsto \pi \circ i_{a,b}$ where $i_{a,b}$ is the inversion of a and b .

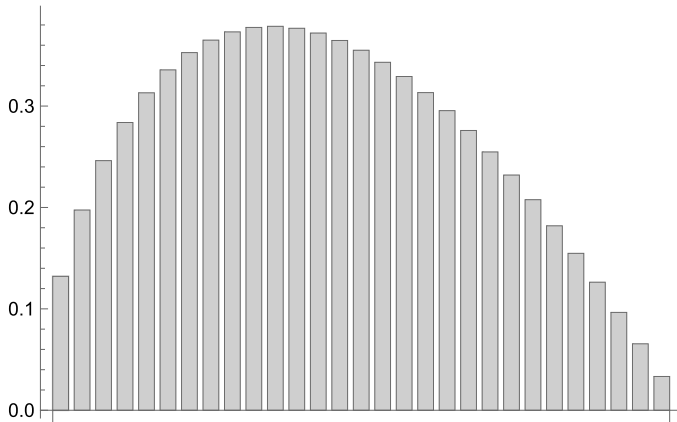


Figure 5.4: The values of $P(r)$, $1 \leq r \leq 29$, in the case of 30 applicants.

Remark 5.4.5. For large numbers n , it might be quite difficult to find the number $r < n$ for which the function in (5.52) becomes maximal. Here assertion (5.29) may be helpful. Using

$$\frac{r}{n} \sum_{k=r}^{n-1} \frac{1}{k} = \frac{r}{n} \sum_{k=1}^{n-1} \frac{1}{k} - \frac{r}{n} \sum_{k=1}^{r-1} \frac{1}{k} \approx \frac{r}{n} (\ln n - \ln r) = \frac{r}{n} \ln \left(\frac{n}{r} \right) = -\frac{r}{n} \ln \left(\frac{r}{n} \right),$$

it follows that for large n the optimal choice of r is $\frac{r}{n} \sim x_0$ where $x \mapsto -x \ln x$ becomes maximal at x_0 . Methods from Calculus imply $x_0 = 1/e \approx 0.367879$. Thus, a rough choice of the optimal r is 37% of n . In the literature, this is quite often called the **37%-rule**. If as above $n = 30$, then $30/e \approx 11.0364$ while 37% of $n = 30$ gives 11.1. Thus, this also leads to $r = 11$ as the optimal choice. Check Figure 5.5 to see that for large n one has $P(r) \approx -(r/n) \ln(r/n)$.

5.4.4 Two-envelope paradox

We finally present a famous paradox in Probability Theory called the “two-envelope paradox” or the “envelope exchange paradox.” Imagine you may choose one of two indistinguishable envelopes, both containing a certain amount of money. You do not know how much money is in the envelopes, but you have the information that one of the two envelopes contains twice as much as the other. Having chosen an envelope at will, you inspect it, and find x dollars. Hence, the other unopened envelope contains either $2x$ or $x/2$ dollars, depending on whether the chosen envelope was that with the smaller or larger sum.

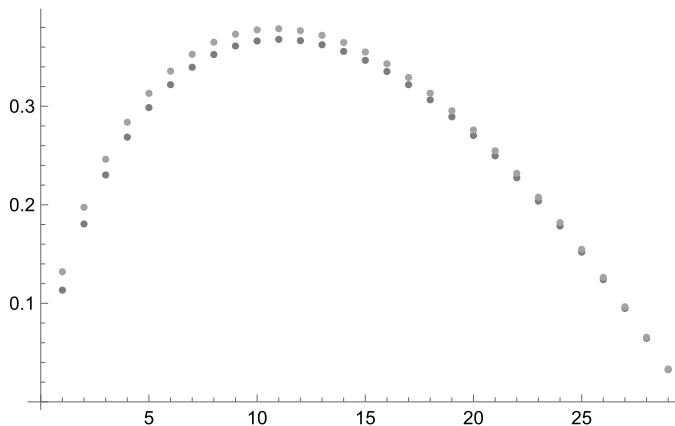


Figure 5.5: The values of $P(r)$ for 30 applicants. The upper dots are the “correct” values while the lower ones are the values of the approximation $r \mapsto -(r/30) \ln(r/30)$.

After that you are given the chance to swap envelopes. Should you use this opportunity? If you do not swap, then you keep the x dollars. Otherwise, you either double your amount or halve it, both with probability $1/2$. Thus, if E is the expected amount after switching, it follows that

$$E = (2x) \cdot \frac{1}{2} + \left(\frac{x}{2}\right) \cdot \frac{1}{2} = \frac{5}{4} \cdot x.$$

Consequently, on average, by switching you gain $x/4$ dollars. Imagine, for example, the chosen envelope contains \$100. Then by switching one either loses \$50 or one wins \$100, both with probability $1/2$. Of course, this contradicts the common sense. But what is wrong?

First, there is a misinterpretation of the observed amount. The observed x is a random value, not the expected value of the money you get. Say for some $c > 0$, the envelopes contain c and $2c$ dollars. Denote by X the money you get without switching, then it follows that

$$\mathbb{P}\{X = c\} = \mathbb{P}\{X = 2c\} = \frac{1}{2}.$$

And after switching the new random variable \tilde{X} also satisfies

$$\mathbb{P}\{\tilde{X} = c\} = \mathbb{P}\{X = 2c\} = \frac{1}{2} \quad \text{and} \quad \mathbb{P}\{\tilde{X} = 2c\} = \mathbb{P}\{X = c\} = \frac{1}{2}.$$

Hence, it follows that $\mathbb{E}X = \mathbb{E}\tilde{X} = 3c/2$, and on average there is no advantage by switching, exactly as one expects.

But there is still another missing information in the scenario. In which way are the sums in the envelopes chosen? Are these, as assumed above, always (in each experiment)

fixed amounts c and $2c$? Or is there a positive random variable Y such that the envelopes contain Y and $2Y$ dollars? In other words, before you take one of the two envelopes at random, the included sums are chosen by another independent random experiment described by a random variable Y . But, and this suggests the formulation of the problem, thereby it is impossible to do it in a way such that all possible amounts of integers (or positive real numbers) are equally likely.

Let us explain this (random) setting with an example. Assume the master of ceremonies rolls a die and, depending on the observed number $k \in \{1, \dots, 6\}$, he puts 2^k dollars into one envelope and 2^{k+1} into the other. In the above setting, the random variable Y satisfies

$$\mathbb{P}\{Y = 2^k\} = \frac{1}{6}, \quad k = 1, \dots, 6, \quad (5.53)$$

and if X denotes the obtained amount, then¹¹

$$\mathbb{P}\{X = 2\} = \frac{1}{2} \cdot \mathbb{P}\{Y = 1\} = \frac{1}{12}, \quad \mathbb{P}\{X = 128\} = \frac{1}{2} \cdot \mathbb{P}\{Y = 6\} = \frac{1}{12}$$

and

$$\mathbb{P}\{X = 2^k\} = \frac{1}{6} \cdot \mathbb{P}\{X = 2^k | Y = k - 1\} + \frac{1}{6} \cdot \mathbb{P}\{X = 2^k | Y = k\} = \frac{1}{6}, \quad k = 2, \dots, 6.$$

So we get

$$\mathbb{E}X = 2 \cdot \frac{1}{12} + \frac{1}{6} \cdot \sum_{k=2}^6 2^k + 128 \cdot \frac{1}{12} = 31.5. \quad (5.54)$$

If \tilde{X} denotes the obtained amount after always switching, then

$$\mathbb{P}\{\tilde{X} = 2^k\} = \mathbb{P}\{X = 2^k\}, \quad k = 1, \dots, 7,$$

hence nothing changes by always swapping.

But what happens if one swaps only in the case that the opened envelope contains a “small” amount? Say, one swaps in the above example if there are less than \$60 in the envelope and otherwise one does not. Then the probability to get \$32 diminishes to $1/12$ while the probability of obtaining \$64 increases to $3/12$. Thus, after eventually swapping, the average of the amount \bar{X} equals

$$\mathbb{E}\bar{X} = 2 \frac{1}{12} + 4 \frac{1}{6} + 8 \frac{1}{6} + 16 \frac{1}{6} + 32 \frac{1}{12} + 64 \frac{3}{12} + 128 \frac{1}{12} = 34.1\bar{6}.$$

¹¹ The possible pairs of included amounts are

(2, 4), (4, 8), (8, 16), (16, 32), (32, 64), and (64, 128).

So we see, this strategy improves the average of the money obtained. Moreover, the optimal case occurs if there are \$32 and \$64 in the envelopes¹² and, furthermore, one had chosen the envelope containing the smaller amount. Then by swapping one gets extra \$32. The probability that this happens equals $1/12$.

But note that this strategy heavily depends on some foreknowledge about the size of the amount in the envelopes. For example, if one decides to switch provided there are less than \$200 in the opened envelope,¹³ then there is no improvement of $\mathbb{E}X$.

Let us finally shortly discuss the case of general (discretely) distributed amounts in the envelopes.¹⁴ So suppose there are certain positive numbers x_1, x_2, \dots and nonnegative p_k s with $\sum_{k=1}^{\infty} p_k = 1$. Choose a random variable Y for which

$$\mathbb{P}\{Y = x_k\} = p_k, \quad k = 1, 2, \dots$$

Put with probability p_k into one envelope x_k and into the other $2x_k$ dollars. After that, choose equally likely one envelope at random.¹⁵ Then we get for the expected amount X that

$$\mathbb{P}\{X = x_k\} = \frac{p_k}{2} \quad \text{and} \quad \mathbb{P}\{X = 2x_k\} = \frac{p_k}{2}.$$

This implies

$$\mathbb{E}X = \sum_{k=1}^{\infty} x_k \frac{p_k}{2} + \sum_{k=1}^{\infty} (2x_k) \frac{p_k}{2} = \frac{3}{2} \sum_{k=1}^{\infty} p_k x_k = \frac{3}{2} \mathbb{E}Y.$$

For example, choosing Y as in eq. (5.53), it follows that

$$\mathbb{E}Y = \frac{1}{6} \sum_{k=1}^6 2^k = 21 \quad \Rightarrow \quad \mathbb{E}X = \frac{3}{2} \cdot \mathbb{E}Y = \frac{3}{2} \cdot 21 = 31.5.$$

This coincides with the result obtained in eq. (5.54).

If, as before, \tilde{X} denotes the obtained amount after switching, then

$$\mathbb{P}\{\tilde{X} = x_k\} = \mathbb{P}\{X = 2x_k\} = \frac{p_k}{2} \quad \text{and} \quad \mathbb{P}\{\tilde{X} = 2x_k\} = \mathbb{P}\{X = x_k\} = \frac{p_k}{2},$$

so nothing has changed and $\mathbb{E}\tilde{X} = \mathbb{E}X = 3 \mathbb{E}Y/2$.

¹² The result of rolling the die was “5.”

¹³ We encourage the reader to evaluate $\mathbb{E}\tilde{X}$ when swapping in the case that there are either less than 20 or less than 10 dollars in the chosen envelope. Find the optimal threshold for switching and nonswitching.

¹⁴ One may also choose continuous distributions of the included amounts, but this is more involved and uses facts not included in the present book.

¹⁵ Note that the previous example fits into this setting. There we had $x_k = 2^k$ as well as $p_1 = \dots = p_6 = 1/6$ and $p_k = 0$ if $k > 6$.

Thus, always swapping does not yield any advantage. But what happens if we use the following strategy: Choose a threshold $N > 0$. If the amount x in the chosen envelope satisfies $x < N$, then switch. Otherwise, if $x \geq N$, do not do so. For simplicity, we answer this question only for special distributions of the amounts.

So suppose that for a certain $k = 0, 1, 2, \dots$ one envelope contains 2^k dollars and the other 2^{k+1} , and that the probability to choose this pair equals p_k where $p_k \geq 0$ and $\sum_{k=0}^{\infty} p_k = 1$. That is, the contents of one envelope is Y , that of the other $2Y$ where

$$\mathbb{P}\{Y = 2^k\} = p_k, \quad k = 0, 1, 2, \dots$$

This leads to

$$\mathbb{E}X = \frac{3}{2} \cdot \mathbb{E}Y = \frac{3}{2} \sum_{k=0}^{\infty} p_k 2^k,$$

which implies that $\mathbb{E}X < \infty$ if and only if $\sum_{k=0}^{\infty} p_k 2^k < \infty$.

If $\mathbb{E}X = \infty$, it does not make sense to ask whether $\mathbb{E}X$ increases or decreases by swapping or nonswapping. So let us assume that the expected gained amount is finite. Choose a threshold $N > 1$. Swap if the amount in the chosen envelope is less than N . Do not swap otherwise. Does this improve the expected value of gained money?

To answer this question take the integer $\ell \geq 0$ for which $2^\ell < N \leq 2^{\ell+1}$. Let \bar{X} be the obtained amount after switching those envelopes where one observes an amount smaller than N . Then the only change of the distribution of X occurs in the case where 2^ℓ and $2^{\ell+1}$ dollars are in the two envelopes. No matter if one were to choose the envelope with the smaller or with the larger amount, in this case one would always get that with $2^{\ell+1}$ dollars. Hence, it follows that

$$\mathbb{E}\bar{X} - \mathbb{E}X = \frac{p_\ell}{2} (2^{\ell+1} - 2^\ell) = p_\ell 2^{\ell-1}.$$

This is the good news. But what is the bad? Since we assumed that the expected value of X exists, it follows that

$$\lim_{\ell \rightarrow \infty} p_\ell 2^\ell = 0.$$

That is, the larger the threshold N , the less the expected advantage by choosing this strategy.

Another way to formulate the result is as follows: choosing the threshold N such that $2^\ell < N \leq 2^{\ell+1}$, by switching one may gain $2^{\ell+1} - 2^\ell = 2^\ell$ dollars, maybe a huge amount. But the likelihood that this happens is $p_\ell/2$, a very small number. Recall that is so if and only if, firstly, there are 2^ℓ and $2^{\ell+1}$ dollars in the envelopes and, secondly, one had chosen the envelope containing the smaller amount.

To illustrate the obtained results, choose $p_k = 2/3^{k+1}$, $k = 0, 1, \dots$. That is,

$$\mathbb{P}\{Y = 2^k\} = \frac{2}{3^{k+1}}, \quad k = 0, 1, 2, \dots$$

In this case,

$$\mathbb{E}Y = \sum_{k=0}^{\infty} \left(\frac{2}{3}\right)^{k+1} = 2 \quad \Rightarrow \quad \mathbb{E}X = \frac{3}{2} \cdot \mathbb{E}Y = 3.$$

If the chosen threshold N satisfies $2^\ell < N \leq 2^{\ell+1}$, then by switching envelopes with small amounts, the expected value increases to

$$\mathbb{E}\bar{X} = 3 + \frac{2^\ell}{3^{\ell+1}}.$$

Note that the maximal expected advantage of $1/3$ occurs if $\ell = 0$.

Remark 5.4.6. There exists an interesting tightly related version of the two envelopes paradox, sometimes called the **two-number problem**. A person writes two different numbers on two slips of paper, one on each, so that you cannot see what is written. Next you choose at random one of these two slips, turn it around and read the number stated there. After that you may decide whether you keep the chosen slip or you better switch and choose the other. At the end, after switching or nonswitching, you lose the game if you have chosen the slip with the smaller number. Otherwise you win. It looks like that your chance of winning is 50 %. But there exists a strategy to increase your chance slightly. Take an arbitrary probability distribution \mathbb{Q} on \mathbb{R} satisfying $\mathbb{Q}([a, b]) > 0$ for all $a < b$. Simulate a random real number z distributed according to \mathbb{Q} . If your number x at the chosen slip satisfies $x < z$, then switch. Otherwise, if $z < x$, keep the chosen slip.

Let us heuristically explain why this strategy improves your chance of winning. Suppose the two numbers on the slips are $a \in \mathbb{R}$ and $b \in \mathbb{R}$ with $a < b$. If the simulated number z satisfies $z < a$, then, no matter which of the two slips you chose, you do not switch. Hence, in this case your chance of winning remains 50 % as it was at the beginning. Similarly, if $z > b$, then you always switch, and your chance of winning remains 50 % as it was before switching. But what happens in the case $a < z < b$? If you have chosen the slip with a on it, you switch and win. Otherwise, if your choice was already the larger number b , you do not switch and you win as well. Due to the assumption about the underlying probability distribution \mathbb{Q} , no matter how big/small $a < b$ are, with probability $\mathbb{Q}([a, b]) > 0$ the simulated number z will satisfy $a < z < b$, a case where the strategy always leads to a win. Putting together all three cases, the chance of winning becomes slightly greater than 50 %. We refer to [Sam04] for a precise presentation. Note that we did not say anything about the rules for the choice of the numbers $a < b$. Recall that there is no probability distribution \mathbb{P} on \mathbb{R} such that $\mathbb{P}(\{a\}) = \mathbb{P}(\{b\})$ for all $a < b$.

Summary: In the previous section, we presented three famous examples in Probability Theory: The “Boy or Girl Paradox,” the “Secretary Problem,” and the “Envelope Exchange Paradox.” We gave full solutions and discussed some generalizations of these classical problems.

5.5 Gambler's ruin

Two players, say player A and his opponent B , play a series of independent games. Player A wins each single game with probability p , hence the success probability for B equals $1 - p$. Here and later on, we always assume $0 < p < 1$ because otherwise either A or B always win. Each time the winner gets \$1 from the loser. At the beginning, A has $a \geq 1$ dollars in his wallet, B possesses $b \geq 1$ dollars. The gamblers decide to play as long as one of them lost all of his money.

The basic question is how likely is it that A and/or B go bankrupt. To answer this question, we use the technique of random walks as presented in Example 4.1.7. There we investigated walks starting at zero. But, of course, this easily extends to walks starting at an arbitrary integer $k \in \mathbb{Z}$.

Definition 5.5.1. Given an independent sequence $(X_i)_{i \geq 1}$ with

$$\mathbb{P}\{X_i = 1\} = p \quad \text{and} \quad \mathbb{P}\{X_i = -1\} = 1 - p, \quad i = 1, 2, \dots$$

Let $S_0 = k$ and $S_n = k + X_1 + \dots + X_n$ if $n \geq 1$. Then $(S_n)_{n \geq 0}$ is a (simple) random walk starting at $k \in \mathbb{Z}$.

In this setting, player A wins if a random walk $(S_n)_{n \geq 0}$ starting at $a \geq 1$ satisfies $S_n = a + b$ for some $n \geq 0$ and, moreover, $S_j > 0$ if $0 \leq j \leq n$. Note that S_n is the amount of money which owns player A after n games. Compare Figure 5.6.

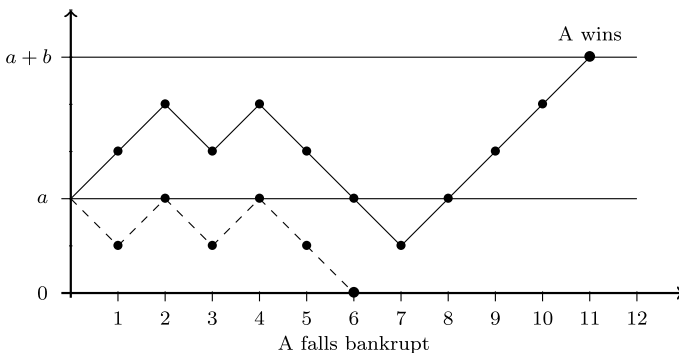


Figure 5.6: Players A and B start their series of games with $a \geq 1$ and $b \geq 1$ dollars, respectively.

Let $0 \leq k \leq a + b$ be an arbitrary integer. Set

$$A_k = \{(S_n)_{n \geq 0} \text{ starts at } k, \text{ and } \exists n \geq 0, S_n = a + b \text{ and } S_j > 0, j \leq n\}.$$

In other words, the event A_k occurs if player A starts with k dollars, at some time he reaches level $a + b$ and before that he does not go bankrupt.

The basic properties of $\mathbb{P}(A_k)$ are as follows:

$$\mathbb{P}(A_0) = 0 \quad \text{and} \quad \mathbb{P}(A_{a+b}) = 1,$$

and if $1 \leq k < a + b$, then

$$\begin{aligned} \mathbb{P}(A_k) &= \mathbb{P}(A_k|X_1 = 1)\mathbb{P}\{X_1 = 1\} + \mathbb{P}(A_k|X_1 = -1)\mathbb{P}\{X_1 = -1\} \\ &= p\mathbb{P}(A_{k+1}) + (1-p)\mathbb{P}(A_{k-1}). \end{aligned}$$

To see the last property, imagine A and B play one (their first) game and after that they start a new series of games where now, depending on the result in the first game, player A either owns $k + 1$ or $k - 1$ dollars.

Letting $x_k = \mathbb{P}(A_k)$, and setting $q = 1 - p$, for any $0 < k < a + b$ we get

$$x_{k+1} = \frac{1}{p}x_k - \frac{q}{p}x_{k-1}, \quad x_0 = 0, \quad \text{and} \quad x_{a+b} = 1. \quad (5.55)$$

So we obtained for the x_k s a linear recurrence formula of second order with two boundary conditions at $k = 0$ and $k = a + b$. The technique to solve such recurrence formulas is well known; see, for example, page 41 in [CL23]. A basic role play the zeroes or roots of the characteristic equation, which in the case of eq. (5.55) is given by

$$z^2 - \frac{1}{p}z + \frac{q}{p} = 0.$$

If $p \neq q$, that is, if $p \neq 1/2$, then this equation has two different roots which are $z_1 = 1$ and $z_2 = q/p$ (recall that $q = 1 - p$). Thus, there are constants c and d such that

$$x_k = c \cdot 1^k + d\left(\frac{q}{p}\right)^k = c + d\left(\frac{q}{p}\right)^k, \quad k = 0, \dots, a + b.$$

The boundary conditions tell us that $c + d = 0$ and $c + d(q/p)^{a+b} = 1$, hence

$$c = \frac{-1}{\left(\frac{q}{p}\right)^{a+b} - 1} \quad \text{and} \quad d = \frac{1}{\left(\frac{q}{p}\right)^{a+b} - 1},$$

leading to

$$\mathbb{P}(A_k) = x_k = \frac{\left(\frac{q}{p}\right)^k - 1}{\left(\frac{q}{p}\right)^{a+b} - 1}, \quad k = 0, \dots, a + b.$$

If $p = q = 1/2$, the characteristic equation becomes

$$z^2 - 2z + 1 = 0$$

with root $z_0 = 1$ of multiplicity 2. In this case, see, for example, page 43 in [CL23], one gets

$$x_k = c1^k + dk1^k = c + dk$$

with certain constants c and d . The boundary conditions imply $c + d \cdot 0 = 0$ as well as $c + d(a + b) = 1$, hence $c = 0$ and $d = 1/(a + b)$. So we finally conclude that

$$\mathbb{P}(A_k) = x_k = \frac{k}{a + b}, \quad k = 0, \dots, a + b.$$

Choosing in both cases $k = a$, we obtain the following result.

Proposition 5.5.2. *Suppose A and B play a series of games where every time A wins one dollar with probability p , hence B wins one dollar with probability $q = 1 - p$. If the initial amounts of money are $a \geq 1$ and $b \geq 1$, respectively, then*

$$\mathbb{P}\{B \text{ goes bankrupt}\} = \mathbb{P}\{A \text{ wins}\} = \begin{cases} \left(\frac{q}{p}\right)^a - 1 & \text{if } p \neq \frac{1}{2}, \\ \frac{q}{a+b} & \text{if } p = \frac{1}{2}. \end{cases}$$

What happens if both players start the series of games with identical amount $a > 0$? The following corollary gives the answer.

Corollary 5.5.3. *Suppose both players A and B start their games with the same amount $a > 0$. As before, p and $q = 1 - p$ are the success probabilities of players A and B , respectively. Then it follows that*

$$\mathbb{P}\{B \text{ goes bankrupt}\} = \mathbb{P}\{A \text{ wins}\} = \frac{1}{1 + \left(\frac{q}{p}\right)^a}.$$

Proof. The result is obviously true if $p = 1/2$, that is, if $p = q$. So let us assume that $p \neq q$. Hence it follows that $x := (q/p)^a \neq 1$. We apply now Proposition 5.5.2 with $b = a$ and obtain

$$\mathbb{P}\{A \text{ wins}\} = \frac{\left(\frac{q}{p}\right)^a - 1}{\left(\frac{q}{p}\right)^{2a} - 1} = \frac{x - 1}{x^2 - 1} = \frac{1}{1 + x} = \frac{1}{1 + \left(\frac{q}{p}\right)^a},$$

as asserted. □

Remark 5.5.4. The previous corollary shows the not very surprising fact that the chances of player A are less than $1/2$ whenever $q > p$, that is, if $0 < p < 1/2$. Moreover, in this case one observes the following: the bigger the initial amount $a > 0$, the less the probability for A to win.

Example 5.5.5. Suppose player A wins a single game with probability $p = 0.49$ and both players A and B start their games with the same amount of a dollars. Then Corollary 5.5.3 applies and we get

$$\mathbb{P}\{A \text{ wins}\} = \frac{1}{1 + \left(\frac{51}{49}\right)^a}.$$

For example, if $a = 50$, this probability equals 0.119175 while for $a = 100$ one gets 0.0179768. See Figure 5.7 for other probabilities with respect to the sums $a = 1, \dots, 100$.

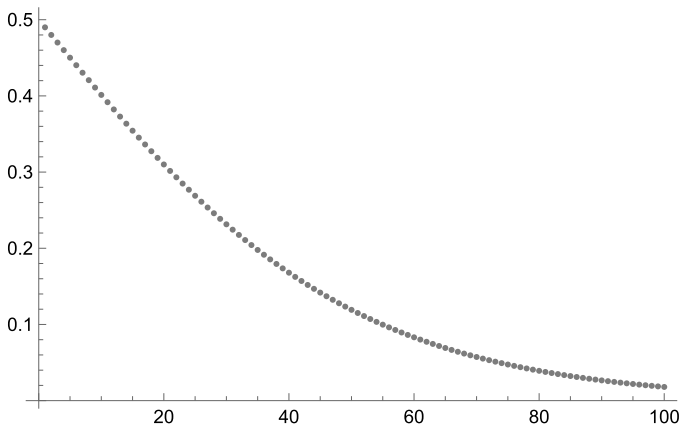


Figure 5.7: The probability that A wins when both players start with a dollars, $a = 1, \dots, 100$. Here player A wins a single game with probability $p = 0.49$.

Example 5.5.6. Let us play roulette where in every game we either win or lose \$1 (for example, put every time \$1 either on red or on black). The chance of winning is $p = 18/37$, hence $q = 19/37$. Say one stops gambling if either one had lost \$10 or if one had won \$100. So, in the previous notation $a = 10$ and $b = 100$. Hence we get

$$\mathbb{P}\{\text{Win } \$100, \text{ starting with } \$10\} = \frac{\left(\frac{19}{18}\right)^{10} - 1}{\left(\frac{19}{18}\right)^{110} - 1} = 0.00187859. \quad (5.56)$$

Does it considerably improve the chance of winning \$100 if one accepts in between a bigger loss? Not really. For example, if one goes bankrupt after loosing \$100, then the chance of winning \$100 equals 0.00446628. Note that this probability equals 1/2 in the case of a fair game. Thus, even the small disadvantage of 1/18 changes the chance of winning dramatically. See Figure 5.8 for the probabilities to win \$100 starting with \$10 up to \$30, respectively.

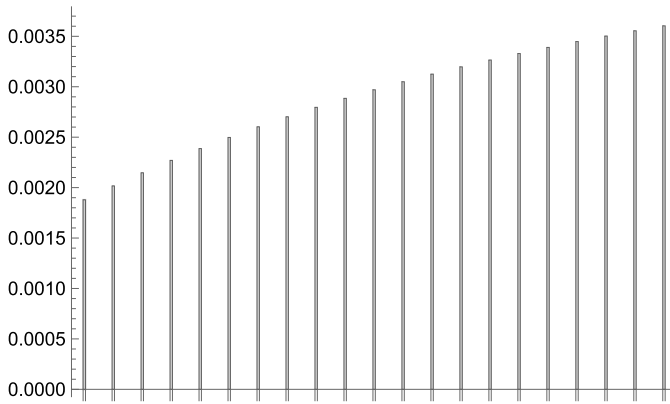


Figure 5.8: The probability to win \$100 starting with $a = 10, \dots, 30$ dollars when playing roulette. Each time one wins or loses \$1.

Another interesting numerical example is $a = b = 10$ and $p = 18/37$. That is, the game is terminated when either one has lost or won \$10. Here Corollary 5.5.3 leads to 0.368031 as probability for winning \$10. In other words, playing roulette 100 times with initial amount \$10, on average in about 37 of the cases you will win \$10, but in 63 of the cases you are going to lose your initial sum.

Let us now come back to the general case of players A and B with success probabilities p and $q = 1 - p$, owning at the beginning a and b dollars, respectively.

How likely is it that B wins? To answer this question, we use Proposition 5.5.2 but turn the tables. Interchange A and B , p and q , as well as a and b . Doing so, we obtain the following:

Proposition 5.5.7. *Suppose A and B play a series of games where every time A wins with probability p , hence B with probability $q = 1 - p$. If the initial amounts of money are a and b , respectively, then*

$$\mathbb{P}\{A \text{ goes bankrupt}\} = \mathbb{P}\{B \text{ wins}\} = \begin{cases} \frac{(\frac{p}{q})^b - 1}{(\frac{p}{q})^{a+b} - 1} & \text{if } p \neq \frac{1}{2}, \\ \frac{b}{a+b} & \text{if } p = \frac{1}{2}. \end{cases}$$

An interesting question remained unanswered until now: is it possible that the series of games between A and B lasts forever? In other words, may it happen that neither A nor B wins?

The following result shows that the answer is negative.

Proposition 5.5.8. *Under the previous assumptions, it follows that*

$$\mathbb{P}\{A \text{ wins}\} + \mathbb{P}\{B \text{ wins}\} = 1.$$

In particular, this implies

$$\mathbb{P}\{\text{The game lasts forever}\} = 0.$$

Proof. If $p = 1/2$, by Propositions 5.5.2 and 5.5.7, one gets

$$\mathbb{P}\{A \text{ wins}\} + \mathbb{P}\{B \text{ wins}\} = \frac{a}{a+b} + \frac{b}{a+b} = 1,$$

completing the proof in this case.

Thus, let us assume now $p \neq q$. To simplify the calculations, set $x = q/p$ and $y = 1/x = p/q$. Note that both numbers are by assumption different from 1. With these notations, Propositions 5.5.2 and 5.5.7 may be written as

$$\mathbb{P}\{A \text{ wins}\} = \frac{x^a - 1}{x^{a+b} - 1} = \frac{y^b - y^{a+b}}{1 - y^{a+b}} \quad \text{and} \quad \mathbb{P}\{B \text{ wins}\} = \frac{y^b - 1}{y^{a+b} - 1}.$$

Consequently, the assertion follows from

$$\frac{y^b - y^{a+b}}{1 - y^{a+b}} + \frac{y^b - 1}{y^{a+b} - 1} = \frac{y^{a+b} - y^b + y^b - 1}{y^{a+b} - 1} = \frac{y^{a+b} - 1}{y^{a+b} - 1} = 1. \quad \square$$

In view of Proposition 5.5.8, the following natural question arises: Let $T_{a,b}$ be the number of rounds that A and B play. That is, given a random walk $(S_n)_{n \geq 0}$ starting at zero, for some $a, b \in \mathbb{N}$ set

$$T_{a,b} = \min\{n \geq 0 : S_n = -a \text{ or } S_n = b\}. \quad (5.57)$$

What is the expected value of $T_{a,b}$? In other words, how long does the series of games last on average. The answer is as follows (for a proof, we refer to [Sti03] or [Fel68]; the basic idea is similar to that used in the proof of Proposition 5.5.2, namely conditioning on the first step which leads to a linear recurrence formula for $\mathbb{E} T_{a,b}$).

Proposition 5.5.9. *Let a, b, p , and $q = 1 - p$ be as before. If $T_{a,b}$ denotes the number of rounds before one of the players goes bankrupt, then*

$$\mathbb{E} T_{a,b} = \begin{cases} \frac{a}{q-p} - \frac{a+b}{q-p} \frac{(\frac{q}{p})^a - 1}{(\frac{q}{p})^{a+b} - 1} & \text{if } p \neq \frac{1}{2}, \\ a \cdot b & \text{if } p = \frac{1}{2}. \end{cases} \quad (5.58)$$

Remark 5.5.10. In the case $b = a$, the first formula in eq. (5.58) simplifies to

$$\mathbb{E} T_{a,a} = \frac{a}{q-p} \left[\frac{(\frac{q}{p})^a - 1}{(\frac{q}{p})^a + 1} \right], \quad p \neq \frac{1}{2}.$$

The proof goes along the same lines as that of Corollary 5.5.3. Furthermore, as can be easily seen, the expected value of $T_{a,a}$ does not change if one interchanges p and q , or,

equivalently, players A and B . This is, of course, because the length of the game does not depend on who players A and B are, provided both start with the same amount of money.

Note that in the case $p \neq 1/2$, it follows that

$$\lim_{a \rightarrow \infty} \frac{1}{a} \mathbb{E} T_{a,a} = \frac{1}{|p - q|}.$$

So in the long run, the expected time of the game is of order $a/|p - q|$. Compare this with the case $p = 1/2$ where the expected time behaves like a^2 .

Example 5.5.11. If the game is fair and both players start either with \$50 or \$100, then on average the gamblers have to play either 2500 or 10,000 rounds before there is a winner. The situation changes drastically if the success probability of one player is diminished to 0.49. That is, the game is “almost” fair. In this case the average number of necessary rounds equals either 1904.13 or 4820.23, respectively.

Let us finally treat a related problem, sometimes called “the monkey at the cliff.” A monkey is standing one step from the edge of a cliff and takes repeated independent steps; forward, with probability p , or backward, with probability $q = 1 - p$. What is the probability that the monkey, sooner or later, will fall off the cliff?

The mathematical formulation is as follows: let $(S_n)_{n \geq 0}$ be a random walk starting at zero jumping with probability p to the right and with probability $q = 1 - p$ to the left. How likely is it that there exists an $n \geq 1$ such that $S_n = 1$. More generally, one may ask for the existence of an $n \geq 1$ with $S_n = b$ for a given integer $b \geq 1$. The answer is as follows:

Proposition 5.5.12. *Let $(S_n)_{n \geq 1}$ be as before. Then for any integer $b \geq 1$ it follows that*

$$\mathbb{P}\{S_n = b \text{ for some } n \geq 1\} = \begin{cases} 1 & \text{if } p \geq \frac{1}{2}, \\ \left(\frac{p}{q}\right)^b & \text{if } p < \frac{1}{2}. \end{cases}$$

Proof. Fix $b \geq 1$. Given $a \geq 1$, define events B_a as follows: B_a occurs if there is an $n \geq 1$ for which $S_n = b$ and, at the same time, $S_j > -a$ if $1 \leq j < n$. Then $A_1 \subseteq A_2 \subseteq \dots$ and, moreover,

$$\{S_n = b \text{ for some } n \geq 1\} = \bigcup_{a=1}^{\infty} B_a.$$

To see this, suppose $S_n = b$ and choose $a \geq 1$ such that $\min_{1 \leq j \leq n} S_j > -a$.

Hence, by the continuity of probability measures from below (see property (6) in Proposition 1.2.1), we obtain

$$\mathbb{P}\{S_n = b \text{ for some } n \geq 1\} = \lim_{a \rightarrow \infty} \mathbb{P}(B_a).$$

Now Proposition 5.5.2 applies and leads to

$$\mathbb{P}\{S_n = b \text{ for some } n \geq 1\} = \lim_{a \rightarrow \infty} \begin{cases} \frac{(\frac{q}{p})^a - 1}{(\frac{q}{p})^{a+b} - 1} & \text{if } p \neq \frac{1}{2}, \\ \frac{a}{a+b} & \text{if } p = \frac{1}{2}. \end{cases}$$

Of course,

$$\lim_{a \rightarrow \infty} \frac{a}{a+b} = 1 \quad \text{and} \quad \lim_{a \rightarrow \infty} \frac{(\frac{q}{p})^a - 1}{(\frac{q}{p})^{a+b} - 1} = 1 \quad \text{if } \frac{q}{p} < 1.$$

Recall that $x^a \rightarrow 0$ provided that $0 < x < 1$.

It remains to investigate the case $q > p$ or, equivalently, $p < 1/2$. As before set $x = q/p > 1$ and $y = 1/x < 1$. Doing so, we get

$$\lim_{a \rightarrow \infty} \frac{x^a - 1}{x^{a+b} - 1} = \lim_{a \rightarrow \infty} \frac{y^b - y^{a+b}}{1 - y^{a+b}} = y^b = \left(\frac{p}{q}\right)^b,$$

which completes the proof in the remaining case. \square

Remark 5.5.13. Proposition 5.5.12 asserts that in the case $p \geq 1/2$, the monkey will fall off the cliff with probability one, even if it is not only 1 but $b > 1$ steps away from the cliff. On the other hand, if $p < 1/2$ and the monkey is b steps away from the cliff, then with probability $1 - (p/q)^b$ the monkey will be safe. Since in this case $p/q < 1$, hence $(p/q)^b \rightarrow 0$ as $b \rightarrow \infty$, the situation of the monkey improves considerably as soon as it is further away from the cliff.

Still another way to formulate Proposition 5.5.12 is as follows. Say player A has an unlimited amount of money while his opponent starts with b dollars. Then A will win with probability one provided his success probability p satisfies $p \geq 1/2$. On the other hand, in the case of $p < 1/2$ his chance of winning equals $(p/q)^b < 1$.

Example 5.5.14. Let us investigate how likely it is to win $b \geq 1$ dollars in a roulette provided one has an unlimited amount of money. As before, every time the chance to win one dollar is $p = 18/37$ while one loses one dollar with probability $q = 19/37$. Hence, in this case we obtain

$$\mathbb{P}\{\text{Win } \$b\} = \left(\frac{18}{19}\right)^b.$$

For example, the chance to win \$100 possessing an infinite amount of money equals

$$\mathbb{P}\{\text{Win } \$100\} = \left(\frac{18}{19}\right)^{100} \approx 0.00448632.$$

Compare this result with eq. (5.56) where we got 0.00187859 for the probability to win \$100 when starting with an initial amount of \$10 or with the probability 0.00446628

when starting with \$100. So one sees, in order to win \$100, it does not make a big difference whether one starts to play with \$100 or with an unlimited amount of money. The result will be the same in both cases: very likely one is going to lose a lot of money. Compare Figure 5.9 for the probabilities to win $b = 1, \dots, 30$ dollars playing roulette possessing an unlimited amount of money.

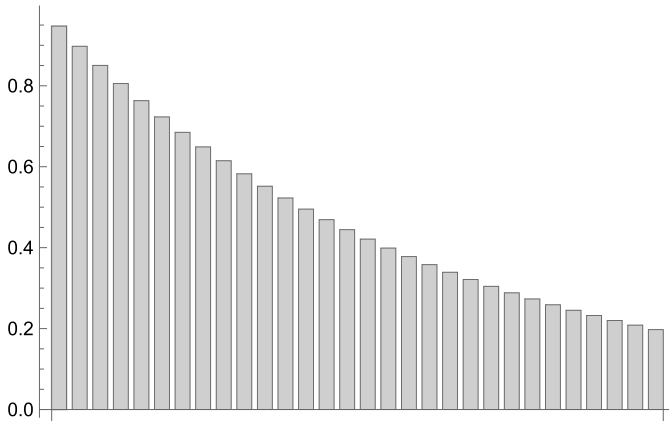


Figure 5.9: The probability to win $b = 1, \dots, 30$ dollars playing roulette possessing an unlimited amount of money.

Suppose now one does not bet \$1 each time, but \$10. How likely is it now to win \$100? The answer is as follows: the likelihood to win \$100 by \$10 steps coincides with that to win \$10 by steps of size \$1. Hence, the probability of this event equals $(18/19)^{10} \approx 0.582357$.

Remark 5.5.15. It might be of interest to compare this with the probability 0.368031 in the case of an initial deposit of \$10. This tells us that it is not unlikely to lose at some time more than \$10 before one finally wins \$10. For example, the chance to win \$10 becomes for the first time greater than $1/2$ if one starts gambling with an amount of \$24. Then the probability to win at some time \$10 without going bankrupt equals 0.503344.

But note that this does not mean that one has an advantage. In the case of success, one wins \$10 while one loses \$24 in the case of failure. Thus, on average there will be a loss, no matter how big the initial amount was.

Remark 5.5.16. The symmetric (fair) case $p = 1/2$ is of special interest. By symmetry, given $b \geq 1$, with probability one there also exists an $n \geq 1$ such that $S_n = -b$. Thus, with probability one, a symmetric random walk attains any value in \mathbb{Z} . See Example 7.2.15 for further asymptotic properties of symmetric walks.

In view of Proposition 5.5.12 and the previous remark, the following question arises: let $(S_n)_{n \geq 1}$ be a random walk starting at zero. Given $b \in \mathbb{N}$, how long does it take on average before the walk reaches level b ? To make it more precise, given $b \geq 1$, let

$$T_b = \begin{cases} \min\{n \geq 0 : S_n = b\} & \text{if there is an } n \geq 0 \text{ with } S_n = b, \\ \infty & \text{otherwise.} \end{cases}$$

Then Proposition 5.5.12 may be rephrased as follows:

$$\mathbb{P}\{T_b < \infty\} = \begin{cases} 1 & \text{if } p \geq \frac{1}{2}, \\ \left(\frac{p}{q}\right)^b & \text{if } p < \frac{1}{2}. \end{cases}$$

Proposition 5.5.17. *Let $(S_n)_{n \geq 1}$ be a random walk starting at zero. Given $b \in \mathbb{N}$,*

$$\mathbb{E} T_b = \begin{cases} \frac{b}{p-q} & \text{if } p > \frac{1}{2}, \\ \infty & \text{if } p \leq \frac{1}{2}. \end{cases}$$

Proof. If $p < 1/2$, then $\mathbb{P}\{T_b = \infty\} = 1 - (p/q)^b > 0$, hence $\mathbb{E} T_b = \infty$ as asserted.

The case $1/2 \leq p$ is more involved and may be found in [Sti03, Section 5.6]. Basic ingredient is the so-called hitting time theorem asserting

$$\mathbb{P}\{T_b = n\} = \frac{b}{n} \mathbb{P}\{S_n = b\} \quad \text{and} \quad \mathbb{E} T_b = b \sum_{n=1}^{\infty} \mathbb{P}\{S_n = b\}, \quad b \geq 1. \quad (5.59)$$

A heuristic proof of Proposition 5.5.17 (without using the hitting time theorem) can be given by using eq. (5.58). Assume we know that

$$\mathbb{E} T_b = \lim_{a \rightarrow \infty} \mathbb{E} T_{a,b}$$

(which is true and can be made precise). Then, if $p > 1/2$, hence $q/p < 1$, it follows that

$$\mathbb{E} T_b = \lim_{a \rightarrow \infty} \left[\frac{a}{q-p} - \frac{a+b}{q-p} \frac{\left(\frac{q}{p}\right)^a - 1}{\left(\frac{q}{p}\right)^{a+b} - 1} \right] = \frac{b}{p-q}.$$

The case $p = 1/2$ is even easier to handle and follows by evaluating the infinite sum in the right-hand formula of (5.59) or from

$$\mathbb{E} T_b = \lim_{a \rightarrow \infty} \mathbb{E} T_{a,b} = \lim_{a \rightarrow \infty} a \cdot b = \infty. \quad \square$$

Remark 5.5.18. Most interesting in Proposition 5.5.17 is the case $p = 1/2$. Assume players A and B play a series of fair games. Player A has an unlimited amount of money while B starts with \$1. Then A wins with probability 1, but on average it takes an arbitrarily long time until B goes bankrupt.

Similarly, in the symmetric case, for sure the monkey will fall off the cliff, but on average it takes a lot of time before this happens.

Final remark: There exist many other interesting results about random walks not included in the present book. For example, what can be said about the behavior of $\max\{S_k : 0 \leq k \leq n\}$? How are the zeroes of a symmetric walk $(S_n)_{n \geq 0}$ distributed? How about recurrence or transience? How many changes of signs exist until a given time n ? Or what happens if, as in our case, the barriers are not absorbing but reflecting? Neither did we treat random walks in the more general setting of jumping particles in \mathbb{Z}^d . We refer to [Fel68], to [Rev13], or to [Sti03] for further reading about this highly interesting topic.

But let us shortly discuss one property of the symmetric walk, which, in our opinion, is very surprising. Let $(S_n)_{n \geq 0}$ be a symmetric random walk starting at zero. If

$$L_n^+ = |\{k \leq n : S_k > 0\}|, \quad (5.60)$$

then L_n^+ is the number of the times until n where the symmetric walk is located in the positive half-space or, equivalently, where player A is ahead of player B provided both players start with an unlimited amount of money. Probably everybody will guess that for large n , player A will be about half of the time ahead B , and during the other half, B will be ahead A . Recall that the walk is symmetric, hence jumps to the right are as likely as those to the left.

The following result shows that this is not so. It is much more likely that most of the time one of the players is ahead of the other. To verify this, one investigates the proportion L_n^+/n of times where the walk is above zero. Here the following holds (see [Sti03, Section 6.8]):

Proposition 5.5.19 (Arcsine law for random walks). *Let L_n^+ be defined by (5.60). Then for any $0 < t < 1$ it holds that*

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \frac{L_n^+}{n} \leq t \right\} = \frac{2}{\pi} \arcsin(\sqrt{t}) = \frac{1}{\pi} \int_0^t \frac{1}{\sqrt{x(1-x)}} dx.$$

Thus, the random variables $(L_n^+/n)_{n \geq 1}$ converge in distribution, that is, in the sense of Definition 7.2.1, to an arcsine distribution. Recall that the arcsine distribution was introduced in Definition 1.6.35.

Corollary 5.5.20. *If $0 \leq a < b \leq 1$, then for sufficiently large n we have*

$$\mathbb{P} \left\{ a \leq \frac{L_n^+}{n} \leq b \right\} \approx \frac{2}{\pi} [\arcsin(\sqrt{b}) - \arcsin(\sqrt{a})].$$

The numerical values in the cases $a = 0.1$, $a = 0.2$, and $a = 0.5$ are

$\mathbb{P}\{0.05 \leq L_n^+/n \leq 0.15\}$	$\mathbb{P}\{0.15 \leq L_n^+/n \leq 0.25\}$	$\mathbb{P}\{0.45 \leq L_n^+/n \leq 0.55\}$
≈ 0.1096	≈ 0.0802	≈ 0.06377

These values, as well as Figure 5.10, tell us that it is much more likely that L_n^+/n is near zero or one than near 0.5. Recall that L_n^+/n is the proportion of those times $k \leq n$ where $S_k > 0$. Hence, the event that L_n^+/n near 0.5 occurs if the walk is about half of the time positive and the other half it is negative.

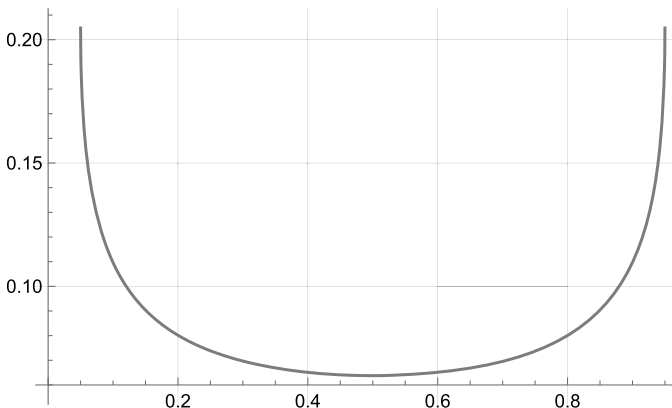


Figure 5.10: The approximate probabilities $\mathbb{P}\{a - 0.05 \leq L_n^+/n \leq a + 0.05\}$ with $a \in [0.05, 0.95]$.

Summary: Two persons A and B play a series of games as long as one of them goes bankrupt. Hereby, in every single game the loser has to pay \$1 to the winner. The describing mathematical model is a random walk starting at zero and with absorbing barriers at $-a$ and b . Here a and b are the initial amounts of A and B , respectively. Equivalently, one may regard a random walk starting at a and with barriers at 0 and $a + b$. Let p be the success probability of player A , thus $q = 1 - p$ is that of player B . Then the basic result asserts

$$\mathbb{P}\{B \text{ goes bankrupt}\} = \mathbb{P}\{A \text{ wins}\} = \begin{cases} \frac{(\frac{q}{p})^a - 1}{(\frac{q}{p})^{a+b} - 1} & \text{if } p \neq \frac{1}{2}, \\ \frac{a}{a+b} & \text{if } p = \frac{1}{2}. \end{cases}$$

5.6 Problems

Problem 5.1.

- Put successively and independently of each other n particles into N boxes. Thereby, each box is equally likely. How many boxes remain empty on average?
Hint: Define random variables X_1, \dots, X_N as follows: set $X_i = 1$ if box i remains empty, and $X_i = 0$ otherwise.
- Fifty persons write randomly (according to the uniform distribution), and independently of each other, one of the 26 letters in the alphabet on a sheet of paper. On average, how many different letters appear?

3. In a factory with $N \geq 1$ employees, every day of the year on which one of the employees has a birthday is a holiday. Let E_N be the expected number of working days, that is, the expected number of days which are not a holiday. For which $N \geq 1$ does $N \cdot E_N$ (the expected total working time) become maximal? Hereby one assumes that all 365 days of the year are equally likely to be birthdays.

Problem 5.2 (A. E. Lawrance, 1969). An urn contains eight white balls and two black. Choose one after another a ball without replacing the chosen one. Let $1 \leq r \leq 9$ be the number of that choice where for the first time a black ball occurs. Which number r is most likely for the appearance of the first black ball? Evaluate the average value over all possible numbers $r \leq 9$.

Answer the same questions in the case that $n \geq 2$ balls are in the urn where 2 are black and $n - 2$ are white.

Problem 5.3 (De Moivre, 1756). A man rolls a fair die six times. He gets an amount of M francs, every time he

1. rolls a “1” or
2. if he rolls at least one “1.”

Evaluate in both cases the expected amount of money he gets.

Problem 5.4. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Given (not necessarily disjoint) events A_1, \dots, A_n in \mathcal{A} and real numbers $\alpha_1, \dots, \alpha_n$, define $X : \Omega \rightarrow \mathbb{R}$ by¹⁶

$$X := \sum_{j=1}^n \alpha_j \mathbb{1}_{A_j}.$$

1. Why is X a random variable?
2. Prove

$$\mathbb{E}X = \sum_{j=1}^n \alpha_j \mathbb{P}(A_j) \quad \text{and} \quad \mathbb{V}X = \sum_{i,j=1}^n \alpha_i \alpha_j [\mathbb{P}(A_i \cap A_j) - \mathbb{P}(A_i)\mathbb{P}(A_j)].$$

How does $\mathbb{V}X$ simplify for independent events A_1, \dots, A_n ?

Problem 5.5. Suppose a fair “die” has k faces labeled by the numbers from 1 to k .

1. How often one has to roll the die on the average before the first “1” shows up?
2. Suppose one rolls the die exactly k times. Let p_k be the probability that “1” appears exactly once and q_k is the probability that “1” shows up at least once. Compute p_k and q_k and determine their behavior as $k \rightarrow \infty$, that is, find $\lim_{k \rightarrow \infty} p_k$ and $\lim_{k \rightarrow \infty} q_k$.

¹⁶ For the definition of indicator functions $\mathbb{1}_{A_i}$, see eq. (3.21).

Problem 5.6.

1. Let X be a random variable with values in $\mathbb{N}_0 = \{0, 1, 2, \dots\}$. Prove that

$$\mathbb{E}X = \sum_{k=1}^{\infty} \mathbb{P}\{X \geq k\}.$$

2. Suppose now that X is continuous with $\mathbb{P}\{X \geq 0\} = 1$. Verify

$$\sum_{k=1}^{\infty} \mathbb{P}\{X \geq k\} \leq \mathbb{E}X \leq 1 + \sum_{k=1}^{\infty} \mathbb{P}\{X \geq k\}.$$

Problem 5.7. Let X be an \mathbb{N}_0 -valued random variable with

$$\mathbb{P}\{X = k\} = q^{-k}, \quad k = 1, 2, \dots$$

for some $q \geq 2$.

(a) Why do we have to suppose $q \geq 2$, although $\sum_{k=1}^{\infty} q^{-k} < \infty$ for $q > 1$?

(b) Determine $\mathbb{P}\{X = 0\}$?

(c) Compute $\mathbb{E}X$ by the formula in Problem 5.6.

(d) Compute $\mathbb{E}X$ directly by $\mathbb{E}X = \sum_{k=0}^{\infty} k \mathbb{P}\{X = k\}$.

Problem 5.8. Two independent random variables X and Y with finite third moment satisfy $\mathbb{E}X = \mathbb{E}Y = 0$. Prove that then

$$\mathbb{E}(X + Y)^3 = \mathbb{E}X^3 + \mathbb{E}Y^3.$$

Problem 5.9. A random variable X is Pois_λ -distributed for some $\lambda > 0$. Evaluate

$$\mathbb{E}\left(\frac{1}{1+X}\right) \quad \text{and} \quad \mathbb{E}\left(\frac{X}{1+X}\right).$$

Problem 5.10. In a lottery, 6 of 49 numbers are randomly chosen. Let X be the largest number of the 6. Show that

$$\mathbb{E}X = \frac{6 \cdot 43!}{49!} \sum_{k=6}^{49} k(k-1)(k-2)(k-3)(k-4)(k-5) = 42.8571.$$

Evaluate $\mathbb{E}X$ if X is the smallest number of the 6 chosen.

Hint: Either one modifies the calculations for the maximal value suitably or one reduces the second problem to the first by an easy algebraic operation.

Problem 5.11. A fair coin is labeled by “0” on one side and with “1” on the other. Toss it four times. Let X be the sum of the first two tosses and Y be the sum of all four. Determine the joint distribution of X and Y . Evaluate $\text{Cov}(X, Y)$, as well as $\rho(X, Y)$.

Problem 5.12. In an urn there are five balls, two labeled by “0” and three by “1.” Choose two balls without replacement. Let X be the number on the first ball and Y that on the second.

1. Determine the distribution of the random vector (X, Y) and its marginal distributions.
2. Compute $\rho(X, Y)$.
3. Which distribution does $X + Y$ possess?

Problem 5.13. Among 40 students there are 30 men and 10 women. Also, 25 of the 30 men and 8 of the 10 women passed an exam successfully. Choose randomly, according to the uniform distribution, one of the 40 students. Let $X = 0$ if the chosen person is a man, and $X = 1$ if it is a woman. Furthermore, set $Y = 0$ if the person failed the exam, and $Y = 1$ if she or he passed.

1. Find the joint distribution of X and Y .
2. Are X and Y independent? If not, evaluate $\text{Cov}(X, Y)$.
3. Are X and Y negatively or positively correlated? What does it express, when X and Y are positively or negatively correlated?

Problem 5.14. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Prove, for any two events A and B in \mathcal{A} , the estimate

$$|\mathbb{P}(A \cap B) - \mathbb{P}(A)\mathbb{P}(B)| \leq \frac{1}{4}.$$

Is it possible to improve the upper bound $\frac{1}{4}$?

Problem 5.15 (Problem of Luca Pacioli in 1494; the first correct solution was found by Blaise Pascal in 1654). Two players, say A and B , are playing a fair game consisting of several rounds. The first player who wins six rounds wins the game and the stakes of 20 taler¹⁷ that have been bet throughout the game. However, one day the game is interrupted and must be stopped. If player A has won five rounds and player B has won three rounds, how should the stakes be divided fairly among the players?

Problem 5.16 (B. Pascal, 1654). Three players, say A , B , and C , play a series of fair games. Whoever first wins three games is the winner. One day the series of games had to be stopped before one of the three players had won three games. Player A still needs one win, B and C still need two wins each. How to distribute the stakes in this case in a fair way?

Problem 5.17. In Example 5.1.49, we computed the average number of necessary purchases to get all n pictures. Let m be an integer with $1 \leq m < n$. How many purchases are necessary on average to possess m of the n pictures?

¹⁷ Former German currency, root of the word “dollar.”

For n even, choose $m = n/2$, and for n odd take $m = (n - 1)/2$. Let M_n be the average number of purchases to get m pictures, that is, to get half of the pictures. Determine

$$\lim_{n \rightarrow \infty} \frac{M_n}{n}.$$

Hint: Use eq. (5.29).

Problem 5.18. Compute $\mathbb{E}|X|^{2n+1}$ for a standard normal distributed X and $n = 0, 1, \dots$

Problem 5.19. Suppose X has the density

$$p(x) = \begin{cases} 0 & \text{if } x < 1, \\ c_\alpha x^\alpha & \text{if } x \geq 1, \end{cases}$$

for some $\alpha < -1$.

1. Determine c_α such that p is a density.
2. For which $n \geq 1$ does X possess an n th moment?

Problem 5.20. Let U be uniform distributed on an interval $[\alpha, \beta]$. Show that for $n \geq 1$,

$$\mathbb{E} U^n = \frac{\beta^n + \alpha\beta^{n-1} + \dots + \alpha^{n-1}\beta + \alpha^n}{n+1}.$$

Problem 5.21. Let X_1, \dots, X_n be random variables with finite second moment and with $\mathbb{E}X_j = 0$. Show that

$$\mathbb{E}[X_1 + \dots + X_n]^2 = \sum_{i,j=1}^n \text{Cov}(X_i, X_j) = \sum_{j=1}^n \mathbb{V}X_j + 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j).$$

Problem 5.22. Show that

$$\mathbb{E}X = \frac{nM}{N}$$

for a hypergeometrically distributed random variable X with

$$\mathbb{P}\{X = m\} = \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad m = 0, \dots, n.$$

Problem 5.23. Let X be $\mathcal{N}(0, 1)$ -distributed. Determine $\mathbb{V}X^3$ and $\mathbb{V}X^4$.

Problem 5.24. Given a nonnegative random variable X , define φ_X from $[0, 1]$ to $[0, 1]$ by $\varphi_X(t) = \mathbb{E} t^X$. Then φ_X is called the **generating function** of X (see [GS20, Section 5.1]).

- (1) Suppose X has values in \mathbb{N}_0 . Show that, if $t \geq 0$, then this “new” definition of the generating function coincides with that given in Problem 4.2.

- (2) Let X_1, \dots, X_n be independent and nonnegative. For $\alpha_j \geq 0, 1 \leq j \leq n$, let

$$X = \alpha_1 X_1 + \dots + \alpha_n X_n.$$

Prove

$$\varphi_X(t) = \varphi_{X_1}(t^{\alpha_1}) \cdots \varphi_{X_n}(t^{\alpha_n}).$$

- (3) Find φ_X for an exponentially distributed X .

Problem 5.25. Complete the arguments stated in Remark 5.4.4. That is, argue why the following questions are equivalent:

1. Toss a fair coin n times and choose after that at random a number $1 \leq j \leq n$. Suppose the j th toss was a “1.” What is the probability that under this condition the observed sequence has k “1”s for some $1 \leq k \leq n$?
2. One tosses a fair coin $n - 1$ times. How likely is the appearance of $k - 1$ “1”s?

6 Normally distributed random vectors

6.1 Representation and density

In Example 3.4.3 we considered a two-dimensional random vector (X_1, X_2) , where X_1 was the height of a randomly chosen person and X_2 was his weight. From experience and in view of the central limit theorem (cf. Section 7.2), it is quite reasonable to assume that X_1 and X_2 are normally distributed. Suppose we are able to determine their expected values and their variances. However, this is not sufficient to describe the experiment. Why? The random variables X_1 and X_2 are surely dependent, and the most interesting problem is to describe their degree of dependence. This cannot be done based only on the knowledge of their distributions. What we really need to know is their joint distribution. Therefore, we not only have to suppose X_1 and X_2 to be normal, but the generated vector (X_1, X_2) has to be as well.

But what does it mean that a random vector is normally distributed? This section is devoted to answer this and related questions.

Let us first recall the univariate case, investigated in Example 4.2.2 and in the subsequent Proposition 4.2.3. The main observation was that a random variable Y is normally distributed if and only if it may be written as

$$Y = aX + \mu \tag{6.1}$$

for some $a \neq 0$, $\mu \in \mathbb{R}$, and a standard normal random variable X .

Let now $\vec{Y} = (Y_1, \dots, Y_n)$ be an n -dimensional random vector. We want to represent it in the same way as Y in eq. (6.1). Consequently, we have to replace X by a multivariate standard normal vector and the function $x \mapsto ax + \mu$ by a suitable mapping from \mathbb{R}^n to \mathbb{R}^n . But which kind of mapping should this be and what is an n -dimensional standard normal vector?

Let us begin by answering the second question. Therefore, recall the definition of the multivariate standard normal distribution $\mathcal{N}(0, 1)^{\otimes n}$ introduced in Definition 1.9.21. This probability measure on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ acts as follows: if $B \in \mathcal{B}(\mathbb{R}^n)$, then its probability equals

$$\begin{aligned} \mathcal{N}(0, 1)^{\otimes n}(B) &= \frac{1}{(2\pi)^{n/2}} \int_B e^{-|x|^2/2} dx \\ &= \frac{1}{(2\pi)^{n/2}} \underbrace{\int \dots \int}_B e^{-(x_1^2 + \dots + x_n^2)/2} dx_n \dots dx_1. \end{aligned}$$

Thus, a random vector \vec{X} should be standard normally distributed whenever its probability distribution is $\mathcal{N}(0, 1)^{\otimes n}$. Let us formulate this as a definition.

Definition 6.1.1. A random vector $\vec{X} = (X_1, \dots, X_n)$ is **standard normally distributed** (or is standard normal) if its probability distribution satisfies $\mathbb{P}_{\vec{X}} = \mathcal{N}(0, 1)^{\otimes n}$.

To make this definition more descriptive, let us state some equivalent properties.

Proposition 6.1.2. For a random vector $\vec{X} = (X_1, \dots, X_n)$, the following are equivalent:

1. \vec{X} is standard normal.
2. If $B \in \mathcal{B}(\mathbb{R}^n)$, then

$$\mathbb{P}\{\vec{X} \in B\} = \frac{1}{(2\pi)^{n/2}} \int_B e^{-|x|^2/2} dx.$$

3. The coordinate mappings X_1, \dots, X_n are (univariate) standard normally distributed and independent. That is, for all $t_j \in \mathbb{R}$, $1 \leq j \leq n$,

$$\begin{aligned} \mathbb{P}\{X_1 \leq t_1, \dots, X_n \leq t_n\} &= \mathbb{P}\{X_1 \leq t_1\} \cdots \mathbb{P}\{X_n \leq t_n\} \\ &= \left(\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{t_1} e^{-x_1^2/2} dx_1 \right) \cdots \left(\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{t_n} e^{-x_n^2/2} dx_n \right). \end{aligned}$$

Proof. Taking into account the definition of $\mathcal{N}(0, 1)^{\otimes n}$, this is an immediate consequence of Propositions 3.6.5 and 3.6.20. Compare also the considerations in Example 3.6.24. \square

Remark 6.1.3. The density of the n -dimensional standard normal distribution possesses nice properties: it attains its maximal value $(2\pi)^{-n/2}$ at zero, it is invariant under rotations of the arguments and its level sets are circles. See Figure 6.1 for the graph of this density in the case $n = 2$.

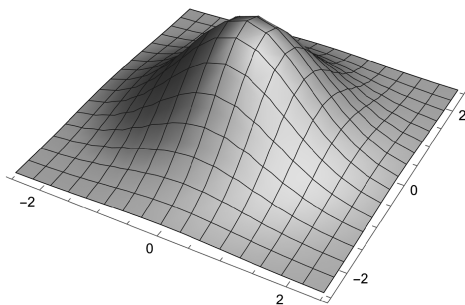


Figure 6.1: The density of a 2-dimensional standard normal vector.

An adequate substitute for $x \mapsto ax + \mu$ in representation (6.1) is still undetermined. Which mappings in \mathbb{R}^n should be considered?

Observe that $x \mapsto ax + \mu$ is affine linear from \mathbb{R} to \mathbb{R} . The counterpart in \mathbb{R}^n is of the form $x \mapsto Ax + \mu$, where A is a linear mapping in \mathbb{R}^n and $\mu \in \mathbb{R}^n$. Linear mappings in \mathbb{R}^n are described by $n \times n$ matrices $A = (\alpha_{ij})_{i,j=1}^n$ and act as follows:

$$Ax = \left(\sum_{j=1}^n \alpha_{1j}x_j, \dots, \sum_{j=1}^n \alpha_{nj}x_j \right), \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n.$$

Consequently, the suitable generalization of $x \mapsto ax + \mu$ is the mapping $x \mapsto Ax + \mu$ with an $n \times n$ matrix A and $\mu \in \mathbb{R}^n$. The condition $a \neq 0$ transfers to $\det(A) \neq 0$ or, equivalently, A has to be *regular*, that is, the generated mapping is one-to-one from \mathbb{R}^n onto \mathbb{R}^n . Here and in the sequel we will use results and notations as presented in Section A.4.

Now we are in position to define normally (distributed) random vectors.

Definition 6.1.4. A random vector \vec{Y} is said to be **normally distributed** (or simply, normal) provided there exists a regular $n \times n$ matrix A and a vector $\mu \in \mathbb{R}^n$ such that

$$\vec{Y} = A\vec{X} + \mu \tag{6.2}$$

for some standard normal \vec{X} .

Remark 6.1.5. Let us reformulate Definition 6.1.4 due to its importance. A random vector $\vec{Y} = (Y_1, \dots, Y_n)$ is normally distributed if and only if there exists a regular matrix $A = (\alpha_{ij})_{i,j=1}^n$ and a vector $\mu = (\mu_1, \dots, \mu_n)$ such that

$$Y_i = \sum_{j=1}^n \alpha_{ij}X_j + \mu_i, \quad 1 \leq i \leq n,$$

with X_1, \dots, X_n independent $\mathcal{N}(0, 1)$ -distributed.

Example 6.1.6. Suppose the three-dimensional random vector $\vec{Y} = (Y_1, Y_2, Y_3)$ is defined by

$$\begin{aligned} Y_1 &= 2X_1 + X_2 - X_3 + 4, & Y_2 &= X_1 - 2X_2 + X_3 - 2, & \text{and} \\ Y_3 &= X_1 - 2X_3 + 5 \end{aligned}$$

with $\mathcal{N}(0, 1)$ -distributed independent X_1, X_2, X_3 . Then \vec{Y} is normally distributed. Observe that it may be represented in the form of eq. (6.2) with A given by

$$A = \begin{pmatrix} 2 & 1 & -1 \\ 1 & -2 & 1 \\ 1 & 0 & -2 \end{pmatrix}$$

and with $\mu = (4, -2, 5)$. Moreover, we have $\det(A) = 9$, hence A is regular.

Remark 6.1.7. If the n -dimensional vector \vec{Y} is represented as $\vec{Y} = A\vec{X} + \mu$ with \vec{X} standard normal, $\mu \in \mathbb{R}^n$ and a *nonregular* $n \times n$ -matrix A , then \vec{Y} may also be regarded as

normal, yet in a more general setting. In this case it follows that $\mathbb{P}\{\vec{Y} \in \text{range}(A) + \mu\} = 1$ where $\text{range}(A)$ is a strict subspace of \mathbb{R}^n . For example, if $Y_1 = X_1 + X_2$ and $Y_2 = -X_1 - X_2$, then $\mathbb{P}\{\vec{Y} \in E\} = 1$ with $E = \{(t, -t) : t \in \mathbb{R}\}$. Thus, $\mathbb{P}_{\vec{Y}}$ is concentrated on the subspace $E \neq \mathbb{R}^2$. Here and in the sequel, we want to exclude such “degenerated” normal vectors by assuming that the generating matrix A is regular.

Given a normal vector \vec{Y} , how do we get the standard normal \vec{X} in representation (6.2)? The next proposition answers this question.

Proposition 6.1.8. *A random vector $\vec{Y} = (Y_1, \dots, Y_n)$ is normal if and only if there exists a regular $n \times n$ matrix $B = (\beta_{ij})_{i,j=1}^n$ and a vector $v = (v_1, \dots, v_n) \in \mathbb{R}^n$ such that the random variables X_i , defined by*

$$X_i := \sum_{j=1}^n \beta_{ij} Y_j + v_i, \quad 1 \leq i \leq n,$$

are independent standard normal.

Proof. This is a direct consequence of the following observation. One has $\vec{Y} = A\vec{X} + \mu$ if and only if \vec{X} may be represented as $\vec{X} = A^{-1}\vec{Y} - A^{-1}\mu$. Therefore, the assertion follows by choosing B and v such that $B = A^{-1}$ and $v = -A^{-1}\mu$. \square

Example 6.1.9. For the random vector \vec{Y} investigated in Example 6.1.6, the generated independent standard normal random variables X_1 , X_2 , and X_3 may be represented as follows:

$$\begin{aligned} X_1 &= \frac{1}{9}(4Y_1 + 2Y_2 - Y_3 + 7), & X_2 &= \frac{1}{9}(Y_1 - Y_2 - Y_3 + 1), \\ X_3 &= \frac{1}{9}(2Y_1 + Y_2 - 5Y_3 - 19). \end{aligned}$$

Suppose $\vec{Y} = A\vec{X} + \mu$ is a normal vector. How can we evaluate its distribution density? To answer this question, we introduce the following function. Let $R > 0$ be an $n \times n$ -matrix and $\mu \in \mathbb{R}^n$. The inverse matrix of R is R^{-1} , and to simplify the notation, set $|R| = \det(R)$. Observe that $R > 0$ implies $|R| > 0$. With these notations, we define a function $p_{\mu,R}$ from \mathbb{R}^n to \mathbb{R} by

$$p_{\mu,R}(x) := \frac{1}{(2\pi)^{n/2}|R|^{1/2}} e^{-(R^{-1}(x-\mu), (x-\mu))/2}, \quad x \in \mathbb{R}^n. \quad (6.3)$$

Note that the expression in the exponent may be written as follows. If $R^{-1} = (\tilde{r}_{ij})_{i,j=1}^n$ is the inverse matrix of R , then one gets

$$\langle R^{-1}(x - \mu), (x - \mu) \rangle / 2 = \frac{1}{2} \sum_{i,j=1}^n \tilde{r}_{ij} (x_i - \mu_i)(x_j - \mu_j).$$

Now we are prepared to answer the above question about the density of \vec{Y} .

Proposition 6.1.10. *Suppose the normal vector \vec{Y} is represented as in eq. (6.2) with regular A and $\mu \in \mathbb{R}^n$. Define the positive matrix R by $R = AA^T$. Then $p_{\mu,R}$, as given in eq. (6.3), is the distribution density of \vec{Y} . In other words, if $B \in \mathcal{B}(\mathbb{R}^n)$, then*

$$\mathbb{P}\{\vec{Y} \in B\} = \frac{1}{(2\pi)^{n/2}|R|^{1/2}} \int_B e^{-\langle R^{-1}(x-\mu), (x-\mu) \rangle / 2} dx.$$

Proof. Because $\vec{Y} = A\vec{X} + \mu$ with \vec{X} standard normal, Proposition 6.1.2 implies

$$\begin{aligned} \mathbb{P}\{\vec{Y} \in B\} &= \mathbb{P}\{A\vec{X} + \mu \in B\} = \mathbb{P}\{\vec{X} \in A^{-1}(B - \mu)\} \\ &= \frac{1}{(2\pi)^{n/2}} \int_{A^{-1}(B-\mu)} e^{-|y|^2/2} dy \end{aligned}$$

for any Borel set $B \subseteq \mathbb{R}^n$. Hereby, $B - \mu$ denotes the set $\{b - \mu : b \in B\}$. In the next step, we change the variables by setting $x = Ay + \mu$. Then $dx = |\det(A)|dy$, where by assumption $\det(A) \neq 0$ and, moreover, we have $y \in A^{-1}(B - \mu)$ if and only if $x \in B$. Therefore, the last integral transforms to

$$\mathbb{P}\{\vec{Y} \in B\} = \frac{1}{(2\pi)^{n/2}} |\det(A)|^{-1} \int_B e^{-|A^{-1}(x-\mu)|^2/2} dx. \quad (6.4)$$

Proposition A.4.1 implies $R > 0$ and, moreover,

$$|R| = \det(R) = \det(AA^T) = \det(A) \cdot \det(A^T) = \det(A)^2.$$

Since $|R| = \det(R) > 0$, this leads to $|R|^{1/2} = |\det(A)|$, that is, to

$$|\det(A)|^{-1} = |R|^{-1/2}. \quad (6.5)$$

Note that

$$|A^{-1}(x - \mu)|^2 = \langle A^{-1}(x - \mu), A^{-1}(x - \mu) \rangle = \langle (A^{-1})^T A^{-1}(x - \mu), (x - \mu) \rangle,$$

which, due to

$$(A^{-1})^T \circ A^{-1} = (A^T)^{-1} \circ A^{-1} = (A \circ A^T)^{-1} = R^{-1},$$

implies

$$|A^{-1}(x - \mu)|^2 = \langle R^{-1}(x - \mu), (x - \mu) \rangle. \quad (6.6)$$

Plugging eqs. (6.5) and (6.6) into eq. (6.4), we get

$$\mathbb{P}\{\vec{Y} \in B\} = \int_B p_{\mu,R}(x) \, dx$$

with $p_{\mu,R}$ as in eq. (6.3). This completes the proof. \square

Remark 6.1.11. How does Proposition 6.1.10 look like for $n = 1$? Here $Y = aX + \mu$, that is, $A = (a)$, and since A has to be regular, this implies $a \neq 0$. Hence we get $R = AA^T = (a^2)$, $R^{-1} = (a^{-2})$, and $|R|^{1/2} = |a|$. Thus, the density of Y is given by

$$p_{\mu,R}(x) = \frac{1}{(2\pi)^{1/2}|R|^{1/2}} e^{-\langle R^{-1}(x-\mu), x-\mu \rangle / 2} = \frac{1}{(2\pi)^{1/2}|a|} e^{-(x-\mu)^2 / 2a^2}, \quad x \in \mathbb{R}.$$

This coincides with the result obtained in Example 4.2.2.

In view of Proposition 6.1.10, we will use the following notation.

Definition 6.1.12. A normal vector \vec{Y} is said to be $\mathcal{N}(\mu, R)$ -distributed if $p_{\mu,R}$ is its density, that is, if

$$\mathbb{P}\{\vec{Y} \in B\} = \frac{1}{(2\pi)^{n/2}|R|^{1/2}} \int_B e^{-\langle R^{-1}(x-\mu), x-\mu \rangle / 2} \, dx.$$

Remark 6.1.13. It follows from Proposition A.4.2 that, given *any* $\mu \in \mathbb{R}^n$ and *any* $R > 0$, there exists a normal vector \vec{Y} that is $\mathcal{N}(\mu, R)$ -distributed. Indeed, write $R > 0$ as $R = AA^T$ and set $\vec{Y} = A\vec{X} + \mu$ with \vec{X} standard normal. Then \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed by Proposition 6.1.10.

{Distributions of \mathbb{R}^n -valued normal vectors} \Leftrightarrow $\{(\mu, R) : \mu \in \mathbb{R}^n, n \times n \text{ matrix } R > 0\}$.



Example 6.1.14. Assume

$$Y_1 = X_1 - X_2 + 3 \quad \text{and} \quad Y_2 = 2X_1 + X_2 - 2$$

for X_1, X_2 independent $\mathcal{N}(0, 1)$ -distributed. Then we get

$$\mu = (3, -2) \quad \text{and} \quad A = \begin{pmatrix} 1 & -1 \\ 2 & 1 \end{pmatrix},$$

which implies

$$R = AA^T = \begin{pmatrix} 1 & -1 \\ 2 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix} = \begin{pmatrix} 2 & 1 \\ 1 & 5 \end{pmatrix}. \quad (6.7)$$

Thus, \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed with $\mu = (3, -2)$ and R as in eq. (6.7).

Which density does \vec{Y} possess? To answer this, we have to compute $\det(R)$ and R^{-1} . One easily gets $\det(R) = 9$. The inverse matrix of R equals

$$R^{-1} = \frac{1}{9} \begin{pmatrix} 5 & -1 \\ -1 & 2 \end{pmatrix}.$$

Therefore, the distribution density $p_{\mu,R}$ of $\vec{Y} = (Y_1, Y_2)$ is given by (see Fig. 6.2)

$$\begin{aligned} p_{\mu,R}(x_1, x_2) &= \frac{1}{6\pi} \exp\left(-\frac{1}{2} \langle R^{-1}(x_1 - 3, x_2 + 2), (x_1 - 3, x_2 + 2) \rangle\right) \\ &= \frac{1}{6\pi} \exp\left(-\frac{1}{18} [5(x_1 - 3)^2 - 2(x_1 - 3)(x_2 + 2) + 2(x_2 + 2)^2]\right). \end{aligned} \quad (6.8)$$

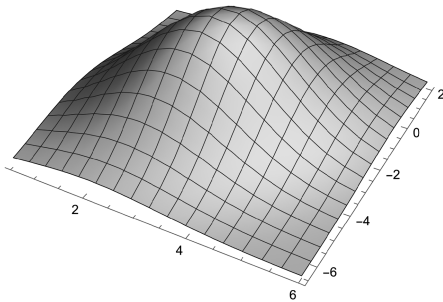


Figure 6.2: The density given by eq. (6.8). It attains its maximal value $1/6\pi$ at the point $(3, -2)$.

For later purposes, we have to name the probability measures on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ appearing as distributions of normal vectors.

Definition 6.1.15. Given $\mu \in \mathbb{R}^n$ and $R > 0$, the probability measure $\mathcal{N}(\mu, R)$ on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ is defined by

$$\mathcal{N}(\mu, R)(B) = \int_B p_{\mu,R}(x) \, dx = \frac{1}{(2\pi)^{n/2} |R|^{1/2}} \int_B e^{-\frac{1}{2} \langle R^{-1}(x-\mu), (x-\mu) \rangle} \, dx.$$

Measure $\mathcal{N}(\mu, R)$ is called a **multivariate normal distribution** with expected value μ and covariance matrix R .

According to Definition 6.1.15, we may now formulate Proposition 6.1.10 as follows:

Proposition 6.1.16. *Let \vec{Y} be a random vector. Then the following are equivalent:*

1. \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed.
2. $\mathbb{P}_{\vec{Y}} = \mathcal{N}(\mu, R)$.
3. *There is a regular $n \times n$ matrix A with $R = AA^T$ such that for some standard normal \vec{X} one has $\vec{Y} = A\vec{X} + \mu$.*

Remark 6.1.17. If \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed, the representing matrix A in (3) is not unique; compare Remark A.4.3. This is already so in the univariate case where an

$\mathcal{N}(\mu, \sigma^2)$ -distributed random variable Y may either be represented as $Y = \sigma X + \mu$ or as $Y = (-\sigma)X' + \mu$ for some standard normal random variables X and X' .

Remark 6.1.18. The case $R = I_n$ (as in Section A.4, we denote the identity matrix in \mathbb{R}^n by I_n) and $\mu = 0$ is of special interest. Because $I_n^{-1} = I_n$ and $\det(I_n) = 1$, we get

$$p_{0, I_n}(x) = \frac{1}{(2\pi)^{n/2}} e^{-|x|^2/2}, \quad x \in \mathbb{R}^n.$$

This tells us that $\mathcal{N}(0, I_n)$ is nothing else as the multivariate standard normal distribution introduced in Definition 1.9.21. Written as formula, this means

$$\mathcal{N}(0, 1)^{\otimes n} = \mathcal{N}(0, I_n).$$

More generally, in view of eq. (1.86), it follows that

$$\mathcal{N}(\mu, \sigma^2)^{\otimes n} = \mathcal{N}(\vec{\mu}, \sigma^2 I_n)$$

where $\vec{\mu} = (\mu, \dots, \mu) \in \mathbb{R}^n$ and $\sigma > 0$. In other words,

$$\mathcal{N}(\mu, \sigma^2)^{\otimes n}(B) = \mathcal{N}(\vec{\mu}, \sigma^2 I_n)(B) = \frac{1}{(2\pi)^{n/2} \sigma^n} \int_B e^{-|x - \vec{\mu}|^2 / 2\sigma^2} dx. \quad (6.9)$$

For later purposes, the next result is of importance.

Proposition 6.1.19. *Suppose a normal vector $\vec{Y} = (Y_1, \dots, Y_n)$ may be written as*

$$\vec{Y} = U\vec{X}$$

with an $\mathcal{N}(0, I_n)$ -distributed (standard normal) \vec{X} and a unitary matrix U . Then its coordinate mappings Y_1, \dots, Y_n are independent standard normal random variables.

Proof. The random vector \vec{Y} is $\mathcal{N}(0, UU^T)$ -distributed. But U is unitary, hence, $UU^T = I_n$ and \vec{Y} is $\mathcal{N}(0, I_n)$ or, equivalently, standard normally distributed. Then the assertion follows by Proposition 6.1.2. \square

Remark 6.1.20. The previous result may be phrased also as follows: If B is a Borel set in \mathbb{R}^n and $B' = U(B) = \{U(x) : x \in B\}$ for some unitary transformation U , then this implies

$$\mathcal{N}(0, I_n)(B) = \mathcal{N}(0, I_n)(B').$$

That is, the n -dimensional standard normal distribution is invariant under unitary transformations as, e. g., rotations or reflections.

Example 6.1.21. For $\theta \in [0, 2\pi)$, define the 2×2 matrix U by

$$U = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}.$$

The matrix U is unitary (it is a rotation by the angle θ) and, due to Proposition 6.1.19, the vector $\vec{Y} = U\vec{X}$ is standard normal. In other words, given independent standard normal X_1 and X_2 , for each $\theta \in [0, 2\pi)$ the random variables

$$Y_1 := \cos \theta X_1 + \sin \theta X_2 \quad \text{and} \quad Y_2 = -\sin \theta X_1 + \cos \theta X_2$$

are independent and standard normally distributed as well. That is, if $a_1 < b_1$ and $a_2 < b_2$, then we obtain

$$\begin{aligned} \mathbb{P}\{a_1 \leq Y_1 \leq b_1, a_2 \leq Y_2 \leq b_2\} &= \frac{1}{2\pi} \left(\int_{a_1}^{b_1} e^{-x_1^2/2} dx_1 \right) \left(\int_{a_2}^{b_2} e^{-x_2^2/2} dx_2 \right) \\ &= \frac{1}{2\pi} \iint_{[a_1, b_1] \times [a_2, b_2]} e^{-|x|^2/2} dx. \end{aligned}$$

6.2 Expected value and covariance matrix

We start with the following definition.

Definition 6.2.1. Let $\vec{Y} = (Y_1, \dots, Y_n)$ be a random vector such that $\mathbb{E}|Y_j| < \infty$ for all $1 \leq j \leq n$. Then the vector

$$\mathbb{E}\vec{Y} := (\mathbb{E}Y_1, \dots, \mathbb{E}Y_n) = (\mu_1, \dots, \mu_n)$$

is called the (multivariate) **expected value** of \vec{Y} .

If $\mathbb{E}Y_j^2 < \infty$, $1 \leq j \leq n$, then the matrix

$$\text{Cov}_{\vec{Y}} := (\text{Cov}(Y_i, Y_j))_{i,j=1}^n = (\mathbb{E}(Y_i - \mu_i)(Y_j - \mu_j))_{i,j=1}^n$$

is said to be the **covariance matrix** of \vec{Y} .

Remark 6.2.2. It is important to notice that both $\mathbb{E}\vec{Y}$ and the covariance matrix $\text{Cov}_{\vec{Y}}$ depend only on the distribution of \vec{Y} . That is, whenever $\mathbb{P}_{\vec{Y}_1} = \mathbb{P}_{\vec{Y}_2}$, then

$$\mathbb{E}\vec{Y}_1 = \mathbb{E}\vec{Y}_2 \quad \text{and} \quad \text{Cov}_{\vec{Y}_1} = \text{Cov}_{\vec{Y}_2}.$$

The next proposition describes the (multivariate) expected value and the covariance matrix of a normally distributed vector.

Proposition 6.2.3. Assume $\vec{Y} = A\vec{X} + \mu$ for some regular matrix A , with \vec{X} standard normal and $\mu \in \mathbb{R}^n$. Define $R = (r_{ij})_{i,j=1}^n$ as $R = AA^T$. Then the following are valid:

- (1) We have $\mathbb{E}\vec{Y} = \mu$ and $\text{Cov}_{\vec{Y}} = (\text{Cov}(Y_i, Y_j))_{i,j=1}^n = R$.
- (2) Given $a \in \mathbb{R}^n$, $a \neq 0$, then $\langle \vec{Y}, a \rangle$ is a normal random variable with expected value $\langle \mu, a \rangle$ and variance $\langle Ra, a \rangle$.

(3) The coordinate mappings Y_i are $\mathcal{N}(\mu_i, r_{ii})$ -distributed, $1 \leq i \leq n$, that is, the marginal distributions of \vec{Y} are the probability measures $\mathcal{N}(\mu_i, r_{ii})$.

Proof. By assumption,

$$Y_i = \sum_{j=1}^n \alpha_{ij} X_j + \mu_i, \quad i = 1, \dots, n, \quad (6.10)$$

hence, the linearity of the expected value and $\mathbb{E}X_j = 0$ imply

$$\mathbb{E}Y_i = \sum_{j=1}^n \alpha_{ij} \mathbb{E}X_j + \mu_i = \mu_i, \quad 1 \leq i \leq n.$$

This proves $\mathbb{E}\vec{Y} = (\mathbb{E}Y_1, \dots, \mathbb{E}Y_n) = \mu$.

Let us now verify the second part of property (1). Using $\mu_j = \mathbb{E}Y_j$, from representation (6.10) we get

$$\begin{aligned} \text{Cov}(Y_i, Y_j) &= \mathbb{E}[(Y_i - \mu_i)(Y_j - \mu_j)] = \mathbb{E}\left(\sum_{k=1}^n \alpha_{ik} X_k\right) \left(\sum_{\ell=1}^n \alpha_{j\ell} X_\ell\right) \\ &= \sum_{k, \ell=1}^n \alpha_{ik} \alpha_{j\ell} \mathbb{E}X_k X_\ell. \end{aligned}$$

The X_j s are independent $\mathcal{N}(0, 1)$ -distributed, hence

$$\mathbb{E}X_k X_\ell = \begin{cases} 1 & \text{if } k = \ell, \\ 0 & \text{if } k \neq \ell, \end{cases}$$

leading to

$$\text{Cov}(Y_i, Y_j) = \sum_{k=1}^n \alpha_{ik} \alpha_{jk} = r_{ij}.$$

To see this, recall that $R = AA^T$, hence $r_{ij} = \sum_{k=1}^n \alpha_{ik} \alpha_{jk}$. This proves $\text{Cov}_{\vec{Y}} = R$, as asserted.

To verify property (2), we first treat a special case, namely that the random vector is standard normally distributed. So suppose that \vec{X} is $\mathcal{N}(0, I_n)$ -distributed. In this case, property (2) asserts the following. For any $b \in \mathbb{R}^n$, $b \neq 0$, we have

$$\langle \vec{X}, b \rangle \text{ is distributed according to } \mathcal{N}(0, |b|^2). \quad (6.11)$$

If $b = (b_1, \dots, b_n)$, then

$$\langle \vec{X}, b \rangle = \sum_{j=1}^n b_j X_j = \sum_{j=1}^n Z_j$$

with $Z_j = b_j X_j$. The random variables Z_1, \dots, Z_n are independent and, moreover, by Proposition 4.2.3, the Z_j s are $\mathcal{N}(0, b_j^2)$ -distributed. Proposition 4.6.11 implies that $\sum_{j=1}^n Z_j$ is distributed according to $\mathcal{N}(0, \sum_{j=1}^n b_j^2)$. In view of $\sum_{j=1}^n b_j^2 = |b|^2$, this proves assertion (6.11).

Let us now turn to the general case. Recall that

$$\vec{Y} = A\vec{X} + \mu$$

and $R = AA^T$. If $a \in \mathbb{R}^n$ is a nonzero vector, then we take the scalar product with respect to a on both sides of the last equation and obtain

$$\langle \vec{Y}, a \rangle = \langle A\vec{X}, a \rangle + \langle \mu, a \rangle = \langle \vec{X}, A^T a \rangle + \langle \mu, a \rangle.$$

An application of statement (6.11) with $b = A^T a$ lets us conclude that $\langle \vec{X}, A^T a \rangle$ is $\mathcal{N}(0, |A^T a|^2)$ -distributed, that is, $\langle \vec{Y}, a \rangle$ is $\mathcal{N}(\langle \mu, a \rangle, |A^T a|^2)$ -distributed. Here we used that A , hence also A^T , is regular, so that $a \neq 0$ yields $b = A^T a \neq 0$, and statement (6.11) applies. Assertion (2) follows now from

$$|A^T a|^2 = \langle A^T a, A^T a \rangle = \langle AA^T a, a \rangle = \langle Ra, a \rangle.$$

Property (3) is an immediate consequence of the second. An application of property (2) to the i th unit vector $e_i = (0, \dots, 0, \underbrace{1}_i, 0, \dots, 0)$ in \mathbb{R}^n leads, on the one hand, to

$$\langle \vec{Y}, e_i \rangle = Y_i, \quad 1 \leq i \leq n,$$

and, on the other hand, to

$$\langle Re_i, e_i \rangle = r_{ii} \quad \text{and} \quad \langle \mu, e_i \rangle = \mu_i, \quad 1 \leq i \leq n.$$

Thus, by property (2), for each $i \leq n$ the random variable Y_i is $\mathcal{N}(\mu_i, r_{ii})$ -distributed. This completes the proof. \square

Corollary 6.2.4. *If \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed, then $\mathbb{E}\vec{Y} = \mu$ and $\text{Cov}_{\vec{Y}} = R$.*

Proof. Choose any regular $n \times n$ matrix \tilde{A} such that $R = \tilde{A}\tilde{A}^T$. The existence of such an \tilde{A} is proved in Proposition A.4.2. Set $\vec{Z} = \tilde{A}\vec{X} + \mu$ for some standard normal vector \vec{X} . Then \vec{Y} and \vec{Z} are both $\mathcal{N}(\mu, R)$ -distributed, hence $\vec{Z} \stackrel{d}{=} \vec{Y}$. Proposition 6.2.3 implies $\mathbb{E}\vec{Z} = \mu$ and $\text{Cov}_{\vec{Z}} = R$. Consequently, by Remark 6.2.2, it follows that

$$\mathbb{E}\vec{Y} = \mathbb{E}\vec{Z} = \mu \quad \text{and} \quad \text{Cov}_{\vec{Y}} = \text{Cov}_{\vec{Z}} = R,$$

which completes the proof. \square

In view of property Corollary 6.2.4, we will use the following notation.

Definition 6.2.5. If \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed, then the parameters μ and R are called the (multivariate) **expected value** and the **covariance matrix** of \vec{Y} , respectively.

Remark 6.2.6. We proved above that, for any normal vector \vec{Y} , the coordinate mappings $Y_i = \langle \vec{Y}, e_i \rangle$ are normal as well. The converse is not valid. There are random vectors \vec{Y} with all random variables $\langle \vec{Y}, e_i \rangle$ normal, $1 \leq i \leq n$, but \vec{Y} is not normal.

In contrast to this remark, the following is valid.

Proposition 6.2.7. *If $\langle \vec{Y}, a \rangle$ is normal for all nonzero $a \in \mathbb{R}^n$, then \vec{Y} is normal as well.*

Idea of the proof. By assumption, for each $a \neq 0$ there are real numbers μ_a and $\sigma_a > 0$ such that $\langle \vec{Y}, a \rangle$ is $\mathcal{N}(\mu_a, \sigma_a^2)$ -distributed. In order to prove the proposition, one has to show that there are a $\mu \in \mathbb{R}^n$ with $\mu_a = \langle \mu, a \rangle$ and an $R > 0$ such that $\sigma_a^2 = \langle Ra, a \rangle$, $a \in \mathbb{R}^n$. The existence of the vector μ easily follows from

$$\mu_{\alpha a + \beta b} = \mathbb{E}\langle \vec{Y}, \alpha a + \beta b \rangle = \alpha \mathbb{E}\langle \vec{Y}, a \rangle + \beta \mathbb{E}\langle \vec{Y}, b \rangle = \alpha \mu_a + \beta \mu_b,$$

using the fact that each linear mapping from \mathbb{R}^n to \mathbb{R} is of the form $a \mapsto \langle a, \mu \rangle$ for a suitable $\mu \in \mathbb{R}^n$.

The existence of an $R > 0$ with $\sigma_a^2 = \langle Ra, a \rangle$ is consequence of a representation theorem for positive quadratic forms on \mathbb{R}^n . To this end, one has to show that $a \mapsto \sigma_a^2$ is a positive quadratic form, which follows by using $\sigma_a^2 = \mathbb{E}\langle \vec{Y}, a \rangle^2$. \square

As we saw above (see Proposition 5.3.10), independent random variables are uncorrelated. On the other hand, Examples 5.3.12 and 5.3.15 showed the existence of uncorrelated variables that are not independent. Thus, in general, the property of being uncorrelated is weaker than that of being independent.

One of the basic features of normal vectors is that for them uncorrelated coordinate mappings are already independent. This somehow explains why in the common speech these properties are synonyms.

Proposition 6.2.8. *Let $\vec{Y} = (Y_1, \dots, Y_n)$ be a normally distributed vector. Then the following are equivalent:*

- (1) Y_1, \dots, Y_n are independent.
- (2) Y_1, \dots, Y_n are uncorrelated.
- (3) The covariance matrix $\text{Cov}_{\vec{Y}}$ is a diagonal matrix.

Proof. The implication (1) \Rightarrow (2) follows by Proposition 5.3.10. If the Y_i s are uncorrelated, then this tells us that $\text{Cov}(Y_i, Y_j) = 0$ whenever $i \neq j$. Thus, $\text{Cov}_{\vec{Y}}$ is a diagonal matrix, which proves (2) \Rightarrow (3).

It remains to verify (3) \Rightarrow (1). Thus assume that \vec{Y} is $\mathcal{N}(\mu, R)$ -distributed, where $R > 0$ is a diagonal matrix. Let r_{11}, \dots, r_{nn} be the entries of R at the diagonal. Define A as diagonal matrix with $r_{11}^{1/2}, \dots, r_{nn}^{1/2}$ on the diagonal. Note that $R > 0$ implies $r_{ii} > 0$, hence A is well defined. Of course, then $AA^T = R$, hence \vec{Y} has the same distribution as

the vector (Z_1, \dots, Z_n) with

$$Z_i = r_{ii}^{1/2} X_i + \mu_i, \quad 1 \leq i \leq n,$$

where X_1, \dots, X_n are independent standard normal. Proposition 4.1.9 lets us conclude that Z_1, \dots, Z_n are independent normal random variables. But since $\vec{Y} \stackrel{d}{=} \vec{Z}$, the random variables Y_1, \dots, Y_n are independent as well.¹ \square

Remark 6.2.9. Another property, being equivalent to those in Proposition 6.2.8, is as follows. The density function of \vec{Y} with independent coordinates equals

$$p_{\mu,R}(x) = \frac{1}{(2\pi)^{n/2} |R|^{1/2}} e^{-\sum_{j=1}^n (x_j - \mu_j)^2 / 2r_{jj}}, \quad x = (x_1, \dots, x_n).$$

Note that $|R| = \det(R) = r_{11} \cdots r_{nn}$.

Finally, we investigate the case of two-dimensional normal vectors more thoroughly. Thus assume $\vec{Y} = (Y_1, Y_2)$ is a normal vector. Then the covariance matrix R is given by

$$R = \begin{pmatrix} \mathbb{V}Y_1 & \text{Cov}(Y_1, Y_2) \\ \text{Cov}(Y_1, Y_2) & \mathbb{V}Y_2 \end{pmatrix}.$$

Let σ_1^2 and σ_2^2 be the variances of Y_1 and Y_2 , respectively, and let $\rho = \rho(Y_1, Y_2)$ be their correlation coefficient.² Because of

$$\text{Cov}(Y_1, Y_2) = (\mathbb{V}Y_1)^{1/2} (\mathbb{V}Y_2)^{1/2} \rho(Y_1, Y_2) = \sigma_1 \sigma_2 \rho,$$

we may rewrite R as

$$R = \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix}.$$

This implies $\det(R) = \sigma_1^2 \sigma_2^2 (1 - \rho^2)$. Since $\sigma_1^2 > 0$, the matrix R is positive if and only if $|\rho| < 1$. The inverse matrix R^{-1} can be computed by Cramer's rule as

$$R^{-1} = \frac{1}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} \begin{pmatrix} \sigma_2^2 & -\rho \sigma_1 \sigma_2 \\ -\rho \sigma_1 \sigma_2 & \sigma_1^2 \end{pmatrix} = \frac{1}{1 - \rho^2} \begin{pmatrix} \frac{1}{\sigma_1^2} & \frac{-\rho}{\sigma_1 \sigma_2} \\ \frac{-\rho}{\sigma_1 \sigma_2} & \frac{1}{\sigma_2^2} \end{pmatrix}.$$

Consequently,

1 Indeed, use the characterization of independent random variables given in Proposition 3.6.5. The condition stated there depends only on the joint distribution.

2 Recall that ρ describes the degree and the way of the dependence between Y_1 and Y_2 . These two random variables are positively correlated if $\rho > 0$, negatively if $\rho < 0$, strongly dependent if ρ is near 1 or -1 , and only weakly dependent in the case that ρ is near zero.

$$\langle R^{-1}x, x \rangle = \frac{1}{1-\rho^2} \left(\frac{x_1^2}{\sigma_1^2} - \frac{2\rho x_1 x_2}{\sigma_1 \sigma_2} + \frac{x_2^2}{\sigma_2^2} \right), \quad x = (x_1, x_2) \in \mathbb{R}^2.$$

If $\mu = (\mu_1, \mu_2) = (\mathbb{E}Y_1, \mathbb{E}Y_2)$ denotes the expected value of \vec{Y} , then for $a_1 < b_1$ and $a_2 < b_2$,

$$\begin{aligned} \mathbb{P}\{a_1 \leq Y_1 \leq b_1, a_2 \leq Y_2 \leq b_2\} &= \frac{1}{2\pi(1-\rho^2)^{1/2}\sigma_1\sigma_2} \\ &\times \int_{a_1}^{b_1} \int_{a_2}^{b_2} \exp\left(-\frac{1}{2(1-\rho^2)} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} - \frac{2\rho(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right]\right) dx_2 dx_1. \end{aligned} \quad (6.12)$$

Compare this with the case of *independent* Y_1 and Y_2 or, equivalently, with the case $\rho = 0$. Here it follows that

$$\begin{aligned} \mathbb{P}\{a_1 \leq Y_1 \leq b_1, a_2 \leq Y_2 \leq b_2\} \\ = \frac{1}{2\pi\sigma_1\sigma_2} \int_{a_1}^{b_1} \int_{a_2}^{b_2} \exp\left(-\frac{1}{2} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right]\right) dx_2 dx_1. \end{aligned} \quad (6.13)$$

It is worthwhile to mention that in both cases (dependent and independent) the marginal distributions are the same, namely $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$. A comparison of eqs. (6.12) and (6.13) shows clearly the influence of the correlation coefficient to the density (see Fig. 6.3).

Summary: An n -dimensional random vector $\vec{Y} = (Y_1, \dots, Y_n)$ is said to be normal (or the Y_j s are called jointly normal) if there are a regular $n \times n$ matrix $A = (a_{ij})_{i,j=1}^n$ and a vector $\mu = (\mu_1, \dots, \mu_n) \in \mathbb{R}^n$ such that with independent $\mathcal{N}(0, 1)$ -distributed X_1, \dots, X_n ,

$$Y_i = \sum_{j=1}^n a_{ij} X_j + \mu_i, \quad 1 \leq i \leq n.$$

Equivalently, \vec{Y} is normal if and only if there are a positive $n \times n$ matrix $R = (r_{ij})_{i,j=1}^n$ and a $\mu \in \mathbb{R}^n$ such that

$$\mathbb{P}\{\vec{Y} \in B\} = \frac{1}{(2\pi)^{n/2} |R|^{1/2}} \int_B e^{-\langle R^{-1}(x-\mu), (x-\mu) \rangle / 2} dx.$$

The matrices A and R are linked by $R = AA^T$. If \vec{Y} is normal, the coordinates Y_1, \dots, Y_n are normal (univariate) random variables (the converse is in general not true) with

$$\mathbb{E}Y_i = \mu_i \quad \text{and} \quad \text{Cov}(Y_i, Y_j) = r_{ij}, \quad 1 \leq i, j \leq n.$$

Let \vec{Y} be a normal vector. Then the Y_j s are pairwise uncorrelated if and only if they are independent:

$$Y_1, \dots, Y_n \text{ independent} \quad \Leftrightarrow \quad \text{Cov}(Y_i, Y_j) = 0, \quad 1 \leq i < j \leq n.$$

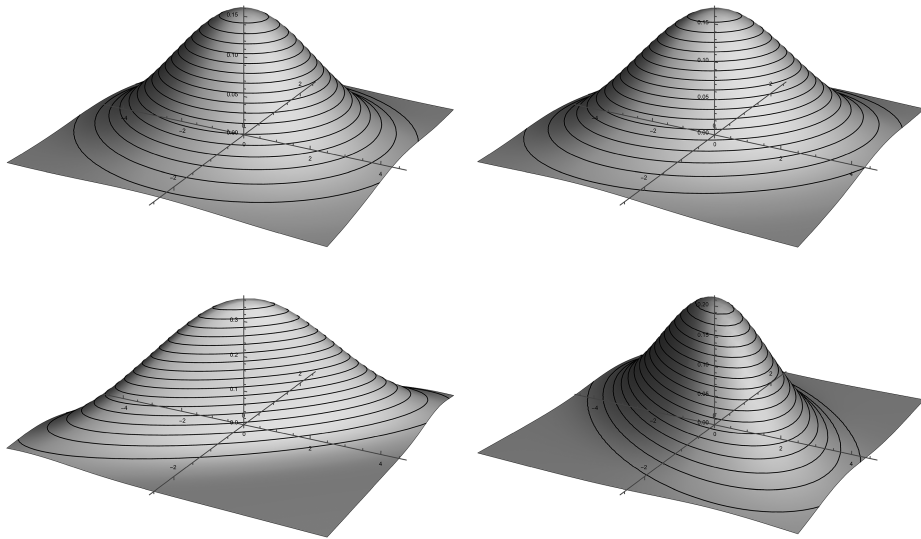


Figure 6.3: The 2-dimensional densities of a normal vector with $\mu_1 = \mu_2 = 0$, $\sigma_1 = 2$, $\sigma_2 = 1$, and $\rho = 0, 0.25, 0.75, -0.5$ from top left to bottom right. Thus, the coordinates of normal vectors with these densities are either independent ($\rho = 0$), weakly ($\rho = 0.25$) or strongly ($\rho = 0.75$) positively correlated. In the last case ($\rho = -0.5$), they are moderately negatively correlated; values in the regions $\{x > 0, y < 0\}$ and $\{x < 0, y > 0\}$ become more likely.

6.3 Problems

Problem 6.1. Let $\vec{Y} = (Y_1, \dots, Y_n)$ be an arbitrary (not necessarily normal) random vector.

1. Show that $\mathbb{E}|Y_j| < \infty$, $1 \leq j \leq n$, if and only if $\mathbb{E}|\vec{Y}| < \infty$. Here $|\vec{Y}|$ denotes the Euclidean distance of \vec{Y} .
2. Let A be an arbitrary $n \times n$ matrix. Prove that

$$\mathbb{E}(A\vec{Y}) = A(\mathbb{E}\vec{Y})$$

provided that $\mathbb{E}|\vec{Y}| < \infty$.

3. Show that $\mathbb{E}|Y_j|^2 < \infty$, $1 \leq j \leq n$, if and only if $\mathbb{E}|\vec{Y}|^2 < \infty$.
4. Suppose $\mathbb{E}|\vec{Y}|^2 < \infty$. Let $\text{Cov}_{\vec{Y}}$ be the covariance matrix of \vec{Y} . Prove that $\text{Cov}_{\vec{Y}}$ is nonnegative definite, that is,

$$\langle \text{Cov}_{\vec{Y}}x, x \rangle \geq 0, \quad x \in \mathbb{R}^n.$$

Problem 6.2. Roll a fair die two times. Let X_1 be the greater of the two rolls and X_2 denotes the smaller one. Evaluate the expected value $\mu \in \mathbb{R}^2$ and the covariance matrix of the random vector $\vec{X} = (X_1, X_2)$.

Problem 6.3. Let X_1 and X_2 be two independent standard normal random variables. Define Y_1 and Y_2 by

$$Y_1 = 2X_1 - 2X_2 + 1 \quad \text{and} \quad Y_2 = 3X_1 + X_2 - 2.$$

1. Find $\mu \in \mathbb{R}^2$ and the positive 2×2 matrix R such that $\vec{Y} = (Y_1, Y_2)$ is $\mathcal{N}(\mu, R)$ -distributed.
2. Determine $\text{Cov}_{\vec{Y}}$ and the correlation coefficient $\rho = \rho(Y_1, Y_2)$. Are Y_1 and Y_2 positively or negatively correlated?
3. Which distribution do $Y_1 + Y_2$ and $Y_1 - Y_2$ possess?
4. Evaluate the distribution density of \vec{Y} .

Problem 6.4. Let $\vec{X} = (X_1, X_2)$ be a two-dimensional standard normal vector. Compute

$$\mathbb{P}\{|\vec{X}| \leq 1\} = \mathbb{P}\{X_1^2 + X_2^2 \leq 1\}.$$

Hint: Compare with the proof of Proposition 1.6.7.

Problem 6.5. Let X_1, \dots, X_{n+m} be a sequence of independent standard normal random variables. For an $n \times n$ matrix $A = (a_{ij})_{i,j=1}^n$ and an $m \times m$ matrix $B = (\beta_{kl})_{k,l=1}^m$, define two normal vectors \vec{Y} and \vec{Z} by

$$Y_i = \sum_{j=1}^n a_{ij} X_j \quad \text{and} \quad Z_k = \sum_{l=1}^m \beta_{kl} X_{l+n},$$

with $1 \leq i \leq n$ and $1 \leq k \leq m$. Let (\vec{Y}, \vec{Z}) be the $(n+m)$ -dimensional vector

$$(\vec{Y}, \vec{Z}) = (Y_1, \dots, Y_n, Z_1, \dots, Z_m).$$

Why is (\vec{Y}, \vec{Z}) normal? Show that the covariance matrix $\text{Cov}_{(\vec{Y}, \vec{Z})}$ is given by

$$\text{Cov}_{(\vec{Y}, \vec{Z})} = \begin{pmatrix} \text{Cov}_{\vec{Y}} & 0 \\ 0 & \text{Cov}_{\vec{Z}} \end{pmatrix}.$$

Problem 6.6. Let X_1, X_2 , and X_3 be three standard normal independent random variables. Define the random vector \vec{Y} by

$$\vec{Y} := (X_1 - 1, X_1 + X_2 - 1, X_1 + X_2 + X_3 - 1).$$

1. Argue why \vec{Y} is normal. Determine its expected value, covariance matrix, and the correlation coefficients $\rho(Y_i, Y_j)$, $1 \leq i < j \leq 3$.
2. Determine the distribution density of \vec{Y} .

Problem 6.7. The random vector $\vec{Y} = (Y_1, \dots, Y_n)$ is $\mathcal{N}(\mu, R)$ -distributed for some $\mu \in \mathbb{R}^n$ and $R > 0$. Determine the distribution of $Y_1 + \dots + Y_n$.

Problem 6.8. Prove the following assertion: If \vec{Y} is $\mathcal{N}(0, R)$ -distributed, then there exist an orthonormal basis $(f_j)_{j=1}^n$ in \mathbb{R}^n , positive numbers $\lambda_1, \dots, \lambda_n$ and independent $\mathcal{N}(0, 1)$ -distributed ξ_1, \dots, ξ_n such that

$$\vec{Y} = \sum_{j=1}^n \lambda_j \xi_j f_j. \quad (6.14)$$

Hint: Use the principal axis transformation for symmetric matrices and the fact that unitary matrices map an orthonormal basis onto an orthonormal basis.

Conclude from eq. (6.14) the following: If \vec{Y} is $\mathcal{N}(0, R)$ -distributed, then there are a_1, \dots, a_n in \mathbb{R}^n such that $\langle \vec{Y}, a_1 \rangle, \dots, \langle \vec{Y}, a_n \rangle$ is a sequence of independent standard normal random variables.

Problem 6.9. The n -dimensional vector \vec{Y} is distributed according to $\mathcal{N}(\mu, R)$. For some regular $n \times n$ matrix S , define \vec{Z} by $\vec{Z} := S\vec{Y}$. Is \vec{Z} normal? If this is so, determine the expected value and the covariance matrix of \vec{Z} .

Problem 6.10. Let $\vec{X} = (X_1, X_2)$ be standard normal. Define random variables Y_1 and Y_2 by

$$Y_1 := \frac{1}{\sqrt{2}}(X_1 + X_2) \quad \text{and} \quad Y_2 := \frac{1}{\sqrt{2}}(X_1 - X_2).$$

Why are Y_1 and Y_2 also independent and standard normal?

7 Limit theorems

Probability Theory does not have the ability to predict the occurrence or nonoccurrence of a single event in a random experiment; besides, this event occurs either with probability one or with probability zero. For example, Probability Theory does not give any information about the next result when rolling a die, it does not predict the numbers appearing next week on the lottery nor is it able to foresee the lifetime of a component in a machine. Such statements are impossible within the theory. The theory is only able to say that some events are more likely and others are less likely. For instance, when rolling a die twice, it is more likely that the sum of both rolls will be “7” than “2.” Nevertheless, next when we roll the die the sum may be “2,” not “7.” The event “the sum is 2” is not impossible, only less likely.

In contrast, Probability Theory provides us with very precise and far-reaching information about the behavior of the results when we execute “many” identical random experiments. As already said, we cannot tell anything about the expected number on a die when we roll it once, but we are able to say a lot about the frequency of the number “6” when rolling a die many times, namely that, on average, this number will appear in one of six cases (provided the die is fair). In this example, certain laws of Probability Theory, which we will present in this section, are operating. These laws are only applicable in the case of many experiments, not in that of a single one.

Limit theorems in Probability Theory belong to the most beautiful and most important assertions within this theory. They are always the highlight of a lecture about advanced Probability Theory. However, their proofs require a longer comprehensive mathematical explanation, which is impossible to give here within the framework of this book. Those who are interested in knowing more about this topic may look into one of the more advanced books, such as [Bil12, Dur19] or [Kho07]. Although the proofs of the limit theorems are mostly quite complicated, they are very important, and their consequences influence our daily lives. Moreover, great parts of Mathematical Statistics are based on these results. Therefore, we decided to state here the crucial assertions without proving most of them. Thus, our main focus is to present the most important limit theorems, to explain them in detail, and to give examples that show how they apply. If possible, we give some hint as to how the results are derived, but mostly we must resign to prove them.

7.1 Laws of large numbers

7.1.1 Chebyshev’s inequality

Our first objective is to prove Chebyshev’s inequality. To do so, we need the following lemma.

Lemma 7.1.1. *Let Y be a nonnegative random variable. Then for each $\lambda > 0$ it follows that*

$$\mathbb{P}\{Y \geq \lambda\} \leq \frac{\mathbb{E}Y}{\lambda}. \quad (7.1)$$

Proof. Let us first treat the case that Y is discrete. Since $Y \geq 0$, its possible values y_1, y_2, \dots are nonnegative real numbers. Therefore, we get

$$\begin{aligned} \mathbb{E}Y &= \sum_{j=1}^{\infty} y_j \mathbb{P}\{Y = y_j\} \geq \sum_{y_j \geq \lambda} y_j \mathbb{P}\{Y = y_j\} \\ &\geq \lambda \sum_{y_j \geq \lambda} \mathbb{P}\{Y = y_j\} = \lambda \mathbb{P}\{Y \geq \lambda\}. \end{aligned}$$

Solving the inequality for $\mathbb{P}\{Y \geq \lambda\}$ proves inequality (7.1).

The proof of estimate (7.1) for continuous Y uses similar methods. If q denotes the distribution density of Y , by $Y \geq 0$ we may suppose $q(y) = 0$ for $y < 0$. Then, as in the discrete case, we conclude that

$$\mathbb{E}Y = \int_0^{\infty} yq(y) \, dy \geq \int_{\lambda}^{\infty} yq(y) \, dy \geq \lambda \int_{\lambda}^{\infty} q(y) \, dy = \lambda \mathbb{P}\{Y \geq \lambda\}.$$

From this inequality, (7.1) follows directly. \square

Remark 7.1.2. Sometimes it is useful to apply inequality (7.1) in a slightly modified way. For example, if $Y \geq 0$ and $\alpha > 0$, then one derives

$$\mathbb{P}\{Y \geq \lambda\} = \mathbb{P}\{Y^\alpha \geq \lambda^\alpha\} \leq \frac{\mathbb{E}Y^\alpha}{\lambda^\alpha}.$$

Or, if Y is real valued, then for $\lambda \in \mathbb{R}$ we obtain

$$\mathbb{P}\{Y \leq \lambda\} = \mathbb{P}\{e^{-Y} \geq e^{-\lambda}\} \leq \frac{\mathbb{E}e^{-Y}}{e^{-\lambda}} = e^\lambda \mathbb{E}e^{-Y}.$$

Now we are in a position to state and to prove Chebyshev's inequality.

Proposition 7.1.3 (Chebyshev's inequality). *Let X be a random variable with finite second moment. Then, if $c > 0$, it follows that*

$$\mathbb{P}\{|X - \mathbb{E}X| \geq c\} \leq \frac{\mathbb{V}X}{c^2}. \quad (7.2)$$

Proof. Setting $Y := |X - \mathbb{E}X|^2$, we have $Y \geq 0$ and $\mathbb{E}Y = \mathbb{V}X$. Now apply inequality (7.1) to Y with $\lambda = c^2$. This leads to

$$\mathbb{P}\{|X - \mathbb{E}X| \geq c\} = \mathbb{P}\{|X - \mathbb{E}X|^2 \geq c^2\} = \mathbb{P}\{Y \geq c^2\} \leq \frac{\mathbb{E}Y}{c^2} = \frac{\mathbb{V}X}{c^2},$$

and estimate (7.2) is proven. \square

Interpretation: Inequality (7.2) quantifies the interpretation of $\mathbb{V}X$ as a measure for the dispersion of X . The smaller the $\mathbb{V}X$, the less the probability that the values of X are far away from its expected value $\mathbb{E}X$.

Remark 7.1.4. Another way to formulate inequality (7.2) is as follows. If $\kappa > 0$, then

$$\mathbb{P}\{|X - \mathbb{E}X| \geq \kappa (\mathbb{V}X)^{1/2}\} \leq \frac{1}{\kappa^2}.$$

To see this, apply inequality (7.2) with $c = \kappa (\mathbb{V}X)^{1/2}$.

Example 7.1.5. Roll a fair die n times. If, for example, $A = \{6\}$, we are interested in the relative frequency $r_n(A)$ of the occurrence of A . Recall that this frequency was defined in eq. (1.1). Moreover, we claimed in this section that $\lim_{n \rightarrow \infty} r_n(A) = \mathbb{P}(A) = \frac{1}{6}$. Is it possible to estimate the probability for $|r_n(A) - \frac{1}{6}|$ being bigger than some given $c > 0$?

Answer: Define the random variable X as the absolute frequency of the occurrence of A , that is, we have $X = k$ for some $k = 0, \dots, n$ provided that A occurred exactly k times. Then X is binomial distributed with parameters n and $p = 1/6$. To see this, define “success” as appearance of “6.” Consequently, the relative frequency can be represented as $r_n(A) = \frac{X}{n}$. An application of eqs. (5.8) and (5.36) gives

$$\mathbb{E} r_n(A) = \frac{1}{n} \mathbb{E}X = \frac{np}{n} = p = \frac{1}{6} \quad \text{and} \quad \mathbb{V} r_n(A) = \frac{np(1-p)}{n^2} = \frac{5}{36n}.$$

Thus, inequality (7.2) leads to

$$\mathbb{P}\left\{\left|r_n(A) - \frac{1}{6}\right| \geq c\right\} \leq \frac{5}{36c^2n}.$$

If, for example, $n = 10^3$, and if we choose $c = 1/36$, then Chebyshev’s inequality yields

$$\mathbb{P}\left\{\frac{5}{36} < r_{10^3}(A) < \frac{7}{36}\right\} \geq 1 - \frac{9}{50} = 0.82.$$

For the absolute frequency, this means

$$\mathbb{P}\{139 \leq a_{10^3}(A) \leq 194\} \geq 0.82.$$

Let us interpret the result. Suppose we roll a fair die 1000 times. Then, with a probability of at least 82 %, the frequency of “6” will be between 139 and 194.

Let us present a second quite similar example.

Example 7.1.6. Roll a fair die n times and let S_n be the sum of the n results. Then $S_n = X_1 + \dots + X_n$, where X_1, \dots, X_n are uniformly distributed on $\{1, \dots, 6\}$ and independent. By Example 5.2.17, we know that

$$\mathbb{E}S_n = \mathbb{E}X_1 + \dots + \mathbb{E}X_n = \frac{7n}{2} \quad \text{and} \quad \mathbb{V}S_n = \mathbb{V}X_1 + \dots + \mathbb{V}X_n = \frac{35n}{12},$$

hence

$$\mathbb{E}\left(\frac{S_n}{n}\right) = \frac{7}{2} \quad \text{and} \quad \mathbb{V}\left(\frac{S_n}{n}\right) = \frac{35}{12n}.$$

An application of inequality (7.2) leads then to

$$\mathbb{P}\left\{\left|\frac{S_n}{n} - \frac{7}{2}\right| \geq c\right\} \leq \frac{35}{12nc^2}.$$

For example, if $n = 10^3$ and c is chosen as $c = 0.1$, then

$$\mathbb{P}\left\{3.4 < \frac{S_{10^3}}{10^3} < 3.6\right\} \geq 0.709.$$

The interpretation of this result is as in the previous example. With a probability larger than 70 % the sum of 1000 rolls of a fair die will be a number between 3400 and 3600.

This looks like to be a pretty rough estimate and, indeed, this is so. Sharper bounds follow by using the central limit theorem as we will see in Example 7.2.14.

Summary: Let X be a random variable with finite second moment. Then Chebyshev's inequality asserts that for any $c > 0$,

$$\mathbb{P}\{|X - \mathbb{E}X| \geq c\} \leq \frac{\mathbb{V}X}{c^2}.$$

This inequality clarifies once more the role of the variance $\mathbb{V}X$. The smaller the $\mathbb{V}(X)$, the more likely the random observations are concentrated around $\mathbb{E}X$.

7.1.2 Infinite sequences of independent random variables*

Whenever one wants to describe the limit behavior of random variables or random events, one needs a model for the infinite performance of random experiments. Otherwise, we cannot investigate limits or other related quantities. This is comparable with similar investigations in Calculus. In order to analyze limits, infinite sequences are necessary, not finite ones. Thus, for the examination of limits of random variables we need an infinite sequence X_1, X_2, \dots of random variables, which are, on the one hand, independent in the sense of Definition 4.3.4 and, on the other hand, possess some given probability distributions.

Example 7.1.7. In order to describe the infinite tossing of a fair coin, we need independent random variables X_1, X_2, \dots such that $\mathbb{P}\{X_j = 0\} = \mathbb{P}\{X_j = 1\} = \frac{1}{2}$. Or, similarly, for a model of rolling a die infinitely often, we need infinitely many independent random variables all uniformly distributed on $\{1, \dots, 6\}$.

In Proposition 4.3.3, we presented the construction of independent $(X_j)_{j=1}^{\infty}$ distributed according to $B_{1,1/2}$. This technique can be extended to more general sequences of random variables, but it is quite complicated. Another, much smarter way is to use so-called infinite product measures.¹ Their existence follows by a deep theorem due to A. N. Kolmogorov. As a consequence, one gets the following result, which cannot be proven within the framework of this book. We refer to [Kho07, Chapter 5, § 2] or [Ros06, Theorem 7.1.1] for proofs.

Proposition 7.1.8. *Let $\mathbb{P}_1, \mathbb{P}_2, \dots$ be arbitrary probability measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Then there are a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and an infinite sequence of random variables $X_j : \Omega \rightarrow \mathbb{R}$ such that the following hold:*

1. *The probability distribution of X_j is $\mathbb{P}_j, j = 1, 2, \dots$. That is, for all $j \geq 1$ and all $B \in \mathcal{B}(\mathbb{R})$, it follows that*

$$\mathbb{P}\{X_j \in B\} = \mathbb{P}_j(B).$$

2. *The random variables X_1, X_2, \dots are independent in the sense of Definition 4.3.4. This says, for all $n \geq 1$ and all $B_j \in \mathcal{B}(\mathbb{R})$, it follows that*

$$\begin{aligned} \mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} &= \mathbb{P}\{X_1 \in B_1\} \cdots \mathbb{P}\{X_n \in B_n\} \\ &= \mathbb{P}_1(B_1) \cdots \mathbb{P}_n(B_n). \end{aligned}$$

Of special interest is the case $\mathbb{P}_1 = \mathbb{P}_2 = \dots = \mathbb{P}_0$ for a certain probability measure \mathbb{P}_0 on \mathbb{R} . Then the previous proposition implies the following.

Corollary 7.1.9. *Given an arbitrary probability measure \mathbb{P}_0 on $\mathcal{B}(\mathbb{R})$, there are random variables X_1, X_2, \dots , defined on some probability space $(\Omega, \mathcal{A}, \mathbb{P})$, such that for all $n \geq 1$ and all $B_j \in \mathcal{B}(\mathbb{R})$,*

$$\mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} = \mathbb{P}_0(B_1) \cdots \mathbb{P}_0(B_n).$$

Example 7.1.10. Choosing as \mathbb{P}_0 the uniform distribution on $[0, 1]$, the previous corollary ensures the existence of (independent) random variables X_1, X_2, \dots such that for all $n \geq 1$ and all $0 \leq a_j < b_j \leq 1$,

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n\} = \prod_{j=1}^n (b_j - a_j).$$

The sequence X_1, X_2, \dots models the independent choosing of infinitely many numbers uniformly distributed in $[0, 1]$.

¹ Compare with Proposition 3.6.7 for the construction of *finitely many* independent random variables possessing given distributions. There we used *finite* product measures to obtain independent random variables possessing given distributions $\mathbb{P}_1, \dots, \mathbb{P}_n$.

Remark 7.1.11. One may ask whether the kind of independence in Definition 4.3.4 suffices for later purposes. Recall, we only require X_1, \dots, X_n to be independent for all (finite) $n \geq 1$. Maybe one would expect a condition that involves the whole infinite sequence, not only a finite part of it. The answer is that such a condition for the whole sequence is a consequence of Definition 4.3.4. Namely, if B_1, B_2, \dots are Borel sets in \mathbb{R} , then, by the continuity of probability measures from above, it follows that

$$\begin{aligned} \mathbb{P}\{X_1 \in B_1, X_2 \in B_2, \dots\} &= \lim_{n \rightarrow \infty} \mathbb{P}\{X_1 \in B_1, \dots, X_n \in B_n\} \\ &= \lim_{n \rightarrow \infty} \mathbb{P}\{X_1 \in B_1\} \cdots \mathbb{P}\{X_n \in B_n\} \\ &= \lim_{n \rightarrow \infty} \prod_{j=1}^n \mathbb{P}\{X_j \in B_j\} = \prod_{j=1}^{\infty} \mathbb{P}\{X_j \in B_j\}. \end{aligned}$$

In particular, if $a_j < b_j, j = 1, 2, \dots$, this implies

$$\mathbb{P}\{a_1 \leq X_1 \leq b_1, a_2 \leq X_2 \leq b_2, \dots\} = \prod_{j=1}^{\infty} \mathbb{P}\{a_j \leq X_j \leq b_j\}. \quad (7.3)$$

Example 7.1.12. Let X_1, X_2, \dots be a sequence of independent E_λ -distributed random variables for some $\lambda > 0$. Given real numbers $a_j > 0$, we ask for the probability of

$$\mathbb{P}\{X_1 \leq a_1, X_2 \leq a_2, \dots\}.$$

Answer: If we apply eq. (7.3) with $a_j = 0$ and with $b_j = a_j$, then we get

$$\mathbb{P}\{X_1 \leq a_1, X_2 \leq a_2, \dots\} = \prod_{j=1}^{\infty} \mathbb{P}\{X_j \leq a_j\} = \prod_{j=1}^{\infty} [1 - e^{-\lambda a_j}].$$

Of special interest are sequences $(a_j)_{j \geq 1}$ such that the infinite product converges, that is, for these sequences $(a_j)_{j \geq 1}$ we have $\prod_{j=1}^{\infty} [1 - e^{-\lambda a_j}] > 0$. This happens if and only if

$$\ln \left(\prod_{j=1}^{\infty} [1 - e^{-\lambda a_j}] \right) = \sum_{j=1}^{\infty} \ln [1 - e^{-\lambda a_j}] > -\infty. \quad (7.4)$$

Because of

$$\lim_{x \rightarrow 0} \frac{\ln(1-x)}{-x} = 1,$$

by the limit comparison test for infinite series, condition (7.4) holds if and only if

$$\sum_{j=1}^{\infty} e^{-\lambda a_j} < \infty.$$

If, for example, $a_j = c \cdot \ln(j + 1)$ for some $c > 0$, then

$$\sum_{j=1}^{\infty} e^{-\lambda a_j} = \sum_{j=1}^{\infty} \frac{1}{(j+1)^{\lambda c}}.$$

This sum is known to be finite if and only if $\lambda c > 1$, that is, if $c > 1/\lambda$.

Another way to formulate this observation is as follows. One has

$$\mathbb{P}\left\{\sup_{j \geq 1} \frac{X_j}{\ln(j+1)} \leq c\right\} = \mathbb{P}\{X_j \leq c \ln(j+1), \forall j \geq 1\} = \prod_{j=1}^{\infty} \left(1 - \frac{1}{(j+1)^{\lambda c}}\right),$$

and this probability is positive if and only if $c > 1/\lambda$.

Summary: An infinite sequence X_1, X_2, \dots of random variables is said to be independent if for each $n \geq 1$ the finite sequence X_1, \dots, X_n is independent. In particular, this implies for all $a_i < b_i$ that

$$\mathbb{P}\{a_i \leq X_i \leq b_i, i = 1, 2, \dots\} = \prod_{i=1}^{\infty} \mathbb{P}\{a_i \leq X_i \leq b_i\}.$$

Given probability measures $\mathbb{P}_1, \mathbb{P}_2, \dots$ on \mathbb{R} , there are independent random variables X_1, X_2, \dots on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ so that $\mathbb{P}_{X_j} = \mathbb{P}_j, j = 1, 2, \dots$. That is, for all Borel sets B in $\mathcal{B}(\mathbb{R})$ and $j = 1, 2, \dots$, we have

$$\mathbb{P}\{\omega \in \Omega : X_j(\omega) \in B\} = \mathbb{P}\{X_j \in B\} = \mathbb{P}_j(B).$$

For example, there are infinitely many independent random variables X_j such that

$$\mathbb{P}\{X_j = 1\} = \dots = \mathbb{P}\{X_j = 6\} = \frac{1}{6}.$$

These X_j s may serve as model for rolling a die infinitely often.

7.1.3 Borel–Cantelli lemma*

The aim of this section is to present one of the most useful tools for the investigation of the limit behavior of infinite sequences of random variables and events. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and let A_1, A_2, \dots be a sequence of events in \mathcal{A} . Then two typical questions arise. What is the probability that there exists some $n \in \mathbb{N}$ such that all events A_m with $m \geq n$ occur? The other related question asks for the probability that infinitely many of the events A_n occur.

To explain why these questions are of interest, let us once more regard Example 4.1.7 of the random walk. Here S_n denotes the integer where the particle is located after n random jumps. For example, letting $A_n := \{\omega \in \Omega : S_n(\omega) > 0\}$, then the existence of an $n \in \mathbb{N}$ such that A_m occurs for all $m \geq n$ says that the particle from a certain (random)

moment attains only positive numbers and never goes back to the negative ones. Or, if we investigate the events $B_n := \{\omega \in \Omega : S_n(\omega) = 0\}$, then the B_n s occur infinitely often if and only if the particle returns to zero infinitely often. Equivalently, there are (random) $n_1 < n_2 < \dots$ with $S_{n_j}(\omega) = 0$.

To formulate the two previous questions more precisely, let us introduce the following two events.

Definition 7.1.13. Let A_1, A_2, \dots be subsets of Ω . Then

$$\liminf_{n \rightarrow \infty} A_n := \bigcup_{n=1}^{\infty} \bigcap_{m=n}^{\infty} A_m \quad \text{and} \quad \limsup_{n \rightarrow \infty} A_n := \bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} A_m$$

are called the **lower** and **upper limit** of the A_n s.

Remark 7.1.14. Let us characterize when the lower and the upper limit occur.

1. An element $\omega \in \Omega$ belongs to $\liminf_{n \rightarrow \infty} A_n$ if and only if there is an $n \in \mathbb{N}$ such that $\omega \in \bigcap_{m=n}^{\infty} A_m$, that is, if it is an element of A_m for $m \geq n$. In other words, the lower limit occurs if there is an $n \in \mathbb{N}$ such that after n the events A_m always occur. Therefore, we say that $\liminf_{n \rightarrow \infty} A_n$ occurs if the A_n s **finally always** (abbreviated as f. a.) occur. Thus,

$$\mathbb{P}\{\omega \in \Omega : \exists n \text{ such that } \omega \in A_m, m \geq n\} = \mathbb{P}\left(\liminf_{n \rightarrow \infty} A_n\right).$$

2. An element $\omega \in \Omega$ belongs to $\limsup_{n \rightarrow \infty} A_n$ if and only if for each $n \in \mathbb{N}$ there is an $m \geq n$ such that $\omega \in A_m$. But this is nothing else as saying that the number of A_n s with $\omega \in A_n$ is infinite. Therefore, the upper limit consists of those elements for which we have **infinitely often** (abbreviated as i. o.) $\omega \in A_n$. Note that also these events may be different for different ω . Thus,

$$\mathbb{P}\{\omega \in \Omega : \omega \in A_n \text{ for infinitely many } n\} = \mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right).$$

Example 7.1.15. Suppose a fair coin is labeled on one side with “0” and on the other side with “1.” We toss it infinitely often. Let A_n occur if the n th toss is “1.” Then $\liminf_{n \rightarrow \infty} A_n$ occurs if after a certain number of tosses “1” always shows up. On the other hand, $\limsup_{n \rightarrow \infty} A_n$ occurs if and only if the number “1” appears infinitely often. The subsequent results imply that the probability of the lower limit of these A_n s equals zero, while with probability one they will occur infinitely often.

Let us formulate and prove some easy properties of the lower and upper limit.

² Note that this n is random, that is, it may depend on the chosen $\omega \in \Omega$.

Proposition 7.1.16. *If A_1, A_2, \dots are subsets of Ω , then*

$$(1) \quad \liminf_{n \rightarrow \infty} A_n \subseteq \limsup_{n \rightarrow \infty} A_n,$$

$$(2) \quad \left(\limsup_{n \rightarrow \infty} A_n \right)^c = \liminf_{n \rightarrow \infty} A_n^c \quad \text{and} \quad \left(\liminf_{n \rightarrow \infty} A_n \right)^c = \limsup_{n \rightarrow \infty} A_n^c.$$

Proof. We prove these properties in the interpretation of the lower and upper limit given in Remark 7.1.14.

Suppose that $\omega \in \liminf_{n \rightarrow \infty} A_n$. Then for some $n \geq 1$ it follows that $\omega \in A_m$, $m \geq n$. Of course, then the number of events with $\omega \in A_n$ is infinite, which implies $\omega \in \limsup_{n \rightarrow \infty} A_n$. This proves (1).

Observe that we have $\omega \notin \limsup_{n \rightarrow \infty} A_n$ if and only if $\omega \in A_n$ for only finitely many $n \in \mathbb{N}$. Equivalently, there is an $n \geq 1$ such that whenever $m \geq n$, then $\omega \notin A_m$, or, that $\omega \in A_m^c$. In other words, this happens if and only if $\omega \in \liminf_{n \rightarrow \infty} A_n^c$. This proves the left-hand identity in (2). The right-hand one follows by the same arguments. One may also prove this by applying the left-hand identity with A_n^c . \square

Before we can formulate the main result in this section, we have to define when an infinite sequence of events is independent.

Definition 7.1.17. A sequence of events A_1, A_2, \dots in \mathcal{A} is said to be **independent** provided that for all $n \geq 1$ the events A_1, \dots, A_n are independent in the sense of Definition 2.2.12. An equivalent formulation is as follows: given $m \geq 1$ and indices $i_1 < i_2 < \dots < i_m$, then this implies

$$\mathbb{P}(A_{i_1} \cap \dots \cap A_{i_m}) = \mathbb{P}(A_{i_1}) \cdots \mathbb{P}(A_{i_m}).$$

Remark 7.1.18. Using the method for the proof of eq. (7.3), one may deduce the following “infinite” version of independence. For independent A_1, A_2, \dots follows that

$$\mathbb{P}\left(\bigcap_{n=1}^{\infty} A_n\right) = \prod_{n=1}^{\infty} \mathbb{P}(A_n).$$

Remark 7.1.19. According to Proposition 3.6.9, the independence of random variables and events are linked as follows. The random variables X_1, X_2, \dots are independent in the sense of Definition 4.3.4 if and only if for all Borel sets B_1, B_2, \dots in \mathbb{R} the preimages $X_1^{-1}(B_1), X_2^{-1}(B_2), \dots$ are independent events as introduced in Definition 7.1.17.

Now we are in the position to state and prove the main result of this section.

Proposition 7.1.20 (Borel–Cantelli lemma). *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and let $A_n \in \mathcal{A}$, $n = 1, 2, \dots$*

1. *If $\sum_{n=1}^{\infty} \mathbb{P}(A_n) < \infty$, then*

$$\mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 0. \tag{7.5}$$

2. For independent A_1, A_2, \dots , the following is valid. If $\sum_{n=1}^{\infty} \mathbb{P}(A_n) = \infty$, then

$$\mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 1.$$

Proof. We start with proving the first assertion. Thus, take arbitrary subsets $A_n \in \mathcal{A}$ satisfying $\sum_{n=1}^{\infty} \mathbb{P}(A_n) < \infty$. Write

$$\limsup_{n \rightarrow \infty} A_n = \bigcap_{n=1}^{\infty} B_n$$

with $B_n := \bigcup_{m=n}^{\infty} A_m$. Since $B_1 \supseteq B_2 \supseteq \dots$, property (7) in Proposition 1.2.1 applies and, together with (5) in the same proposition, leads to

$$\mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = \lim_{n \rightarrow \infty} \mathbb{P}(B_n) \leq \liminf_{n \rightarrow \infty} \sum_{m=n}^{\infty} \mathbb{P}(A_m). \quad (7.6)$$

If a_1, a_2, \dots are nonnegative numbers with $\sum_{n=1}^{\infty} a_n < \infty$, then it is known that $\sum_{m=n}^{\infty} a_m \rightarrow 0$ as $n \rightarrow \infty$. Applying this observation to $a_n = \mathbb{P}(A_n)$, assertion (7.5) is a direct consequence of estimate (7.6). Thus, the first part is proven.

To prove the second assertion, we investigate the probability of the complementary event. Here we have

$$\left(\limsup_{n \rightarrow \infty} A_n\right)^c = \bigcup_{n=1}^{\infty} \bigcap_{m=n}^{\infty} A_m^c.$$

An application of (5) in Proposition 1.2.1 implies

$$\mathbb{P}\left(\left(\limsup_{n \rightarrow \infty} A_n\right)^c\right) \leq \sum_{n=1}^{\infty} \mathbb{P}\left(\bigcap_{m=n}^{\infty} A_m^c\right). \quad (7.7)$$

Fix $n \in \mathbb{N}$ and for $k \geq n$ set $B_k := \bigcap_{m=n}^k A_m^c$. Then $B_n \supseteq B_{n+1} \supseteq \dots$, hence by property (7) in Proposition 1.2.1 it follows that

$$\mathbb{P}\left(\bigcap_{m=n}^{\infty} A_m^c\right) = \mathbb{P}\left(\bigcap_{k=n}^{\infty} B_k\right) = \lim_{k \rightarrow \infty} \mathbb{P}(B_k) = \lim_{k \rightarrow \infty} \prod_{m=n}^k (1 - \mathbb{P}(A_m)).$$

Here in the last step we used that, due to Proposition 2.2.15, the events A_1^c, A_2^c, \dots are independent as well. Next we apply the elementary inequality

$$1 - x \leq e^{-x}, \quad 0 \leq x \leq 1,$$

for $x = \mathbb{P}(A_m)$ and, because of $\sum_{m=n}^{\infty} \mathbb{P}(A_m) = \infty$, we arrive at

$$\mathbb{P}\left(\bigcap_{m=n}^{\infty} A_m^c\right) \leq \limsup_{k \rightarrow \infty} \exp\left(-\sum_{m=n}^k \mathbb{P}(A_m)\right) = 0.$$

Plugging this into estimate (7.7) finally yields

$$\mathbb{P}\left(\left(\limsup_{n \rightarrow \infty} A_n\right)^c\right) = 0, \quad \text{hence } \mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 1,$$

as asserted. \square

Remark 7.1.21. The second assertion in Proposition 7.1.20 remains valid under the weaker condition of pairwise independence. But then the proof becomes more complicated (see Examples 6.4 and 6.5 in [Bil12] or Lemma 11.1 in [Bau96]).

Corollary 7.1.22. *Let $A_n \in \mathcal{A}$ be independent events. Then the following are equivalent:*

$$\begin{aligned} \mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 0 &\Leftrightarrow \sum_{n=1}^{\infty} \mathbb{P}(A_n) < \infty, \\ \mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 1 &\Leftrightarrow \sum_{n=1}^{\infty} \mathbb{P}(A_n) = \infty. \end{aligned}$$

Example 7.1.23. Suppose we play a series of independent games where the success probability in the n th game is p_n for some given p_1, p_2, \dots . How likely is it to win infinitely many of the games?

Answer: Let the event A_n occur if one wins game n . By the choice of the p_n s, it follows that $\mathbb{P}(A_n) = p_n$. Thus, the Borel–Cantelli lemma asserts that one wins with probability 1 infinitely many games if and only if $\sum_{n=1}^{\infty} p_n = \infty$. On the other hand, if $\sum_{n=1}^{\infty} p_n < \infty$, then with probability 1 one loses all games after a finite random number of games. So, for example, if $p_n = 1/n$, then with probability 1 one wins infinitely often although the success probability becomes smaller and smaller. On the contrary, in the case $p_n = 1/n^2$ there will be an $N \in \mathbb{N}$ so that one loses all games after the N th one.

Example 7.1.24. Let $(U_n)_{n \geq 1}$ be a sequence of independent random variables, uniformly distributed on $[0, 1]$. Given positive real numbers $(\alpha_n)_{n \geq 1}$, we define events A_n by setting $A_n := \{U_n \leq \alpha_n\}$. Since the U_n s are independent, so are the events A_n , and Corollary 7.1.22 applies. Because of $\mathbb{P}(A_n) = \alpha_n$, this leads to

$$\mathbb{P}\{U_n \leq \alpha_n \text{ i. o.}\} = \begin{cases} 0 & \text{if } \sum_{n=1}^{\infty} \alpha_n < \infty, \\ 1 & \text{if } \sum_{n=1}^{\infty} \alpha_n = \infty, \end{cases}$$

or, equivalently, to

$$\mathbb{P}\{U_n > \alpha_n \text{ f. a.}\} = \begin{cases} 0 & \text{if } \sum_{n=1}^{\infty} \alpha_n = \infty, \\ 1 & \text{if } \sum_{n=1}^{\infty} \alpha_n < \infty. \end{cases}$$

For example, we have

$$\mathbb{P}\{U_n \leq 1/n \text{ i. o.}\} = 1 \quad \text{and} \quad \mathbb{P}\{U_n \leq 1/n^2 \text{ i. o.}\} = 0.$$

Example 7.1.25. Let $(X_n)_{n \geq 1}$ be a sequence of independent $\mathcal{N}(0, 1)$ -distributed random variables and let $c_n > 0$. What probability does the event to observe $\{|X_n| \geq c_n\}$ infinitely often possess?

Answer: It holds that

$$\sum_{n=1}^{\infty} \mathbb{P}\{|X_n| \geq c_n\} = \frac{2}{\sqrt{2\pi}} \sum_{n=1}^{\infty} \int_{c_n}^{\infty} e^{-x^2/2} dx = \frac{2}{\sqrt{2\pi}} \sum_{n=1}^{\infty} \varphi(c_n),$$

where

$$\varphi(t) := \int_t^{\infty} e^{-x^2/2} dx, \quad t \in \mathbb{R}.$$

Setting $\psi(t) := t^{-1}e^{-t^2/2}$, $t > 0$, one obtains

$$\varphi'(t) = -e^{-t^2/2} \quad \text{and} \quad \psi'(t) = -\left(1 + \frac{1}{t^2}\right)e^{-t^2/2},$$

hence l'Hôpital's rule implies

$$\lim_{t \rightarrow \infty} \frac{\varphi'(t)}{\psi'(t)} = 1, \quad \text{thus} \quad \lim_{t \rightarrow \infty} \frac{\varphi(t)}{\psi(t)} = 1.$$

The limit comparison test for infinite series tells us that $\sum_{n=1}^{\infty} \varphi(c_n) < \infty$ if and only if $\sum_{n=1}^{\infty} \psi(c_n) < \infty$. Thus, by the definition of ψ , the following are equivalent:

$$\sum_{n=1}^{\infty} \mathbb{P}\{|X_n| \geq c_n\} < \infty \quad \iff \quad \sum_{n=1}^{\infty} \frac{e^{-c_n^2/2}}{c_n} < \infty.$$

In other words, we have

$$\mathbb{P}\{|X_n| \geq c_n \text{ i. o.}\} = \begin{cases} 0 & \text{if } \sum_{n=1}^{\infty} \frac{e^{-c_n^2/2}}{c_n} < \infty, \\ 1 & \text{if } \sum_{n=1}^{\infty} \frac{e^{-c_n^2/2}}{c_n} = \infty. \end{cases}$$

For example, if $c_n = c\sqrt{\ln n}$ for some $c > 0$, then

$$\sum_{n=1}^{\infty} \frac{e^{-c_n^2/2}}{c_n} = \frac{1}{c} \sum_{n=1}^{\infty} \frac{1}{n^{c^2/2} \sqrt{\ln n}} < \infty$$

if and only if $c > \sqrt{2}$. In particular, this yields the following interesting fact:

$$\mathbb{P}\{|X_n| \geq \sqrt{2 \ln n} \text{ i. o.}\} = 1,$$

while for each $c > 2$,

$$\mathbb{P}\{|X_n| \geq \sqrt{c \ln n} \text{ i. o.}\} = 0.$$

From this, we derive

$$\mathbb{P}\left\{\omega \in \Omega : \limsup_{n \rightarrow \infty} \frac{|X_n(\omega)|}{\sqrt{\ln n}} = \sqrt{2}\right\} = 1.$$

Example 7.1.26. In a lottery, 6 of 49 numbers are randomly chosen. Find the probability to have infinitely often the six chosen numbers on your lottery ticket.

Answer: Let A_n be the event to have in the n th drawing the six chosen numbers on the ticket. We saw (see Example 1.4.3) that

$$\mathbb{P}(A_n) = \frac{1}{\binom{49}{6}} := \delta > 0.$$

Consequently, it follows that $\sum_{n=1}^{\infty} \mathbb{P}(A_n) = \infty$ and, since the A_n s are independent, Proposition 7.1.20 implies

$$\mathbb{P}\{\text{The } A_n\text{s occur i. o.}\} = 1.$$

Therefore, the event to win infinitely often has probability 1. One does only not play long enough!

Remark 7.1.27. Corollary 7.1.22 shows in particular that for independent A_n s either

$$\mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 0 \quad \text{or} \quad \mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 1.$$

Because of Proposition 7.1.16, the same is valid for the lower limit. Here the so-called 0–1 laws operate, which, roughly speaking, assert the following. Whenever the occurrence or nonoccurrence of an event is independent of the first finitely many results, then such events occur either with probability 0 or 1. For example, the occurrence or nonoccurrence of the lower or upper limit is completely independent of what had happened during the first n results, $n \geq 1$.

Summary: Let A_1, A_2, \dots be a sequence of events in a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. A basic question in Probability Theory is how likely the occurrence of infinitely many of the events A_n is. This question is answered by the Borel–Cantelli Lemma. It asserts

$$\sum_{n=1}^{\infty} \mathbb{P}(A_n) < \infty \quad \Rightarrow \quad \mathbb{P}\{\text{The } A_n\text{s occur i. o.}\} = 0.$$

Conversely, if the A_n s are independent events, then it follows that

$$\sum_{n=1}^{\infty} \mathbb{P}(A_n) = \infty \quad \Rightarrow \quad \mathbb{P}\{\text{The } A_n\text{s occur i. o.}\} = 1.$$

In particular, if the events A_1, A_2, \dots are independent, then the A_n s occur infinitely often either with probability 0 or 1.

7.1.4 Weak law of large numbers

Given random variables X_1, X_2, \dots , let

$$S_n := X_1 + \dots + X_n \tag{7.8}$$

be the sum of the first n values. One of the most important questions in Probability Theory is that about the behavior of S_n as $n \rightarrow \infty$. Suppose we play a series of games and X_j denotes the loss or the gain in game $j \geq 1$. Then S_n is nothing else than the total loss or gain after n games. Also recall the random walk presented in Example 4.1.7. Set $X_j = -1$ if in step j the particle jumps to the left, and $X_j = 1$ otherwise. Then S_n is the point in \mathbb{Z} where the particle is located after n jumps.

Let us come back to the general case. We are given arbitrary independent and identically distributed random variables X_1, X_2, \dots . Recall that “identically distributed” says that they all possess the same probability distribution. Set $S_n = X_1 + \dots + X_n$. The first result gives some information about the behavior of the arithmetic mean S_n/n as $n \rightarrow \infty$.

Proposition 7.1.28 (Weak law of large numbers). *Let X_1, X_2, \dots be independent identically distributed random variables with (common) expected value $\mu \in \mathbb{R}$. If $\varepsilon > 0$, then it follows that*

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{\left|\frac{S_n}{n} - \mu\right| \geq \varepsilon\right\} = 0.$$

Proof. We prove the result only with an additional condition, namely that X_1 and hence all X_j possess a finite second moment. The result remains true without this condition, but then its proof becomes significantly more complicated.

From (3) in Proposition 5.1.38, we derive

$$\mathbb{E}\left(\frac{S_n}{n}\right) = \frac{\mathbb{E}S_n}{n} = \frac{\mathbb{E}(X_1 + \dots + X_n)}{n} = \frac{\mathbb{E}X_1 + \dots + \mathbb{E}X_n}{n} = \frac{n\mu}{n} = \mu.$$

Furthermore, by the independence of the X_j s, property (iv) in Proposition 5.2.15 also gives

$$\mathbb{V}\left(\frac{S_n}{n}\right) = \frac{\mathbb{V}S_n}{n^2} = \frac{\mathbb{V}X_1 + \dots + \mathbb{V}X_n}{n^2} = \frac{\mathbb{V}X_1}{n}.$$

Consequently, inequality (7.2) implies

$$\mathbb{P}\left\{\left|\frac{S_n}{n} - \mu\right| \geq \varepsilon\right\} \leq \frac{\mathbb{V}(S_n/n)}{\varepsilon^2} = \frac{\mathbb{V}X_1}{n\varepsilon^2},$$

and the desired assertion follows from

$$\limsup_{n \rightarrow \infty} \mathbb{P}\left\{\left|\frac{S_n}{n} - \mu\right| \geq \varepsilon\right\} \leq \lim_{n \rightarrow \infty} \frac{\mathbb{V}X_1}{n\varepsilon^2} = 0. \quad \square$$

Remark 7.1.29. The type of convergence appearing in Proposition 7.1.28 is usually called **convergence in probability**. More precisely, given random variables Y_1, Y_2, \dots , they converge in probability to some random variable Y provided that for each $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}\{|Y_n - Y| \geq \varepsilon\} = 0.$$

Hence, in this language the weak law of large numbers asserts that S_n/n converges in probability to a random variable Y , which is the constant μ .

Interpretation of Proposition 7.1.28. Fix $\varepsilon > 0$ and define events A_n , $n \geq 1$, by

$$A_n := \left\{ \omega \in \Omega : \left| \frac{S_n(\omega)}{n} - \mu \right| < \varepsilon \right\}.$$

Then Proposition 7.1.28 implies $\lim_{n \rightarrow \infty} \mathbb{P}(A_n) = 1$. Hence, given $\delta > 0$, there is an $n_0 = n_0(\varepsilon, \delta)$ such that $\mathbb{P}(A_n) \geq 1 - \delta$ whenever $n \geq n_0$. In other words, if n is sufficiently large, then with high probability (recall, μ is the expected value of the X_j s)

$$\mu - \varepsilon \leq \frac{1}{n} \sum_{j=1}^n X_j \leq \mu + \varepsilon.$$

This confirms once more the interpretation of the expected value as (approximate) arithmetic mean of the observed values, provided that we execute the same experiment arbitrarily often and the results do not depend on each other.

7.1.5 Strong law of large numbers

Proposition 7.1.28 does not imply $S_n/n \rightarrow \mu$ in the usual sense. It only asserts the convergence of S_n/n in probability, which, in general, does not imply pointwise convergence. The following theorem due to A. N. Kolmogorov shows that, nevertheless, a strong type of convergence takes place. The proof of this result is much more complicated than that of Proposition 7.1.28. Therefore, we cannot present it in the scope of this book, and we refer to [Dur19, Section 2.4] or [Ros06, Chapter 5] for a proof.

Proposition 7.1.30 (Strong law of large numbers). *Let X_1, X_2, \dots be a sequence of independent identically distributed random variables with expected value $\mu = \mathbb{E}X_1$. If S_n is defined by eq. (7.8), then*

$$\mathbb{P}\left\{\omega \in \Omega : \lim_{n \rightarrow \infty} \frac{S_n(\omega)}{n} = \mu\right\} = 1.$$

Remark 7.1.31. Given random variables Y_1, Y_2, \dots and Y , one says that the Y_n s converge to Y **almost surely**, if

$$\mathbb{P}\left\{\lim_{n \rightarrow \infty} Y_n = Y\right\} = \mathbb{P}\left\{\omega \in \Omega : \lim_{n \rightarrow \infty} Y_n(\omega) = Y(\omega)\right\} = 1.$$

Thus, Proposition 7.1.30 asserts that S_n/n converges almost surely to a random variable Y , which is a constant μ .

Remark 7.1.32. Proposition 7.1.30 allows the following interpretation. There exists a subset Ω_0 in the sample space Ω with $\mathbb{P}(\Omega_0) = 1$ such that for all $\omega \in \Omega_0$ and all $\varepsilon > 0$, there is an $n_0 = n_0(\varepsilon, \omega)$ with

$$\left|\frac{S_n(\omega)}{n} - \mu\right| < \varepsilon$$

whenever $n \geq n_0$.

In other words, with probability one the following happens: given $\varepsilon > 0$, there is a certain n_0 depending on ω , hence random, such that for $n \geq n_0$ the arithmetic mean S_n/n is in an ε -neighborhood of μ and never leaves it again.

Let us emphasize once more that S_n/n is random, hence S_n/n may attain different values for a different series of experiments. Nevertheless, starting from a certain point, which may be different for different experiments, the arithmetic mean of the first n results will be in $(\mu - \varepsilon, \mu + \varepsilon)$.

When we introduced probability measures in Section 1.1.3, we claimed that the number $\mathbb{P}(A)$ may be regarded as the limit of the relative frequencies of the occurrence of the event A . As the first consequence of Proposition 7.1.30, we show that this is indeed true.

Proposition 7.1.33. *Suppose a random experiment is described by a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Execute this experiment arbitrarily often. Given an event $A \in \mathcal{A}$, let $r_n(A)$ be the relative frequency of A in n trials as defined in eq. (1.1). Then almost surely*

$$\lim_{n \rightarrow \infty} r_n(A) = \mathbb{P}(A).$$

Proof. Define random variables X_1, X_2, \dots as follows. Set $X_j = 1$ if A occurs in trial j , while $X_j = 0$ otherwise. Since the experiments are executed independently of each other,

the X_j s are independent as well. Moreover, we execute every time exactly the same experiment, hence the X_j s are also identically distributed.

By the definition of the X_j s,

$$\frac{S_n}{n} = r_n(A).$$

Thus, it remains to evaluate $\mu = \mathbb{E}X_j$. To this end observe that the X_j s are $B_{1,p}$ -distributed with success probability $p = \mathbb{P}(A)$. Recall that $X_j = 1$ if and only if A occurs in experiment j , and since the experiment is described by $(\Omega, \mathcal{A}, \mathbb{P})$, the probability for X_j being 1 is $\mathbb{P}(A)$. Consequently, $\mathbb{E}X_j = \mathbb{P}(A)$.

Proposition 7.1.30 now implies that almost surely

$$\lim_{n \rightarrow \infty} r_n(A) = \lim_{n \rightarrow \infty} \frac{S_n}{n} = \mathbb{E}X_1 = \mathbb{P}(A).$$

This completes the proof. \square

What happens in the case when the X_j s do not possess an expected value? Does then S_n/n converge nevertheless? If this is so, could we take this limit as a “generalized” expected value? The next proposition shows that such an approach does not work. For a proof, see [Dur19, Theorem 2.4.5]; see also [Eri73] for further reading.

Proposition 7.1.34. *Let X_1, X_2, \dots be independent and identically distributed with $\mathbb{E}|X_1| = \infty$. Then it follows that*

$$\mathbb{P}\left\{\omega \in \Omega : \frac{S_n(\omega)}{n} \text{ diverges}\right\} = 1.$$

For example, if we take an independent sequence $(X_j)_{j \geq 1}$ of Cauchy distributed random variables, then their arithmetic means S_n/n will diverge almost surely.

Remark 7.1.35. Why does one need a *weak* law of large numbers when there exists a *strong* one? This question is justified and, in fact, in the situation described in this book the weak law is a consequence of the strong one, thus, it is not necessarily needed.

The situation is different if one investigates independent, but not necessarily identically distributed, random variables. Then there are sequences X_1, X_2, \dots satisfying the weak law but not the strong one.³

Let us state two applications of Proposition 7.1.30, one taken from Numerical Mathematics, the other from Number Theory.

³ In the case of nonidentically distributed X_j s, one investigates if $\frac{1}{n} \sum_{j=1}^n (X_j - \mathbb{E}X_j)$ converges to zero either in probability (weak law) or almost surely (strong law).

Example 7.1.36 (Monte Carlo method for integrals). Suppose we are given a quite “complicated” function $f : [0, 1]^n \rightarrow \mathbb{R}$. The task is to find the numerical value of

$$\int_{[0,1]^n} f(x) \, dx = \int_0^1 \dots \int_0^1 f(x_1, \dots, x_n) \, dx_n \cdots dx_1.$$

For large n , this can be a highly nontrivial problem. One way to overcome this difficulty is to use a probabilistic approach that is based on the strong law of large numbers.

To this end, choose an independent sequence $\vec{U}_1, \vec{U}_2, \dots$ of random vectors uniformly distributed on $[0, 1]^n$. For example, such a sequence can be constructed as follows. Take independent U_1, U_2, \dots uniformly distributed on⁴ $[0, 1]$ and build random vectors by $\vec{U}_1 = (U_1, \dots, U_n)$, $\vec{U}_2 = (U_{n+1}, \dots, U_{2n})$, and so on.

Proposition 7.1.37. *As above, let $\vec{U}_1, \vec{U}_2, \dots$ be independent random vectors uniformly distributed on $[0, 1]^n$. Given an integrable function $f : [0, 1]^n \rightarrow \mathbb{R}$, with probability one,*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N f(\vec{U}_j) = \int_{[0,1]^n} f(x) \, dx.$$

Proof. Set $X_j := f(\vec{U}_j)$, $j = 1, 2, \dots$. By construction, the X_j s are independent and identically distributed random variables. Proposition 3.6.20 implies (compare also with Example 3.6.23) that the distribution densities of the random vectors \vec{U}_j are given by

$$p(x) = \begin{cases} 1 & \text{if } x \in [0, 1]^n, \\ 0 & \text{if } x \notin [0, 1]^n. \end{cases}$$

As already mentioned in Remark 5.3.5, formula (5.42), stated for a function of two variables, also holds for functions of n variables, $n \geq 1$ arbitrary. This implies

$$\mathbb{E}X_1 = \mathbb{E}f(\vec{U}_1) = \int_{\mathbb{R}^n} f(x) p(x) \, dx = \int_{[0,1]^n} f(x) \, dx.$$

Thus, Proposition 7.1.30 applies and leads to

$$\mathbb{P} \left\{ \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N f(\vec{U}_j) = \int_{[0,1]^n} f(x) \, dx \right\} = \mathbb{P} \left\{ \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N X_j = \mathbb{E}X_1 \right\} = 1,$$

as asserted. □

⁴ Use the methods developed in Section 4.4 to construct such U_j s.

Remark 7.1.38. The numerical application of the preceding proposition is as follows. Choose independent numbers $u_i^{(j)}$, $1 \leq i \leq n$, $1 \leq j \leq N$, uniformly distributed on $[0, 1]$ and set

$$R_N(f) := \frac{1}{N} \sum_{j=1}^N f(u_1^{(j)}, \dots, u_n^{(j)}).$$

Proposition 7.1.37 asserts that $R_N(f)$ converges almost surely to $\int_{[0,1]^n} f(x) dx$. Thus, if $N \geq 1$ is large, then $R_N(f)$ may be taken as approximate value for $\int_{[0,1]^n} f(x) dx$.

If we apply Proposition 7.1.37 to the indicator function of a Borel set $B \subseteq [0, 1]^n$, that is, we choose $f = \mathbb{1}_B$ with $\mathbb{1}_B$ as in Definition 3.6.16, then with probability 1 it follows that

$$\text{vol}_n(B) = \int_{[0,1]^n} \mathbb{1}_B(x) dx = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N \mathbb{1}_B(\vec{U}_j) = \lim_{N \rightarrow \infty} \frac{|\{j \leq N : \vec{U}_j \in B\}|}{N}.$$

This provides us with a method to determine the volume $\text{vol}_n(B)$, even for quite “complicated” Borel sets $B \subseteq \mathbb{R}^n$.

Example 7.1.39. A way to approximate $\pi/4$ by the described method is as follows: draw a quadrant Q of a circle with radius 1 inside a square S of side length 1. Next choose independently points u_1, u_2, \dots, u_n uniformly distributed in $[0, 1]^2$. Then, as $n \rightarrow \infty$,

$$\frac{|\{j \leq n : u_j \in Q\}|}{|\{j \leq n : u_j \in S\}|}$$

converges to $\text{vol}_2(Q)/\text{vol}_2(S) = \pi/4$. See Fig. 7.1.

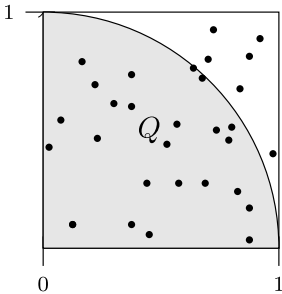


Figure 7.1: As $n \rightarrow \infty$, the proportion of the randomly chosen points inside the quadrant Q converges to $\pi/4$. In the above figure, there are 26 of the 32 points inside Q . This gives 0.8125 as an approximation of $\pi/4 \approx 0.7854$.

Example 7.1.40 (Normal numbers). As we saw in Section 4.3.1, each $x \in [0, 1)$ admits a representation as binary fraction $x = 0.x_1x_2\dots$ with $x_j \in \{0, 1\}$. Take some fixed $x \in [0, 1)$

with binary representation $x = 0.x_1x_2\dots$. Then one may ask whether in the binary representation of x one of the numbers 0 or 1 occurs more frequently than the other. Or do both numbers possess the same frequency, at least on average?

To investigate this question, for $n \in \mathbb{N}$ set

$$a_n^0(x) := |\{k \leq n : x_k = 0\}| \quad \text{and} \quad a_n^1(x) := |\{k \leq n : x_k = 1\}|, \quad x = 0.x_1x_2\dots$$

Thus, $a_n^0(x)$ is the frequency of the number 0 among the first n positions in the representation of x .

Definition 7.1.41. An $x \in [0, 1)$ is said to be **normal** (with respect to base 2) if

$$\lim_{n \rightarrow \infty} \frac{a_n^0(x)}{n} = \lim_{n \rightarrow \infty} \frac{a_n^1(x)}{n} = \frac{1}{2}.$$

In other words, a number $x \in [0, 1)$ is normal with respect to base 2 if, on average, in its binary representation the frequency of 0, and hence also of 1, equals $1/2$. Are there many normal numbers as, for example, $x = 0.0101010\dots$, or are there maybe only a few? The next proposition gives the answer.

Proposition 7.1.42. Let \mathbb{P} be the uniform distribution on $[0, 1)$. Then there is a subset $M \subseteq [0, 1)$ with $\mathbb{P}(M) = 1$ such that all $x \in M$ are normal with respect to base 2.

Proof. Define random variables $X_k : [0, 1) \rightarrow \mathbb{R}$, $k = 0, 1, \dots$, by $X_k(x) := x_k$ whenever $x = 0.x_1x_2\dots$. Proposition 4.3.3 tells us that the X_k s are independent with $\mathbb{P}\{X_k = 0\} = 1/2$ and $\mathbb{P}\{X_k = 1\} = 1/2$. Recall that the underlying probability measure \mathbb{P} on $[0, 1)$ is the uniform distribution. By the definition of the X_k s it follows that

$$S_n(x) := X_1(x) + \dots + X_n(x) = |\{k \leq n : X_k(x) = 1\}| = a_n^1(x).$$

Since $\mathbb{E}X_1 = 1/2$, Proposition 7.1.30 implies the existence of a subset $M \subseteq [0, 1)$ with $\mathbb{P}(M) = 1$ such that for $x \in M$ it follows that

$$\lim_{n \rightarrow \infty} \frac{a_n^1(x)}{n} = \lim_{n \rightarrow \infty} \frac{S_n(x)}{n} = \mathbb{E}X_1 = \frac{1}{2}.$$

Since $a_n^0(x) = n - a_n^1(x)$, this completes the proof. \square

Remark 7.1.43. The previous considerations do not depend on the fact that the base of the representation was 2. It extends easily to representations with respect to any base $b \geq 2$. Here, the definition of normal numbers has to be extended slightly. Fix $b \geq 2$. Each $x \in [0, 1)$ admits the representation $x = 0.x_1x_2\dots$ where $x_j \in \{0, \dots, b-1\}$ provided that $x = \sum_{k=1}^{\infty} \frac{x_k}{b^k}$. To make this representation unique, we do not allow representations $x = 0.x_1x_2\dots$ where for some $k_0 \in \mathbb{N}$ we have $x_k = b-1$ whenever $k \geq k_0$.

Then a number x is said to be normal with respect to the base $b \geq 2$ if for all $\ell = 0, \dots, b-1$,

$$\lim_{n \rightarrow \infty} \frac{|\{j \leq n : x_j = \ell\}|}{n} = \frac{1}{b}, \quad x = 0.x_1x_2\dots$$

Similar methods as used in the proof of Proposition 7.1.42 show that there is a set $M_b \subseteq [0, 1]$ with $\mathbb{P}(M_b) = 1$ such that all $x \in M_b$ are normal with respect to base b . Letting $M = \bigcap_{b=2}^{\infty} M_b$, then property (5) (Boole's inequality) in Proposition 1.2.1 easily gives $\mathbb{P}(M) = 1$. Numbers $x \in M$ are **completely normal**, which says that they are normal for any base $b \geq 2$. Again we see that with respect to the uniform distribution on $[0, 1]$ almost all numbers are completely normal.

Summary: Laws of large numbers are among the most important assertions in Probability Theory.⁵ Verbally said, these laws assert the following: if x_1, x_2, \dots are the random results of identical experiments, obtained independently of each other, then the sequence of arithmetic means $(x_1 + \dots + x_n)/n$ converges in a weak, as well as in a strong sense, to the expected value of the x_j s, provided the expected value exists.

In particular, these laws justify regarding the probability $\mathbb{P}(A)$ of an event A as the limit of the relative frequencies $r_n(A)$ of its occurrence, as we claimed in Section 1.1.3. We emphasize once more, the arithmetic mean, as well as the relative frequency, is random, thus may be different in different trials, but the expected value and $\mathbb{P}(A)$ are both fixed nonrandom real numbers.

7.2 Central limit theorem

Why does the normal distribution play such an important role in Probability Theory and why are so many observed random phenomena normally distributed? The reason for this is the central limit theorem, which we are going to present in this section.

Consider a sequence of independent and identically distributed random variables $(X_j)_{j \geq 1}$ with finite second moment. As in eq. (7.8), let S_n be the sum of X_1, \dots, X_n . For example, if X_j is the loss or gain in the j th game, then S_n is the total loss or gain after n games. Which probability distribution does S_n possess? Theoretically, this can be evaluated by the convolution formulas stated in Section 4.5. But practically, this is mostly impossible; imagine, we want to determine the distribution of the sum of 100 rolls with a fair die. Therefore, one is very interested in asymptotic statements about the distribution of S_n .

To get a clue about possible asymptotic distributions of S_n , take independent $B_{1,p}$ -distributed X_j s. In this case, the distribution of S_n is known to be $B_{n,p}$.

For example, if $p = 0.4$ and $n = 30$, then $\mathbb{P}\{S_n = k\} = B_{n,p}(\{k\})$, $k = 0, \dots, 30$, may be described in Fig. 7.2.

⁵ Sometimes it is said that the strong law of large numbers is one of the three pearls in Probability Theory. The two other are the central limit theorem and the so-called law of iterated logarithm, shortly discussed in Remark 7.2.16.

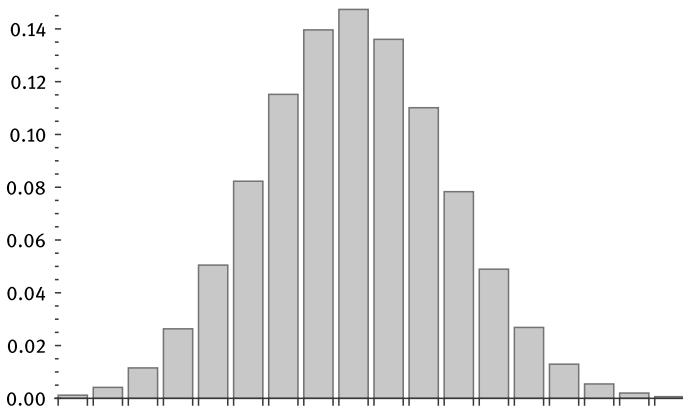


Figure 7.2: Probability mass function of $B_{n,p}$, $n = 30$ and $p = 0.4$.

The peak of the diagram occurs at $k = 12$, which is the expected value of S_{30} . Enlarging the number of trials leads to a shift of the peak to the right. At the same time, the height of the peak becomes smaller.

The shape of the diagram in Fig. 7.2 lets us suggest that sums of independent, identically distributed random variables are “almost” normally distributed. If this is so, which expected value and which variance will the approximating normal distribution possess?

Let us investigate this question in the general setting. Thus, we are given a sequence $(X_j)_{j \geq 1}$ of independent identically distributed random variables with finite second moment and with $\mu = \mathbb{E}X_1$ and $\sigma^2 = \mathbb{V}X_1 > 0$. If, as before, $S_n = X_1 + \cdots + X_n$, then

$$\mathbb{E}S_n = n\mu \quad \text{and} \quad \mathbb{V}S_n = n\sigma^2.$$

Consequently, if we conjecture that S_n is “approximately” normally distributed, then the normalized sum $(S_n - n\mu)/\sigma\sqrt{n}$ should be “approximately” $\mathcal{N}(0, 1)$ -distributed. Recall that Propositions 5.1.38 and 5.2.15 imply

$$\mathbb{E}\left(\frac{S_n - n\mu}{\sigma\sqrt{n}}\right) = 0 \quad \text{and} \quad \mathbb{V}\left(\frac{S_n - n\mu}{\sigma\sqrt{n}}\right) = 1.$$

The question about the possible limit of the normalized sums $(S_n - n\mu)/\sigma\sqrt{n}$ remained open for long time. In 1718 Abraham de Moivre investigated the limit behavior for a special case of binomial distributed random variables. As limit he found some infinite series, not a concrete function. In 1808 the American scientist and mathematician Robert Adrain published a paper where for the first time the normal distribution occurred. A year later, independently of the former work, Carl Friedrich Gauß used the normal distribution for error estimates. In 1812 Pierre-Simon Laplace proved that the normalized sums of independent binomial distributed random variables approximate the normal distribution. Later on, Andrei Andreyevich Markov, Aleksandr Mikhailovich Lyapunov,

Jarl Waldemar Lindeberg, Paul Lévy, and other mathematicians continued the work of De Moivre and Laplace. In particular, they showed that the normal distribution occurs *always* as a limit, not only for binomial distributed random variables. The only assumption is that the random variables possess a finite second moment. We refer to the very interesting book [Fis11] for further reading about the history of normal approximation.

It remains the question in which sense does $(S_n - n\mu)/\sigma\sqrt{n}$ converge to the standard normal distribution. To answer this, we have to introduce the concept of the *convergence in distribution*.

Definition 7.2.1. Let Y_1, Y_2, \dots and Y be random variables with distribution functions F_1, F_2, \dots and F , respectively. The sequence $(Y_n)_{n \geq 1}$ **converges to Y in distribution** provided that

$$\lim_{n \rightarrow \infty} F_n(t) = F(t) \quad \text{for all } t \in \mathbb{R} \quad \text{at which } F \text{ is continuous.} \quad (7.9)$$

In this case, one writes $Y_n \xrightarrow{\mathcal{D}} Y$.

Remark 7.2.2. An alternative way to formulate property (7.9) is as follows:

$$\lim_{n \rightarrow \infty} \mathbb{P}\{Y_n \leq t\} = \mathbb{P}\{Y \leq t\} \quad \text{for all } t \in \mathbb{R} \text{ with } \mathbb{P}\{Y = t\} = 0.$$

Without proof, we state two other characterizations of convergence in distribution.

Proposition 7.2.3. One has $Y_n \xrightarrow{\mathcal{D}} Y$ if and only if for all bounded continuous functions $f : \mathbb{R} \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \mathbb{E}f(Y_n) = \mathbb{E}f(Y).$$

Furthermore, this is also equivalent to

$$\limsup_{n \rightarrow \infty} \mathbb{P}\{Y_n \in A\} \leq \mathbb{P}\{Y \in A\}$$

for all closed subsets $A \subseteq \mathbb{R}$, or also to

$$\liminf_{n \rightarrow \infty} \mathbb{P}\{Y_n \in G\} \geq \mathbb{P}\{Y \in G\}$$

whenever $G \subseteq \mathbb{R}$ is an open set.

Remark 7.2.4. Note that in general $Y_n \xrightarrow{\mathcal{D}} Y$ does *not* imply

$$\lim_{n \rightarrow \infty} \mathbb{P}\{Y_n \in B\} = \mathbb{P}\{Y \in B\}$$

for all Borel sets $B \subseteq \mathbb{R}$. For example, if $\mathbb{P}\{Y_n = 1/n\} = 1$ and $\mathbb{P}\{Y = 0\} = 1$, then the Y_n s converge to Y in distribution (check this), but if $B = (0, 1)$, then $\mathbb{P}\{Y_n \in B\}$ does not converge to $\mathbb{P}\{Y \in B\}$.

If the distribution function of Y is continuous, that is, we have $\mathbb{P}\{Y = t\} = 0$ for all $t \in \mathbb{R}$, then $Y_n \xrightarrow{\mathcal{D}} Y$ is equivalent to $\lim_{n \rightarrow \infty} F_n(t) = F(t)$ for all $t \in \mathbb{R}$. Besides, in this case, the type of convergence is stronger as the next proposition shows.

Proposition 7.2.5. *Let Y_1, Y_2, \dots and Y be random variables with $\mathbb{P}\{Y = t\} = 0$ for all $t \in \mathbb{R}$. Then $Y_n \xrightarrow{\mathcal{D}} Y$ implies that the distribution functions converge uniformly:*

$$\lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} |\mathbb{P}\{Y_n \leq t\} - \mathbb{P}\{Y \leq t\}| = 0, \quad \text{hence also}$$

$$\lim_{n \rightarrow \infty} \sup_{a < b} |\mathbb{P}\{a \leq Y_n \leq b\} - \mathbb{P}\{a \leq Y \leq b\}| = 0.$$

We have now all notations and definitions that are necessary to formulate the central limit theorem. Mostly, this theorem is proved via properties of the so-called characteristic functions (see Chapter 3 of [Dur19] for such a proof). For alternative proofs using properties of moment generating functions, we refer to [Rss14] and [Gha19].

Proposition 7.2.6 (Central limit theorem). *Let $(X_j)_{j \geq 1}$ be a sequence of independent identically distributed random variables with finite second moment. Let μ be the expected value of the X_j s and let $\sigma^2 > 0$ be their variance. Then for the sums $S_n = X_1 + \dots + X_n$ it follows that*

$$\frac{S_n - n\mu}{\sigma\sqrt{n}} \xrightarrow{\mathcal{D}} Z. \quad (7.10)$$

Here Z is an $\mathcal{N}(0, 1)$ -distributed random variable.

Since the limit Z in statement (7.10) is a continuous random variable, Proposition 7.2.5 applies, and the central limit theorem may also be formulated as follows.

Proposition 7.2.7. *Suppose $(X_j)_{j \geq 1}$ and S_n are as in Proposition 7.2.6. Then it follows that*

$$\lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq t \right\} - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx \right| = 0 \quad \text{and} \quad (7.11)$$

$$\lim_{n \rightarrow \infty} \sup_{a < b} \left| \mathbb{P}\left\{ a \leq \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq b \right\} - \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx \right| = 0. \quad (7.12)$$

Remark 7.2.8. Recall that Φ denotes the distribution function of the standard normal distribution as introduced in eq. (1.70). Thus, another way to write eq. (7.11) is as follows: if $F_n(t) = \mathbb{P}\{S_n \leq t\}$ denotes the distribution function of S_n , then

$$\sup_{t \in \mathbb{R}} |F_n(\sigma\sqrt{n}t + n\mu) - \Phi(t)| = \sup_{t \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq t \right\} - \Phi(t) \right| \xrightarrow{n \rightarrow \infty} 0. \quad (7.13)$$

Example 7.2.9. Suppose X_1, X_2, \dots are independent Pois_1 -distributed. Hence, their sum $S_n = X_1 + \dots + X_n$ is a Pois_n -distributed random variable, and

$$F_n(t) = \mathbb{P}\{S_n \leq t\} = \sum_{0 \leq k \leq t} \frac{n^k}{k!} e^{-n}.$$

Since $\mu = \sigma = \lambda = 1$, eq. (7.13) tells us that (compare Figure 7.3)

$$F_n(\sqrt{n}t + n) = \sum_{0 \leq k \leq \sqrt{n}t + n} \frac{n^k}{k!} e^{-n} \xrightarrow{n \rightarrow \infty} \Phi(t). \quad (7.14)$$

Moreover, the convergence takes place uniformly in $t \in \mathbb{R}$.

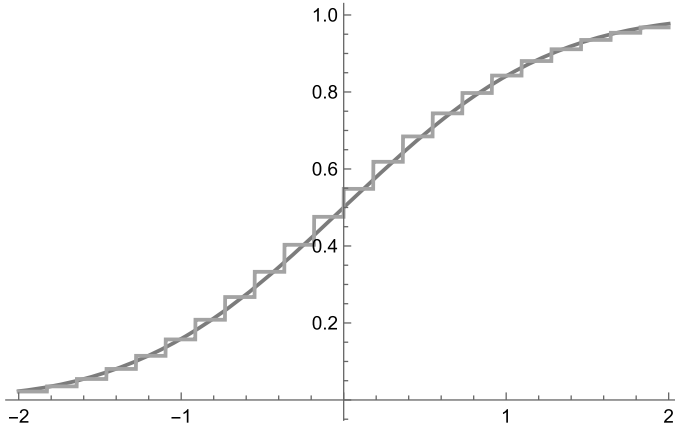


Figure 7.3: The approximation of $\Phi(t)$ in eq. (7.14) with $n = 30$.

Our next objective is another reformulation of eq. (7.12). If we set $a' = a\sigma\sqrt{n} + n\mu$ and $b' = b\sigma\sqrt{n} + n\mu$, then these numbers depend on $n \in \mathbb{N}$. But since the convergence in eq. (7.12) is uniform, we may replace a and b by a' and b' , respectively and obtain

$$\lim_{n \rightarrow \infty} \sup_{a' < b'} \left| \mathbb{P}\{a' \leq S_n \leq b'\} - \mathbb{P}\left\{\frac{a' - n\mu}{\sigma\sqrt{n}} \leq Z \leq \frac{b' - n\mu}{\sigma\sqrt{n}}\right\} \right| = 0. \quad (7.15)$$

Here, as before, Z denotes a standard normally distributed random variable. For a final reformulation, set

$$Z_n := \sigma\sqrt{n}Z + n\mu.$$

Then eq. (7.15) is equivalent to

$$\lim_{n \rightarrow \infty} \sup_{a' < b'} |\mathbb{P}\{a' \leq S_n \leq b'\} - \mathbb{P}\{a' \leq Z_n \leq b'\}| = 0. \quad (7.16)$$

By Proposition 4.2.3, the random variables Z_n are $\mathcal{N}(n\mu, n\sigma^2)$ -distributed, which allows us to interpret eq. (7.15), or eq. (7.16), as follows. If $\mu = \mathbb{E}X_j$ and $\sigma^2 = \mathbb{V}X_j$, then for large n , the sum S_n is “approximately” $\mathcal{N}(n\mu, n\sigma^2)$ -distributed.

In other words, for $-\infty \leq a < b \leq \infty$ it follows that



$$\mathbb{P}\{a \leq S_n \leq b\} \approx \Phi\left(\frac{b - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{a - n\mu}{\sigma\sqrt{n}}\right).$$

Interpretation: We emphasize once more that the central limit theorem is valid for *all* sequences of independent identically distributed random variables possessing a second moment. For example, it is true for X_j s that are binomial distributed, for X_j s being exponentially distributed, and so on. Thus, no matter how the random variables with second moment are distributed, all their normalized sums possess the same limit, the normal distribution. This explains the outstanding role of the normal distribution.

The deeper reason for this phenomenon is that S_n may be viewed as the superposition of many “small” independent errors or perturbations, all of the same kind.⁶ Although each perturbation is distributed according to \mathbb{P}_{X_1} , the independent superposition of the perturbations leads to the fact that the final result is approximately normally distributed. This explains why so many random phenomena may be described by normally distributed random variables.

Remark 7.2.10 (Continuity correction). A slight technical problem arises in the case of discrete random variables X_j . Then the S_n s are discrete as well, hence their distribution functions F_n have jumps. If these noncontinuous functions F_n approximate the continuous function Φ , then certain errors occur at the points where the jumps of F_n are (compare Figure 7.4).

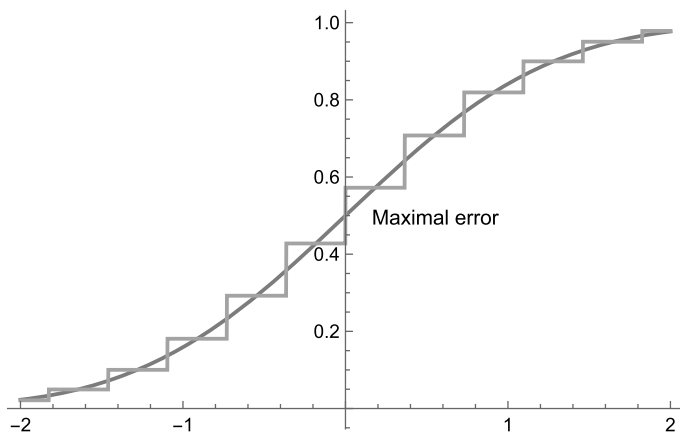


Figure 7.4: A sequence of noncontinuous functions (here a sequence of distribution functions of binomial random variables) approximates the continuous function Φ .

⁶ The central limit theorem also holds for not necessarily identically distributed random variables provided that all “errors” become uniformly small. That is, one has to exclude that certain errors are dominating the others.

To understand the problem, assume that the X_j s possess values in \mathbb{Z} , then S_n is also \mathbb{Z} -valued, hence for any $0 \leq h < 1$, and all integers $k < \ell$, it follows that (see Figure 7.5)

$$\mathbb{P}\{k \leq S_n \leq \ell\} = \mathbb{P}\{k - h \leq S_n \leq \ell + h\}.$$

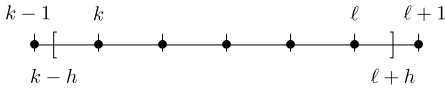


Figure 7.5: If S_n has values in \mathbb{Z} , then for any $0 \leq h < 1$ one has $k \leq S_n \leq \ell$ if and only if $k - h \leq S_n \leq \ell + h$. Thus, both events possess the same probability.

Consequently, for each such number h , the value

$$\Phi\left(\frac{\ell + h - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{k - h - n\mu}{\sigma\sqrt{n}}\right)$$

may be taken as normal approximation of the above probability. Which number $h < 1$ should be chosen?

To answer this question, observe the following. If $k < m < \ell$, then

$$\mathbb{P}\{k \leq S_n \leq \ell\} = \mathbb{P}\{k \leq S_n \leq m\} + \mathbb{P}\{m + 1 \leq S_n \leq \ell\},$$

which, after choosing h in $[0, 1)$, is approximated by

$$\Phi\left(\frac{\ell + h - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{m + 1 - h - n\mu}{\sigma\sqrt{n}}\right) + \Phi\left(\frac{m + h - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{k - h - n\mu}{\sigma\sqrt{n}}\right).$$

Thus, in order to get neither an overlap nor a gap between $m + 1 - h - n\mu$ and $m + h - n\mu$, it is customary to choose $h = 0.5$. This leads to the following definition.

Definition 7.2.11. Suppose X_1, X_2, \dots are independent identically distributed with values in \mathbb{Z} . Then the corrected normal approximation is given by

$$\mathbb{P}\{k \leq S_n \leq \ell\} \approx \Phi\left(\frac{\ell + 0.5 - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{k - 0.5 - n\mu}{\sigma\sqrt{n}}\right).$$

It is called the **continuity correction** or **histogram correction** for the normal approximation. In a similar way, one corrects the approximation for infinite intervals by

$$\mathbb{P}\{S_n \leq \ell\} \approx \Phi\left(\frac{\ell + 0.5 - n\mu}{\sigma\sqrt{n}}\right)$$

and by

$$\mathbb{P}\{S_n \geq k\} \approx 1 - \Phi\left(\frac{k - 0.5 - n\mu}{\sigma\sqrt{n}}\right) = \Phi\left(\frac{n\mu - k + 0.5}{\sigma\sqrt{n}}\right). \quad (7.17)$$

The next result tells us that the continuity correction is only needed for small values of $n \geq 1$.

Proposition 7.2.12. For all $x \in \mathbb{R}$ and $h \in \mathbb{R}$, it follows that

$$\left| \Phi\left(\frac{x+h-n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{x-n\mu}{\sigma\sqrt{n}}\right) \right| \leq \frac{|h|}{\sigma\sqrt{2\pi n}}.$$

Proof. The mean value theorem of Calculus implies the existence of an intermediate value ξ in $(\frac{x-|h|-n\mu}{\sigma\sqrt{n}}, \frac{x+|h|-n\mu}{\sigma\sqrt{n}})$ such that

$$\left| \Phi\left(\frac{x+h-n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{x-n\mu}{\sigma\sqrt{n}}\right) \right| = |h| \frac{\Phi'(\xi)}{\sigma\sqrt{n}}.$$

Using

$$\Phi'(\xi) = \frac{1}{\sqrt{2\pi}} e^{-\xi^2/2} \leq \frac{1}{\sqrt{2\pi}},$$

this proves the asserted estimate. \square

Remark 7.2.13. An application of Proposition 7.2.12 with $x = k$ and/or $x = \ell$, and with $h = \pm 0.5$, shows that the improvement by the continuity correction is at most of order $n^{-1/2}$. Thus, it is no longer needed for large n .

Example 7.2.14. Roll a fair die n times. Let S_n be the sum of the n rolls. In view of eq. (7.16), this sum S_n is approximately $\mathcal{N}(\frac{7n}{2}, \frac{35n}{12})$ -distributed. In other words, it follows that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{a \leq \frac{S_n - 7n/2}{\sqrt{35n/12}} \leq b\right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx = \Phi(b) - \Phi(a).$$

Moreover, this convergence takes place uniformly for all $a < b$. Therefore, at least for large n , the right-hand side of the last equation may be taken as an approximate value of the left-hand one.

At first, we consider an example with a small number of trials. We roll a die three times and ask for the probability of the event $\{7 \leq S_3 \leq 8\}$. Let us compare the exact value

$$\mathbb{P}\{7 \leq S_3 \leq 8\} = \frac{1}{6} = 0.1\bar{6}$$

with that we get by applying the central limit theorem. Without continuity correction, the approximate value is

$$\Phi\left(\frac{8 - 21/2}{\sqrt{3 \cdot 35/12}}\right) - \Phi\left(\frac{7 - 21/2}{\sqrt{3 \cdot 35/12}}\right) \approx 0.08065,$$

while an application of the continuity correction leads to

$$\Phi\left(\frac{8 + 0.5 - 21/2}{\sqrt{3 \cdot 35/12}}\right) - \Phi\left(\frac{7 - 0.5 - 21/2}{\sqrt{3 \cdot 35/12}}\right) \approx 0.16133.$$

We see an improvement using the continuity correction.

Next we treat an example with large n . Let us investigate once more Example 7.1.6, but this time from the point of view of the central limit theorem. Choose again $n = 10^3$, $a = -\frac{100\sqrt{12}}{\sqrt{35000}}$, and $b = \frac{100\sqrt{12}}{\sqrt{35000}}$. Then it follows that

$$\mathbb{P}\{3400 \leq S_n \leq 3600\} \approx \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx \approx 0.93592.$$

As we see, the use of the central limit theorem improves considerably the bound 0.709 obtained by Chebyshev's inequality.

Example 7.2.15. The aim of this example is to apply the central limit theorem for the investigation of the asymptotic behavior of random walks as they were introduced in Example 4.1.7 and Section 5.5. Recall that we suppose $S_0 = 0$ and, if $n \geq 1$, then $S_n = X_1 + \dots + X_n$, where the X_i s are independent and attain the values -1 and 1 with probability $1 - p$ and p , respectively. In eqs. (5.26) and (5.37), we got

$$\mathbb{E}S_n = n(2p - 1) \quad \text{and} \quad \mathbb{V}S_n = 4np(1 - p).$$

Consequently, the central limit theorem leads in this case to the following: for all real numbers $a < b$, we have

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{a \leq \frac{S_n - n(2p - 1)}{2\sqrt{np(1 - p)}} \leq b\right\} = \Phi(b) - \Phi(a).$$

In the case of a symmetric walk, that is, $p = 1/2$, this simplifies to

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{a \leq \frac{S_n}{\sqrt{n}} \leq b\right\} = \Phi(b) - \Phi(a) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

For instance, if $a = -2$ and $b = 2$, then it follows that

$$\lim_{n \rightarrow \infty} \mathbb{P}\{-2\sqrt{n} \leq S_n \leq 2\sqrt{n}\} = \frac{1}{\sqrt{2\pi}} \int_{-2}^2 e^{-x^2/2} dx \approx 0.954.$$

Keep in mind that the possible values of S_n are between $-n$ and n . But in reality, if n is large enough, then with probability greater than 0.95, the value of S_n will be in the much smaller interval $[-2\sqrt{n}, 2\sqrt{n}]$ (see Figure 7.6 for a graphically presentation of this fact).

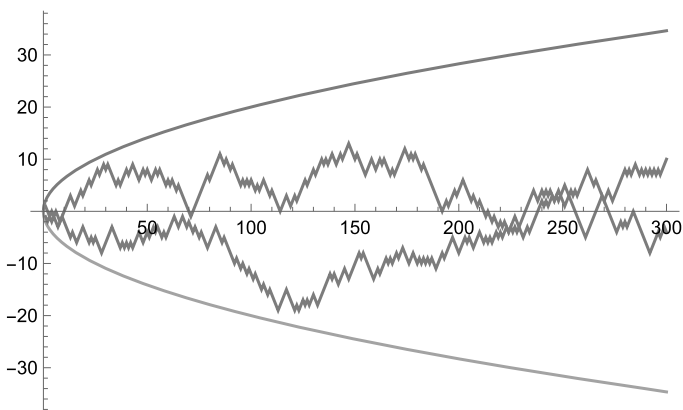


Figure 7.6: Two independent symmetric random walks compared with $t \mapsto \pm 2\sqrt{t}$.

On the other hand, if we ask for the probability that S_n is between $-\sqrt{n}$ and \sqrt{n} , we obtain

$$\lim_{n \rightarrow \infty} \mathbb{P}\{-\sqrt{n} \leq S_n \leq \sqrt{n}\} = \frac{1}{\sqrt{2\pi}} \int_{-1}^1 e^{-x^2/2} dx \approx 0.6827.$$

Maybe more impressive than the previous statements is the following fact: for any $c > 0$, it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}\{S_n \geq c\sqrt{n}\} = 1 - \Phi(c) = \frac{1}{\sqrt{2\pi}} \int_c^{\infty} e^{-t^2/2} dt.$$

Remark 7.2.16. More precise statements about the asymptotic behavior of symmetric random walks are available. For example, the central limit theorem implies

$$\mathbb{P}\left\{\limsup_{n \rightarrow \infty} \frac{S_n}{\sqrt{n}} \geq c\right\} > 0 \tag{7.18}$$

for any $c > 0$. A zero–one law tells us that the probability in (7.18) is not only positive, but equals 1. This leads to

$$\mathbb{P}\left\{\limsup_{n \rightarrow \infty} \frac{S_n}{\sqrt{n}} = \infty\right\} = 1 \quad \text{and} \quad \mathbb{P}\left\{\liminf_{n \rightarrow \infty} \frac{S_n}{\sqrt{n}} = -\infty\right\} = 1.$$

Thus, \sqrt{n} is not the right scaling factor for S_n . Some other, bigger, sequence is needed.

On the other hand, a scaling of S_n by n is also not appropriate. Why? Observe that the strong law of large numbers asserts that

$$\mathbb{P}\left\{\lim_{n \rightarrow \infty} \frac{S_n}{n} = 0\right\} = 1.$$

Hence, an appropriate scaling of the S_n s should be a sequence lying between \sqrt{n} and n . Surprisingly, the correct sequence of normalization is $\sqrt{2n \log \log n}$. The **law of iterated logarithm** due to Hartman and Wintner (see, e. g., [Bil12, Theorem 9.5], or [Kle20]) implies

$$\mathbb{P}\left\{\limsup_{n \rightarrow \infty} \frac{S_n}{\sqrt{2n \log \log n}} = 1\right\} = \mathbb{P}\left\{\liminf_{n \rightarrow \infty} \frac{S_n}{\sqrt{2n \log \log n}} = -1\right\} = 1.$$

Consequently, for any $\varepsilon > 0$ one gets

$$\mathbb{P}\{S_n \geq (1 - \varepsilon)\sqrt{2n \log \log n} \text{ i. o.}\} = 1,$$

while

$$\mathbb{P}\{S_n \geq (1 + \varepsilon)\sqrt{2n \log \log n} \text{ i. o.}\} = 0.$$

Example 7.2.17 (Round-off errors). Many calculations in a bank, for instance, of interest, lead to amounts that are not integral in cents. In this case the bank rounds the calculated value either up or down, whether the remainder is larger or smaller than 0.5 cent. For example, if the calculations lead to \$12.837, then the bank transfers \$12.84. Thus, in this case, the bank loses 0.3 cent. This seems to be a small amount, but if, for example, the bank performs 10^6 calculations per day, the total loss or gain could sum up to an amount of \$5000.00. But does this really happen?

Answer: Theoretically, the rounding procedure could lead to huge losses or gains of the bank. But, as the central limit theorem shows, in reality such a scenario is extremely unlikely. To make this more precise, we use the following model. Let X_j be the loss or gain (in cents) of the bank in calculation j . Then the X_j are independent and uniformly distributed on $[-0.5, 0.5]$. Thus, the total loss or gain after n calculations equals $S_n = X_1 + \dots + X_n$. By Propositions 5.1.27 and 5.2.25, we know that

$$\mu = \mathbb{E}X_1 = 0 \quad \text{and} \quad \sigma^2 = \mathbb{V}X_1 = \frac{1}{12},$$

hence, if $a < b$, the central limit theorem implies

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{\frac{a\sqrt{n}}{\sqrt{12}} \leq S_n \leq \frac{b\sqrt{n}}{\sqrt{12}}\right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

For example, if $n = 10^6$, then taking $a = \sqrt{12}$ and $b = \infty$, this leads to

$$\mathbb{P}\{S_n \geq \$10\} = \mathbb{P}\{S_n \geq 10^3 \text{ cents}\} \approx \frac{1}{\sqrt{2\pi}} \int_{\sqrt{12}}^{\infty} e^{-x^2/2} dx \approx 0.00026603,$$

which is an extremely small probability. By symmetry, it also follows that

$$\mathbb{P}\{S_n \leq -\$10\} \approx \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\sqrt{12}} e^{-x^2/2} dx \approx 0.00026603.$$

In a similar way, one obtains

$$\begin{aligned} \mathbb{P}\{S_n \geq \$1\} &\approx 0.364517, & \mathbb{P}\{S_n \geq \$2\} &\approx 0.244211, \\ \mathbb{P}\{S_n \geq \$5\} &\approx 0.0416323 & \text{and } \mathbb{P}\{S_n \geq \$20\} &\approx 2.1311 \times 10^{-12}. \end{aligned}$$

This shows that even for many calculations, in our case 10^6 , the probability for a loss or gain of more than \$5 is very unlikely. Recall that theoretically an amount of \$5000.00 would be possible.

Example 7.2.18. Suppose n people choose independently of each other an integer in $\{0, \dots, 9\}$. Thereby, each of the 10 numbers is equally likely. Of course, the expected value of the chosen numbers is $\mu = 9/2$. Moreover, the variance of a single choice can be evaluated by

$$\sigma^2 = \frac{1}{10} \sum_{k=0}^9 (k - 9/2)^2 = \frac{33}{4}.$$

Let S_n be the sum of the chosen n numbers. Then the strong law of large numbers yields that with probability 1,

$$\lim_{n \rightarrow \infty} \frac{S_n}{n} = \frac{9}{2}.$$

We ask now how far or near we may expect S_n/n to $9/2$, of course, depending on n .

To answer this question, we apply the central limit theorem. It asserts that, given $a < b$,

$$\mathbb{P}\left\{a \leq 2 \frac{S_n - 9n/2}{\sqrt{33n}} \leq b\right\} \approx \Phi(b) - \Phi(a).$$

Equivalently, this is

$$\mathbb{P}\left\{\frac{a\sqrt{33}}{2\sqrt{n}} \leq \frac{S_n}{n} - \frac{9}{2} \leq \frac{b\sqrt{33}}{2\sqrt{n}}\right\} \approx \Phi(b) - \Phi(a).$$

Setting $a = -2c/\sqrt{33}$ and $b = 2c/\sqrt{33}$, an application of Corollary 4.2.5 leads to

$$\mathbb{P}\left\{\left|\frac{S_n}{n} - \frac{9}{2}\right| \leq \frac{c}{\sqrt{n}}\right\} \approx \Phi\left(\frac{2c}{\sqrt{33}}\right) - \Phi\left(-\frac{2c}{\sqrt{33}}\right) \geq 0.9973$$

provided that $2c/\sqrt{33} \geq 3$, that is, if $c > 8.17$.

We refer to Problem 7.8 for a general approach to the question treated in this example.

Special cases of the central limit theorem

Binomial distributed random variables. In 1738 De Moivre, and later on in 1812 Laplace, investigated the normal approximation of binomial distributed⁷ random variables. This was the starting point for the investigation of general central limit theorems. Let us state their result.

Proposition 7.2.19 (De Moivre–Laplace theorem). *Let the X_j s be independent $B_{1,p}$ -distributed random variables. Then their sums $S_n = X_1 + \dots + X_n$ satisfy*

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ a \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq b \right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx. \quad (7.19)$$

Proof. Recall that for a $B_{1,p}$ -distributed random variable X , we have $\mu = EX = p$ and $\sigma^2 = \mathbb{V}X = p(1-p)$. Consequently, Proposition 7.2.7 applies and leads to eq. (7.19). \square

Remark 7.2.20. By Corollary 4.6.2, we know that $S_n = X_1 + \dots + X_n$ is $B_{n,p}$ -distributed. Consequently, eq. (7.19) may also be written as

$$\lim_{n \rightarrow \infty} \sum_{k \in I_{n,a,b}} \binom{n}{k} p^k (1-p)^{n-k} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx,$$

where

$$I_{n,a,b} := \left\{ k \geq 0 : a \leq \frac{k - np}{\sqrt{np(1-p)}} \leq b \right\}.$$

Another way to formulate the De Moivre–Laplace theorem is as follows. For “large” n , S_n is approximative $\mathcal{N}(np, np(1-p))$ -distributed. That is, if $0 \leq \ell < m \leq n$, then

$$\sum_{k=\ell}^m \binom{n}{k} p^k (1-p)^{n-k} \approx \Phi \left(\frac{m - np}{\sqrt{np(1-p)}} \right) - \Phi \left(\frac{\ell - np}{\sqrt{np(1-p)}} \right). \quad (7.20)$$

Since the sums S_n are integer-valued, the continuity correction should be applied for n small, that is, on the right-hand side of eq. (7.20) the numbers m and ℓ should be replaced by $m + 0.5$ and by $\ell - 0.5$, respectively.

Example 7.2.21. Play a series of games with success probability $0 < p < 1$. Let $\alpha \in (0, 1)$ be a given security probability, and $m \in \mathbb{N}$ is some integer. How many games does one

⁷ De Moivre investigated sums of $B_{1,1/2}$ -distributed random variables while Laplace treated $B_{1,p}$ -distributed ones for general $0 \leq p \leq 1$.

have to play in order to have with probability greater than or equal to $1 - \alpha$ at least m successes?

Answer: Define random variables X_j by setting $X_j = 1$ when winning game j , while $X_j = 0$ in the case of losing it. Then the X_j s are independent and $B_{1,p}$ -distributed. Hence, if $S_n = X_1 + \cdots + X_n$, then the above question may be formulated as follows. What is the smallest $n \in \mathbb{N}$ for which

$$\mathbb{P}\{S_n \geq m\} \geq 1 - \alpha? \quad (7.21)$$

By Corollary 4.6.2, the sum S_n is $B_{n,p}$ -distributed and, therefore, the estimate in (7.21) transforms to

$$\sum_{k=m}^n \binom{n}{k} p^k (1-p)^{n-k} \geq 1 - \alpha. \quad (7.22)$$

Thus, the “exact” answer to the above question is as follows. Choose the minimal $n \geq 1$ for which estimate (7.22) is valid.

Remark 7.2.22. For large m , it may be a difficult task to determine the minimal n satisfying estimate (7.22). Therefore, one looks for an “approximate” approach via Proposition 7.2.19. Rewriting estimate (7.21) as

$$\mathbb{P}\left\{ \frac{S_n - np}{\sqrt{np(1-p)}} \geq \frac{m - np}{\sqrt{np(1-p)}} \right\} \geq 1 - \alpha,$$

an “approximate” condition for n is

$$1 - \alpha \leq 1 - \Phi\left(\frac{m - np}{\sqrt{np(1-p)}}\right) = \Phi\left(\frac{np - m}{\sqrt{np(1-p)}}\right).$$

Given $\beta \in (0,1)$, let us define⁸ z_β by $\Phi(z_\beta) = \beta$. Consequently, an approximate solution of the above question is to choose the minimal $n \geq 1$ satisfying

$$\frac{np - m}{\sqrt{np(1-p)}} \geq z_{1-\alpha}. \quad (7.23)$$

For “small” n , we have to modify the previous approach slightly. Here we have to use the continuity correction. In view of eq. (7.17), the condition is now

$$1 - \alpha \leq \Phi\left(\frac{np - m + 0.5}{\sqrt{np(1-p)}}\right),$$

⁸ Later on, in Proposition 8.4.3, these numbers z_β will play an important role; compare also with Definition 8.4.8.

leading to

$$\frac{np - m + 0.5}{\sqrt{np(1-p)}} \geq z_{1-\alpha}. \quad (7.24)$$

Let us explain Remark 7.2.22 with the help of a concrete example.

Example 7.2.23. Find the minimal $n \geq 1$ such that, rolling a fair die n times, one observes with probability greater than or equal to 0.9 at least 100 times the number 6?

For the “exact” answer choose the minimal $n \geq 1$ satisfying

$$\sum_{k=100}^n \binom{n}{k} \left(\frac{1}{6}\right)^k \left(\frac{5}{6}\right)^{n-k} \geq 0.9.$$

Numerical calculations give that the left-hand side equals 0.897721 if $n = 670$, and it is 0.900691 if $n = 671$. Thus, in order to observe, with probability greater than 0.9, the number 6 at least 100 times, one has to roll the die at least 671 times.

Let us compare this result with that we obtained by the approximation approach. First, we approximate S_n directly, that is, without applying the continuity correction. Here estimate (7.23) says that we have to look for the minimal $n \geq 1$ satisfying

$$\frac{\frac{n}{6} - m}{\sqrt{\frac{1}{6} \cdot \frac{5}{6} \cdot n}} = \frac{n - 600}{\sqrt{5n}} \geq z_{0.9} = 1.28155. \quad (7.25)$$

Since

$$\frac{665 - 600}{\sqrt{5 \cdot 665}} = 1.12724 \quad \text{and} \quad \frac{666 - 600}{\sqrt{5 \cdot 666}} = 1.4373,$$

the smallest n satisfying estimate (7.25) is 666.

Applying the continuity correction, by estimate (7.24), condition (7.25) has to be replaced by

$$\frac{\frac{n}{6} - m + 0.5}{\sqrt{\frac{1}{6} \cdot \frac{5}{6} \cdot n}} = \frac{n - 600 + 3}{\sqrt{5n}} \geq z_{0.9} = 1.28155.$$

The left-hand side equals 1.27757 for $n = 671$ and 1.29387 if $n = 672$. Consequently, this type of approximation gives the (more precise) value $n = 672$ for the minimal number of necessary rolls of the die.

Poisson distributed random variables. Let X_1, X_2, \dots be independent and Pois_λ -distributed. By Propositions 5.1.16 and 5.2.22, we know

$$\mathbb{E}X_1 = \lambda \quad \text{and} \quad \mathbb{V}X_1 = \lambda.$$

Thus, in this case Proposition 7.2.7 reads as follows.

Proposition 7.2.24. Let $(X_j)_{j \geq 1}$ be independent Pois_λ -distributed random variables. Then the sums $S_n = X_1 + \cdots + X_n$ satisfy

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ a \leq \frac{S_n - n\lambda}{\sqrt{n\lambda}} \leq b \right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx. \quad (7.26)$$

Remark 7.2.25. By Proposition 4.6.5, the sum S_n is $\text{Pois}_{\lambda n}$ -distributed, hence eq. (7.26) transforms to

$$\lim_{n \rightarrow \infty} \sum_{k \in J_{n,a,b}} \frac{(\lambda n)^k}{k!} e^{-\lambda n} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx, \quad (7.27)$$

where

$$J_{n,a,b} := \left\{ k \in \mathbb{N}_0 : a \leq \frac{k - n\lambda}{\sqrt{n\lambda}} \leq b \right\}.$$

Another way to express this is as follows. If $0 \leq \ell < m < \infty$, then

$$\sum_{k=\ell}^m \frac{(\lambda n)^k}{k!} e^{-\lambda n} \approx \Phi\left(\frac{m - n\lambda}{\sqrt{n\lambda}}\right) - \Phi\left(\frac{\ell - n\lambda}{\sqrt{n\lambda}}\right).$$

Remark 7.2.26. Choosing in eq. (7.27) the numbers as $a = -\infty$, $b = 0$, and $\lambda = 1$, we get

$$\lim_{n \rightarrow \infty} e^{-n} \sum_{k=0}^n \frac{n^k}{k!} = \frac{1}{2},$$

which is interesting in its own right. Taking $a = -\infty$ and $b_n = \sqrt{n}$ yields

$$\lim_{n \rightarrow \infty} \left| e^{-n} \sum_{k=0}^{2n} \frac{n^k}{k!} - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{b_n} e^{-x^2/2} dx \right| = 0,$$

hence, because of

$$\lim_{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\sqrt{n}} e^{-x^2/2} dx = 1,$$

we obtain

$$\lim_{n \rightarrow \infty} e^{-n} \sum_{k=0}^{2n} \frac{n^k}{k!} = 1.$$

Gamma distributed random variables. Finally, we investigate sums of gamma distributed random variables. Here the central limit theorem leads to the following result.

Proposition 7.2.27. Let X_1, X_2, \dots be independent $\Gamma_{\alpha, \beta}$ -distributed random variables. Then their sums $S_n = X_1 + \dots + X_n$ satisfy

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ a \leq \frac{S_n - n\alpha\beta}{\alpha\sqrt{n\beta}} \leq b \right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx. \quad (7.28)$$

Proof. Propositions 5.1.28 and 5.2.26 tell us that the expected value and the variance of the X_j s are given by $\mu = \mathbb{E}X_1 = \alpha\beta$ and $\sigma^2 = \mathbb{V}X_1 = \alpha^2\beta$. Therefore, eq. (7.28) follows by an application of Proposition 7.2.7. \square

Remark 7.2.28. Note that Proposition 4.6.4 implies that S_n is $\Gamma_{\alpha, n\beta}$ -distributed. Thus, setting

$$I_{n, a, b} := \left\{ x \geq 0 : a \leq \frac{x - n\alpha\beta}{\alpha\sqrt{n\beta}} \leq b \right\},$$

eq. (7.28) leads to

$$\lim_{n \rightarrow \infty} \frac{1}{\alpha^{n\beta} \Gamma(n\beta)} \int_{I_{n, a, b}} x^{n\beta-1} e^{-x/\alpha} dx = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

Another way to express this is as follows. If $0 \leq a < b$, then

$$\frac{1}{\alpha^{n\beta} \Gamma(n\beta)} \int_a^b x^{n\beta-1} e^{-x/\alpha} dx \approx \Phi \left(\frac{b - n\alpha\beta}{\alpha\sqrt{n\beta}} \right) - \Phi \left(\frac{a - n\alpha\beta}{\alpha\sqrt{n\beta}} \right).$$

Two cases of Proposition 7.2.27, or Remark 7.2.28, are of special interest.

(a) For $n \geq 1$, let S_n be a χ_n^2 -distributed random variable. Then it follows that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ a \leq \frac{S_n - n}{\sqrt{2n}} \leq b \right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

Another way to express this as follows: if the S_n s are χ_n^2 -distributed, then for $t \in \mathbb{R}$ and n sufficiently large,

$$\mathbb{P}\{S_n \leq \sqrt{2n}t + n\} \approx \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx = \Phi(t).$$

In Fig. 7.7, one sees how good the approximation of Φ is even for small n .

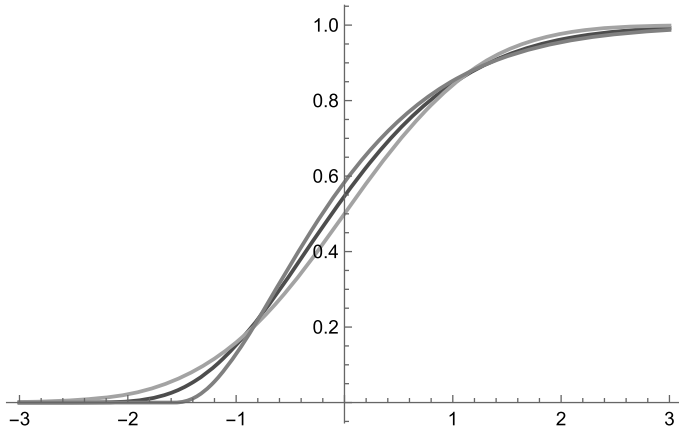


Figure 7.7: The normalized and shifted distribution function $\mathbb{P}\{S_n \leq \sqrt{2n}t + n\}$ of a χ_n^2 -distributed random variable S_n . We chose $n = 5$ (upper graph at zero) and $n = 16$ (middle graph), in comparison with the approximated $\Phi(t)$ (lower graph).

(b) If S_n is distributed according to the Erlang distribution $E_{\lambda,n}$, then we get

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{a \leq \frac{\lambda S_n - n}{\sqrt{n}} \leq b\right\} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

For $\lambda = 1$, this implies (set $a = -\infty$ and $b = 0$) that

$$\lim_{n \rightarrow \infty} \frac{1}{\Gamma(n)} \int_0^n x^{n-1} e^{-x} dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^0 e^{-x^2/2} dx = \frac{1}{2}.$$

Additional remarks

(1) We play a series of the same game. Suppose in each game we may lose or win a certain amount of money. A natural condition for these games (among friends) is whether it should be fair. But *what does it mean that a series of games is fair?* Is this the case

- (i) if the average loss or gain in each single game is zero, or
- (ii) if the probability that, after n games, the total loss or gain is positive, tends to $1/2$ as n tends to infinity?

The mathematical formulation of the previous question is as follows. Let X_1, X_2, \dots denote the win or loss in the first game, the second, and so on. Then the X_j 's are independent identically distributed random variables. The above question reads now as follows. Is the game fair

- (i) if the expected value $\mu = \mathbb{E}X_1$ satisfies $\mu = 0$, or

(ii) if the sum $S_n := X_1 + \cdots + X_n$ fulfills

$$\lim_{n \rightarrow \infty} \mathbb{P}\{S_n \leq 0\} = \lim_{n \rightarrow \infty} \mathbb{P}\{S_n \geq 0\} = \frac{1}{2} ? \quad (7.29)$$

In the sequel, we have to exclude the trivial case $\mathbb{P}\{X_j = 0\} = 1$, that is, in each game one neither wins, nor loses some money. Of course, then eq. (7.29) does not hold.

At a first glance, one might believe that the two possible definitions of fairness describe the same fact. But this is not so as one may see in an example in [Fel68, Chapter X, Section 4]. There one finds a sequence of independent random variables X_1, X_2, \dots with $\mathbb{E}X_1 = 0$, however,

$$\lim_{n \rightarrow \infty} \mathbb{P}\{S_n \leq 0\} = 1.$$

In particular, this tells us that, in general, condition (i) does not imply condition (ii).

The next result clarifies the relation between these two definitions of fairness in the case that the random variables possess a finite second moment.

Proposition 7.2.29. *Let X_1, X_2, \dots be independent and identically distributed with expected value μ . Assume $\mathbb{P}\{X_j = 0\} < 1$.*

1. *Then eq. (7.29) always implies $\mu = 0$. That is, a fair game in the sense of (ii) also satisfies condition (i).*
2. *Conversely, if $\mathbb{E}|X_1|^2 < \infty$, then (ii) is a consequence of (i). Hence, assuming the existence of a second moment, conditions (i) and (ii) are equivalent.*

Proof. We prove the contraposition of the first statement. Thus, suppose that (i) does not hold, that is, we have $\mu \neq 0$. Without losing generality, we may assume $\mu > 0$. Otherwise, investigate $-X_1, -X_2, \dots$. An application of Proposition 7.1.28 with $\varepsilon = \mu/2$ yields

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{\left|\frac{S_n}{n} - \mu\right| \leq \frac{\mu}{2}\right\} = 1. \quad (7.30)$$

Since $\left|\frac{S_n}{n} - \mu\right| \leq \mu/2$ implies $\frac{S_n}{n} \geq \mu/2$, hence $S_n \geq 0$, it follows that

$$\mathbb{P}\left\{\left|\frac{S_n}{n} - \mu\right| \leq \frac{\mu}{2}\right\} \leq \mathbb{P}\{S_n \geq 0\}.$$

Consequently, from eq. (7.30) we derive

$$\lim_{n \rightarrow \infty} \mathbb{P}\{S_n \geq 0\} = 1,$$

hence eq. (7.29) cannot be valid. This proves the first part of the proposition.

We prove now the second assertion. Thus, suppose $\mu = 0$ as well as the existence of the variance $\sigma^2 = \mathbb{V}X_1$. Note that $\sigma^2 > 0$. Why? If a random variable X satisfies $\mathbb{E}X = 0$

and $\mathbb{V}X = 0$, then necessarily $\mathbb{P}\{X = 0\} = 1$. But, since we assumed $\mathbb{P}\{X_1 = 0\} < 1$, we cannot have $\sigma^2 = \mathbb{V}X_1 = 0$.

Thus, Proposition 7.2.7 applies and leads to

$$\lim_{n \rightarrow \infty} \mathbb{P}\{S_n \geq 0\} = \lim_{n \rightarrow \infty} \left\{ \frac{S_n}{\sigma\sqrt{n}} \geq 0 \right\} = \frac{1}{\sqrt{2\pi}} \int_0^{\infty} e^{-x^2/2} dx = \frac{1}{2}.$$

The proof for $\mathbb{P}\{S_n \leq 0\} \rightarrow 1/2$ follows in the same way, thus eq. (7.29) is valid. This completes the proof. \square

(2) *How fast does $\frac{S_n - n\mu}{\sigma\sqrt{n}}$ converge to a normally distributed random variable?* Before we answer this question, we have to determine how this speed is measured. In view of Proposition 7.2.7, we use the following quantity depending on $n \geq 1$:

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq t \right\} - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx \right|.$$

Doing so, the following classical result holds (see [Dur19, Section 3.4.4], for a proof).

Proposition 7.2.30 (Berry–Esséen theorem). *Let X_1, X_2, \dots be independent identically distributed random variables with finite third moment, that is, with $\mathbb{E}|X_1|^3 < \infty$. If $\mu = \mathbb{E}X_1$ and $\sigma^2 = \mathbb{V}X_1 > 0$, then it follows that*

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq t \right\} - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx \right| \leq C \cdot \frac{\mathbb{E}|X_1|^3}{\sigma^3} n^{-1/2}. \quad (7.31)$$

Here $C > 0$ denotes a universal constant.

Remark 7.2.31. The order $n^{-1/2}$ in estimate (7.31) is optimal and cannot be improved. This can be seen by the following example. Take independent random variables X_1, X_2, \dots with $\mathbb{P}\{X_j = -1\} = \mathbb{P}\{X_j = 1\} = 1/2$. Hence, in this case $\mu = 0$ and $\sigma^2 = 1$. Then one has

$$\liminf_{n \rightarrow \infty} n^{1/2} \sup_{t \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{S_n}{\sqrt{n}} \leq t \right\} - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx \right| > 0. \quad (7.32)$$

Assertion (7.32) is a consequence of the fact that, if n is even, then the function $t \mapsto \mathbb{P}\{\frac{S_n}{\sqrt{n}} \leq t\}$ has a jump of order $n^{-1/2}$ at zero. This follows by the calculations in Example 4.1.7. On the other hand, $t \mapsto \Phi(t)$ is continuous, hence the maximal difference between these two functions is at least half of the height of the jump.

Remark 7.2.32. The exact value of the constant $C > 0$ appearing in estimate (7.31) is, in spite of intensive investigations, still unknown. At present, the best-known estimates are $0.40973 < C < 0.4748$.

Summary: The central limit theorem belongs to the most important mathematical results. It explains why so many random observations in nature, community or business, etc., are distributed according to the normal distribution.

The precise statement of the central limit theorem is as follows: if X_1, X_2, \dots are independent identically distributed random variables with expected value μ and variance $\sigma^2 > 0$, then their sum $S_n = X_1 + \dots + X_n$ is approximative $\mathcal{N}(n\mu, n\sigma^2)$ -distributed. After normalizing the sum S_n in the right way, this says that

$$\mathbb{P}\left\{a \leq \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq b\right\} \xrightarrow{n \rightarrow \infty} \Phi(b) - \Phi(a) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

In case of integer valued X_j s and small n the continuity correction improves the approximation of S_n by the normal distribution as follows: if $k \leq \ell$ are integers, then one uses as approximation

$$\mathbb{P}\{k \leq S_n \leq \ell\} \approx \Phi\left(\frac{\ell + 0.5 - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{k - 0.5 - n\mu}{\sigma\sqrt{n}}\right).$$

7.3 Problems

Problem 7.1. Let A_1, A_2, \dots and B_1, B_2, \dots be two sequences of events in a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Prove that

$$\limsup_{n \rightarrow \infty} (A_n \cup B_n) = \limsup_{n \rightarrow \infty} (A_n) \cup \limsup_{n \rightarrow \infty} (B_n).$$

Is this also valid for the intersection? That is, does one have

$$\limsup_{n \rightarrow \infty} (A_n \cap B_n) = \limsup_{n \rightarrow \infty} (A_n) \cap \limsup_{n \rightarrow \infty} (B_n)?$$

Problem 7.2. Let A_1, A_2, \dots be a sequence of subsets in Ω . Show that

$$\mathbb{1}_{\{\liminf A_n\}} = \liminf_{n \rightarrow \infty} \mathbb{1}_{A_n} \quad \text{and} \quad \mathbb{1}_{\{\limsup A_n\}} = \limsup_{n \rightarrow \infty} \mathbb{1}_{A_n}.$$

Here $\mathbb{1}_A$ denotes the indicator function of a set A as introduced in Definition 3.6.16.

Problem 7.3. Let $(X_n)_{n \geq 1}$ be a sequence of independent E_λ -distributed random variables. Characterize sequences $(c_n)_{n \geq 1}$ of positive real numbers for which

$$\mathbb{P}\{X_n \geq c_n \text{ i. o.}\} = 1?$$

Problem 7.4. Let $f : [0, 1] \rightarrow \mathbb{R}$ be a continuous function. Its **Bernstein polynomial** B_n^f of degree n is defined by

$$B_n^f(x) := \sum_{k=0}^n f\left(\frac{k}{n}\right) \binom{n}{k} x^k (1-x)^{n-k}, \quad 0 \leq x \leq 1.$$

Show that Proposition 7.1.30 implies the following. If \mathbb{P} is the uniform distribution on $[0, 1]$, then

$$\mathbb{P}\left\{x \in [0, 1] : \lim_{n \rightarrow \infty} B_n^f(x) = f(x)\right\} = 1.$$

Remark: Using methods from Calculus, one may even show the uniform convergence, that is,

$$\lim_{n \rightarrow \infty} \sup_{0 \leq x \leq 1} |B_n^f(x) - f(x)| = 0.$$

Problem 7.5. Roll a fair die 180 times. What is the probability that the number “6” occurs at most 25 times. Determine this probability by the following three methods:

- Directly via the binomial distribution.
- Approximately by virtue of the central limit theorem.
- Approximately by applying the continuity correction.

Problem 7.6. Toss a fair coin 16 times. Compute the probability to observe exactly eight times “heads” by the following methods:

- Directly via the binomial distribution.
- Approximately by applying the continuity correction.

Why does one not get a reasonable result using the normal approximation directly, that is, without continuity correction?

Problem 7.7. Let X_1, X_2, \dots be a sequence of independent G_p -distributed random variables, that is, for some $0 < p < 1$ one has

$$\mathbb{P}\{X_j = k\} = p(1-p)^{k-1}, \quad k = 1, 2, \dots$$

1. What does the central limit theorem tell us in this case about the behavior of the sums $S_n = X_1 + \dots + X_n$?
2. For two real numbers $a < b$, set

$$I_{n,a,b} := \left\{k \geq 0 : a \leq \frac{pk - n(1-p)}{\sqrt{n(1-p)}} \leq b\right\}.$$

Argue why

$$\lim_{n \rightarrow \infty} \sum_{k \in I_{n,a,b}} \binom{-n}{k} p^n (1-p)^k = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

Hint: Use Corollary 4.6 and investigate $S_n - n$.

Problem 7.8. Extend the question treated in Example 7.2.18 to the general setting. That is, given independent, identically distributed random variables X_1, X_2, \dots with expected value μ and variance σ^2 , find $c > 0$ depending on σ^2 such that for sufficiently large n it follows that

$$\mathbb{P}\left\{\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| \leq \frac{c}{\sqrt{n}}\right\} \geq \Phi(3) - \Phi(-3) \approx 0.9973.$$

8 Mathematical statistics

8.1 Statistical models

8.1.1 Nonparametric statistical models

The main objective of Probability Theory is to describe and analyze random experiments by means of a suitable probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Here it is always assumed that the probability space is known, in particular, that the describing probability measure, \mathbb{P} , is identified.

**Probability Theory:**

Description of a random experiment and its properties by a probability space. The distribution of the outcomes is assumed to be *known*.

Mathematical Statistics deals mainly with the reverse question: one executes an experiment, that is, one draws a sample (e. g., one takes a series of measurements of an item or one questions several people), and, on the basis of the observed sample, one wants to derive as much information as possible about the (unknown) underlying probability measure \mathbb{P} . Sometimes the precise knowledge of \mathbb{P} is not needed; it may suffice to know a certain parameter of \mathbb{P} .

**Mathematical Statistics:**

As a result of a statistical experiment, a (random) sample is observed. On its basis, conclusions are drawn about the unknown underlying probability distribution.

Let us state the mathematical formulation of the task: first, we mention that it is standard practice in Mathematical Statistics to denote the describing probability space by $(\mathcal{X}, \mathcal{F}, \mathbb{P})$. As before, \mathcal{X} is the sample space (the set that contains all possible outcomes of the experiment), and \mathcal{F} is a suitable σ -field of events. The probability measure \mathbb{P} describes the experiment, that is, $\mathbb{P}(A)$ is the probability of observing a sample belonging to A , but recall that \mathbb{P} is unknown.

Based on theoretical considerations or on long-time experience, quite often we are able to restrict the entirety of probability measures in question. Mathematically, this means that we choose a set \mathbf{P} of probability measures on $(\mathcal{X}, \mathcal{F})$ which contains what we believe to be the “correct” \mathbb{P} . Thereby, it is not impossible that \mathbf{P} is the set of *all* probability measures, but for most statistical methods it is very advantageous to take \mathbf{P} as small as possible. On the other hand, the set \mathbf{P} cannot be chosen too small, because we have to be sure that the “correct” \mathbb{P} is really contained in \mathbf{P} . Otherwise, the obtained results are either false or imprecise.

Definition 8.1.1. A subset \mathbf{P} of probability measures on $(\mathcal{X}, \mathcal{F})$ is called a **distribution assumption**, that is, one assumes that the underlying (unknown) \mathbb{P} belongs to the collection \mathbf{P} .

After having fixed the distribution assumption \mathbf{P} , one now regards only probability measures $\mathbb{P} \in \mathbf{P}$ or, equivalently, measures not in \mathbf{P} are discarded.

To get information about the unknown probability measure, one performs a statistical experiment or analyzes some given data. In both cases, the result is a random **sample** $x \in \mathcal{X}$. The task of Mathematical Statistics is to get information about $\mathbb{P} \in \mathbf{P}$, based on the observed sample $x \in \mathcal{X}$. A suitable way to describe the problem is as follows.

Definition 8.1.2. A (nonparametric) **statistical model** is a collection of probability spaces $(\mathcal{X}, \mathcal{F}, \mathbb{P})$ with $\mathbb{P} \in \mathbf{P}$. Here, \mathcal{X} and \mathcal{F} are fixed, and \mathbb{P} varies through the distribution assumption \mathbf{P} . One writes for the model

$$(\mathcal{X}, \mathcal{F}, \mathbb{P})_{\mathbb{P} \in \mathbf{P}} \quad \text{or} \quad \{(\mathcal{X}, \mathcal{F}, \mathbb{P}) : \mathbb{P} \in \mathbf{P}\}.$$

Let us illustrate the previous definition with two examples.

Example 8.1.3. In an urn there are white and black balls of an unknown ratio. Let $\theta \in [0, 1]$ be the (unknown) proportion of white balls, hence $1 - \theta$ is that of the black ones. In order to get some information about θ , one randomly chooses n balls with replacement. The result of this experiment, or the sample, is a number $k \in \{0, \dots, n\}$, the frequency of observed white balls. Thus, the sample space is $\mathcal{X} = \{0, \dots, n\}$ and as σ -field we may choose, as always for finite sample spaces, the powerset $\mathcal{P}(\mathcal{X})$. The possible probability measures describing this experiment are binomial distributions $B_{n, \theta}$ with $0 \leq \theta \leq 1$. Consequently, the distribution assumption is

$$\mathbf{P} = \{B_{n, \theta} : \theta \in [0, 1]\}.$$

Summing up, the statistical model describing the experiment is

$$(\mathcal{X}, \mathcal{P}(\mathcal{X}), \mathbb{P})_{\mathbb{P} \in \mathbf{P}} \quad \text{where } \mathcal{X} = \{0, \dots, n\} \quad \text{and} \quad \mathbf{P} = \{B_{n, \theta} : 0 \leq \theta \leq 1\}.$$

Next, we consider an important example from quality control.

Example 8.1.4. A buyer obtains from a trader a delivery of N machines. Among them $M \leq N$ are defective. The buyer does not know the value of M . To determine it, he randomly chooses n machines from the delivery and checks them. The result, or the sample, is the number $0 \leq m \leq n$ of defective machines among the n tested.

Thus, the sample space is $\mathcal{X} = \{0, \dots, n\}$, $\mathcal{F} = \mathcal{P}(\mathcal{X})$, and the probability measures in question are hypergeometric ones. Therefore, the distribution assumption is

$$\mathbf{P} = \{H_{N, M, n} : M = 0, \dots, N\},$$

where $H_{N, M, n}$ denotes the hypergeometric distribution with parameters N , M , and n , as introduced in Definition 1.4.32.

Before we proceed further, we consider a particularly interesting case of statistical model, which describes the *n-fold independent repetition* of a single experiment. To explain this model, let us investigate the following easy example.

Example 8.1.5. We are given a die that looks biased. To check this, we roll it n times and record the sequence of numbers appearing in each of the trials. Thus, our sample space is $\mathcal{X} = \{1, \dots, 6\}^n$, and the observed sample is $x = (x_1, \dots, x_n)$, with $1 \leq x_k \leq 6$. Let $\theta_1, \dots, \theta_6$ be the probabilities for 1 to 6. Then we want to check whether or not $\theta_1 = \dots = \theta_6 = \frac{1}{6}$, that is, whether \mathbb{P}_0 given by $\mathbb{P}_0(\{k\}) = \theta_k$, $1 \leq k \leq 6$, is the uniform distribution. What are the possible probability measures on $(\mathcal{X}, \mathcal{P}(\mathcal{X}))$ describing the statistical experiment? Since the results of different rolls are independent, the describing measure \mathbb{P} is of the form $\mathbb{P} = \mathbb{P}_0^{\otimes n}$ with

$$\mathbb{P}_0^{\otimes n}(\{x\}) = \mathbb{P}_0(\{x_1\}) \cdots \mathbb{P}_0(\{x_n\}) = \theta_1^{m_1} \cdots \theta_6^{m_6}, \quad x = (x_1, \dots, x_n),$$

and where the m_k s denote the frequency of the number $1 \leq k \leq 6$ in the sequence x . Consequently, the natural distribution assumption is

$$\mathbf{P} = \{\mathbb{P}_0^{\otimes n} : \mathbb{P}_0 \text{ probability measure on } \{1, \dots, 6\}\}.$$

Suppose we are given a probability space $(\mathcal{X}_0, \mathcal{F}_0, \mathbb{P}_0)$ with unknown $\mathbb{P}_0 \in \mathbf{P}_0$. Here, \mathbf{P}_0 denotes a set of probability measures on $(\mathcal{X}_0, \mathcal{F}_0)$, hopefully containing the “correct” \mathbb{P}_0 . We call $(\mathcal{X}_0, \mathcal{F}_0, \mathbf{P}_0)_{\mathbb{P}_0 \in \mathbf{P}_0}$ the *initial model*. In Example 8.1.5, the initial model is $\mathcal{X}_0 = \{1, \dots, 6\}$, while \mathbf{P}_0 is the set of all probability measures on $(\mathcal{X}_0, \mathcal{P}(\mathcal{F}_0))$.

In order to determine \mathbb{P}_0 , we execute n independent trials according to \mathbb{P}_0 . The result, or the observed sample, is a vector $x = (x_1, \dots, x_n)$ with $x_i \in \mathcal{X}_0$. Consequently, the natural sample space is $\mathcal{X} = \mathcal{X}_0^n$.

Which statistical model does this experiment describe? To answer this question, let us recall the basic results in Section 1.9, where exactly those problems have been investigated. As σ -field \mathcal{F} , we choose the n -fold product σ -field of \mathcal{F}_0 , that is,

$$\mathcal{F} = \underbrace{\mathcal{F}_0 \otimes \cdots \otimes \mathcal{F}_0}_{n \text{ times}},$$

and the describing probability measure \mathbb{P} is of the form $\mathbb{P}_0^{\otimes n}$, that is, it is the n -fold product of \mathbb{P}_0 . Recall that, according to Definition 1.9.5, the product $\mathbb{P}_0^{\otimes n}$ is the unique probability measure on $(\mathcal{X}, \mathcal{F})$ satisfying

$$\mathbb{P}_0^{\otimes n}(A_1 \times \cdots \times A_n) = \mathbb{P}_0(A_1) \cdots \mathbb{P}_0(A_n),$$

whenever $A_j \in \mathcal{F}_0$. Since we assumed $\mathbb{P}_0 \in \mathbf{P}_0$, the possible probability measures are $\mathbb{P}_0^{\otimes n}$ with $\mathbb{P}_0 \in \mathbf{P}_0$.

Let us summarize what we obtained until now.

Definition 8.1.6. The statistical model for the n -fold independent repetition of an experiment, determined by the initial model $(\mathcal{X}_0, \mathcal{F}_0, \mathbb{P}_0)_{\mathbb{P}_0 \in \mathbf{P}_0}$, is given by

$$(\mathcal{X}, \mathcal{F}, \mathbb{P}_0^{\otimes n})_{\mathbb{P}_0 \in \mathbf{P}_0}$$

where $\mathcal{X} = \mathcal{X}_0^n$, \mathcal{F} denotes the n -fold product σ -field of \mathcal{F}_0 , and $\mathbb{P}_0^{\otimes n}$ is the n -fold product measure of \mathbb{P}_0 .

Remark 8.1.7. Of course, the main goal in the model of n -fold repetition is to get some knowledge about \mathbb{P}_0 . To obtain the desired information, we perform n independent trials, each time observing a value distributed according to \mathbb{P}_0 . Altogether, the sample is a vector $x = (x_1, \dots, x_n)$, which is now distributed according to $\mathbb{P}_0^{\otimes n}$.

The two following examples explain Definition 8.1.6.

Example 8.1.8. A coin is labeled on one side with “0” and on the other side with “1.” There is some evidence that the coin is biased. To check this, let us execute the following statistical experiment: toss the coin n times and record the sequence of zeroes and ones. Thus, the observed sample is an $x = (x_1, \dots, x_n)$, with each x_k being either “0” or “1.”

Our initial model is given by $\mathcal{X}_0 = \{0, 1\}$ and $\mathbb{P}_0 = B_{1, \theta}$ for a certain (unknown) $\theta \in [0, 1]$. Then the experiment is described by $\mathcal{X} = \{0, 1\}^n$ and $\mathbf{P} = \{B_{1, \theta}^{\otimes n} : 0 \leq \theta \leq 1\}$. Note that

$$B_{1, \theta}^{\otimes n}(\{x\}) = \theta^k (1 - \theta)^{n-k}, \quad k = x_1 + \dots + x_n.$$

Example 8.1.9. A company produces a new type of light bulb with an unknown distribution of the lifetime. To determine it, n light bulbs are switched on at the same time. Let $t = (t_1, \dots, t_n)$ be the times when the bulbs burn out. Then our sample is the vector $t \in (0, \infty)^n$.

From long-time experience, one knows the lifetime of each light bulb is exponentially distributed. Thus, the initial model is $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathbf{P}_0)$ with $\mathbf{P}_0 = \{E_\lambda : \lambda > 0\}$. Consequently, the experiment of testing n light bulbs is described by the model

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathbb{P}_0^{\otimes n})_{\mathbb{P}_0 \in \mathbf{P}_0} = (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), E_\lambda^{\otimes n})_{\lambda > 0},$$

where $\mathbf{P}_0 = \{E_\lambda : \lambda > 0\}$. Recall that $E_\lambda^{\otimes n}$ is the probability measure on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ with density $p(t_1, \dots, t_n) = \lambda^n e^{-\lambda(t_1 + \dots + t_n)}$ for $t_j \geq 0$.

Summary: The basic problem in Mathematical Statistics is as follows: One observes a random sample x belonging to a sample space \mathcal{X} and, depending on this observed x , one wants to get as much information as possible about the underlying unknown probability measure \mathbb{P} . The statistical model to describe this task is a triple

$$(\mathcal{X}, \mathcal{F}, \mathbb{P})_{\mathbb{P} \in \mathbf{P}} \quad \text{or} \quad \{(\mathcal{X}, \mathcal{F}, \mathbb{P}) : \mathbb{P} \in \mathbf{P}\}.$$

Here \mathcal{X} is the sample space, \mathcal{F} denotes a σ -field of subsets of \mathcal{X} , mostly $\mathcal{P}(\mathcal{X})$ or $\mathcal{B}(\mathbb{R})$, and \mathbf{P} is a certain set of probability measures defined on \mathcal{F} . The collection \mathbf{P} of probability measures is said to be the distribution assumption. One conjectures or knows by theoretical considerations that the underlying unknown probability measure \mathbb{P} belongs to \mathbf{P} .

8.1.2 Parametric statistical models

In all of our previous examples, there was a parameter that parametrized the probability measures in \mathbf{P} in natural way. In Example 8.1.3, this is the parameter $\theta \in [0, 1]$, in Example 8.1.4, the probability measures are parametrized by $M \in \{0, \dots, N\}$, in Example 8.1.8 the parameter is also $\theta \in [0, 1]$, and, finally, in Example 8.1.9 the natural parameter is $\lambda > 0$. Therefore, from now on, we assume that there is a parameter set Θ such that \mathbf{P} may be represented as

$$\mathbf{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}.$$

Definition 8.1.10. A **parametric statistical model** is defined as

$$(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$$

with **parameter set** Θ . Equivalently, we suppose that the distribution assumption \mathbf{P} , appearing in Definition 8.1.2, may be represented as $\mathbf{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$.

In this notation, the parameter sets in Examples 8.1.3, 8.1.4, 8.1.8, and 8.1.9 are $\Theta = [0, 1]$, $\Theta = \{0, \dots, N\}$, $\Theta = [0, 1]$, and $\Theta = (0, \infty)$, respectively.

Remark 8.1.11. It is worthwhile mentioning that the parameter can be quite general; for example, it can be a vector $\theta = (\theta_1, \dots, \theta_k)$, so that in fact there are k unknown parameters θ_j , combined into a single vector θ . For instance, in Example 8.1.5, the unknown parameters are $\theta_1, \dots, \theta_6$, thus, the parameter set is given by

$$\Theta = \{\theta = (\theta_1, \dots, \theta_6) : \theta_k \geq 0, \theta_1 + \dots + \theta_6 = 1\}.$$

Let us present two further examples with slightly more complicated parameter sets.

Example 8.1.12. We are given an item of unknown length. It is measured by an instrument of an unidentified precision. We assume that the instrument is unbiased, that is, on average, it shows the correct value. In view of the central limit theorem, we may suppose that the measurements are distributed according to a normal distribution $\mathcal{N}(\mu, \sigma^2)$. Here μ is the “correct” length of the item, and $\sigma > 0$ reflects the precision of the measuring instrument. A small $\sigma > 0$ says that the instrument is quite precise, while a large $\sigma > 0$ corresponds to an inaccurate instrument. Consequently, by the distribution assumption the initial model is given as

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{N}(\mu, \sigma^2))_{\mu \in \mathbb{R}, \sigma^2 > 0}.$$

In order to determine μ (and maybe also σ), we measure the item n times by the same method. As a result, we obtain a random sample $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. Thus, our model describing this experiment is

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}.$$

Because of eq. (6.9), the model may also be written as

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\vec{\mu}, \sigma^2 I_n))_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}$$

with $\vec{\mu} = (\mu, \dots, \mu) \in \mathbb{R}^n$, and with diagonal matrix $\sigma^2 I_n$. The unknown parameter is (μ, σ^2) , taken from the parameter set $\mathbb{R} \times (0, \infty)$.

Example 8.1.13. Suppose now we have two different items of lengths μ_1 and μ_2 . We take m measurements of the first item and n of the second. Thereby, we use different instruments with maybe different degrees of precision. All measurements are taken independently of each other. As a result, we get a vector $(x, y) \in \mathbb{R}^{m+n}$, where $x = (x_1, \dots, x_m)$ are the values of the first m measurements and $y = (y_1, \dots, y_n)$ those of the other n . As before we assume that the x_i s are distributed according to $\mathcal{N}(\mu_1, \sigma_1^2)$, and the y_j s according to $\mathcal{N}(\mu_2, \sigma_2^2)$. We neither know μ_1 and μ_2 nor σ_1^2 and σ_2^2 . Thus, the sample space is \mathbb{R}^{m+n} and the vectors (x, y) are distributed according to $\mathcal{N}((\vec{\mu}_1, \vec{\mu}_2), R_{\sigma_1^2, \sigma_2^2})$ with diagonal matrix $R_{\sigma_1^2, \sigma_2^2}$ having σ_1^2 as its first m entries and σ_2^2 as the remaining n .

Note that by Definition 1.9.5,

$$\mathcal{N}((\vec{\mu}_1, \vec{\mu}_2), R_{\sigma_1^2, \sigma_2^2}) = \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n}.$$

This is valid because, if $A \in \mathcal{B}(\mathbb{R}^m)$ and $B \in \mathcal{B}(\mathbb{R}^n)$, then it follows that

$$\mathcal{N}((\vec{\mu}_1, \vec{\mu}_2), R_{\sigma_1^2, \sigma_2^2})(A \times B) = \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m}(A) \cdot \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n}(B).$$

The parameter set in this example is given as $\mathbb{R}^2 \times (0, \infty)^2$, hence the statistical model may be written as

$$(\mathbb{R}^{m+n}, \mathcal{B}(\mathbb{R}^{m+n}), \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n})_{(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2) \in \mathbb{R}^2 \times (0, \infty)^2}.$$

Summary: A parametric statistical model is a statistical model represented as $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$. Equivalently, the distribution assumption may be written as $\mathbf{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$ with a suitable index set Θ . A typical example is $\mathcal{X} = \{0, \dots, n\}$, $\mathcal{F} = \mathcal{P}(\mathcal{X})$, $\Theta = [0, 1]$, and $\mathbb{P}_\theta = B_{n, \theta}$.

8.2 Statistical hypothesis testing

8.2.1 Hypotheses and tests

We start with a parametric statistical model $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$. Suppose the parameter set Θ is split up into disjoint subsets Θ_0 and Θ_1 . The aim of a test is to decide, on the basis of the observed sample, whether or not the “true” parameter θ belongs to Θ_0 or to Θ_1 .

Let us explain the problem with two examples.

Example 8.2.1. Consider once more the situation described in Example 8.1.4. Assume there exists a critical value $M_0 \leq N$ such that the buyer accepts the delivery if the number M of defective machines satisfies $M \leq M_0$. Otherwise, if $M > M_0$, the buyer rejects it and sends the machines back to the trader. In this example the parameter set is $\Theta = \{0, \dots, N\}$. Letting $\Theta_0 = \{0, \dots, M_0\}$ and $\Theta_1 = \{M_0 + 1, \dots, N\}$, the question about acceptance or rejection of the delivery is equivalent to whether $M \in \Theta_0$ or $M \in \Theta_1$. Assume now the buyer checked n of the N machines and found m defective. On the basis of this observation, the buyer has to decide about acceptance or rejection, or, equivalently, about $M \in \Theta_0$ or $M \in \Theta_1$.

Example 8.2.2. Let us consider once more Example 8.1.13. There we had two measuring instruments, both being unbiased. Consequently, the expected values μ_1 and μ_2 are the correct lengths of the two items. The parameter set was $\Theta = \mathbb{R}^2 \times (0, \infty)^2$. Suppose we conjecture that both items are of equal length, that is, we conjecture $\mu_1 = \mu_2$. Letting

$$\Theta_0 := \{(\mu, \mu, \sigma_1^2, \sigma_2^2) : \mu \in \mathbb{R}, \sigma_1^2, \sigma_2^2 > 0\}$$

and $\Theta_1 = \Theta \setminus \Theta_0$, to prove or disprove the conjecture, we have to check whether $(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$ belongs to Θ_0 or Θ_1 .

On the other hand, if we want to know whether or not the first item is smaller than the second, then we have to choose

$$\Theta_0 := \{(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2) : -\infty < \mu_1 \leq \mu_2 < \infty, \sigma_1^2, \sigma_2^2 > 0\}$$

and to check whether or not $(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$ belongs to Θ_0 .

An exact mathematical formulation of the previous problems is as follows.

Definition 8.2.3. Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model and suppose $\Theta = \Theta_0 \cup \Theta_1$ with $\Theta_0 \cap \Theta_1 = \emptyset$.

Then the **hypothesis** or, more precisely, **null hypothesis** \mathbb{H}_0 says that for the “correct” $\theta \in \Theta$ one has $\theta \in \Theta_0$. This is expressed by writing $\mathbb{H}_0 : \theta \in \Theta_0$.

The **alternative hypothesis** \mathbb{H}_1 says $\theta \in \Theta_1$. Thus, $\mathbb{H}_1 : \theta \in \Theta_1$, and we have to check

$$\mathbb{H}_0 : \theta \in \Theta_0 \quad \text{against} \quad \mathbb{H}_1 : \theta \in \Theta_1 .$$

After the hypothesis is set, one executes a statistical experiment. Here the order is important: first, one has to set the hypothesis, then test it, not vice versa. If the hypothesis is chosen on the basis of the observed results, then, of course, the sample will confirm it.

Say the result of the experiment is some sample $x \in \mathcal{X}$. One of the fundamental problems in Mathematical Statistics is to decide, on the basis of the observed sample, about acceptance or rejection of \mathbb{H}_0 . The mathematical formulation of the problem is as follows.

Definition 8.2.4. A (hypothesis) **test** \mathbf{T} for checking \mathbb{H}_0 (against \mathbb{H}_1) is a disjoint partition $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ of the sample space \mathcal{X} . The set \mathcal{X}_0 is called the **region of acceptance** while \mathcal{X}_1 is said to be the **critical region**, sometimes also called **critical section** or **region of rejection**. By mathematical reasoning, we have to assume $\mathcal{X}_0 \in \mathcal{F}$, which of course implies $\mathcal{X}_1 \in \mathcal{F}$ as well.

Remark 8.2.5. A hypothesis test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ operates as follows: if the statistical experiment leads to a sample $x \in \mathcal{X}_1$, then we reject \mathbb{H}_0 . But, if we get an $x \in \mathcal{X}_0$, then this does not contradict the hypothesis, and for now we may furthermore work with it.

Important comment: If we observe an $x \in \mathcal{X}_0$, then this does *not* say that \mathbb{H}_0 is correct. It only asserts that we failed to reject it or that there is a lack of evidence against it. Let us illustrate the procedure with Example 8.2.1.

Example 8.2.6. By the choice of Θ_0 and Θ_1 , the hypothesis \mathbb{H}_0 is given by

$$\mathbb{H}_0 : 0 \leq M \leq M_0, \quad \text{hence } \mathbb{H}_1 : M_0 < M \leq N.$$

To test \mathbb{H}_0 against \mathbb{H}_1 , the sample space $\mathcal{X} = \{0, \dots, n\}$ is split up into the two regions $\mathcal{X}_0 := \{0, \dots, m_0\}$ and $\mathcal{X}_1 := \{m_0 + 1, \dots, n\}$ with some (for now) arbitrary number $m_0 \in \{0, \dots, n\}$. If among the checked n machines m are defective with some $m > m_0$, then $m \in \mathcal{X}_1$, hence one rejects \mathbb{H}_0 . In this case the buyer refuses to take the delivery and sends it back to the trader. On the other hand, if $m \leq m_0$, then $m \in \mathcal{X}_0$, which does not contradict \mathbb{H}_0 , and the buyer will accept the delivery and pay for it. Of course, the key question is how to choose the value m_0 in a proper way.

Remark 8.2.7. Sometimes tests are also defined as mappings $\varphi : \mathcal{X} \rightarrow \{0, 1\}$. The link between these two approaches is immediately clear. Starting with φ , the hypothesis test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is constructed by $\mathcal{X}_0 = \{x \in \mathcal{X} : \varphi(x) = 0\}$ and $\mathcal{X}_1 = \{x \in \mathcal{X} : \varphi(x) = 1\}$. Conversely, if $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is a given test, then set $\varphi(x) = 0$ if $x \in \mathcal{X}_0$ and $\varphi(x) = 1$ for $x \in \mathcal{X}_1$. The advantage of this approach is that it allows us to define the so-called **randomized tests**. Here $\varphi : \mathcal{X} \rightarrow [0, 1]$. Then, as before, $\mathcal{X}_0 = \{x \in \mathcal{X} : \varphi(x) = 0\}$ and $\mathcal{X}_1 = \{x \in \mathcal{X} : \varphi(x) = 1\}$. If $0 < \varphi(x) < 1$, then

$$\varphi(x) = \mathbb{P}\{\text{reject } \mathbb{H}_0 \text{ if } x \text{ is observed}\}.$$

That is, for certain observations $x \in \mathcal{X}$, an additional random experiment (e. g., tossing a coin) decides whether we accept or reject \mathbb{H}_0 . Randomized tests are useful in the case of finite or countably infinite sample spaces.

When applying a test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ to check the null hypothesis $\mathbb{H}_0 : \theta \in \Theta_0$, two different types of errors may occur.

Definition 8.2.8. An **error of the first kind** or **type I error** occurs if \mathbb{H}_0 is true but one observes a sample $x \in \mathcal{X}_1$, hence rejects \mathbb{H}_0 .

! Type I error = incorrect rejection of a true null hypothesis

In other words, a type I error happens if the “true” θ is in Θ_0 , but we observe an $x \in \mathcal{X}_1$.

Definition 8.2.9. An **error of the second kind** or **type II error** occurs if \mathbb{H}_0 is false, but the observed sample lies in \mathcal{X}_0 , hence we do not reject the false hypothesis \mathbb{H}_0 .

! Type II error = failure to reject a false null hypothesis

Consequently, a type II error occurs if the “true” θ is in Θ_1 , but the observed sample is an element of the region of acceptance \mathcal{X}_0 .

Example 8.2.10. In the context of Example 8.2.6, a type I error occurs if the delivery was good, but among the checked machines more than m_0 were defective, so that the buyer rejected the delivery. Since the trader was not able to sell a proper delivery, this error is also called the **risk of the trader**.

On the other hand, a type II error occurs if the delivery is not in good order, but among the checked machines only a few were defective (at most m_0). Thus, the buyer accepted the bad delivery and paid for it. Therefore, this type of error is also called the **risk of the buyer**.

Summary: Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model and suppose $\Theta = \Theta_0 \cup \Theta_1$ with $\Theta_0 \cap \Theta_1 = \emptyset$. Then the hypotheses \mathbb{H}_0 and \mathbb{H}_1 are

$$\mathbb{H}_0 : \theta \in \Theta_0 \quad \text{against} \quad \mathbb{H}_1 : \theta \in \Theta_1.$$

A (hypothesis) test \mathbf{T} for checking \mathbb{H}_0 (against \mathbb{H}_1) is a disjoint partition $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ of the sample space \mathcal{X} . The set \mathcal{X}_0 is called the region of acceptance while \mathcal{X}_1 is said to be the critical region. If the observed sample $x \in \mathcal{X}_1$, one rejects \mathbb{H}_0 . If $x \in \mathcal{X}_0$, we cannot reject \mathbb{H}_0 and have to work with it furthermore.

Type I error = incorrect rejection of a true null hypothesis $\Leftrightarrow x \in \mathcal{X}_1$ and $\theta \in \Theta_0$,

Type II error = failure to reject a false null hypothesis $\Leftrightarrow x \in \mathcal{X}_0$ and $\theta \in \Theta_1$.

Verbally said:¹ A type I error is the mistaken rejection of a null hypothesis that is actually true; for example, “due to a false witness report, an innocent person is convicted.” A type II error is the failure to reject a null hypothesis that is actually false; for example, “due to the lack of evidence, a guilty person is not convicted.”

8.2.2 Power function and significance tests

The power of a test is described by its power function defined as follows.

Definition 8.2.11. Let $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ be a test for $\mathbb{H}_0 : \theta \in \Theta_0$ against $\mathbb{H}_1 : \theta \in \Theta_1$. The function $\beta_{\mathbf{T}}$ from Θ to $[0, 1]$ defined as

$$\beta_{\mathbf{T}}(\theta) := \mathbb{P}_{\theta}(\mathcal{X}_1)$$

is called the **power function** of the test \mathbf{T} .

Remark 8.2.12. If $\theta \in \Theta_0$, that is, if \mathbb{H}_0 is true, then $\beta_{\mathbf{T}}(\theta) = \mathbb{P}_{\theta}(\mathcal{X}_1)$ is the probability that \mathcal{X}_1 occurs or, equivalently, that a type I error happens.

On the contrary, if $\theta \in \Theta_1$, that is, \mathbb{H}_0 is false, then $1 - \beta_{\mathbf{T}}(\theta) = \mathbb{P}_{\theta}(\mathcal{X}_0)$ is the probability that \mathcal{X}_0 occurs or, equivalently, that a type II error appears.

Thus, a “good” test should satisfy the following conditions: the power function $\beta_{\mathbf{T}}$ attains small values on Θ_0 and/or $1 - \beta_{\mathbf{T}}$ has small values on Θ_1 . Then the probabilities for the occurrence of type I and/or type II errors are not too big.²

Example 8.2.13. What is the power function of the test presented in Example 8.2.6? Recall that $\Theta = \{0, \dots, N\}$ and $\mathcal{X}_1 = \{m_0 + 1, \dots, n\}$. Hence, $\beta_{\mathbf{T}}$ maps $\{0, \dots, N\}$ to $[0, 1]$ in the following way:

$$\beta_{\mathbf{T}}(M) = H_{N,M,n}(\mathcal{X}_1) = \sum_{m=m_0+1}^n \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}. \quad (8.1)$$

If the hypotheses are

$$\mathbb{H}_0 : 0 \leq M \leq M_0 \quad \text{against} \quad \mathbb{H}_1 : M_0 < M \leq N,$$

then the maximal probability for a type I error is given by

$$\max_{0 \leq M \leq M_0} \beta_{\mathbf{T}}(M) = \max_{0 \leq M \leq M_0} \sum_{m=m_0+1}^n \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad (8.2)$$

¹ See [Fsh71].

² In the literature, the power function is sometimes defined in a slightly different way. If $\theta \in \Theta_0$, then it is as in our Definition 8.2.11 while for $\theta \in \Theta_1$ one defines it as $1 - \beta_{\mathbf{T}}(\theta)$. Moreover, for $1 - \beta_{\mathbf{T}}$ one finds the notion of the **operation characteristics** or **oc-function**.

while the maximal probability for a type II error equals

$$\max_{M_0 < M \leq N} (1 - \beta_T(M)) = \max_{M_0 < M \leq N} \sum_{m=0}^{m_0} \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}. \tag{8.3}$$

Let us give a concrete example for the power function in eq. (8.1). Suppose the trader submits a delivery of 30 machines. Hence, $N = 30$ and $\Theta = \{0, \dots, 30\}$. The buyer chooses randomly 10 machines and tests them. That is, $n = 10$. Assume, the test is $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ where $\mathcal{X}_0 = \{0, 1, 2, 3, 4\}$, hence $\mathcal{X}_1 = \{5, 6, 7, 8, 9, 10\}$. Thus, $m_0 = 4$. In other words, if there are at most 4 defective machines among the tested 10, then the buyer accepts the delivery. Otherwise, he rejects it.

Then the power function of the test \mathbf{T} equals

$$\beta_T(M) = \sum_{m=5}^{10} \frac{\binom{M}{m} \binom{30-M}{10-m}}{\binom{30}{10}}, \quad M = 0, \dots, 30. \tag{8.4}$$

Note that

$$\beta_T(0) = \dots = \beta_T(4) = 0 \quad \text{while} \quad \beta_T(25) = \dots = \beta_T(30) = 1.$$

Why is this so? First, if there are only 4 or less defective machines among the delivered, it is impossible to observe 5 or more defective machines among the 10 chosen. On the other hand, if there are 25 or more defective machines among 30, then at most 5 machines are nondefective. Hence, among the 10 chosen there are at least 5 defective.

Some other interesting values of β_T are (see also Figure 8.1)

M	5	6	7	8	9	10	11
$\beta_T(M)$	0.001768	0.008842	0.02564	0.05632	0.1037	0.1687	0.2500

These values of β_T tell us the following: if, for example, there are 11 defective machines in the delivery, then there is a 25% chance to observe in the sample 5 or more defective. Equivalently, the probability to observe at most 4 defective machines is still 75%. Thus, in this case the likelihood of a type II error is rather big. The situation changes drastically for larger M . If, for example, the number of defective machines is 18, among 100 trials the sample will on average 88.2 times contain 5 or more defective items.

Remark 8.2.14. Formulas (8.2) and (8.3) already illustrate the **dilemma of hypothesis testing**. To minimize the type I error, one has to choose m_0 as large as possible. But increasing m_0 enlarges the type II error.

This dilemma occurs always in the theory of hypothesis testing. In order to minimize the probability of a type I error, the critical region \mathcal{X}_1 has to be chosen as small as possible. But making \mathcal{X}_1 smaller enlarges \mathcal{X}_0 , hence the probability for the occurrence

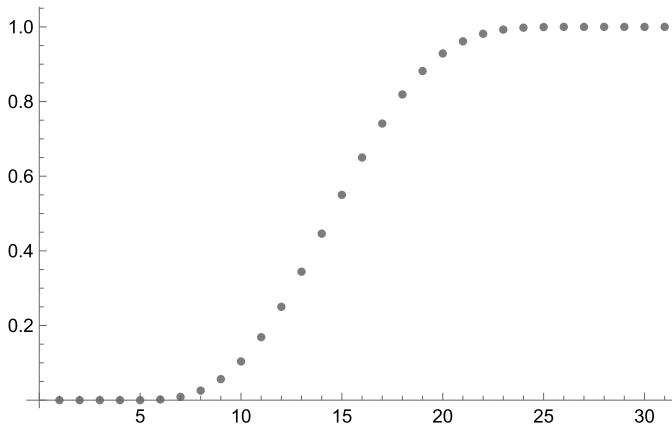


Figure 8.1: The power function β_T in eq. (8.4).

of a type II error increases. In the extreme case, if $\mathcal{X}_1 = \emptyset$, hence $\mathcal{X}_0 = \mathcal{X}$, then a type I error cannot occur at all. In the context of Example 8.2.6 that means the buyer accepts all deliveries and the trader takes no risk.

On the other hand, to minimize the occurrence of a type II error, the region of acceptance \mathcal{X}_0 has to be as small as possible. In the extreme case, if we choose $\mathcal{X}_0 = \emptyset$, then a type II error cannot occur because we always reject the hypothesis. In the context of Example 8.2.6, this says the buyer rejects all deliveries. In this way he avoids buying any delivery of bad quality, but he also never gets a proper one. Thus the buyer takes no risk.

It is pretty clear that both extreme cases presented above are absurd. Therefore, one has to find a suitable compromise. The approach for such a compromise is as follows: in the first step, one chooses tests where the probability of a type I error is bounded from above. And in the second step, among all these tests satisfying this bound, one takes that which minimizes the probability of a type II error. More precisely, we will investigate tests satisfying the following condition.

Definition 8.2.15. Suppose we are given a number $\alpha \in (0, 1)$, the so-called **significance level**. A test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ for testing the hypothesis $\mathbb{H}_0 : \theta \in \Theta_0$ against $\mathbb{H}_1 : \theta \in \Theta_1$ is said to be an **α -significance test** (or simply **α -test**), provided the probability for the occurrence of a type I error is bounded by α . That is, the test has to satisfy

$$\sup_{\theta \in \Theta_0} \beta_T(\theta) = \sup_{\theta \in \Theta_0} \mathbb{P}_\theta(\mathcal{X}_1) \leq \alpha.$$

Interpretation: The significance level α is assumed to be small. Typical choices are $\alpha = 0.1$ or $\alpha = 0.01$. Let \mathbf{T} be an α -significance test and assume that \mathbb{H}_0 is true. If we observe now a sample in the critical region \mathcal{X}_1 , then an event occurred with probability less than or equal to α , that is, a very unlikely event has been observed. Therefore, we can be very

sure that this could not have happened provided \mathbb{H}_0 had been true, and we reject this hypothesis. The probability that we made a mistake is less than or equal to the chosen $\alpha > 0$, hence very small.

Recall that α -significance tests admit no bound for the probability of a type II error. Therefore, we look for those α -significance tests that minimize the probability for a type II error.

Definition 8.2.16. Let \mathbf{T}_1 and \mathbf{T}_2 be two α -significance tests for checking \mathbb{H}_0 against \mathbb{H}_1 . If their power functions satisfy

$$\beta_{\mathbf{T}_1}(\theta) \geq \beta_{\mathbf{T}_2}(\theta), \quad \theta \in \Theta_1,$$

then we say that \mathbf{T}_1 is (uniformly) **more powerful** than \mathbf{T}_2 .

A (uniformly) **most powerful** α -test \mathbf{T} is that which is more powerful than all other α -tests.

Remark 8.2.17. Note that $\beta_{\mathbf{T}_1}(\theta) \geq \beta_{\mathbf{T}_2}(\theta)$ implies $1 - \beta_{\mathbf{T}_1}(\theta) \leq 1 - \beta_{\mathbf{T}_2}(\theta)$, hence if \mathbf{T}_1 is more powerful than \mathbf{T}_2 , then, according to Remark 8.2.12, the probability for the occurrence of a type II error is smaller for \mathbf{T}_1 than it is for \mathbf{T}_2 . Therefore, a most powerful α -test is that which minimizes the probability of occurrence of a type II error.

Remark 8.2.18. The question about existence and uniqueness of most powerful α -tests is treated in the Neyman–Pearson lemma and its consequences. We will not discuss that problem here; instead, we will construct most powerful tests in concrete situations. See [CB02, Chapter 8.3.2] for a detailed discussion of the Neyman–Pearson lemma and its consequences.

We start with the construction of such tests in the hypergeometric case. Here we have the following.

Proposition 8.2.19. *If the statistical model equals $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{N,M,n})_{M=0,\dots,N}$ with $\mathcal{X} = \{0, \dots, n\}$, then the most powerful α -test for testing $M \leq M_0$ against $M > M_0$ is given by $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$, where $\mathcal{X}_0 = \{0, \dots, m_0\}$, and m_0 is defined by*

$$\begin{aligned} m_0 &:= \max \left\{ k \leq n : \sum_{m=k}^n \frac{\binom{M_0}{m} \binom{N-M_0}{n-m}}{\binom{N}{n}} > \alpha \right\} \\ &= \min \left\{ k \leq n : \sum_{m=k+1}^n \frac{\binom{M_0}{m} \binom{N-M_0}{n-m}}{\binom{N}{n}} \leq \alpha \right\}. \end{aligned}$$

Proof. The proof of Proposition 8.2.19 needs the following lemma.

Lemma 8.2.20. *The power function, defined by eq. (8.1), is a nondecreasing function on the set $\{0, \dots, N\}$.*

Proof. Suppose we get a delivery of N machines containing M defective. Now there are not only defective machines within the delivery, but also $\tilde{M} - M$ false ones for some $\tilde{M} \geq M$. We take a sample of size n and test these machines. Let X be the number of

defective machines and let \tilde{X} be the number of machines that are either defective or false. Of course, we have $X \leq \tilde{X}$ implying $\mathbb{P}(X > m_0) \leq \mathbb{P}(\tilde{X} > m_0)$. Note that X is $H_{N,M,n}$ -distributed while \tilde{X} is distributed according to $H_{N,\tilde{M},n}$. These observations lead to

$$\begin{aligned}\beta_{\mathbf{T}}(M) &= H_{N,M,n}(\{m_0 + 1, \dots, n\}) = \mathbb{P}\{X > m_0\} \leq \mathbb{P}\{\tilde{X} > m_0\} \\ &= H_{N,\tilde{M},n}(\{m_0 + 1, \dots, n\}) = \beta_{\mathbf{T}}(\tilde{M}).\end{aligned}$$

This being true for all $M \leq \tilde{M}$ proves that $\beta_{\mathbf{T}}$ is nondecreasing. \square

Let us come back to the proof of Proposition 8.2.19. Set $\mathcal{X}_0 := \{0, \dots, m_0\}$, thus $\mathcal{X}_1 = \{m_0 + 1, \dots, n\}$ for some (at the moment arbitrary) $m_0 \leq n$. Because of Lemma 8.2.20, the test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test if and only if it satisfies

$$\sum_{m=m_0+1}^n \frac{\binom{M_0}{m} \binom{N-M_0}{n-m}}{\binom{N}{n}} = H_{N,M_0,n}(\mathcal{X}_1) = \sup_{M \leq M_0} H_{N,M,n}(\mathcal{X}_1) \leq \alpha.$$

To minimize the probability for the occurrence of a type II error, we have to choose \mathcal{X}_1 as large as possible or, equivalently, m_0 as small as possible, that is, if we replace m_0 by $m_0 - 1$, then the new test is no longer an α -test. Thus, in order that \mathbf{T} is an α -test that minimizes the probability for a type II error, the number m_0 has to be chosen such that

$$\sum_{m=m_0+1}^n \frac{\binom{M_0}{m} \binom{N-M_0}{n-m}}{\binom{N}{n}} \leq \alpha \quad \text{and} \quad \sum_{m=m_0}^n \frac{\binom{M_0}{m} \binom{N-M_0}{n-m}}{\binom{N}{n}} > \alpha.$$

This completes the proof. \square

Example 8.2.21. A buyer gets a delivery of 100 machines. In the case that there are strictly more than 10 defective machines in the delivery, he will reject it. Thus, his hypothesis is $\mathbb{H}_0 : M \leq 10$. In order to test \mathbb{H}_0 , he chooses 15 machines and checks them. Let m be the number of defective machines among the checked. For which m does he reject the delivery with a significance level $\alpha = 0.01$?

Answer: We have $N = 100$, $M_0 = 10$, and $n = 15$. Since $\alpha = 0.01$, from

$$\sum_{m=5}^{15} \frac{\binom{10}{m} \binom{90}{15-m}}{\binom{100}{15}} = 0.0063 \dots < \alpha \quad \text{and} \quad \sum_{m=4}^{15} \frac{\binom{10}{m} \binom{90}{15-m}}{\binom{100}{15}} = 0.04 \dots > \alpha,$$

it follows that the optimal choice is $m_0 = 4$. Consequently, we have $\mathcal{X}_0 = \{0, \dots, 4\}$, thus, $\mathcal{X}_1 = \{5, \dots, 15\}$. If there are 5 or even more defective machines among the tested 15, then the buyer should reject the delivery. The probability that his decision is wrong is less than or equal to 0.01.

What can be said about the probability for a type II error? For this test, we have

$$\beta_{\mathbf{T}}(M) = \sum_{m=5}^{15} \frac{\binom{M}{m} \binom{100-M}{15-m}}{\binom{100}{15}}, \quad (8.5)$$

hence

$$1 - \beta_{\mathbf{T}}(M) = \sum_{m=0}^4 \frac{\binom{M}{m} \binom{100-M}{15-m}}{\binom{100}{15}}.$$

Since $\beta_{\mathbf{T}}$ is nondecreasing, $1 - \beta_{\mathbf{T}}$ is nonincreasing, and the probability for a type II error becomes maximal for $M = 11$. Recall that $\Theta_0 = \{0, \dots, 10\}$ and, therefore, $\Theta_1 = \{11, \dots, 100\}$. Thus, an upper bound for the probability of a type II error is given by

$$1 - \beta_{\mathbf{T}}(M) \leq 1 - \beta_{\mathbf{T}}(11) = \sum_{m=0}^4 \frac{\binom{11}{m} \binom{89}{15-m}}{\binom{100}{15}} = 0.989471, \quad M = 11, \dots, 100.$$

This tells us that even in the case of most powerful tests the likelihood for a type II error may be quite large. Even if the number of defective machines is big, this error may occur with higher probability. For example, we have $1 - \beta_{\mathbf{T}}(20) = 0.853089$ and $1 - \beta_{\mathbf{T}}(40) = 0.197057$. See also Figure 8.2.

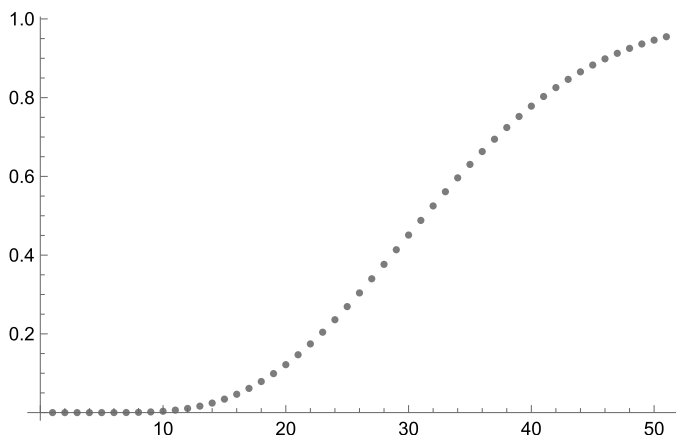


Figure 8.2: The power function $\beta_{\mathbf{T}}$ defined by eq. (8.5).

If one is willing to take a greater risk and chooses $\alpha = 0.1$, then, since

$$\sum_{m=4}^{15} \frac{\binom{10}{m} \binom{90}{15-m}}{\binom{100}{15}} = 0.0063 \dots < \alpha \quad \text{and} \quad \sum_{m=3}^{15} \frac{\binom{10}{m} \binom{90}{15-m}}{\binom{100}{15}} = 0.1705 \dots > \alpha,$$

one may take $\mathcal{X}_0 = \{0, 1, 2, 3\}$ as the region of acceptance. Thus, in this case the buyer will reject the delivery if there are 4 defective machines among the 15 tested. The probability that this is a wrong decision is less than 0.1, but greater than 0.01.

Remark 8.2.22 (Important!). An α -significance test provides us with quite precise information when rejecting the hypothesis \mathbb{H}_0 . In contrast, when we observe a sample $x \in \mathcal{X}_0$,

the only information we get is that we failed to reject \mathbb{H}_0 , thus, we must continue to regard it as true. Consequently, whenever fixing the null hypothesis, we have to fix it in a way that either a type I error has the most serious consequences or that we can attain the most information by rejecting \mathbb{H}_0 . Let us explain this with two examples.

Example 8.2.23. A certain type of food sometimes contains a special kind of poison. Suppose there are μ milligrams of poison in one kilogram of the food. If $\mu > \mu_0$, then eating this becomes dangerous while for $\mu \leq \mu_0$ it is unproblematic. How do we successfully choose the hypothesis when testing some sample of the food? We could take either $\mathbb{H}_0 : \mu > \mu_0$ or $\mathbb{H}_0 : \mu \leq \mu_0$. Which is the right choice?

Answer: The correct choice is $\mathbb{H}_0 : \mu > \mu_0$. Why? If we reject \mathbb{H}_0 , then we can be very sure that the food is not poisoned and may be eaten. The probability that someone will be poisoned is less than α . A type II error occurs if the food is harmless, but we discard it because our test tells us that it is poisoned. That results in a loss for the company that produced it, but no one will suffer from poisoning. If we were to choose $\mathbb{H}_0 : \mu \leq \mu_0$, then a type II error would occur if \mathbb{H}_0 is false, that is, the food is poisoned, but our test says that it is eatable. Of course, this error is much more serious, and we have no control in regards to its probability.

Example 8.2.24. Suppose the height of 18-year-old males in the US is normally distributed with expected value μ and variance $\sigma^2 > 0$. We want to know whether the average height is above or below 6 feet. There is strong evidence that we will have $\mu \leq 6$, but we cannot prove this. To do so, we execute a statistical experiment and choose randomly n males of age 18 and measure their height. Which hypothesis should be checked? If we take $\mathbb{H}_0 : \mu \leq 6$, then it is very likely that our experiment will lead to a result that does not contradict this hypothesis, resulting in a small amount of information gained. But, if we work with the hypothesis $\mathbb{H}_0 : \mu > 6$, then a rejection of this hypothesis tells us that \mathbb{H}_0 is very likely wrong, and we may say the conjecture is true with high probability, namely that we have $\mu \leq 6$. Here the probability that our conclusion is wrong is very small.

Summary: The power function of a test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ for $\mathbb{H}_0 : \theta \in \Theta_0$ against $\mathbb{H}_1 : \theta \in \Theta_1$ is defined by

$$\beta_{\mathbf{T}}(\theta) = \mathbb{P}_{\theta}(\mathcal{X}_1), \quad \theta \in \Theta.$$

If \mathbf{T} is a “good” test, then the power function $\beta_{\mathbf{T}}$ should be small on Θ_0 and near to one on Θ_1 . Given $\alpha > 0$, an α -significance test \mathbf{T} satisfies

$$\sup_{\theta \in \Theta_0} \beta_{\mathbf{T}}(\theta) = \sup_{\theta \in \Theta_0} \mathbb{P}_{\theta}(\mathcal{X}_1) \leq \alpha.$$

That is, if \mathbf{T} is an α -test, the probability of a type I error is bounded by α .

8.3 Tests for binomial distributed populations

Because of their importance, we present tests for binomial distributed populations in a separate section. The starting point is the problem described in Examples 8.1.3 and 8.1.8. In a single experiment, we may observe either “0” or “1,” but we do not know the probabilities for the occurrence of these events. To obtain some information about the unknown probabilities, we execute n independent trials and record how often “1” occurs. This number is $B_{n,\theta}$ -distributed for some $0 \leq \theta \leq 1$. Hence, the describing statistical model is given by

$$(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{\theta \in [0,1]} \quad \text{where } \mathcal{X} = \{0, \dots, n\}. \quad (8.6)$$

Two-sided tests: We want to check whether the unknown parameter θ satisfies $\theta = \theta_0$ or $\theta \neq \theta_0$ for some given $\theta_0 \in [0, 1]$. Thus, $\Theta_0 = \{\theta_0\}$ and $\Theta_1 = [0, 1] \setminus \{\theta_0\}$. In other words, the null and the alternative hypothesis are

$$\mathbb{H}_0 : \theta = \theta_0 \quad \text{and} \quad \mathbb{H}_1 : \theta \neq \theta_0,$$

respectively.

To construct a suitable α -significance test for checking \mathbb{H}_0 , we introduce two numbers n_0 and n_1 as follows. Note that these numbers depend on θ_0 and, of course, also on α . The numbers are defined by

$$\begin{aligned} n_0 &:= \min \left\{ k \leq n : \sum_{j=0}^k \binom{n}{j} \theta_0^j (1 - \theta_0)^{n-j} > \alpha/2 \right\} \\ &= \max \left\{ k \leq n : \sum_{j=0}^{k-1} \binom{n}{j} \theta_0^j (1 - \theta_0)^{n-j} \leq \alpha/2 \right\} \end{aligned} \quad (8.7)$$

and

$$\begin{aligned} n_1 &:= \max \left\{ k \leq n : \sum_{j=k}^n \binom{n}{j} \theta_0^j (1 - \theta_0)^{n-j} > \alpha/2 \right\} \\ &= \min \left\{ k \leq n : \sum_{j=k+1}^n \binom{n}{j} \theta_0^j (1 - \theta_0)^{n-j} \leq \alpha/2 \right\}. \end{aligned} \quad (8.8)$$

Proposition 8.3.1. *Consider the statistical model (8.6) and let $0 < \alpha < 1$ be a significance level. The hypothesis test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ with*

$$\mathcal{X}_0 := \{n_0, n_0 + 1, \dots, n_1 - 1, n_1\} \quad \text{and} \quad \mathcal{X}_1 = \{0, \dots, n_0 - 1\} \cup \{n_1 + 1, \dots, n\} \quad (8.9)$$

is an α -significance test to check $\mathbb{H}_0 : \theta = \theta_0$ against $\mathbb{H}_1 : \theta \neq \theta_0$. Here n_0 and n_1 are defined as in eqs. (8.7) and (8.8).

Proof. Since Θ_0 consists only of the point $\{\theta_0\}$, an arbitrary test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test if and only if $B_{n, \theta_0}(\mathcal{X}_1) \leq \alpha$. Now let \mathbf{T} be as in the formulation of the proposition. By the definition of the numbers n_0 and n_1 , we obtain

$$B_{n, \theta_0}(\mathcal{X}_1) = \sum_{j=0}^{n_0-1} \binom{n}{j} \theta_0^j (1-\theta_0)^{n-j} + \sum_{j=n_1+1}^n \binom{n}{j} \theta_0^j (1-\theta_0)^{n-j} \leq \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha,$$

that is, as claimed, the test $\mathbf{T} := (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test. \square

Remark 8.3.2. In this test the critical region \mathcal{X}_1 consists of two parts or tails. Therefore, this type of hypothesis test is called a **two-sided test**.

Remark 8.3.3. By the choice of n_0 and n_1 , the regions \mathcal{X}_0 and \mathcal{X}_1 in eq. (8.9) are optimal in the following sense: If $\tilde{\mathbf{T}} = (\tilde{\mathcal{X}}_0, \tilde{\mathcal{X}}_1)$ with

$$\tilde{\mathcal{X}}_0 := \{n_0 + 1, \dots, n_1 - 1\} \quad \text{and} \quad \tilde{\mathcal{X}}_1 := \{0, \dots, n_0\} \cup \{n_1, \dots, n\},$$

then

$$B_{n, \theta_0}(\tilde{\mathcal{X}}_1) = \sum_{j=0}^{n_0} \binom{n}{j} \theta_0^j (1-\theta_0)^{n-j} + \sum_{j=n_1}^n \binom{n}{j} \theta_0^j (1-\theta_0)^{n-j} > \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha.$$

Hence, $\tilde{\mathbf{T}}$ is no longer an α -significance test. But note that we cannot exclude that the tests with either

$$\tilde{\mathcal{X}}_0 := \{n_0, \dots, n_1 - 1\} \quad \text{or} \quad \tilde{\mathcal{X}}_1 := \{n_0 + 1, \dots, n_1\}$$

are still α -significance tests.

Example 8.3.4. In an urn there is an unknown number of white and black balls. Let $\theta \in [0, 1]$ be the proportion of white balls. We conjecture that there are as many white as black balls in the urn. That is, the null hypothesis is $H_0 : \theta = 0.5$. To test this hypothesis, we choose one after another 100 balls with replacement. In order to determine n_0 and n_1 in this situation, let φ be defined as

$$\varphi(k) := \sum_{j=0}^k \binom{100}{j} \cdot \left(\frac{1}{2}\right)^{100} = B_{100, 0.5}(\{0, \dots, k\}).$$

Numerical calculations give

$$\begin{aligned} \varphi(36) &= 0.00331856, & \varphi(37) &= 0.00601649, & \varphi(38) &= 0.0104894, \\ \varphi(39) &= 0.0176001, & \varphi(40) &= 0.028444, & \varphi(41) &= 0.044313, \\ \varphi(42) &= 0.0666053, & \varphi(43) &= 0.096674, & \varphi(44) &= 0.135627, \\ \varphi(45) &= 0.184101, & \varphi(46) &= 0.242059, & \varphi(47) &= 0.30865, \\ \varphi(48) &= 0.382177, & \varphi(49) &= 0.460205. \end{aligned}$$

If the significance level is chosen as $\alpha = 0.1$, we see that $\varphi(41) \leq 0.05$, but $\varphi(42) > 0.05$. Hence, by the definition of n_0 in eq. (8.7), it follows that $n_0 = 42$. Either by symmetry or by similar calculations, for n_1 defined in eq. (8.8), we get $n_1 = 58$. Consequently, the regions of acceptance and rejection are given by

$$\mathcal{X}_0 = \{42, 43, \dots, 57, 58\} \quad \text{and} \quad \mathcal{X}_1 = \{0, \dots, 41\} \cup \{59, \dots, 100\}.$$

For example, if we observe during 100 trials k white balls for some $k < 42$ or some $k > 58$, then we may be quite sure that our null hypothesis is wrong, that is, the numbers of white and black balls are significantly different. This assertion is 90 % sure. The power function of this test (see Fig. 8.3) is given by

$$\beta_{\mathbf{T}}(\theta) = \sum_{k=0}^{41} \binom{100}{k} \theta^k (1-\theta)^{100-k} + \sum_{k=59}^{100} \binom{100}{k} \theta^k (1-\theta)^{100-k}, \quad 0 < \theta < 1. \quad (8.10)$$

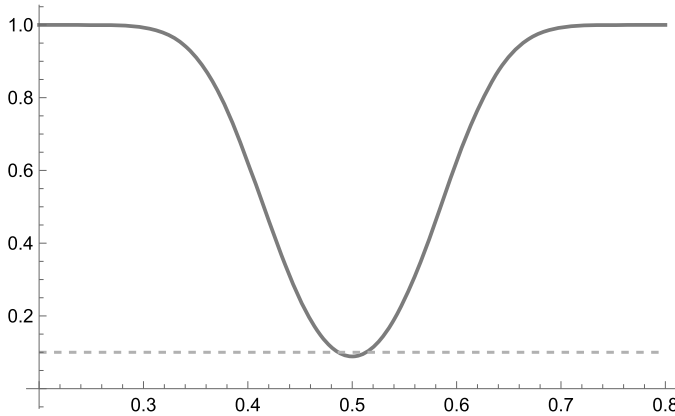


Figure 8.3: The power function $\beta_{\mathbf{T}}$ in eq. (8.10) with significance level $\alpha = 0.1$.

If we want to be more certain about the conclusion, we have to choose a smaller significance level. For example, if we take $\alpha = 0.01$, the values of φ imply $n_0 = 37$ and $n_1 = 63$, hence in this case we conclude that

$$\mathcal{X}_0 = \{37, 38, \dots, 62, 63\} \quad \text{and} \quad \mathcal{X}_1 = \{0, \dots, 36\} \cup \{64, \dots, 100\}.$$

Again we see that a smaller bound for the probability of a type I error leads to an enlargement of \mathcal{X}_0 , thus, to an increase of the chance for a type II error.

One-sided tests: Now the null hypothesis is $\mathbb{H}_0 : \theta \leq \theta_0$ for some $\theta_0 \in [0, 1]$. In the context of Example 8.1.3, we claim that the proportion of white balls in the urn does not exceed θ_0 . For instance, if $\theta_0 = 1/2$, then we want to test whether or not the number of white balls is less than or equal to that of black.

Before we present a most powerful test for this situation, let us define a number n_0 depending on θ_0 and on the significance level $0 < \alpha < 1$, namely

$$\begin{aligned} n_0 &= \max \left\{ k \leq n : \sum_{j=k}^n \binom{n}{j} \theta_0^j (1 - \theta_0)^{n-j} > \alpha \right\} \\ &= \min \left\{ k \leq n : \sum_{j=k+1}^n \binom{n}{j} \theta_0^j (1 - \theta_0)^{n-j} \leq \alpha \right\}. \end{aligned} \quad (8.11)$$

Now we are in a position to state the most powerful one-sided α -test for a binomial distributed population.

Proposition 8.3.5. *Suppose $\mathcal{X} = \{0, \dots, n\}$, and let $(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{\theta \in [0,1]}$ be the statistical model describing a binomial distributed population. Given $0 < \alpha < 1$, define n_0 by (8.11) and set $\mathcal{X}_0 = \{0, \dots, n_0\}$, hence $\mathcal{X}_1 = \{n_0 + 1, \dots, n\}$. Then $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is the most powerful α -test to check the null hypothesis $\mathbb{H}_0 : \theta \leq \theta_0$ against $\mathbb{H}_1 : \theta > \theta_0$.*

Proof. Fixing an arbitrary $0 \leq m \leq n$, we define the region of acceptance \mathcal{X}_0 of a test \mathbf{T} by $\mathcal{X}_0 = \{0, \dots, m\}$. Its power function is given by

$$\beta_{\mathbf{T}}(\theta) = B_{n,\theta}(\mathcal{X}_1) = \sum_{j=m+1}^n \binom{n}{j} \theta^j (1 - \theta)^{n-j}, \quad 0 \leq \theta \leq 1. \quad (8.12)$$

To proceed further, we need the following lemma.

Lemma 8.3.6. *The power function (8.12) is nondecreasing in $[0, 1]$.*

Proof. Suppose in an urn there are white, red, and black balls. Their proportions are θ_1 , $\theta_2 - \theta_1$ and $1 - \theta_2$ for some $0 \leq \theta_1 \leq \theta_2 \leq 1$. Choose n balls with replacement. Let X be the number of chosen white balls, and Y the number of balls that were either white or red. Then X is B_{n,θ_1} -distributed, while Y is distributed according to B_{n,θ_2} . Moreover, $X \leq Y$, hence it follows that $\mathbb{P}(X > m) \leq \mathbb{P}(Y > m)$, which leads to

$$\begin{aligned} \beta_{\mathbf{T}}(\theta_1) &= B_{n,\theta_1}(\{m+1, \dots, n\}) = \mathbb{P}(X > m) \leq \mathbb{P}(Y > m) \\ &= B_{n,\theta_2}(\{m+1, \dots, n\}) = \beta_{\mathbf{T}}(\theta_2). \end{aligned}$$

This being true for all $\theta_1 \leq \theta_2$ completes the proof of the lemma. \square

An application of Lemma 8.3.6 implies that the above test \mathbf{T} is an α -significance test if and only if

$$\sum_{j=m+1}^n \binom{n}{j} \theta_0^j (1-\theta_0)^{n-j} = \beta_{\mathbf{T}}(\theta_0) = \sup_{\theta \leq \theta_0} \beta_{\mathbf{T}}(\theta) \leq \alpha.$$

In order to minimize the probability of a type II error, we have to choose \mathcal{X}_0 as small as possible. That is, if we replace m by $m-1$, the modified test is no longer an α -test. Thus, the optimal choice is $m = n_0$ where n_0 is defined by eq. (8.11). This completes the proof of Proposition 8.3.5. \square

Example 8.3.7. Let us come back to the problem investigated in Example 8.1.3. Our null hypothesis is $H_0 : \theta \leq 1/2$, that is, we claim that at most half of the balls are white. To test H_0 , we choose 100 balls and record their color. Let k be the number of observed white balls. For which k must we reject H_0 with a confidence of 90%?

Answer: Since

$$\sum_{k=56}^{100} \binom{100}{k} 2^{-100} = 0.135627 \quad \text{and} \quad \sum_{k=57}^{100} \binom{100}{k} 2^{-100} = 0.096674,$$

for $\alpha = 0.1$ the number n_0 in eq. (8.11) equals $n_0 = 56$. Consequently, the region of acceptance for the best 0.1-test is given by $\mathcal{X}_0 = \{0, \dots, 56\}$. Thus, whenever there are 57 or more white balls among the chosen 100, the hypothesis has to be rejected. The probability for a wrong decision is less than or equal to 0.1. The power function of this test \mathbf{T} is given by (compare Figure 8.4)

$$\beta_{\mathbf{T}}(\theta) = \sum_{k=57}^{100} \binom{100}{k} \theta^k (1-\theta)^{100-k}, \quad 0 < \theta < 1. \quad (8.13)$$

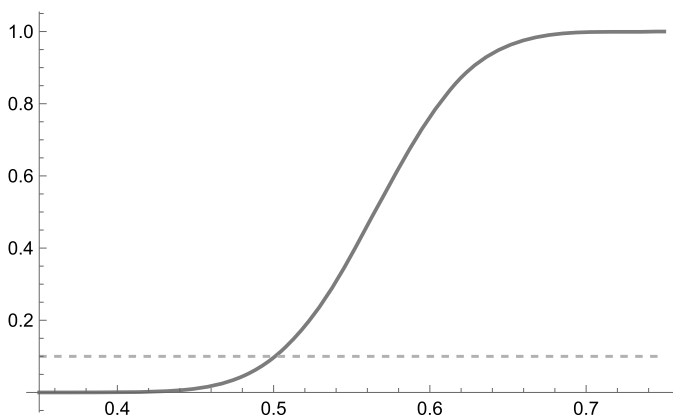


Figure 8.4: The power function $\beta_{\mathbf{T}}$ in eq. (8.13) with significance level $\alpha = 0.1$.

Making the significance level smaller, for example, taking $\alpha = 0.01$, since

$$\sum_{k=62}^{100} \binom{100}{k} 2^{-100} = 0.0104894 \quad \text{and} \quad \sum_{k=63}^{100} \binom{100}{k} 2^{-100} = 0.00601649,$$

we obtain $n_0 = 62$. Hence, if the number of white balls is 63 or larger, a rejection of \mathbb{H}_0 is 99 % sure.

Remark 8.3.8. Example 8.3.7 emphasizes once more the *dilemma* of hypothesis testing. The price one pays for higher confidence when rejecting \mathbb{H}_0 is the increase of the likelihood of a type II error. For instance, replacing $\alpha = 0.1$ by $\alpha = 0.01$ in the previous example leads to an enlargement of \mathcal{X}_0 from $\{0, \dots, 56\}$ to $\{0, \dots, 62\}$. Thus, if we observe 60 white balls, we reject \mathbb{H}_0 in the former case, but we cannot reject it in the latter one. This once more stresses the fact that an observation of an $x \in \mathcal{X}_0$ does not guarantee that \mathbb{H}_0 is true. It only means that the observed sample does not allow us to reject the hypothesis with high probability.

Summary: The model for testing binomial distributed populations is $(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{\theta \in [0,1]}$ where $\mathcal{X} = \{0, \dots, n\}$. Given $\theta_0 \in [0, 1]$, the hypotheses in the two-sided case are $\mathbb{H}_0 : \theta = \theta_0$ against $\mathbb{H}_1 : \theta \neq \theta_0$. If

$$n_0 = \min \left\{ k \leq n : B_{n,\theta_0}(\{0, \dots, k\}) > \frac{\alpha}{2} \right\} \quad \text{and} \quad n_1 = \max \left\{ k \leq n : B_{n,\theta_0}(\{k, \dots, n\}) > \frac{\alpha}{2} \right\},$$

then $\mathcal{X}_0 = \{n_0, \dots, n_1\}$ is the region of acceptance of an α -test checking \mathbb{H}_0 against \mathbb{H}_1 .

In the one-sided case $\mathbb{H}_0 : \theta \leq \theta_0$ against $\mathbb{H}_1 : \theta > \theta_0$ choose $\mathcal{X}_0 = \{0, \dots, n_0\}$ where now

$$n_0 = \max \{ k \leq n : B_{n,\theta_0}(\{k, \dots, n\}) > \alpha \}.$$

8.4 Tests for normally distributed populations

In this section we always assume $\mathcal{X} = \mathbb{R}^n$. That is, our samples are vectors $x = (x_1, \dots, x_n)$ with $x_j \in \mathbb{R}$. Given a sample $x \in \mathbb{R}^n$, we derive from it the following quantities that will soon play a crucial role.

Definition 8.4.1. If $x = (x_1, \dots, x_n) \in \mathbb{R}^n$, then we set

$$\bar{x} := \frac{1}{n} \sum_{j=1}^n x_j, \quad s_x^2 := \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2, \quad \text{and} \quad \sigma_x^2 := \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2. \quad (8.14)$$

The number \bar{x} is said to be the **sample mean** of x , while s_x^2 and σ_x^2 are said to be the **unbiased sample variance** and the **(biased) sample variance** of the vector x , respectively.

Analogously, if $X = (X_1, \dots, X_n)$ is an n -dimensional random vector,³ then we define the corresponding expressions pointwise. For instance, we have

$$\bar{X}(\omega) := \frac{1}{n} \sum_{j=1}^n X_j(\omega) \quad \text{and} \quad s_X^2(\omega) := \frac{1}{n-1} \sum_{j=1}^n (X_j(\omega) - \bar{X}(\omega))^2.$$

8.4.1 Fisher's theorem

We are going to prove important properties of normally distributed populations. They turn out to be the basis for all hypothesis tests in the normally distributed case. The starting point is a crucial lemma going back to Ronald Aylmer Fisher (1890–1962).

Lemma 8.4.2 (Fisher's lemma). *Let Y_1, \dots, Y_n be independent $\mathcal{N}(0, 1)$ -distributed random variables and let $B = (\beta_{ij})_{i,j=1}^n$ be a unitary $n \times n$ matrix. The random variables Z_1, \dots, Z_n are defined as*

$$Z_i := \sum_{j=1}^n \beta_{ij} Y_j, \quad 1 \leq i \leq n.$$

They possess the following properties:

- (i) *The variables Z_1, \dots, Z_n are also independent and $\mathcal{N}(0, 1)$ -distributed.*
- (ii) *For $m < n$, let the (random) quadratic form Q on \mathbb{R}^n be defined by*

$$Q := \sum_{j=1}^n Y_j^2 - \sum_{i=1}^m Z_i^2.$$

Then Q is independent of all Z_1, \dots, Z_m and, moreover, distributed according to χ_{n-m}^2 .

Proof. Assertion (i) was already proven in Proposition 6.1.19.

Let us verify (ii). The matrix B is unitary, thus it preserves the length of vectors in \mathbb{R}^n . Applying this to $Y = (Y_1, \dots, Y_n)$ and $Z = BY$ gives

$$\sum_{i=1}^n Z_i^2 = |Z|_2^2 = |BY|_2^2 = |Y|_2^2 = \sum_{j=1}^n Y_j^2,$$

which leads to

$$Q = Z_{m+1}^2 + \dots + Z_n^2. \tag{8.15}$$

³ To simplify the notation, now and later on, we denote random vectors by X , not by \bar{X} as we did before. This should not lead to confusion. For example, $\bar{\bar{X}}$ does not look very nice.

By virtue of (i), the random variables Z_1, \dots, Z_n are independent, hence by eq. (8.15) and Remark 4.1.10 the quadratic form Q is independent of Z_1, \dots, Z_m .

Recall that Z_{m+1}, \dots, Z_n are independent $\mathcal{N}(0, 1)$ -distributed. Thus, in view of eq. (8.15), Proposition 4.6.10 implies that Q is χ_{n-m}^2 -distributed. Observe that Q is the sum of $n - m$ squares. \square

Now we are in a position to state and prove one of the most important results in Mathematical Statistics.

Proposition 8.4.3 (Fisher's theorem). *Suppose X_1, \dots, X_n are independent and distributed according to $\mathcal{N}(\mu, \sigma^2)$ for some $\mu \in \mathbb{R}$ and some $\sigma^2 > 0$. Then the following are valid:*

$$\sqrt{n} \frac{\bar{X} - \mu}{\sigma} \text{ is } \mathcal{N}(0, 1)\text{-distributed;} \quad (8.16)$$

$$(n - 1) \frac{s_X^2}{\sigma^2} \text{ is } \chi_{n-1}^2\text{-distributed;} \quad (8.17)$$

$$\sqrt{n} \frac{\bar{X} - \mu}{s_X} \text{ is } t_{n-1}\text{-distributed, where } s_X := +\sqrt{s_X^2}. \quad (8.18)$$

Furthermore, \bar{X} and s_X^2 are independent random variables.⁴

Proof. Let us begin with the proof of assertion (8.16). Since the X_j s are independent and $\mathcal{N}(\mu, \sigma^2)$ -distributed, by Proposition 4.6.11 their sum $X_1 + \dots + X_n$ possesses an $\mathcal{N}(n\mu, n\sigma^2)$ distribution. Consequently, an application of Proposition 4.2.3 implies that \bar{X} is $\mathcal{N}(\mu, \sigma^2/n)$ -distributed, hence, another application of Proposition 4.2.3 tells us that $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ is standard normal. This completes the proof of statement (8.16).

We turn now to the verification of the remaining assertions. Letting

$$Y_j := \frac{X_j - \mu}{\sigma}, \quad 1 \leq j \leq n, \quad (8.19)$$

the random variables Y_1, \dots, Y_n are independent $\mathcal{N}(0, 1)$ -distributed. Moreover, their (unbiased) sample variance may be calculated by

$$\begin{aligned} s_Y^2 &= \frac{1}{n-1} \sum_{j=1}^n (Y_j - \bar{Y})^2 = \frac{1}{n-1} \left\{ \sum_{j=1}^n Y_j^2 - 2\bar{Y} \sum_{j=1}^n Y_j + n\bar{Y}^2 \right\} \\ &= \frac{1}{n-1} \left\{ \sum_{j=1}^n Y_j^2 - 2n\bar{Y}^2 + n\bar{Y}^2 \right\} = \frac{1}{n-1} \left\{ \sum_{j=1}^n Y_j^2 - (n\bar{Y})^2 \right\}. \end{aligned} \quad (8.20)$$

⁴ Recall that

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j \quad \text{and} \quad s_X^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2.$$

To proceed further, set $b_1 := (n^{-1/2}, \dots, n^{-1/2})$, and note that b_1 is a normalized n -dimensional vector, that is, we have $|b_1|_2 = 1$. Let $E \subseteq \mathbb{R}^n$ be the $(n - 1)$ -dimensional subspace consisting of elements that are perpendicular to b_1 . Choosing an orthonormal basis b_2, \dots, b_n in E , by the choice of E , the vectors b_1, \dots, b_n form an orthonormal basis in \mathbb{R}^n . If $b_i = (\beta_{i1}, \dots, \beta_{in})$, $1 \leq i \leq n$, let B be the $n \times n$ -matrix with entries β_{ij} , that is, the vectors b_1, \dots, b_n are the rows of B . Since $(b_i)_{i=1}^n$ are orthonormal, B is unitary.

As in Lemma 8.4.2, define Z_1, \dots, Z_n by

$$Z_i := \sum_{j=1}^n \beta_{ij} Y_j, \quad 1 \leq i \leq n,$$

and the quadratic form Q (with $m = 1$) as

$$Q := \sum_{j=1}^n Y_j^2 - Z_1^2.$$

Because of Lemma 8.4.2, the quadratic form Q is χ_{n-1}^2 -distributed and, furthermore, it is independent of Z_1 . By the choice of B and b_1 ,

$$\beta_{11} = \dots = \beta_{1n} = n^{-1/2},$$

hence $Z_1 = n^{1/2} \bar{Y}$ and, due to eq. (8.20), this leads to

$$Q = \sum_{j=1}^n Y_j^2 - (n^{1/2} \bar{Y})^2 = (n - 1) s_Y^2.$$

This observation implies $(n - 1) s_Y^2$ is χ_{n-1}^2 -distributed and, moreover, $(n - 1) s_Y^2$ and Z_1 are independent, thus also s_Y^2 and Z_1 .

The choice of the Y_j in eq. (8.19) immediately implies $\bar{Y} = \frac{\bar{X} - \mu}{\sigma}$, hence

$$(n - 1) s_Y^2 = \sum_{j=1}^n (Y_j - \bar{Y})^2 = \sum_{j=1}^n \left(\frac{X_j - \mu}{\sigma} - \frac{\bar{X} - \mu}{\sigma} \right)^2 = \frac{s_X^2}{\sigma^2} (n - 1),$$

which proves assertion (8.17).

Recall that $Z_1 = n^{1/2} \bar{Y} = n^{1/2} \frac{\bar{X} - \mu}{\sigma}$, which leads to $\bar{X} = n^{-1/2} \sigma Z_1 + \mu$. Thus, because of Proposition 4.1.9, the independence of $s_Y^2 = s_X^2 / \sigma^2$ and Z_1 implies that s_X^2 and \bar{X} are independent as well.

It remains to prove statement (8.18). We already know that $V := \sqrt{n} \frac{\bar{X} - \mu}{\sigma}$ is standard normal, and $W := (n - 1) s_X^2 / \sigma^2$ is χ_{n-1}^2 -distributed. Since they are independent, by Proposition 4.6, applied with $n - 1$, we get

$$\sqrt{n} \frac{\bar{X} - \mu}{s_X} = \frac{V}{\sqrt{\frac{1}{n-1} W}} \text{ is } t_{n-1}\text{-distributed.}$$

This implies assertion (8.18) and completes the proof of the proposition. \square

Remark 8.4.4. It is important to mention that the random variables X_1, \dots, X_n satisfy the assumptions of Proposition 8.4.3 if and only if the vector (X_1, \dots, X_n) is $\mathcal{N}(\mu, \sigma^2)^{\otimes n}$ -distributed or, equivalently, if its probability distribution is $\mathcal{N}(\bar{\mu}, \sigma^2 I_n)$.

Summary: Let X_1, \dots, X_n be independent and $\mathcal{N}(\mu, \sigma^2)$ -distributed. Then

$$\sqrt{n} \frac{\bar{X} - \mu}{\sigma} \sim \mathcal{N}(0, 1), \quad (n-1) \frac{s_X^2}{\sigma^2} \sim \chi_{n-1}^2, \quad \sqrt{n} \frac{\bar{X} - \mu}{s_X} \sim t_{n-1}.$$

Furthermore, \bar{X} and s_X^2 are independent random variables.

8.4.2 Quantiles

Let X be a (real-valued) random variable. Given a number $0 < \beta < 1$, a $u_\beta \in \mathbb{R}$ is said to be a β -quantile of X provided that

$$\mathbb{P}\{X \leq u_\beta\} \geq \beta \quad \text{and} \quad \mathbb{P}\{X \geq u_\beta\} \geq 1 - \beta.$$

Another way to write this is

$$\mathbb{P}_X((-\infty, u_\beta]) \geq \beta \quad \text{and} \quad \mathbb{P}_X([u_\beta, \infty)) \geq 1 - \beta.$$

In particular, this implies that the quantile only depends on the distribution of a random variable, not on the way it is defined. Since

$$\mathbb{P}\{X \geq u_\beta\} = 1 - \mathbb{P}\{X < u_\beta\},$$

the condition for the quantile may also be formulated as

$$\mathbb{P}\{X \leq u_\beta\} \geq \beta \quad \text{and} \quad \mathbb{P}\{X < u_\beta\} \leq \beta.$$

Example 8.4.5. Suppose that $\mathbb{P}\{X = 0\} = \mathbb{P}\{X = 1\} = \frac{1}{2}$. Then there is no β -quantile in the case $\beta \neq \frac{1}{2}$. Why? This is due to the fact that the distribution function $t \mapsto \mathbb{P}\{X \leq t\}$ attains only the values $0, \frac{1}{2}$, and 1 . If $\beta = 1/2$, then every number $u \in [0, 1)$ satisfies

$$\mathbb{P}\{X \leq u\} \geq \frac{1}{2} \quad \text{and} \quad \mathbb{P}\{X < u\} \leq \frac{1}{2},$$

hence each $u \in [0, 1)$ is a $(1/2)$ -quantile of X .

This example tells us two facts: quantiles do not always exist and, moreover, if they exist, then they need not be unique.

The situation becomes completely different if X possesses a positive distribution density.

Proposition 8.4.6. *Suppose that there this a positive density p such that*

$$F_X(t) = \mathbb{P}\{X \leq t\} = \int_{-\infty}^t p(x)dx, \quad t \in \mathbb{R}.$$

Then for each $\beta \in (0, 1)$, there is a unique β -quantile u_β . That is, there is a unique $u_\beta \in \mathbb{R}$ for which

$$F_X(u_\beta) = \mathbb{P}\{X \leq u_\beta\} = \int_{-\infty}^{u_\beta} p(x)dx = \beta. \quad (8.21)$$

Proof. Under these assumptions about X , its distribution function F_X is a one-to-one mapping from \mathbb{R} onto $(0, 1)$. Hence, its inverse function F_X^{-1} exists and $u_\beta = F_X^{-1}(\beta)$ is the unique number satisfying (8.21). \square

Remark 8.4.7. Of course, Proposition 8.4.6 remains valid if there are $a \in \mathbb{R}$ and/or $b \in \mathbb{R}$ such that the density p satisfies $p(x) = 0$ if $x < a$ and/or $p(x) = 0$ if $x > b$. In this case the quantile u_β satisfies either $u_\beta > a$ or $u_\beta < b$, respectively.

Let us now introduce some quantiles which will play an important later on. The first quantiles we consider are those of the standard normal distribution.

Definition 8.4.8. Let Φ be the distribution function of $\mathcal{N}(0, 1)$, as it was introduced in Definition 1.6.2. For a given $\beta \in (0, 1)$, the β -quantile z_β of the standard normal distribution is the unique real number satisfying

$$\Phi(z_\beta) = \beta \quad \text{or, equivalently,} \quad z_\beta = \Phi^{-1}(\beta).$$

Another way to define is as follows. Let X be a standard normal random variable. Then z_β is the unique real number such that

$$\mathbb{P}\{X \leq z_\beta\} = \beta.$$

The following properties of z_β will be used later on. Compare also Figure 8.5 for the assertions.

Proposition 8.4.9. *Let X be standard normally distributed. Then the following are valid:*

1. *We have $z_{1/2} = 0$, $z_\beta < 0$ for $0 < \beta < 1/2$, and $z_\beta > 0$ for $1/2 < \beta < 1$.*
2. *For all $0 < \beta < 1$, it follows that $\mathbb{P}\{X \geq z_\beta\} = 1 - \beta$.*
3. *If $0 < \beta < 1$, then $z_{1-\beta} = -z_\beta$.*
4. *For $0 < \alpha < 1$ we have $\mathbb{P}\{|X| \geq z_{1-\alpha/2}\} = \alpha$.*

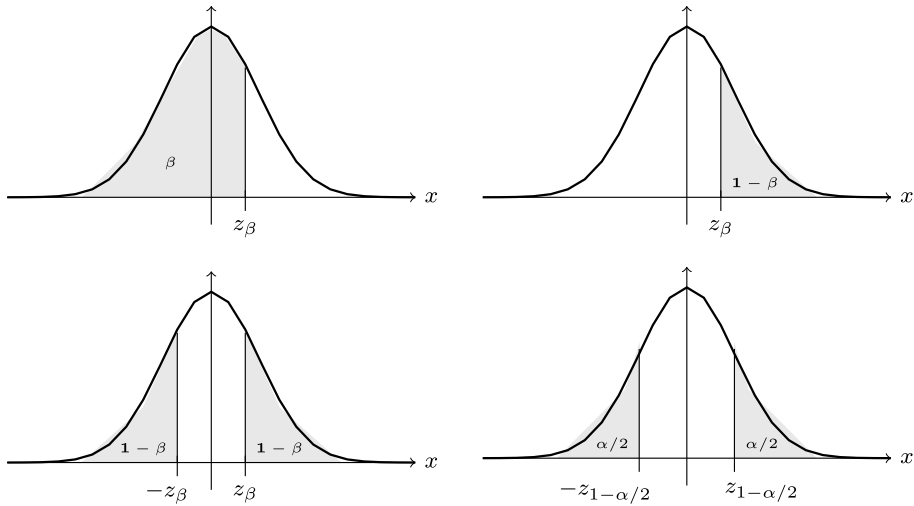


Figure 8.5: The assertions of Proposition 8.4.9.

Proof. The first property easily follows from $\Phi(0) = 1/2$, hence $\Phi(t) > 1/2$ if and only if $t > 0$.

Let X be standard normal. Then $\mathbb{P}\{X \geq z_\beta\} = 1 - \mathbb{P}\{X \leq z_\beta\} = 1 - \beta$, which proves the second assertion.

Since $-X$ is standard normal as well, by property 2 it follows that

$$\mathbb{P}\{X \leq -z_\beta\} = \mathbb{P}\{-X \geq z_\beta\} = \mathbb{P}\{X \geq z_\beta\} = 1 - \beta = \mathbb{P}\{X \leq z_{1-\beta}\},$$

hence $z_{1-\beta} = -z_\beta$ as asserted.

To prove the fourth assertion, note that properties 2 and 3 imply

$$\begin{aligned} \mathbb{P}\{|X| \geq z_{1-\alpha/2}\} &= \mathbb{P}\{X \leq -z_{1-\alpha/2} \text{ or } X \geq z_{1-\alpha/2}\} \\ &= \mathbb{P}\{X \leq -z_{1-\alpha/2}\} + \mathbb{P}\{X \geq z_{1-\alpha/2}\} \\ &= \mathbb{P}\{X \leq z_{\alpha/2}\} + \mathbb{P}\{X \geq z_{1-\alpha/2}\} = \alpha/2 + \alpha/2 = \alpha. \end{aligned}$$

Here we used $1 - \alpha/2 > 1/2$ implying $z_{1-\alpha/2} > 0$, hence the events $\{X \leq -z_{1-\alpha/2}\}$ and $\{X \geq z_{1-\alpha/2}\}$ are disjoint. \square

To get an impression about the size of the quantiles z_β , let us state a few of them.

β	0.999	0.995	0.99	0.95	0.9	0.8	0.75
z_β	3.0902	2.5758	2.3263	1.6449	1.2816	0.8416	0.6745

The values for small $\beta > 0$ follow by $z_{1-\beta} = -z_\beta$. So, for example,

$$z_{0,1} = -z_{0,9} = -1.2816.$$

The next quantiles, needed later on, are those of a χ_n^2 distribution.

Definition 8.4.10. Let X be distributed according to χ_n^2 and let $0 < \beta < 1$. The unique (positive) number $\chi_{n;\beta}^2$ satisfying

$$\mathbb{P}\{X \leq \chi_{n;\beta}^2\} = \beta$$

is called the β -quantile of the χ_n^2 distribution.

Two other, equivalent, ways to introduce these quantiles are as follows:

1. If X, \dots, X_n are independent standard normal, then

$$\mathbb{P}\{X_1^2 + \dots + X_n^2 \leq \chi_{n;\beta}^2\} = \beta.$$

2. The quantile $\chi_{n;\beta}^2$ satisfies

$$\frac{1}{2^{n/2}\Gamma(n/2)} \int_0^{\chi_{n;\beta}^2} x^{n/2-1} e^{-x/2} dx = \beta.$$

For later purposes, we mention also the following property. If $0 < \alpha < 1$, then for any χ_n^2 -distributed random variable X (see Figure 8.6),

$$\mathbb{P}\{X \notin [\chi_{n;\alpha/2}^2, \chi_{n;1-\alpha/2}^2]\} = \alpha. \tag{8.22}$$

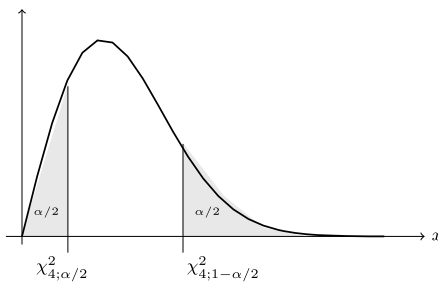


Figure 8.6: Graphic presentation of formula (8.22) for χ_4^2 .

In a similar way, we define now the quantiles of Student's t_n and of Fisher's $F_{m,n}$ distributions. For their descriptions, we refer to Definitions 4.7.6 and 4.7.13, respectively.

Definition 8.4.11. Let X be t_n -distributed and let Y be distributed according to $F_{m,n}$. For $\beta \in (0, 1)$ the β -quantiles $t_{n;\beta}$ and $F_{m,n;\beta}$ of the t_n and $F_{m,n}$ distributions are the unique numbers satisfying

$$\mathbb{P}\{X \leq t_{n;\beta}\} = \beta \quad \text{and} \quad \mathbb{P}\{Y \leq F_{m,n;\beta}\} = \beta.$$

Remark 8.4.12. Let X be t_n -distributed. Then $-X$ is t_n -distributed as well, hence $\mathbb{P}\{X \leq s\} = \mathbb{P}\{-X \leq s\}$ for $s \in \mathbb{R}$. Therefore, as in the case of the normal distribution, we get $-t_{n;\beta} = t_{n;1-\beta}$, and also

$$\mathbb{P}\{|X| > t_{n;1-\alpha/2}\} = \mathbb{P}\{|X| \geq t_{n;1-\alpha/2}\} = \alpha. \quad (8.23)$$

Remark 8.4.13. Another possibility to introduce the β -quantile of the t_n distribution is as follows: if X and X_1, \dots, X_n are independent standard normal variables, then

$$\mathbb{P}\left\{\frac{X}{\sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}} \leq t_{n;\beta}\right\} = \beta.$$

Similarly, one may characterize the quantiles of the $F_{m,n}$ distribution in the following way. Let X and Y be independent and distributed according to χ_m^2 and χ_n^2 , respectively. Then the β -quantile $F_{m,n;\beta}$ is the unique number satisfying

$$\mathbb{P}\left\{\frac{X/m}{Y/n} \leq F_{m,n;\beta}\right\} = \beta.$$

Note that the quantiles $F_{m,n;\beta}$ are positive numbers while the $t_{n;\beta}$ are negative if $0 < \beta < 1/2$ and positive in the case $1/2 < \beta < 1$.

If $s > 0$, then

$$\mathbb{P}\left\{\frac{X/m}{Y/n} \leq s\right\} = \mathbb{P}\left\{\frac{Y/n}{X/m} \geq \frac{1}{s}\right\} = 1 - \mathbb{P}\left\{\frac{Y/n}{X/m} \leq \frac{1}{s}\right\},$$

which immediately implies

$$F_{m,n;\beta} = \frac{1}{F_{n,m;1-\beta}}.$$

Summary: The following β -quantiles will play an important role later on:

$$\begin{aligned} X \sim \mathcal{N}(0, 1) &\Rightarrow \mathbb{P}\{X \leq z_\beta\} = \beta, & X \sim \chi_n^2 &\Rightarrow \mathbb{P}\{X \leq \chi_{n;\beta}\} = \beta, \\ X \sim t_n &\Rightarrow \mathbb{P}\{X \leq t_{n;\beta}\} = \beta, & X \sim F_{m,n} &\Rightarrow \mathbb{P}\{X \leq F_{m,n;\beta}\} = \beta. \end{aligned}$$

8.4.3 Z-tests or Gauss tests

Suppose we have an item of unknown length. In order to get some information about its length, we measure the item n times with an instrument of known accuracy. As sample we get a vector $x = (x_1, \dots, x_n)$, where x_j is the value obtained in the j th measurement. These measurements were executed independently, thus, we may assume that the x_j s are independent $\mathcal{N}(\mu, \sigma_0^2)$ -distributed with *known* $\sigma_0^2 > 0$ and *unknown* length $\mu \in \mathbb{R}$. Therefore, the describing statistical model is

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma_0^2)^{\otimes n})_{\mu \in \mathbb{R}} = (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\vec{\mu}, \sigma_0^2 I_n))_{\mu \in \mathbb{R}}.$$

From the hypothesis, two types of test apply in this case. We start with the so-called **one-sided Z-test** (also called one-sided Gauss test). Here the null hypothesis is $\mathbb{H}_0 : \mu \leq \mu_0$, where $\mu_0 \in \mathbb{R}$ is a given real number. Consequently, the alternative hypothesis is $\mathbb{H}_1 : \mu > \mu_0$, that is, $\Theta_0 = (-\infty, \mu_0]$ while $\Theta_1 = (\mu_0, \infty)$. In the above context, this says that we claim that the length of the item is less than or equal to a given μ_0 , and to check this we measure the item n times.

Proposition 8.4.14. *Let $\alpha \in (0, 1)$ be a given significance level. Then $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ with⁵*

$$\mathcal{X}_0 := \{x \in \mathbb{R}^n : \bar{x} \leq \mu_0 + n^{-1/2} \sigma_0 z_{1-\alpha}\}$$

and with

$$\mathcal{X}_1 := \{x \in \mathbb{R}^n : \bar{x} > \mu_0 + n^{-1/2} \sigma_0 z_{1-\alpha}\}$$

is an α -significance test to check $\mathbb{H}_0 : \mu \leq \mu_0$ against $\mathbb{H}_1 : \mu > \mu_0$. Here $z_{1-\alpha}$ denotes the $(1 - \alpha)$ -quantile introduced in Definition 8.4.8.

Proof. The assertion of Proposition 8.4.14 says that

$$\sup_{\mu \leq \mu_0} \mathbb{P}_\mu(\mathcal{X}_1) = \sup_{\mu \leq \mu_0} \mathcal{N}(\mu, \sigma_0^2)^{\otimes n}(\mathcal{X}_1) \leq \alpha.$$

To verify this, let us choose an arbitrary $\mu \leq \mu_0$ and define $S : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$S(x) := \sqrt{n} \frac{\bar{x} - \mu}{\sigma_0}, \quad x \in \mathbb{R}^n. \quad (8.24)$$

Regard S as a random variable on the probability space $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma_0^2)^{\otimes n})$. We claim that S is a standard normally distributed random variable. This fact is crucial. Therefore, let us give a more detailed reasoning.

⁵ Recall that \bar{x} denotes the arithmetic mean of an vector $x = (x_1, \dots, x_n)$ in \mathbb{R}^n .

Define random variables X_j on the probability space $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma_0^2)^{\otimes n})$ by $X_j(x) = x_j$, where $x = (x_1, \dots, x_n)$. Then the random vector $X = (X_1, \dots, X_n)$ is the identity on \mathbb{R}^n , hence $\mathcal{N}(\mu, \sigma_0^2)^{\otimes n}$ -distributed. In view of Remark 8.4.4 and since

$$S(x) = \sqrt{n} \frac{\bar{X}(x) - \mu}{\sigma_0},$$

assertion (8.16) applies for S , that is, it is $\mathcal{N}(0, 1)$ -distributed. Consequently,

$$\mathcal{N}(\mu, \sigma_0^2)^{\otimes n} \{x \in \mathbb{R}^n : S(x) > z_{1-\alpha}\} = \alpha. \quad (8.25)$$

Since $\mu \leq \mu_0$, we have

$$\begin{aligned} \mathcal{X}_1 &= \{x \in \mathbb{R}^n : \bar{x} > \mu_0 + n^{-1/2} \sigma_0 z_{1-\alpha}\} \\ &\subseteq \{x \in \mathbb{R}^n : \bar{x} > \mu + n^{-1/2} \sigma_0 z_{1-\alpha}\} = \{x \in \mathbb{R}^n : S(x) > z_{1-\alpha}\}, \end{aligned}$$

hence, by eq. (8.25), it follows that

$$\mathcal{N}(\mu, \sigma_0^2)^{\otimes n}(\mathcal{X}_1) \leq \mathcal{N}(\mu, \sigma_0^2)^{\otimes n} \{x \in \mathbb{R}^n : S(x) > z_{1-\alpha}\} = \alpha.$$

This completes the proof. \square

How does the power function of the Z-test in Proposition 8.4.14 look like? If S is as in eq. (8.24), then, according to Definition 8.2.11,

$$\begin{aligned} \beta_T(\mu) &= \mathcal{N}(\mu, \sigma_0^2)^{\otimes n}(\mathcal{X}_1) = \mathcal{N}(\mu, \sigma_0^2)^{\otimes n} \left\{ x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{\sigma_0} > z_{1-\alpha} \right\} \\ &= \mathcal{N}(\mu, \sigma_0^2)^{\otimes n} \left\{ x \in \mathbb{R}^n : S(x) > z_{1-\alpha} + (\mu_0 - \mu) \frac{\sqrt{n}}{\sigma_0} \right\} \\ &= 1 - \Phi \left(z_{1-\alpha} + (\mu_0 - \mu) \frac{\sqrt{n}}{\sigma_0} \right) = \Phi \left(z_\alpha + (\mu - \mu_0) \frac{\sqrt{n}}{\sigma_0} \right). \end{aligned}$$

In particular, β_T is increasing on \mathbb{R} with $\beta_T(\mu_0) = \alpha$. Moreover, we see that $\beta_T(\mu) < \alpha$ if $\mu < \mu_0$, and $\beta_T(\mu) > \alpha$ for $\mu > \mu_0$. See Figure 8.7 for an example of the power function.

While the critical region of a one-sided Z-test is an interval, in the case of the **two-sided Z-test** it is the union of two intervals. Here the null hypothesis is $\mathbb{H}_0 : \mu = \mu_0$, hence the alternative hypothesis is given as $\mathbb{H}_1 : \mu \neq \mu_0$.

Proposition 8.4.15. *The test $T = (\mathcal{X}_0, \mathcal{X}_1)$, where*

$$\begin{aligned} \mathcal{X}_0 &:= \{x \in \mathbb{R}^n : \mu_0 - n^{-1/2} \sigma_0 z_{1-\alpha/2} \leq \bar{x} \leq \mu_0 + n^{-1/2} \sigma_0 z_{1-\alpha/2}\} \\ &= \left\{ x \in \mathbb{R}^n : \sqrt{n} \frac{|\bar{x} - \mu_0|}{\sigma_0} \leq z_{1-\alpha/2} \right\} \end{aligned}$$

and $\mathcal{X}_1 = \mathbb{R}^n \setminus \mathcal{X}_0$, is an α -significance test for $\mathbb{H}_0 : \mu = \mu_0$ against $\mathbb{H}_1 : \mu \neq \mu_0$.

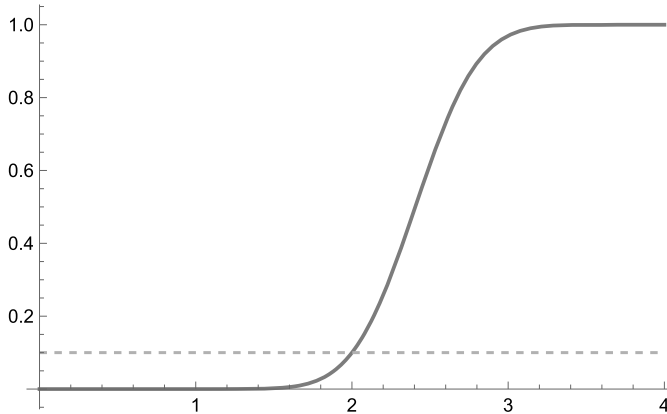


Figure 8.7: Power function of the one-sided Z-test \mathbf{T} with $\alpha = 0.1$, $\mu_0 = 2$, $\sigma_0 = 1$, and $n = 10$.

Proof. Since here $\Theta_0 = \{\mu_0\}$, the proof becomes easier than in the one-sided case. We only have to verify that

$$\mathcal{N}(\mu_0, \sigma_0^2)^{\otimes n}(\mathcal{X}_1) \leq \alpha. \quad (8.26)$$

Regarding S , defined by

$$S(x) := \sqrt{n} \frac{\bar{x} - \mu_0}{\sigma_0},$$

as a random variable on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu_0, \sigma_0^2)^{\otimes n})$, by the same arguments as in the previous proof, it is standard normally distributed. Thus, using assertion (4) of Proposition 8.4.9, we obtain

$$\mathcal{N}(\mu_0, \sigma_0^2)^{\otimes n}(\mathcal{X}_1) = \mathcal{N}(\mu_0, \sigma_0^2)^{\otimes n}\{x \in \mathbb{R}^n : |S(x)| > z_{1-\alpha/2}\} = \alpha.$$

Of course, this completes the proof. \square

Recall that $\beta_{\mathbf{T}}(\mu)$ is the probability to observe an $x = (x_1, \dots, x_n)$ in the critical region \mathcal{X}_1 provided μ is the “true” mean value. As can be seen in Fig. 8.8, this probability is small (equal α) if $\mu = \mu_0$ and becomes rapidly big for μ different of the suggested μ_0 .

Remark 8.4.16 (Important!). How to apply the one- or two-sided Z-test in a concrete situation? If the hypothesis is either $\mathbb{H}_0 : \mu \leq \mu_0$ or $\mathbb{H}_0 : \mu = \mu_0$, then the regions \mathcal{X}_1 of rejection are either

$$\left\{x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{\sigma_0} > z_{1-\alpha}\right\} \quad \text{or} \quad \left\{x \in \mathbb{R}^n : \sqrt{n} \frac{|\bar{x} - \mu_0|}{\sigma_0} > z_{1-\alpha/2}\right\}.$$

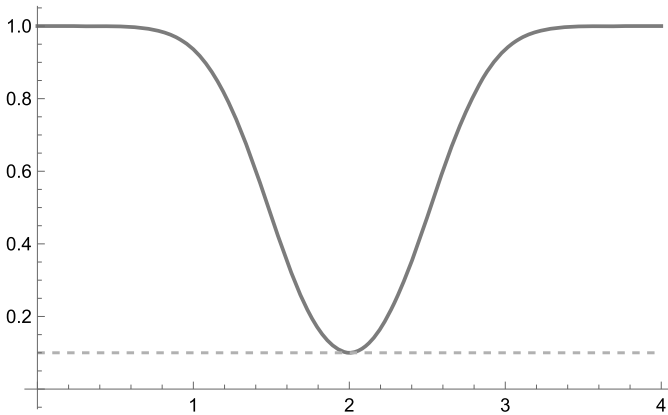


Figure 8.8: Power function of the two-sided Z-test \mathbf{T} with $\alpha = 0.1$, $\mu_0 = 2$, $\sigma_0 = 1$, and $n = 10$.

By the definition of the quantile, the former is equivalent to

$$\Phi\left(\sqrt{n} \frac{\bar{x} - \mu_0}{\sigma_0}\right) > 1 - \alpha \Leftrightarrow \Phi\left(\sqrt{n} \frac{\mu_0 - \bar{x}}{\sigma_0}\right) < \alpha.$$

Similarly, in the two-sided test, one has $x \in \mathcal{X}_1$ if and only if either

$$\Phi\left(\sqrt{n} \frac{\bar{x} - \mu_0}{\sigma_0}\right) < \alpha/2 \quad \text{or} \quad \Phi\left(\sqrt{n} \frac{\mu_0 - \bar{x}}{\sigma_0}\right) < \alpha/2.$$

Suppose now we observed n values $x = (x_1, \dots, x_n)$ which are independent and $\mathcal{N}(\mu, \sigma_0^2)$ -distributed for some unknown $\mu \in \mathbb{R}$. If

$$\Phi\left(\sqrt{n} \frac{\mu_0 - \bar{x}}{\sigma_0}\right) < \alpha \Rightarrow \text{reject } \mathbb{H}_0 : \mu \leq \mu_0.$$

Similarly, we have to reject $\mathbb{H}_0 : \mu = \mu_0$ if either

$$\Phi\left(\sqrt{n} \frac{\bar{x} - \mu_0}{\sigma_0}\right) < \alpha/2 \quad \text{or} \quad \Phi\left(\sqrt{n} \frac{\mu_0 - \bar{x}}{\sigma_0}\right) < \alpha/2.$$

In both cases (one- and two-sided test) the probability for an erroneous decision is bounded by $\alpha > 0$.

Example 8.4.17. Suppose the hypothesis is $\mathbb{H}_0 : \mu \leq \mu_0$ and our calculations lead to

$$\sqrt{n} \frac{\mu_0 - \bar{x}}{\sigma_0} = -2.13.$$

Since $\Phi(-2.13) \approx 0.0165858$, we may reject \mathbb{H}_0 with significance level α whenever $\alpha > 0.0165858$. But we cannot reject it with smaller risk. For example, if $\alpha = 0.01$, then

the result does not contradict the hypothesis. There is a great likelihood that \mathbb{H}_0 is wrong, but if we want to be very sure that this is so, we cannot derive this from the obtained result.

8.4.4 t-tests

The problem is similar to that considered in the case of the Z-test. But there is one important difference. We do no longer assume that the variance is known, which will be so in most cases. Therefore, this test is more realistic than the Z-test.

The starting point is the statistical model

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}.$$

Observe that the unknown parameter is now a vector $(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$. We begin by investigating the **one-sided t-test**. Given some $\mu_0 \in \mathbb{R}$, the null hypothesis is as before, that is, we have $\mathbb{H}_0 : \mu \leq \mu_0$. In the general setting, this means $\Theta_0 = (-\infty, \mu_0] \times (0, \infty)$, while $\Theta_1 = (\mu_0, \infty) \times (0, \infty)$.

To formulate the next result, let us shortly recall the following notations. If s_x^2 denotes the unbiased sample variance, as defined in eq. (8.14), set $s_x := +\sqrt{s_x^2}$. Furthermore, $t_{n-1;1-\alpha}$ denotes the $(1 - \alpha)$ -quantile of the t_{n-1} -distribution, as introduced in Definition 8.4.11.

Proposition 8.4.18. *Given $\alpha \in (0, 1)$, the regions \mathcal{X}_0 and \mathcal{X}_1 in \mathbb{R}^n are defined by*

$$\mathcal{X}_0 := \left\{ x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{s_x} \leq t_{n-1;1-\alpha} \right\}$$

and $\mathcal{X}_1 = \mathbb{R}^n \setminus \mathcal{X}_0$. With this choice of \mathcal{X}_0 and \mathcal{X}_1 , the test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test for $\mathbb{H}_0 : \mu \leq \mu_0$ against $\mathbb{H}_1 : \mu > \mu_0$.

Proof. Given $\mu \leq \mu_0$, define the random variable S on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})$ as

$$S(x) := \sqrt{n} \frac{\bar{x} - \mu}{s_x}, \quad x \in \mathbb{R}^n.$$

Property (8.18) implies that S is t_{n-1} -distributed, hence by the definition of the quantile $t_{n-1;1-\alpha}$, it follows that

$$\mathcal{N}(\mu, \sigma^2)^{\otimes n} \{x \in \mathbb{R}^n : S(x) > t_{n-1;1-\alpha}\} = \alpha.$$

From $\mu \leq \mu_0$, we easily derive

$$\mathcal{X}_1 \subseteq \{x \in \mathbb{R}^n : S(x) > t_{n-1;1-\alpha}\},$$

thus, as asserted,

$$\sup_{\mu \leq \mu_0} \mathcal{N}(\mu, \sigma^2)^{\otimes n}(\mathcal{X}_1) \leq \mathcal{N}(\mu, \sigma^2)^{\otimes n}\{x \in \mathbb{R}^n : S(x) > t_{n-1, 1-\alpha}\} = \alpha. \quad \square$$

As in the case of the Z-test, the null hypothesis of the **two-sided** t-test is $\mathbb{H}_0 : \mu = \mu_0$ for some $\mu_0 \in \mathbb{R}$. Again, we do not assume that the variance is known.

A two-sided t-test with significance level α may be constructed as follows.

Proposition 8.4.19. *Given $\alpha \in (0, 1)$, define regions \mathcal{X}_0 and \mathcal{X}_1 in \mathbb{R}^n by*

$$\mathcal{X}_0 := \left\{ x \in \mathbb{R}^n : \sqrt{n} \left| \frac{\bar{x} - \mu_0}{s_x} \right| \leq t_{n-1, 1-\alpha/2} \right\}$$

and $\mathcal{X}_1 = \mathbb{R}^n \setminus \mathcal{X}_0$. Then $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test for $\mathbb{H}_0 : \mu = \mu_0$ against $\mathbb{H}_1 : \mu \neq \mu_0$.

Proposition 8.4.19 is proven by similar methods, as we have used for the proofs of Propositions 8.4.15 and 8.4.18. Therefore, we decline to prove it here.

Example 8.4.20. We claim a certain workpiece has a length of 22 inches. Thus, the null hypothesis is $\mathbb{H}_0 : \mu = 22$. To check \mathbb{H}_0 , we measure the piece 10 times under the same conditions. The 10 values we obtained are (in inches)

22.17, 22.11, 22.10, 22.14, 22.02, 21.95, 22.02, 22.08, 21.98, 22.15

Do these values allow us to reject the hypothesis or do they confirm it? We have

$$\bar{x} = 22.072 \quad \text{and} \quad s_x = 0.07554248, \quad \text{hence} \quad \sqrt{10} \frac{\bar{x} - 22}{s_x} = 3.013986.$$

If we choose the significance level $\alpha = 0.05$, we have to investigate the quantile $t_{9, 0.975}$, which equals $t_{9, 0.975} = 2.26$. This lets us conclude the observed vector $x = (x_1, \dots, x_{10})$ belongs to \mathcal{X}_1 , and we may reject \mathbb{H}_0 . Consequently, with a confidence of 95% we may say, $\mu \neq 22$.

Another way to argue is as follows: Let $S(x) = \frac{\bar{x} - 22}{s_x} = 3.013986$. Then we have $S(x) \leq t_{9, 0.975}$ if and only if for a t_9 -distributed X it follows that

$$F_X(S(x)) = \mathbb{P}\{X \leq S(x)\} \leq 0.975.$$

But in our case $F_X(S(x)) = 0.992687$. So, we also get by this argument that \mathbb{H}_0 has to be rejected.

Remark 8.4.21. If we plug these 10 values, together with $\mu_0 = 22$, into a mathematical program, the result will be a number $\alpha_0 = 0.00128927$. What does this number tell us? It says the following. If we have chosen a significance level $\alpha > \alpha_0$, then we have to

reject \mathbb{H}_0 . But, if the chosen α satisfies $\alpha < \alpha_0$, then we fail to reject \mathbb{H}_0 . In our case we had $\alpha = 0.05 > 0.00128927 = \alpha_0$, hence we may reject \mathbb{H}_0 .

Thus, the price we pay for choosing a higher certainty and taking $\alpha < \alpha_0$ is that we are no longer able to reject the hypothesis \mathbb{H}_0 . It is as in the daily life: if one wants to be 99.9%-sure of not having a car accident, the best way is to avoid driving a car.

8.4.5 χ^2 -tests for the variance

The aim of this section is to get some information about the (unknown) variance of a normal distribution. Again we have to distinguish between the following two cases. The expected value is known or, otherwise, the expected value is unknown.

Let us start with the former case, that is, we assume that the *expected value is known* to be some $\mu_0 \in \mathbb{R}$. Then the statistical model is

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu_0, \sigma^2)^{\otimes n})_{\sigma^2 > 0}.$$

In the **one-sided χ^2 -test**, the null hypothesis is $\mathbb{H}_0 : \sigma^2 \leq \sigma_0^2$, for some given $\sigma_0^2 > 0$, while in the **two-sided χ^2 -test** we claim that $\mathbb{H}_0 : \sigma^2 = \sigma_0^2$.

Proposition 8.4.22. *In the one-sided setting, an α -significance χ^2 -test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is given by*

$$\mathcal{X}_0 := \left\{ x \in \mathbb{R}^n : \sum_{j=1}^n \frac{(x_j - \mu_0)^2}{\sigma_0^2} \leq \chi_{n;1-\alpha}^2 \right\}.$$

For the two-sided case, choose

$$\mathcal{X}_0 := \left\{ x \in \mathbb{R}^n : \chi_{n;\alpha/2}^2 \leq \sum_{j=1}^n \frac{(x_j - \mu_0)^2}{\sigma_0^2} \leq \chi_{n;1-\alpha/2}^2 \right\}, \quad (8.27)$$

to obtain an α -significance test. In both cases, the critical region is $\mathcal{X}_1 := \mathbb{R}^n \setminus \mathcal{X}_0$.

Proof. We prove the assertion only in the (slightly more difficult) one-sided case. For an arbitrarily chosen $\sigma^2 \leq \sigma_0^2$, let $\mathcal{N}(\mu_0, \sigma^2)^{\otimes n}$ be the underlying probability measure. We define now the random variables $X_j : \mathbb{R}^n \rightarrow \mathbb{R}$ as $X_j(x) = x_j$ for $x = (x_1, \dots, x_n)$. Then the X_j s are independent $\mathcal{N}(\mu_0, \sigma^2)$ -distributed. The normalization $Y_j := \frac{X_j - \mu_0}{\sigma}$ leads to independent standard normal Y_j s. Thus, if

$$S := \sum_{j=1}^n \frac{(X_j - \mu_0)^2}{\sigma^2} = \sum_{j=1}^n Y_j^2,$$

then, by Proposition 4.6.10, the random variable S is χ_n^2 -distributed. By the definition of quantile, we arrive at

$$\mathcal{N}(\mu_0, \sigma^2)^{\otimes n} \{x \in \mathbb{R}^n : S(x) > \chi_{n;1-\alpha}\} = \alpha.$$

Since $\sigma^2 \leq \sigma_0^2$, it follows that

$$\mathcal{X}_1 \subseteq \{x \in \mathbb{R}^n : S(x) > \chi_{n;1-\alpha}\},$$

hence $\mathcal{N}(\mu_0, \sigma^2)^{\otimes n}(\mathcal{X}_1) \leq \alpha$. This proves, as asserted, that $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test. \square

Let us now turn to the case where the *expected value is unknown*. Here the statistical model is given by

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}.$$

In the *one-sided case*, the null hypothesis is $\mathbb{H}_0 : \theta \leq \theta_0$. Thus, the parameter set $\Theta = \mathbb{R} \times (0, \infty)$ splits into $\Theta = \Theta_0 \cup \Theta_1$ with

$$\Theta_0 = \mathbb{R} \times (0, \sigma_0^2] \quad \text{and} \quad \Theta_1 = \mathbb{R} \times (\sigma_0^2, \infty).$$

In the *two-sided case*, the null hypothesis is $\mathbb{H}_0 : \theta = \theta_0$. Hence, in this case we have

$$\Theta_0 = \mathbb{R} \times \{\sigma_0^2\} \quad \text{and} \quad \Theta_1 = \mathbb{R} \times [(0, \sigma_0^2) \cup (\sigma_0^2, \infty)].$$

Proposition 8.4.23. *In the one-sided case, an α -significance test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is given by*

$$\mathcal{X}_0 := \left\{ x \in \mathbb{R}^n : (n-1) \frac{S_x^2}{\sigma_0^2} \leq \chi_{n-1;1-\alpha}^2 \right\}.$$

In the two-sided case, choose the region of acceptance as

$$\mathcal{X}_0 := \left\{ x \in \mathbb{R}^n : \chi_{n-1; \alpha/2}^2 \leq (n-1) \frac{S_x^2}{\sigma_0^2} \leq \chi_{n-1; 1-\alpha/2}^2 \right\} \quad (8.28)$$

to get an α -significance test. Again, the critical regions are given by $\mathcal{X}_1 := \mathbb{R}^n \setminus \mathcal{X}_0$.

Proof. The proof is very similar to that of Proposition 8.4.22, but with some important difference. Here we have to set

$$S(x) := (n-1) \frac{S_x^2}{\sigma^2}, \quad x \in \mathbb{R}^n.$$

Then property (8.17) applies, and it lets us conclude that S is χ_{n-1}^2 -distributed, provided that $\mathcal{N}(\mu, \sigma^2)^{\otimes n}$ is the true probability measure. After that observation the proof is completed as that of Proposition 8.4.22. \square

Summary: The most important *one-sample* α -significance tests for normally distributed populations are

Name	Parameters	Hypotheses	Critical region \mathcal{X}_1
One-sided Z-test	$\sigma^2 > 0$ known	$H_0 : \mu \leq \mu_0$ $H_1 : \mu > \mu_0$	$\{x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{\sigma} > z_{1-\alpha}\}$
Two-sided Z-test	$\sigma^2 > 0$ known	$H_0 : \mu = \mu_0$ $H_1 : \mu \neq \mu_0$	$\{x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{\sigma} > z_{1-\alpha/2}\}$
One-sided t-test	$\sigma^2 > 0$ unknown	$H_0 : \mu \leq \mu_0$ $H_1 : \mu > \mu_0$	$\{x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{s_x} > t_{n-1; 1-\alpha}\}$
Two-sided t-test	$\sigma^2 > 0$ unknown	$H_0 : \mu = \mu_0$ $H_1 : \mu \neq \mu_0$	$\{x \in \mathbb{R}^n : \sqrt{n} \frac{\bar{x} - \mu_0}{s_x} > t_{n-1; 1-\alpha/2}\}$
One-sided χ^2 -test with known mean value	$\mu \in \mathbb{R}$ known	$H_0 : \sigma^2 \leq \sigma_0^2$ $H_1 : \sigma^2 > \sigma_0^2$	$\{x \in \mathbb{R}^n : \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma_0^2} > \chi_{n-1; 1-\alpha}^2\}$
Two-sided χ^2 -test with known mean value	$\mu \in \mathbb{R}$ known	$H_0 : \sigma^2 = \sigma_0^2$ $H_1 : \sigma^2 \neq \sigma_0^2$	$\{x \in \mathbb{R}^n : \sum \frac{(x_i - \mu)^2}{\sigma_0^2} < \chi_{n; \alpha/2}^2\} \cup$ $\{x \in \mathbb{R}^n : \sum \frac{(x_i - \mu)^2}{\sigma_0^2} > \chi_{n; 1-\alpha/2}^2\}$
One-sided χ^2 -test with unknown mean value	$\mu \in \mathbb{R}$ unknown	$H_0 : \sigma^2 \leq \sigma_0^2$ $H_1 : \sigma^2 > \sigma_0^2$	$\{x \in \mathbb{R}^n : \sum \frac{(x_i - \bar{x})^2}{\sigma_0^2} > \chi_{n-1; 1-\alpha}^2\}$
Two-sided χ^2 -test with unknown mean value	$\mu \in \mathbb{R}$ unknown	$H_0 : \sigma^2 = \sigma_0^2$ $H_1 : \sigma^2 \neq \sigma_0^2$	$\{x \in \mathbb{R}^n : \sum \frac{(x_i - \bar{x})^2}{\sigma_0^2} < \chi_{n-1; \alpha/2}^2\} \cup$ $\{x \in \mathbb{R}^n : \sum \frac{(x_i - \bar{x})^2}{\sigma_0^2} > \chi_{n-1; 1-\alpha/2}^2\}$

8.4.6 Two-sample Z-tests

The two-sample Z-test compares the parameters of two different populations. Suppose we are given two different series of data, say $x = (x_1, \dots, x_m)$ and $y = (y_1, \dots, y_n)$, which were obtained independently by executing m experiments of the first kind and n experiments of the second. Combine both series to a single vector $(x, y) \in \mathbb{R}^{m+n}$.

A typical example for the described situation is as follows. A farmer grows grain on two different lots. On one lot he added fertilizer; on the other he did not. Now he wants to figure out whether or not adding fertilizer influenced the amount of grain gathered. Therefore, he measures the amount of grain on the first lot at m different spots and that on the second lot at n spots. The aim is to compare the mean values in both series of experiments.

We suppose that the samples x_1, \dots, x_m of the first population are independent and $N(\mu_1, \sigma_1^2)$ -distributed, while the y_1, \dots, y_n of the second population are independent and $N(\mu_2, \sigma_2^2)$ -distributed. Typical questions are as follows. Do we have $\mu_1 = \mu_2$

or, maybe, only $\mu_1 \leq \mu_2$? One may also ask whether or not $\sigma_1^2 = \sigma_2^2$ or, maybe, only $\sigma_1^2 \leq \sigma_2^2$.

To apply the two-sample Z-test, one has to suppose that the variances σ_1^2 and σ_2^2 are known. This reduces the number of parameters from 4 to 2, namely to μ_1 and μ_2 in \mathbb{R} . Thus, the describing statistical model is given by

$$(\mathbb{R}^{m+n}, \mathcal{B}(\mathbb{R}^{m+n}), \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n})_{(\mu_1, \mu_2) \in \mathbb{R}^2}. \quad (8.29)$$

Recall that $\mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n}$ denotes the multivariate normal distribution with expected value $(\underbrace{\mu_1, \dots, \mu_1}_m, \underbrace{\mu_2, \dots, \mu_2}_n)$ and covariance matrix $R = (r_{ij})_{i,j=1}^{m+n}$, where $r_{ii} = \sigma_1^2$ if $1 \leq i \leq m$, and $r_{ii} = \sigma_2^2$ if $m < i \leq m+n$. Furthermore, $r_{ij} = 0$ if $i \neq j$.

Proposition 8.4.24. *The statistical model is that in (8.29). To test $\mathbb{H}_0 : \mu_1 \leq \mu_2$ against $\mathbb{H}_1 : \mu_1 > \mu_2$, set*

$$\mathcal{X}_0 := \left\{ (x, y) \in \mathbb{R}^{m+n} : \sqrt{\frac{mn}{n\sigma_1^2 + m\sigma_2^2}} (\bar{x} - \bar{y}) \leq z_{1-\alpha} \right\}$$

and $\mathcal{X}_1 = \mathbb{R}^{m+n} \setminus \mathcal{X}_0$. Then the test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test for checking \mathbb{H}_0 against \mathbb{H}_1 . To test $\mathbb{H}_0 : \mu_1 = \mu_2$ against $\mathbb{H}_1 : \mu_1 \neq \mu_2$, let

$$\mathcal{X}_0 := \left\{ (x, y) \in \mathbb{R}^{m+n} : \sqrt{\frac{mn}{n\sigma_1^2 + m\sigma_2^2}} |\bar{x} - \bar{y}| \leq z_{1-\alpha/2} \right\}$$

and $\mathcal{X}_1 = \mathbb{R}^{m+n} \setminus \mathcal{X}_0$. Then the test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test for checking \mathbb{H}_0 against \mathbb{H}_1 .

Proof. Since the proof of the two-sided case is very similar to that of the one-sided, we only prove the first assertion. Thus, let us assume that \mathbb{H}_0 is valid, that is, we have $\mu_1 \leq \mu_2$. Then we have to verify that

$$\mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n}(\mathcal{X}_1) \leq \alpha. \quad (8.30)$$

To prove this, we investigate the random variables X_i and Y_j defined as $X_i(x, y) = x_i$ and $Y_j(x, y) = y_j$. Since the underlying probability space is

$$(\mathbb{R}^{m+n}, \mathcal{B}(\mathbb{R}^{m+n}), \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n}),$$

these random variables are independent and distributed according to $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$, respectively. Consequently, \bar{X} is $\mathcal{N}(\mu_1, \frac{\sigma_1^2}{m})$ -distributed, while \bar{Y} is distributed according to $\mathcal{N}(\mu_2, \frac{\sigma_2^2}{n})$. By the construction, \bar{X} and \bar{Y} are independent as well, and, moreover, since $-\bar{Y}$ is $\mathcal{N}(-\mu_2, \frac{\sigma_2^2}{n})$ -distributed, we conclude that the distribution of $\bar{X} - \bar{Y}$ equals $\mathcal{N}(\mu_1 - \mu_2, \frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n})$. Therefore, the mapping $S : \mathbb{R}^{m+n} \rightarrow \mathbb{R}$ defined by

$$S(x, y) := \left(\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n} \right)^{-1/2} [(\bar{X}(x, y) - \bar{Y}(x, y)) - (\mu_1 - \mu_2)]$$

is standard normal. By the definition of the quantile, this leads to

$$\mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n} \{ (x, y) \in \mathbb{R}^{m+n} : S(x, y) > z_{1-\alpha} \} = \alpha. \quad (8.31)$$

Since we assumed \mathbb{H}_0 to be correct, that is, we suppose $\mu_1 \leq \mu_2$, it follows that

$$S(x, y) \geq \left(\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n} \right)^{-1/2} [\bar{X}(x, y) - \bar{Y}(x, y)] = \sqrt{\frac{mn}{n\sigma_1^2 + m\sigma_2^2}} [\bar{X}(x, y) - \bar{Y}(x, y)].$$

Hence

$$\mathcal{X}_1 \subseteq \{ (x, y) \in \mathbb{R}^{m+n} : S(x, y) > z_{1-\alpha} \},$$

which by eq. (8.31) implies estimate (8.30). This completes the proof of this part of the proposition. \square

8.4.7 Two-sample t-tests

The situation is similar as in the two-sample Z-test, yet with one important difference. The variances σ_1^2 and σ_2^2 of the two populations are no longer known. Instead, we have to assume that they coincide, that is, we suppose

$$\sigma_1^2 = \sigma_2^2 := \sigma^2.$$

Therefore, there are three unknown parameters, the expected values μ_1 , μ_2 , and the common variance σ^2 . Thus, the statistical model describing this situation is given by

$$(\mathbb{R}^{m+n}, \mathcal{B}(\mathbb{R}^{m+n}), \mathcal{N}(\mu_1, \sigma^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma^2)^{\otimes n})_{(\mu_1, \mu_2, \sigma^2) \in \mathbb{R}^2 \times (0, \infty)}. \quad (8.32)$$

To simplify the formulation of the next statement, introduce $T : \mathbb{R}^{m+n} \rightarrow \mathbb{R}$ as

$$T(x, y) := \sqrt{\frac{(m+n-2)mn}{m+n}} \frac{\bar{x} - \bar{y}}{\sqrt{(m-1)s_x^2 + (n-1)s_y^2}}, \quad (x, y) \in \mathbb{R}^{m+n}. \quad (8.33)$$

Proposition 8.4.25. *Let the statistical model be as in (8.32). If*

$$\mathcal{X}_0 := \{ (x, y) \in \mathbb{R}^{m+n} : T(x, y) \leq t_{m+n-2; 1-\alpha} \}$$

and $\mathcal{X}_1 = \mathbb{R}^{m+n} \setminus \mathcal{X}_0$, then $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ is an α -significance test for $\mathbb{H}_0 : \mu_1 \leq \mu_2$ against $\mathbb{H}_1 : \mu_1 > \mu_2$.

On the other hand, the test $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ with

$$\mathcal{X}_0 := \{(x, y) \in \mathbb{R}^{m+n} : |T(x, y)| \leq t_{m+n-2; 1-\alpha/2}\}$$

and $\mathcal{X}_1 = \mathbb{R}^{m+n} \setminus \mathcal{X}_0$ is an α -significance test for $\mathbb{H}_0 : \mu_1 = \mu_2$ against $\mathbb{H}_1 : \mu_1 \neq \mu_2$.

Proof. This time we prove the two-sided case, that is, the null hypothesis is given by $\mathbb{H}_0 : \mu_1 = \mu_2$.

Let the random vectors $X = (X_1, \dots, X_m)$ and $Y = (Y_1, \dots, Y_n)$ on \mathbb{R}^{m+n} be defined with X_i s and Y_j s as in the proof of Proposition 8.4.24, that is, we have $X(x, y) = x$ and $Y(x, y) = y$. Then by Proposition 4.1.9 and Remark 4.1.10, the unbiased sample variances

$$s_X^2 = \frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})^2 \quad \text{and} \quad s_Y^2 = \frac{1}{n-1} \sum_{j=1}^n (Y_j - \bar{Y})^2$$

are independent as well. Furthermore, by virtue of statement (8.17), the random variables

$$(m-1) \frac{s_X^2}{\sigma^2} \quad \text{and} \quad (n-1) \frac{s_Y^2}{\sigma^2}$$

are distributed according to χ_{m-1}^2 and χ_{n-1}^2 , respectively. Proposition 4.6.9 implies that

$$S_{(X,Y)}^2 := \frac{1}{\sigma^2} \{(m-1)s_X^2 + (n-1)s_Y^2\}$$

is χ_{m+n-2}^2 -distributed. Since s_X^2 and \bar{X} , as well as s_Y^2 and \bar{Y} , are independent, by Proposition 8.4.3, this is also so for $S_{(X,Y)}^2$ and $\bar{X} - \bar{Y}$. As in the proof of Proposition 8.4.24, it follows that $\bar{X} - \bar{Y}$ is distributed according to $\mathcal{N}(\mu_1 - \mu_2, \frac{\sigma^2}{m} + \frac{\sigma^2}{n})$. Assume now that \mathbb{H}_0 is true, that is, we have $\mu_1 = \mu_2$. Then the last observation implies that $\frac{\sqrt{mn}}{\sigma\sqrt{m+n}}(\bar{X} - \bar{Y})$ is a standard normally distributed random variable and, furthermore, independent of $S_{(X,Y)}^2$. Thus, by Proposition 4.6, the distribution of the quotient

$$Z := \frac{\sqrt{mn}}{\sqrt{m+n-2}} \frac{\frac{\sqrt{mn}}{\sigma\sqrt{m+n}}(\bar{X} - \bar{Y})}{S_{(X,Y)}},$$

where $S_{(X,Y)} := +\sqrt{S_{(X,Y)}^2}$, is t_{m+n-2} -distributed. If T is as in eq. (8.33), then it is not difficult to prove that $Z = T(X, Y)$. Therefore, by the definition of X and Y , the mapping T is a t_{m+n-2} -distributed random variable on \mathbb{R}^{m+n} , endowed with the probability measure $\mathbb{P}_{\mu_1, \mu_2, \sigma^2} = \mathcal{N}(\mu_1, \sigma^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma^2)^{\otimes n}$. By eq. (8.23), this implies

$$\mathbb{P}_{\mu_1, \mu_2, \sigma^2}(\mathcal{X}_1) = \mathbb{P}_{\mu_1, \mu_2, \sigma^2} \{(x, y) \in \mathbb{R}^{m+n} : |T(x, y)| > t_{m+n-2; 1-\alpha/2}\} = \alpha,$$

as asserted. \square

8.4.8 F-tests

In this final section about tests, we compare the variances of two normally distributed sample series. Since the proofs of the assertions follow the schemes presented in the previous propositions, we decline to verify them here. We only mention the facts that play a crucial role during the proofs.

1. If X_1, \dots, X_m and Y_1, \dots, Y_n are independent and distributed according to $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$, then

$$V := \frac{1}{\sigma_1^2} \sum_{i=1}^m (X_i - \mu_1)^2 \quad \text{and} \quad W := \frac{1}{\sigma_2^2} \sum_{j=1}^n (Y_j - \mu_2)^2$$

are χ_m^2 and χ_n^2 -distributed and independent. Consequently, the quotient $\frac{V/m}{W/n}$ is $F_{m,n}$ -distributed.

2. For X_1, \dots, X_m and Y_1, \dots, Y_n independent and standard normal, the random variables

$$(m-1) \frac{s_X^2}{\sigma_1^2} \quad \text{and} \quad (n-1) \frac{s_Y^2}{\sigma_2^2}$$

are independent and distributed according to χ_{m-1}^2 and χ_{n-1}^2 , respectively. Thus, assuming $\sigma_1 = \sigma_2$, the quotient s_X^2/s_Y^2 possesses an $F_{m-1, n-1}$ -distribution.

When applying an F-test, as before, two different cases have to be considered.

- (K) The expected values μ_1 and μ_2 of the two populations are *known*. Then the statistical model is given by

$$(\mathbb{R}^{m+n}, \mathcal{B}(\mathbb{R}^{m+n}), \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n})_{(\sigma_1^2, \sigma_2^2) \in (0, \infty)^2}.$$

- (U) The expected values are *unknown*. This case is described by the statistical model

$$(\mathbb{R}^{m+n}, \mathcal{B}(\mathbb{R}^{m+n}), \mathcal{N}(\mu_1, \sigma_1^2)^{\otimes m} \otimes \mathcal{N}(\mu_2, \sigma_2^2)^{\otimes n})_{(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2) \in \mathbb{R}^2 \times (0, \infty)^2}.$$

In both cases, the null hypothesis may either be $\mathbb{H}_0 : \sigma_1^2 \leq \sigma_2^2$ in the one-sided case or $\mathbb{H}_0 : \sigma_1^2 = \sigma_2^2$ in the two-sided one. The regions of acceptance in each of the four different cases are given by the following subsets of \mathbb{R}^{m+n} , and always $\mathcal{X}_1 = \mathbb{R}^{m+n} \setminus \mathcal{X}_0$.

Case 1: $\mathbb{H}_0 : \sigma_1^2 \leq \sigma_2^2$ and μ_1, μ_2 are known. Then

$$\mathcal{X}_0 := \left\{ (x, y) \in \mathbb{R}^{m+n} : \frac{\frac{1}{m} \sum_{i=1}^m (x_i - \mu_1)^2}{\frac{1}{n} \sum_{j=1}^n (y_j - \mu_2)^2} \leq F_{m,n; 1-\alpha} \right\}.$$

Case 2: $\mathbb{H}_0 : \sigma_1^2 = \sigma_2^2$ and μ_1, μ_2 are known. Then

$$\mathcal{X}_0 := \left\{ (x, y) \in \mathbb{R}^{m+n} : F_{m,n;a/2} \leq \frac{\frac{1}{m} \sum_{i=1}^m (x_i - \mu_1)^2}{\frac{1}{n} \sum_{j=1}^n (y_j - \mu_2)^2} \leq F_{m,n;1-a/2} \right\}.$$

Case 3: $\mathbb{H}_0 : \sigma_1^2 \leq \sigma_2^2$ and μ_1, μ_2 are unknown. Then

$$\mathcal{X}_0 := \left\{ (x, y) \in \mathbb{R}^{m+n} : \frac{s_x^2}{s_y^2} \leq F_{m-1, n-1; 1-a} \right\}.$$

Case 4: $\mathbb{H}_0 : \sigma_1^2 = \sigma_2^2$ and μ_1, μ_2 are unknown. Then

$$\mathcal{X}_0 := \left\{ (x, y) \in \mathbb{R}^{m+n} : F_{m-1, n-1; a/2} \leq \frac{s_x^2}{s_y^2} \leq F_{m-1, n-1; 1-a/2} \right\}.$$

Example 8.4.26. Suppose there exist two different methods to measure certain items. Some evidence lets us suggest that method 2 is more precise than method 1. That is, we believe that the variance of the measurements by method 2 is smaller than the one by method 1.

To check this we measure some given item 39 times by method 1 and a maybe different item 28 times by method 2. As result we get $x = (x_1, \dots, x_{39})$ values obtained by method 1 and another 28 values $y = (y_1, \dots, y_{28})$ by method 2. Thus, in order to apply an F-test, we have $m = 39$ and $n = 28$.

Assume the unbiased variances of the samples x and y are

$$s_x^2 = 109.63 \quad \text{and} \quad s_y^2 = 65.99, \quad \text{hence} \quad \frac{s_x^2}{s_y^2} = 1.66.$$

Let σ_1^2 and σ_2^2 be the unknown variances of the x_i s and y_j s, respectively. To obtain as much information as possible, let us choose as hypotheses

$$\mathbb{H}_0 : \sigma_1^2 \leq \sigma_2^2 \quad \text{and} \quad \mathbb{H}_1 : \sigma_1^2 > \sigma_2^2. \quad (8.34)$$

Take $\alpha = 0.1$ as a significance level. If $s_x^2/s_y^2 > F_{38,27;0.9}$, then an application of the one-sided F-test implies that we may reject \mathbb{H}_0 . Note that $m - 1 = 38$ and $n - 1 = 27$.

Let X be an $F_{38,27}$ -distributed random variable. Then

$$s_x^2/s_y^2 > F_{38,27;0.9} \quad \Leftrightarrow \quad \mathbb{P}\left\{X \leq \frac{s_x^2}{s_y^2}\right\} > 0.9.$$

Tables or mathematical programs give

$$\mathbb{P}\left\{X \leq \frac{s_x^2}{s_y^2}\right\} = \mathbb{P}\{X \leq 1.66\} = 0.91402 > 0.9 = 1 - \alpha.$$

So we may reject H_0 and conclude that $\sigma_2^2 < \sigma_1^2$. That is, with probability greater than 0.9 we may say that method 2 is more precise than method 1.

Let us one more time emphasize the importance of the choice of the hypotheses in (8.34). If we were to choose $H_0 : \sigma_2^2 \leq \sigma_1^2$, then our test would confirm H_0 , but we could not say that H_0 is valid with a likelihood of at least 90 %.

Summary: The most important *two-sample* α -significance tests for normally distributed populations are

Name	Parameter	Hypotheses	Critical region \mathcal{X}_1
One-sided Z-test	σ_1^2, σ_2^2 known	$H_0 : \mu_1 \leq \mu_2$ $H_1 : \mu_1 > \mu_2$	$\{(x, y) \in \mathbb{R}^{m+n} : \frac{\sqrt{mn}}{\sqrt{n\sigma_1^2 + m\sigma_2^2}} \cdot (\bar{x} - \bar{y}) > z_{1-a}\}$
Two-sided Z-test	σ_1^2, σ_2^2 known	$H_0 : \mu_1 = \mu_2$ $H_1 : \mu_1 \neq \mu_2$	$\{(x, y) \in \mathbb{R}^{m+n} : \frac{\sqrt{mn}}{\sqrt{n\sigma_1^2 + m\sigma_2^2}} \cdot \bar{x} - \bar{y} > z_{1-a/2}\}$
One-sided t-test	$\sigma_1^2 = \sigma_2^2 > 0$ unknown	$H_0 : \mu_1 \leq \mu_2$ $H_1 : \mu_1 > \mu_2$	$\{(x, y) \in \mathbb{R}^{m+n} : \frac{\sqrt{(m+n-2)mn}}{\sqrt{m+n}} \times \frac{\bar{x} - \bar{y}}{\sqrt{(m-1)s_x^2 + (n-1)s_y^2}} > t_{m+n-2; 1-a}\}$
Two-sided t-test	$\sigma_1^2 = \sigma_2^2 > 0$ unknown	$H_0 : \mu_1 = \mu_2$ $H_1 : \mu_1 \neq \mu_2$	$\{(x, y) \in \mathbb{R}^{m+n} : \frac{\sqrt{(m+n-2)mn}}{\sqrt{m+n}} \times \frac{ \bar{x} - \bar{y} }{\sqrt{(m-1)s_x^2 + (n-1)s_y^2}} > t_{m+n-2; 1-a/2}\}$
One-sided F-test	μ_1, μ_2 known	$H_0 : \sigma_1^2 \leq \sigma_2^2$ $H_1 : \sigma_1^2 > \sigma_2^2$	$\{(x, y) \in \mathbb{R}^{m+n} : \frac{\frac{1}{m} \sum_{i=1}^m (x_i - \mu_1)^2}{\frac{1}{n} \sum_{j=1}^n (y_j - \mu_2)^2} > F_{m, n; 1-a}\}$
Two-sided F-test	μ_1, μ_2 known	$H_0 : \sigma_1^2 = \sigma_2^2$ $H_1 : \sigma_1^2 \neq \sigma_2^2$	$\{(x, y) \in \mathbb{R}^{m+n} : \frac{\frac{1}{m} \sum_{i=1}^m (x_i - \mu_1)^2}{\frac{1}{n} \sum_{j=1}^n (y_j - \mu_2)^2} > F_{m, n; 1-a/2}\}$ $\cup \{(x, y) \in \mathbb{R}^{m+n} : \frac{\frac{1}{n} \sum_{j=1}^n (y_j - \mu_2)^2}{\frac{1}{m} \sum_{i=1}^m (x_i - \mu_1)^2} > F_{n, m; 1-a/2}\}$
One-sided F-test	$\mu_1, \mu_2 \in \mathbb{R}$ unknown	$H_0 : \sigma_1^2 \leq \sigma_2^2$ $H_1 : \sigma_1^2 > \sigma_2^2$	$\{(x, y) \in \mathbb{R}^{m+n} : s_x^2 / s_y^2 > F_{m-1, n-1; 1-a}\}$
Two-sided F-test	$\mu_1, \mu_2 \in \mathbb{R}$ unknown	$H_0 : \sigma_1^2 = \sigma_2^2$ $H_1 : \sigma_1^2 \neq \sigma_2^2$	$\{(x, y) \in \mathbb{R}^{m+n} : s_x^2 / s_y^2 > F_{m-1, n-1; 1-a/2}\}$ $\cup \{(x, y) \in \mathbb{R}^{m+n} : s_y^2 / s_x^2 > F_{n-1, m-1; 1-a/2}\}$

8.5 Point estimators

Starting point is a parametric statistical model $(\mathcal{X}, \mathcal{F}, P_\theta)_{\theta \in \Theta}$. Assume we execute a statistical experiment and observe a sample $x \in \mathcal{X}$. The aim of this section is to show how this observation leads to a “good” estimate of the unknown parameter $\theta \in \Theta$.

Example 8.5.1. Suppose the statistical model is $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma_0^2)^{\otimes n})_{\mu \in \mathbb{R}}$ for some known $\sigma_0^2 > 0$. Thus, the unknown parameter is the expected value $\mu \in \mathbb{R}$. To estimate it, we execute n independent measurements and get $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. Knowing this

vector x , what is a “good” estimate for μ ? An intuitive approach is to define the point estimator $\hat{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}$ as

$$\hat{\mu}(x) = \frac{1}{n} \sum_{j=1}^n x_j = \bar{x}, \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n.$$

In other words, if the observed sample is x , then we take its sample mean $\hat{\mu}(x) = \bar{x}$ as an estimate for μ . An immediate question is whether $\hat{\mu}$ is a “good” estimator for μ . Or do there exist maybe “better” (more precise) estimators for μ ?

Before we investigate such and similar questions, the problem has to be generalized slightly. Sometimes it happens that we are not interested in the concrete value of the parameter $\theta \in \Theta$. We only want to know the value $\gamma(\theta)$ derived from θ . Thus, for some function $\gamma : \Theta \rightarrow \mathbb{R}$ we want to find a “good” estimator $\hat{\gamma} : \mathcal{X} \rightarrow \mathbb{R}$ for $\gamma(\theta)$. In other words, if we observe a sample $x \in \mathcal{X}$, then we take $\hat{\gamma}(x)$ as an estimate for the (unknown) value $\gamma(\theta)$. However, in most cases the function γ is not needed. That is, here we have $\gamma(\theta) = \theta$, and we look for a good estimator $\hat{\theta} : \mathcal{X} \rightarrow \Theta$ for θ .

Let us state an example where a nontrivial function γ plays a role.

Example 8.5.2. Let $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}$ be the statistical model. Thus, the unknown parameter is the two-dimensional vector (μ, σ^2) . But, in fact, we are only interested in μ , not in the pair (μ, σ^2) . That is, if

$$\gamma(\mu, \sigma^2) := \mu, \quad (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty),$$

then we want to find an estimate for $\gamma(\mu, \sigma^2)$.

Analogously, if we only want an estimate for σ^2 , then we choose γ as

$$\gamma(\mu, \sigma^2) := \sigma^2, \quad (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty).$$

After these preliminary considerations, we state now the precise definition of an estimator.

Definition 8.5.3. Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model and let $\gamma : \Theta \rightarrow \mathbb{R}$ be a function of the parameter. A mapping $\hat{\gamma} : \mathcal{X} \rightarrow \mathbb{R}$ is said to be a **point estimator** (or simply **estimator**) for $\gamma(\theta)$ if, given $t \in \mathbb{R}$, the set $\{x \in \mathcal{X} : \hat{\gamma}(x) \leq t\}$ belongs to the σ -field \mathcal{F} . In other words, $\hat{\gamma}$ is a random variable defined on \mathcal{X} .

The interpretation of this definition is as follows. If one observes the sample $x \in \mathcal{X}$, then $\hat{\gamma}(x)$ is an estimate for $\gamma(\theta)$. For example, if one measures a workpiece four times and gets 22.03, 21.87, 22.11, and 22.15 inches as results, then using the estimator $\hat{\mu}$ in Example 8.5.2, the estimate for the mean value equals 22.04 inches.

8.5.1 Maximum likelihood estimation

Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model. There exist several methods to construct “good” point estimators for the unknown parameter θ . In this section we present the probably most important of these methods, the so-called **maximum likelihood principle**.

To understand this principle, the following easy example may be helpful.

Example 8.5.4. Suppose the parameter set consists of two elements, say $\Theta = \{0, 1\}$. Moreover, also the sample space \mathcal{X} has cardinality two, that is, $\mathcal{X} = \{a, b\}$. Then the problem is as follows. Depending on the observation a or b , we have to choose either 0 or 1 as an estimate for θ .

For example, let us assume that $\mathbb{P}_0(\{a\}) = 1/4$, hence $\mathbb{P}_0(\{b\}) = 3/4$, and $\mathbb{P}_1(\{a\}) = \mathbb{P}_1(\{b\}) = 1/2$. Say, an experiment has outcome “ a .” What would be a good estimate for θ in this case? Should we take “0” or “1”? The answer is that we should choose “1.” Why? Because the sample “ a ” fits \mathbb{P}_1 better than \mathbb{P}_0 . By the same argument, we should take “0” as an estimate if we observe “ b .” Thus, the point estimator for θ is given by $\hat{\theta}(a) = 1$ and $\hat{\theta}(b) = 0$.

Example 8.5.5. Let us transform the previous example into one of daily life. Say your friend is planning to visit you. He will either arrive by train or by car. If he comes by train, he will be on time with probability $3/4$, and by car with probability $1/2$. Say he arrived on time. What would be your estimate for his choice? Do you guess he came by train or do you conjecture that he used the car? Justify your answer. What if your friend arrived late?

Which property characterizes the estimator $\hat{\theta}$ in Example 8.5.4? To answer this question, fix $x \in \mathcal{X}$ and look at the function

$$\theta \mapsto \mathbb{P}_\theta(\{x\}), \quad \theta \in \Theta. \quad (8.35)$$

If $x = a$, this function becomes maximal for $\theta = 1$, while for $x = b$ it attains its maximal value at $\theta = 0$. Consequently, the estimator $\hat{\theta}$ could also be defined as follows. For each fixed $x \in \mathcal{X}$, choose as an estimate the $\theta \in \Theta$ for which the function (8.35) becomes maximal. But this is exactly the approach of the maximum likelihood principle.

In order to describe this principle in the general setting, we have to introduce the notion of the likelihood function. Let us first assume that the sample space \mathcal{X} consists of *at most countably many elements*.

Definition 8.5.6. The function p from $\Theta \times \mathcal{X}$ to \mathbb{R} defined as

$$p(\theta, x) := \mathbb{P}_\theta(\{x\}), \quad \theta \in \Theta, \quad x \in \mathcal{X},$$

is called the **likelihood function** of the statistical model $(\mathcal{X}, \mathcal{P}(\mathcal{X}), \mathbb{P}_\theta)_{\theta \in \Theta}$.

We come now to the case where all probability measures \mathbb{P}_θ are *continuous*. Thus, we assume that the statistical model is $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathbb{P}_\theta)_{\theta \in \Theta}$ and, moreover, each \mathbb{P}_θ is continuous, that is, it has a density, mapping \mathbb{R}^n to \mathbb{R} . This density is not only a function of $x \in \mathbb{R}^n$, it also depends on the probability measure \mathbb{P}_θ , hence on $\theta \in \Theta$. Therefore, we denote the densities by $p(\theta, x)$. In other words, for each $\theta \in \Theta$ and each box $Q \subseteq \mathbb{R}^n$ as in eq. (1.73) we have

$$\mathbb{P}_\theta(Q) = \int_Q p(\theta, x) dx = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} p(\theta, x_1, \dots, x_n) dx_n \cdots dx_1. \quad (8.36)$$

Definition 8.5.7. The function $p : \Theta \times \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying eq. (8.36) for all boxes Q and all $\theta \in \Theta$ is said to be the **likelihood function** of the statistical model $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathbb{P}_\theta)_{\theta \in \Theta}$.

For a better understanding of Definitions 8.5.6 and 8.5.7, let us give some examples of likelihood functions.

1. First take $(\mathcal{X}, \mathcal{P}(\mathcal{X}), \mathcal{B}_{n,\theta})_{0 \leq \theta \leq 1}$ with $\mathcal{X} = \{0, \dots, n\}$ from Section 8.3. Then its likelihood function equals

$$p(\theta, k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}, \quad \theta \in [0, 1], \quad k \in \{0, \dots, n\}. \quad (8.37)$$

2. Consider the statistical model $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{N,M,n})_{M=0, \dots, N}$ investigated in Example 8.1.4. Then its likelihood function is given by

$$p(M, m) = \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad M = 0, \dots, N, \quad m = 0, \dots, n. \quad (8.38)$$

3. The likelihood function of the model $(\mathbb{N}_0^n, \mathcal{P}(\mathbb{N}_0^n), \text{Pois}_\lambda^{\otimes n})_{\lambda > 0}$ investigated in Example 8.5.22 is

$$p(\lambda, k_1, \dots, k_n) = \frac{\lambda^{k_1 + \dots + k_n}}{k_1! \cdots k_n!} e^{-\lambda n}, \quad \lambda > 0, \quad k_j \in \mathbb{N}_0. \quad (8.39)$$

4. The likelihood function of $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}$ from Example 8.1.12 can be calculated by

$$p(\mu, \sigma^2, x) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left(-\frac{|x - \vec{\mu}|^2}{2\sigma^2}\right), \quad \mu \in \mathbb{R}, \quad \sigma^2 > 0. \quad (8.40)$$

Here, as before, let $\vec{\mu} = (\mu, \dots, \mu)$.

5. The likelihood function of $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), E_\lambda^{\otimes n})_{\lambda > 0}$ from Example 8.1.9 may be represented as

$$p(\lambda, t_1, \dots, t_n) = \begin{cases} \lambda^n e^{-\lambda(t_1 + \dots + t_n)} & \text{if } t_j \geq 0, \lambda > 0, \\ 0 & \text{otherwise.} \end{cases} \quad (8.41)$$

Definition 8.5.8. Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model with likelihood function $p : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$. An estimator $\hat{\theta} : \mathcal{X} \rightarrow \Theta$ is said to be a **maximum likelihood estimator (MLE)** for $\theta \in \Theta$ provided that, for each $x \in \mathcal{X}$, the following is satisfied:

$$p(\hat{\theta}(x), x) = \max_{\theta \in \Theta} p(\theta, x)$$

Remark 8.5.9. Another way to define the MLE is as follows:⁶

$$\hat{\theta}(x) = \arg \max_{\theta \in \Theta} p(\theta, x), \quad x \in \mathcal{X}.$$

How does one find the MLE for concrete statistical models? One observation is that the logarithm is an increasing function. Thus, the likelihood function $p(\cdot, x)$ becomes maximal at a certain parameter $\theta \in \Theta$ if its logarithm $\ln p(\cdot, x)$ does.

Definition 8.5.10. Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a statistical model and let $p : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$ be its likelihood function. Suppose $p(\theta, x) > 0$ for all (θ, x) . Then the function

$$L(\theta, x) := \ln p(\theta, x), \quad \theta \in \Theta, \quad x \in \mathcal{X},$$

is called the **log-likelihood function** of the model.

Thus, $\hat{\theta}$ is an MLE if and only if

$$\hat{\theta}(x) = \arg \max_{\theta \in \Theta} L(\theta, x), \quad x \in \mathcal{X},$$

or, equivalently, if

$$L(\hat{\theta}(x), x) = \max_{\theta \in \Theta} L(\theta, x).$$

Example 8.5.11. If p is the likelihood function in eq. (8.37), then the log-likelihood function equals

$$L(\theta, k) = c + k \ln \theta + (n - k) \ln(1 - \theta), \quad 0 \leq \theta \leq 1, \quad k = 0, \dots, n. \quad (8.42)$$

Here $c \in \mathbb{R}$ denotes a certain constant independent of θ .

Example 8.5.12. The log-likelihood function of p in eq. (8.41) is well defined for $\lambda > 0$ and $t_j \geq 0$. For those λ s and t_j s, it is given by

$$L(\lambda, t_1, \dots, t_n) = n \ln \lambda - \lambda(t_1 + \dots + t_n).$$

⁶ If f is a real-valued function with domain A , then $x = \arg \max_{y \in A} f(y)$ if $x \in A$ and $f(x) \geq f(y)$ for all $y \in A$. In other words, x is one of the points in the domain A where f attains its maximal value.

To proceed further, we assume now that the parameter set Θ is a subset of \mathbb{R}^k for some $k \geq 1$. That is, each parameter θ consists of k unknown components, that is, it may be written as $\theta = (\theta_1, \dots, \theta_k)$ with $\theta_j \in \mathbb{R}$. Furthermore, suppose that for each fixed $x \in \mathcal{X}$ the log-likelihood function $L(\cdot, x)$ is continuously differentiable⁷ on Θ . Then points $\theta^* \in \Theta$ where $L(\cdot, x)$ becomes maximal must satisfy

$$\left. \frac{\partial}{\partial \theta_i} L(\theta, x) \right|_{\theta=\theta^*} = 0, \quad i = 1, \dots, k. \quad (8.43)$$

In particular, this is true for the MLE $\hat{\theta}(x)$. If for each $x \in \mathcal{X}$, the log-likelihood function $L(\cdot, x)$ is continuously differentiable on $\Theta \subseteq \mathbb{R}^k$, then the MLE $\hat{\theta}$ satisfies

$$\left. \frac{\partial}{\partial \theta_i} L(\theta, x) \right|_{\theta=\hat{\theta}(x)} = 0, \quad i = 1, \dots, k.$$



Example 8.5.13. Let us determine the MLE for the log-likelihood function in eq. (8.42). Here we have $\Theta = [0, 1] \subseteq \mathbb{R}$, hence the MLE $\hat{\theta} : \{0, \dots, n\} \rightarrow [0, 1]$ has to satisfy

$$\frac{\partial}{\partial \theta} L(\hat{\theta}(k), k) = \frac{k}{\hat{\theta}(k)} - \frac{n-k}{1-\hat{\theta}(k)} = 0.$$

This easily gives $\hat{\theta}(k) = \frac{k}{n}$, that is, the MLE in this case is defined by

$$\hat{\theta}(k) = \frac{k}{n}, \quad k = 0, \dots, n.$$

Let us interpret this result. In an urn there are white and black balls of unknown proportion. Let θ be the proportion of white balls. To estimate θ , draw n balls out of the urn, with replacement. Assume k of the chosen balls are white. Then $\hat{\theta}(k) = \frac{k}{n}$ is the MLE for the unknown proportion θ of white balls.

Example 8.5.14. The logarithm of the likelihood function p in eq. (8.40) equals

$$L(\mu, \sigma^2, x) = L(\mu, \sigma^2, x_1, \dots, x_n) = c - \frac{n}{2} \cdot \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2$$

with some constant $c \in \mathbb{R}$, independent of μ and of σ^2 . Thus, here $\Theta \subseteq \mathbb{R}^2$, hence, if $\theta^* = (\mu^*, \sigma^{2*})$ denotes the pair satisfying eq. (8.43), then

$$\left. \frac{\partial}{\partial \mu} L(\mu, \sigma^2, x) \right|_{(\mu, \sigma^2) = (\mu^*, \sigma^{2*})} = 0 \quad \text{and} \quad \left. \frac{\partial}{\partial \sigma^2} L(\mu, \sigma^2, x) \right|_{(\mu, \sigma^2) = (\mu^*, \sigma^{2*})} = 0.$$

⁷ The partial derivatives exist and are continuous.

Now

$$\frac{\partial}{\partial \mu} L(\mu, \sigma^2, x) = \frac{1}{\sigma^2} \sum_{j=1}^n (x_j - \mu) = \frac{1}{\sigma^2} \left[\sum_{j=1}^n x_j - n\mu \right],$$

which implies $\mu^* = \frac{1}{n} \sum_{j=1}^n x_j = \bar{x}$.

The derivative of L with respect to σ^2 , taken at $\mu^* = \bar{x}$, equals

$$\frac{\partial}{\partial \sigma^2} L(\bar{x}, \sigma^2, x) = -\frac{n}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{\sigma^4} \sum_{j=1}^n (x_j - \bar{x})^2.$$

It becomes zero at σ^{2*} satisfying

$$\sigma^{2*} = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 = \sigma_x^2,$$

where σ_x^2 was defined in eq. (8.14). Combining these observations, we see that the only pair $\theta^* = (\mu^*, \sigma^{2*})$ satisfying eq. (8.43) is given by (\bar{x}, σ_x^2) . Consequently, as MLE for $\theta = (\mu, \sigma^2)$ we obtain

$$\hat{\mu}(x) = \bar{x} \quad \text{and} \quad \widehat{\sigma^2}(x) = \sigma_x^2, \quad x \in \mathbb{R}^n.$$

Remark 8.5.15. Similar calculations as in the previous examples show that the MLE for the likelihood functions in eqs. (8.39) and (8.41) coincide with

$$\hat{\lambda}(k_1, \dots, k_n) = \frac{1}{n} \sum_{i=1}^n k_i \quad \text{and} \quad \hat{\lambda}(t_1, \dots, t_n) = \frac{1}{\frac{1}{n} \sum_{i=1}^n t_i}.$$

Finally, we present two likelihood functions where we have to determine their maximal values directly. Note that the above approach via the log-likelihood function does not apply if the parameter set Θ is either finite or countably infinite. In this case a derivative of $L(\cdot, x)$ does not make sense, hence we cannot determine points where it vanishes.

The first problem is that discussed in Remark 1.4.33. A retailer gets a delivery of N machines. Among the N machines are M defective. Since M is unknown, the retailer wants a “good” estimate for it. Therefore, he chooses at random n machines and tests them. Suppose he observes m defective machines among the tested. Does this lead to an estimate of the number M of defective machines? The next proposition answers this question.

Proposition 8.5.16. *The statistical model is given by $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{M,N,n})_{M=0,\dots,N}$. Then the MLE \hat{M} for M is of the form*

$$\hat{M}(m) = \begin{cases} \lfloor \frac{m(N+1)}{n} \rfloor & \text{if } m < n, \\ N & \text{if } m = n. \end{cases}$$

Here $\lfloor x \rfloor$ denotes the floor function (integer part) of a real number x . For example, $\lfloor 1.2 \rfloor = 1$ and $\lfloor \pi \rfloor = 3$.

Proof. The likelihood function p was determined in eq. (8.38) as

$$p(M, m) = \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad M = 0, \dots, N, \quad m = 0, \dots, n.$$

First note that $p(M, m) \neq 0$ if and only if $M \in \{m, \dots, N - n + m\}$ and, therefore, it suffices to investigate $p(M, m)$ for M s in this region. Thus, if $M - 1 \geq m$, then easy calculations lead to

$$\frac{p(M, m)}{p(M - 1, m)} = \frac{M}{M - m} \cdot \frac{N - M + 1 - (n - m)}{N - M + 1}. \quad (8.44)$$

By eq. (8.44), it follows that we have $p(M, m) \geq p(M - 1, m)$ if and only if

$$M(N - M + 1 - (n - m)) \geq (M - m)(N - M + 1).$$

Elementary transformations show the last estimate is equivalent to

$$-nM \geq -mN - m,$$

which happens if and only if $M \leq \frac{m(N+1)}{n}$.

Consequently, $M \mapsto p(M, m)$ is nondecreasing on $\{0, \dots, \lfloor \frac{m(N+1)}{n} \rfloor\}$, and it is nonincreasing on $\{\lfloor \frac{m(N+1)}{n} \rfloor, \dots, N\}$. Thus, if $m < n$, then the likelihood function $M \mapsto p(M, m)$ becomes maximal for $M^* = \lfloor \frac{m(N+1)}{n} \rfloor$, and the MLE is given by

$$\hat{M}(m) = \left\lfloor \frac{m(N+1)}{n} \right\rfloor, \quad m = 0, \dots, n - 1.$$

If $m = n$, then $M \mapsto p(M, m)$ is nonincreasing on $\{0, \dots, N\}$, hence in this case the likelihood function attains its maximal value at $M = N$, that is, $\hat{M}(n) = N$. \square

Example 8.5.17. A retailer gets a delivery of 100 TV sets for further selling. He chooses at random 15 sets and tests them. If there is exactly one defective TV set among the 15 tested, then the estimate for the number of defective sets in the delivery is 6. If he observes 2 defective sets, the estimate is 13, for 4 it is 26, and if there are even 6 defective TV sets among the 15 chosen, then the estimate is that 40 sets of the delivery are defective.

Finally, we come back to the question asked in Remark 1.4.36. In order to estimate the number N of fish in a pond, one catches M of them, marks them and puts them back into the pond. After some time one catches fish again, this time n of them. Among them m are marked. Does this number m lead to a “good” estimate of the number of fish in the pond? To describe this problem, we choose as statistical model

$$(\mathcal{X}, \mathbb{P}(\mathcal{X}), H_{N,M,n})_{N=0,1,\dots}$$

where $\mathcal{X} = \{0, \dots, n\}$. Here $H_{N,M,n}$ denotes the hypergeometric probability measure introduced in Definition 1.4.32. Thus, in this case the likelihood function is given by

$$p(N, m) = \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}, \quad N = 0, 1, \dots, \quad m = 0, \dots, n.$$

In the sequel, we have to exclude $m = 0$; in this case, there does not exist a reasonable estimate for N .

Proposition 8.5.18. *If $1 \leq m \leq n$, then the MLE \hat{N} for N is*

$$\hat{N}(m) = \left\lfloor \frac{Mn}{m} \right\rfloor. \quad (8.45)$$

Proof. The proof is quite similar to that of Proposition 8.5.16. Since

$$\frac{p(N, m)}{p(N-1, m)} = \frac{N-M}{N} \cdot \frac{N-n}{N-M-(n-m)},$$

it easily follows that the inequality $p(N, m) \geq p(N-1, m)$ is valid if and only if $N \leq \frac{Mn}{m}$. Therefore, $N \mapsto p(N, m)$ is nondecreasing if $N \leq \lfloor \frac{Mn}{m} \rfloor$ and nonincreasing for the remaining N . This immediately shows that the MLE is given by eq. (8.45). \square

Example 8.5.19. An unknown number of balls are in an urn. In order to estimate this number, we choose 50 balls from the urn and mark them. We put back the marked balls and mix the balls in the urn thoroughly. Then we choose another 30 balls from the urn. If there are 7 marked among the 30, then the estimate for the number of balls in the urn is 214. In the case of two marked balls, the estimate equals 750 while in the case of 16 marked balls we estimate that there are 93 balls in the urn.

Summary: Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model. A function $\hat{\theta} : \mathcal{X} \rightarrow \Theta$ is said to be a point estimator for the unknown parameter θ . That is, observing a sample $x \in \mathcal{X}$, then $\hat{\theta}(x)$ is our estimate for θ . The basic idea of the maximum likelihood estimator (MLE) is as follows: observing an $x \in \mathcal{X}$, one chooses that θ for which the likelihood function $p(\theta, x) = \mathbb{P}_\theta(\{x\})$ becomes maximal. This works well in the case of discrete \mathbb{P}_θ s. In the case of continuous \mathbb{P}_θ s, one asks for the maximum of the densities $p(\theta, x)$ of the \mathbb{P}_θ s. Thus, if $x \in \mathcal{X}$, then $\hat{\theta}(x) = \arg \max_{\theta \in \Theta} p(\theta, x)$ defines the MLE for this model. The logarithm is an increasing function, so that this is equivalent to

$$\hat{\theta}(x) = \arg \max_{\theta \in \Theta} L(\theta, x), \quad x \in \mathcal{X},$$

where the log-likelihood function L is defined by $L(\theta, x) = \ln p(\theta, x)$.

8.5.2 Unbiased estimators

Let us come back to the general setting. We are given a function $\gamma : \Theta \rightarrow \mathbb{R}$ and look for a “good” estimate for $\gamma(\theta)$. If $\hat{\gamma}(x)$ is the estimate, in most cases it will not be the correct value $\gamma(\theta)$. Sometimes the estimate is larger than $\gamma(\theta)$, sometimes one observes an $x \in \mathcal{X}$ for which $\hat{\gamma}(x)$ is smaller than the true value. For example, if the retailer in Example 8.5.17 gets every week a delivery of 100 TV sets, then sometimes his estimate for the number of defective sets will be bigger than the true value, sometimes smaller. Since he only pays for the nondefective sets, sometimes he pays too much, sometimes not enough. Therefore, a crucial condition for a good estimator should be that, on average, it meets the correct value. That is, in the long run, the loss and gain of the retailer should balance. In other words, the estimator should not be biased by a systematic error.

In view of Proposition 7.1.30, this condition for the estimator $\hat{\gamma}$ may be formulated as follows. If $\theta \in \Theta$ is the “true” parameter, then the expected value of $\hat{\gamma}$ should be $\gamma(\theta)$. To make this more precise,⁸ we need the following notation.

Definition 8.5.20. Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a statistical model and let $X : \mathcal{X} \rightarrow \mathbb{R}$ be a random variable. We write $\mathbb{E}_\theta X$ whenever the expected value of X is taken with respect to \mathbb{P}_θ . Similarly, in this case define

$$\mathbb{V}_\theta X = \mathbb{E}_\theta |X - \mathbb{E}_\theta X|^2$$

as variance of X . Of course, we have to assume that the expected value and/or the variance exist.

Remark 8.5.21. If X is discrete with values in $\{t_1, t_2, \dots\}$, then

$$\mathbb{E}_\theta X = \sum_{j=1}^{\infty} t_j \mathbb{P}_\theta\{X = t_j\}.$$

The case of continuous X is slightly more difficult because here we have to describe the density function of X with respect to \mathbb{P}_θ .

To become acquainted with Definition 8.5.20, the two following examples may be helpful. The first deals with the discrete case, while the second with the continuous one.

Example 8.5.22. Suppose the daily number of customers in a shopping center is Poisson distributed with unknown parameter $\lambda > 0$. To estimate this parameter, we record the number of customers on n different days. Thus, the sample we obtain is a vector $\vec{k} = (k_1, \dots, k_n)$ with $k_j \in \mathbb{N}_0$, where k_j is the number of customers on day j . The describing statistical model is given by $(\mathbb{N}_0^n, \mathcal{P}(\mathbb{N}_0^n), \text{Pois}_\lambda^{\otimes n})_{\lambda > 0}$ with distribution Pois_λ . Let $X : \mathbb{N}_0^n \rightarrow \mathbb{R}$ be defined by

⁸ How the expected value is defined? Note that we do not have only one probability measure, but many different ones.

$$X(\vec{k}) = X(k_1, \dots, k_n) := \frac{1}{n} \sum_{j=1}^n k_j, \quad \vec{k} = (k_1, \dots, k_n) \in \mathbb{N}_0^n.$$

Which value does $\mathbb{E}_\lambda X$ possess?

Answer: If we choose $\text{Pois}_\lambda^{\otimes n}$ as probability measure, then all X_j s defined by $X_j(k_1, \dots, k_n) := k_j$ are Pois_λ -distributed (and independent, but this is not needed here). Note that X_j is nothing else as the number of customers at day j . Hence, by Proposition 5.1.16, the expected value of X_j is λ , and since $X = \frac{1}{n} \sum_{j=1}^n X_j$, we finally obtain

$$\mathbb{E}_\lambda X = \mathbb{E}_\lambda \left(\frac{1}{n} \sum_{j=1}^n X_j \right) = \frac{1}{n} \sum_{j=1}^n \mathbb{E}_\lambda X_j = \frac{1}{n} n\lambda = \lambda.$$

Example 8.5.23. Take

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)}$$

as the statistical model. Thus, the parameter is of the form (μ, σ^2) for some $\mu \in \mathbb{R}$ and $\sigma^2 > 0$. Define $X : \mathbb{R}^n \rightarrow \mathbb{R}$ by $X(x) = \bar{x}$. If the underlying measure is $\mathcal{N}(\mu, \sigma^2)^{\otimes n}$, then⁹ X is $\mathcal{N}(\mu, \sigma^2/n)$ -distributed. Consequently, in view of Propositions 5.1.36 and 5.2.29, we obtain

$$\mathbb{E}_{\mu, \sigma^2} X = \mu \quad \text{and} \quad V_{\mu, \sigma^2} X = \frac{\sigma^2}{n}.$$

Using the notation introduced in Definition 8.5.20, the above-mentioned requirement for “good” estimators may now be formulated more precisely.

Definition 8.5.24. An estimator $\hat{\gamma} : \mathcal{X} \rightarrow \mathbb{R}$ is said to be an **unbiased** estimator for $\gamma : \Theta \rightarrow \mathbb{R}$ provided that for each $\theta \in \Theta$,

$$\mathbb{E}_\theta |\hat{\gamma}| < \infty \quad \text{and} \quad \mathbb{E}_\theta \hat{\gamma} = \gamma(\theta).$$

Remark 8.5.25. In view of Proposition 7.1.30, an estimator $\hat{\gamma}$ is unbiased if it possesses the following property: observe N independent samples x^1, \dots, x^N of a statistical experiment. Suppose that $\theta \in \Theta$ is the “true” parameter (according to which the x^j s are distributed). Then

$$\mathbb{P} \left\{ \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N \hat{\gamma}(x^j) = \gamma(\theta) \right\} = 1.$$

Thus, on average, the estimator $\hat{\gamma}$ approximately meets the correct value.

⁹ Compare with the first part of the proof of Proposition 8.4.3.

Example 8.5.26. Let us investigate whether the estimator in Example 8.5.13 is unbiased. The statistical model is $(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{0 \leq \theta \leq 1}$, where $\mathcal{X} = \{0, \dots, n\}$ and the estimator $\hat{\theta}$ acts as

$$\hat{\theta}(k) = \frac{k}{n}, \quad k = 0, \dots, n.$$

Setting $Z := n\hat{\theta}$, then Z is the identity on \mathcal{X} , hence $B_{n,\theta}$ -distributed. Proposition 5.1.13 implies $\mathbb{E}_\theta Z = n\theta$, thus,

$$\mathbb{E}_\theta \hat{\theta} = \mathbb{E}_\theta (Z/n) = \mathbb{E}_\theta Z/n = \theta. \quad (8.46)$$

Equation (8.46) holds for all $\theta \in [0, 1]$, that is, $\hat{\theta}$ is an unbiased estimator for θ .

Example 8.5.27. Next we come back to the problem presented in Example 8.5.22. The number of customers per day is Pois_λ -distributed with an unknown parameter $\lambda > 0$. The data of n days are combined into a vector $\hat{k} = (k_1, \dots, k_n) \in \mathbb{N}_0^n$. Then the parameter $\lambda > 0$ is estimated by $\hat{\lambda}$ defined as

$$\hat{\lambda}(\vec{k}) = \hat{\lambda}(k_1, \dots, k_n) := \frac{1}{n} \sum_{j=1}^n k_j.$$

Is this estimator for λ unbiased?

Answer: Yes, it is unbiased. Observe that $\hat{\lambda}$ coincides with the random variable X investigated in Example 8.5.22. There we proved $\mathbb{E}_\lambda X = \lambda$, hence, if $\lambda > 0$, then we have

$$\mathbb{E}_\lambda \hat{\lambda} = \lambda.$$

Example 8.5.28. We are given certain data x_1, \dots, x_n , which are known to be normally distributed and independent, and where the expected value μ and the variance σ^2 of the underlying probability measure are unknown. Thus, the describing statistical model is

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^n)_{(\mu, \sigma^2) \in \Theta} \quad \text{with } \Theta = \mathbb{R} \times (0, \infty).$$

The aim is to find unbiased estimators for μ and for σ^2 . Let us begin with estimating μ . That is, if γ is defined by $\gamma(\mu, \sigma^2) = \mu$, then we want to construct an unbiased estimator $\hat{\gamma}$ for γ . Let us take the MLE $\hat{\gamma}$ defined as

$$\hat{\gamma}(x) := \bar{x} = \frac{1}{n} \sum_{j=1}^n x_j, \quad x = (x_1, \dots, x_n).$$

Due to the calculations in Example 8.5.23, we obtain

$$\mathbb{E}_{\mu, \sigma^2} \hat{\gamma} = \mu = \gamma(\mu, \sigma^2).$$

This holds for all μ and σ^2 , hence $\hat{\gamma}$ is an unbiased estimator for $\mu = \gamma(\mu, \sigma^2)$.

How to find a suitable estimator for σ^2 ? This time the function γ has to be chosen as $\gamma(\mu, \sigma^2) := \sigma^2$. With s_x^2 defined in eq. (8.14), set

$$\hat{\gamma}(x) := s_x^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2, \quad x \in \mathbb{R}^n.$$

Is this an unbiased estimator for σ^2 ? To answer this, we use property (8.17) of Proposition 8.4.3. It asserts that the random variable $x \mapsto (n-1) \frac{s_x^2}{\sigma^2}$ is χ_{n-1}^2 -distributed, provided it is defined on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\mu, \sigma^2)^{\otimes n})$. Consequently, by Corollary 5.1.32, it follows that

$$\mathbb{E}_{\mu, \sigma^2} \left[(n-1) \frac{s_x^2}{\sigma^2} \right] = n-1.$$

Using the linearity of the expected value, we finally obtain

$$\mathbb{E}_{\mu, \sigma^2} \hat{\gamma} = \mathbb{E}_{\mu, \sigma^2} s_x^2 = \sigma^2.$$

Therefore, $\hat{\gamma}(x) = s_x^2$ is an unbiased¹⁰ estimator for σ^2 .

Remark 8.5.29. Taking the estimator $\hat{\gamma}(x) = s_x^2$ in the previous example, then, in view of $\sigma_x^2 = \frac{n-1}{n} s_x^2$, it follows that

$$\mathbb{E}_{\mu, \sigma^2} \hat{\gamma} = \frac{n-1}{n} \sigma^2.$$

Thus, the estimator $\hat{\gamma}(x) = s_x^2$ is biased. But note that

$$\lim_{n \rightarrow \infty} \frac{n-1}{n} \sigma^2 = \sigma^2,$$

hence, if the sample size n is big, then this estimator is “almost” unbiased. One says in this case the sequence of estimators (in dependence on n) is asymptotically unbiased.

The next example is slightly more involved, but of great interest in application.

Example 8.5.30. The lifetime of light bulbs is supposed to be exponentially distributed with some unknown parameter $\lambda > 0$. To estimate λ , we switch on n light bulbs and record the times t_1, \dots, t_n when they burn out. Thus, the observed sample is a vector $t = (t_1, \dots, t_n)$ in $(0, \infty)^n$. As an estimator for λ we choose

$$\hat{\lambda}(t) := \frac{n}{\sum_{j=1}^n t_j} = 1/\bar{t}.$$

Is this an unbiased estimator for λ ?

¹⁰ This explains why s_x^2 is called the *unbiased* sample variance.

Answer: The statistical model describing this experiment is

$$(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), E_\lambda^{\otimes n})_{\lambda > 0}.$$

If the random variables X_j are defined by $X_j(t) := t_j$, then they are independent and E_λ -distributed. Because of Proposition 4.6.6, their sum $X := \sum_{j=1}^n X_j$ possesses an Erlang distribution with parameters n and λ . An application of eq. (5.24) in Proposition 5.1.38 for $f(x) := \frac{n}{x}$ implies

$$\mathbb{E}_\lambda \hat{\lambda} = \mathbb{E}_\lambda \left(\frac{n}{X} \right) = \int_0^\infty \frac{n}{x} \frac{\lambda^n}{(n-1)!} x^{n-1} e^{-\lambda x} dx.$$

A change of variables $s = \lambda x$ transforms the latter integral into

$$\frac{\lambda n}{(n-1)!} \int_0^\infty s^{n-2} e^{-s} ds = \frac{\lambda n}{(n-1)!} \Gamma(n-1) = \frac{\lambda n}{(n-1)!} \cdot (n-2)! = \lambda \cdot \frac{n}{n-1}.$$

This tells us that $\hat{\lambda}$ is *not* an unbiased estimator for λ . But, as mentioned in Remark 8.5.29 for σ_x^2 , the sequence of estimators is asymptotically unbiased as $n \rightarrow \infty$.

Remark 8.5.31. If we replace the estimator in Example 8.5.30 by

$$\hat{\lambda}(t) := \frac{n-1}{\sum_{j=1}^n t_j} = \frac{1}{\frac{1}{n-1} \sum_{j=1}^n t_j}, \quad t = (t_1, \dots, t_n),$$

then the previous calculations imply

$$\mathbb{E}_\lambda \hat{\lambda} = \frac{n-1}{n} \cdot \lambda \cdot \frac{n}{n-1} = \lambda.$$

Hence, from this small change we get an unbiased estimator $\hat{\lambda}$ for λ .

Observe that the calculations in Example 8.5.30 were only valid for $n \geq 2$. If $n = 1$, then the expected value of $\hat{\lambda}$ does not exist.

Summary: Suppose we want to estimate the value $\gamma(\theta)$ for some function $\gamma : \Theta \rightarrow \mathbb{R}$. Let $\hat{\gamma} : \mathcal{X} \rightarrow \mathbb{R}$ be some point estimator for $\gamma(\theta)$. A basic property of a good estimator $\hat{\gamma}$ is that it should be unbiased. That is, on average, it should give us the correct value $\gamma(\theta)$, no matter which $\theta \in \Theta$ is the right one. In formulas, this means that an estimator $\hat{\gamma}$ is unbiased if given $\theta \in \Theta$, $\mathbb{E}_\theta \hat{\gamma} = \gamma(\theta)$.

8.5.3 Risk function

Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric statistical model. Furthermore, $\gamma : \Theta \rightarrow \mathbb{R}$ is a function of the parameter and $\hat{\gamma} : \mathcal{X} \rightarrow \mathbb{R}$ is an estimator for γ . Suppose $\theta \in \Theta$ is the true

parameter and we observe some $x \in \mathcal{X}$. Then, in general, we will have $\gamma(\theta) \neq \hat{\gamma}(x)$, and the quadratic error $|\gamma(\theta) - \hat{\gamma}(x)|^2$ occurs. Other ways to measure the error are possible and useful, but we restrict ourselves to the quadratic distance. In this way, we get the so-called **loss function** $L : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$ of $\hat{\gamma}$ defined by

$$L(\theta, x) := |\gamma(\theta) - \hat{\gamma}(x)|^2.$$

In other words, if θ is the correct parameter and our sample is $x \in \mathcal{X}$, then, using $\hat{\gamma}$ as the estimator, the (quadratic) error or loss will be $L(\theta, x)$. On average, the (quadratic) loss is evaluated by $\mathbb{E}_\theta |\gamma(\theta) - \hat{\gamma}|^2$.

Definition 8.5.32. The function R describing this average loss of $\hat{\gamma}$ is said to be the **risk function** of the estimator $\hat{\gamma}$. It is defined by

$$R(\theta, \hat{\gamma}) := \mathbb{E}_\theta |\gamma(\theta) - \hat{\gamma}|^2, \quad \theta \in \Theta.$$

Before giving some examples of risk functions, let us rewrite R as follows.

Proposition 8.5.33. *If $\theta \in \Theta$, then it follows that*

$$R(\theta, \hat{\gamma}) = |\gamma(\theta) - \mathbb{E}_\theta \hat{\gamma}|^2 + \mathbb{V}_\theta \hat{\gamma}. \quad (8.47)$$

Proof. The assertion is a consequence of

$$\begin{aligned} R(\theta, \hat{\gamma}) &= \mathbb{E}_\theta |\gamma(\theta) - \hat{\gamma}|^2 = \mathbb{E}_\theta [(\gamma(\theta) - \mathbb{E}_\theta \hat{\gamma}) + (\mathbb{E}_\theta \hat{\gamma} - \hat{\gamma})]^2 \\ &= |\gamma(\theta) - \mathbb{E}_\theta \hat{\gamma}|^2 + 2(\gamma(\theta) - \mathbb{E}_\theta \hat{\gamma}) \mathbb{E}_\theta (\mathbb{E}_\theta \hat{\gamma} - \hat{\gamma}) + \mathbb{V}_\theta \hat{\gamma}. \end{aligned}$$

Because of

$$\mathbb{E}_\theta (\mathbb{E}_\theta \hat{\gamma} - \hat{\gamma}) = \mathbb{E}_\theta \hat{\gamma} - \mathbb{E}_\theta \hat{\gamma} = 0,$$

this implies eq. (8.47). □

Definition 8.5.34. The function $\theta \mapsto |\gamma(\theta) - \mathbb{E}_\theta \hat{\gamma}|^2$, appearing in eq. (8.47), is said to be the **bias** or the **systematic error** of the estimator $\hat{\gamma}$.

Corollary 8.5.35. *A point estimator $\hat{\gamma}$ is unbiased if and only if for all $\theta \in \Theta$ its bias is zero. Moreover, if this is so, then its risk function is given by*

$$R(\theta, \hat{\gamma}) = \mathbb{V}_\theta \hat{\gamma}, \quad \theta \in \Theta.$$

Remark 8.5.36. Another way to formulate eq. (8.47) is as follows. The risk function of an estimator consists of two parts. One part is the systematic error, which does not occur for unbiased estimators. And the second part is given by $\mathbb{V}_\theta \hat{\gamma}$. Thus, the smaller the bias

and/or $\mathbb{V}_\theta \hat{y}$, the smaller the risk to get a wrong estimate for $\gamma(\theta)$, and the better the estimator.

Example 8.5.37. Let us determine the risk functions for the two estimators presented in Example 8.5.28. The estimator \hat{y} for μ was given by $\hat{y}(x) = \bar{x}$. Since this is an unbiased estimator, by Corollary 8.5.35, its risk function is computed as

$$R((\mu, \sigma^2), \hat{y}) = V_{(\mu, \sigma^2)} \hat{y}.$$

The random variable $x \mapsto \bar{x}$ is $\mathcal{N}(\mu, \sigma^2/n)$ -distributed, hence

$$R((\mu, \sigma^2), \hat{y}) = \frac{\sigma^2}{n}.$$

There are two interesting facts about this risk function. First, it does not depend on the parameter μ that we want to estimate. And secondly, if $n \rightarrow \infty$, then the risk tends to zero. In other words, the bigger the sample size, the less the risk for a wrong estimate.

Next we evaluate the risk function of the estimator $\hat{y}(x) = s_x^2$. As we saw in Example 8.5.30, this \hat{y} is also an unbiased estimator for σ^2 , hence

$$R((\mu, \sigma^2), \hat{y}) = V_{(\mu, \sigma^2)} \hat{y}.$$

From eq. (8.17), we know that $\frac{n-1}{\sigma^2} s_x^2$ is χ_{n-1}^2 -distributed, hence Corollary 5.2.28 implies

$$\mathbb{V}_{(\mu, \sigma^2)} \left[\frac{n-1}{\sigma^2} s_x^2 \right] = 2(n-1).$$

From this, one easily derives

$$R((\mu, \sigma^2), \hat{y}) = \mathbb{V}_{(\mu, \sigma^2)} s_x^2 = 2(n-1) \cdot \frac{\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{n-1}.$$

Here, the risk function depends heavily on the parameter σ^2 that we want to estimate. Furthermore, if $n \rightarrow \infty$, then also in this case the risk tends to zero.

Example 8.5.38. Finally, consider the statistical model $(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{0 \leq \theta \leq 1}$, where $\mathcal{X} = \{0, \dots, n\}$. In order to estimate $\theta \in [0, 1]$, we take, as in Example 8.5.26, the estimator $\hat{\theta}(k) = \frac{k}{n}$. There it was shown that the estimator is unbiased, hence, by Corollary 8.5.35, it follows that

$$R(\theta, \hat{\theta}) = \mathbb{V}_\theta \hat{\theta}, \quad 0 \leq \theta \leq 1.$$

If X is the identity on \mathcal{X} , by Proposition 5.2.18, its variance equals $\mathbb{V}_\theta X = n\theta(1-\theta)$. Since $\hat{\theta} = \frac{X}{n}$, this implies

$$R(\theta, \hat{\theta}) = \mathbb{V}_\theta (X/n) = \frac{\mathbb{V}_\theta X}{n^2} = \frac{\theta(1-\theta)}{n}.$$

Consequently, the risk function becomes maximal for $\theta = 1/2$, while for $\theta = 0$ and $\theta = 1$ it vanishes.

We saw in Corollary 8.5.35 that $R(\theta, \hat{\gamma}) = \mathbb{V}_\theta \hat{\gamma}$ for unbiased $\hat{\gamma}$. Thus, for such estimators inequality (7.2) implies

$$\mathbb{P}_\theta\{x \in \mathcal{X} : |\gamma(\theta) - \hat{\gamma}(x)| > c\} \leq \frac{\mathbb{V}_\theta \hat{\gamma}}{c^2},$$

that is, the smaller the $\mathbb{V}_\theta \hat{\gamma}$, the greater the chance to estimate a value near the correct one. This observation leads to the following definition.

Definition 8.5.39. Let $\hat{\gamma}_1$ and $\hat{\gamma}_2$ be two unbiased estimators for $\gamma(\theta)$. Then $\hat{\gamma}_1$ is said to be **uniformly better** than $\hat{\gamma}_2$ provided that

$$\mathbb{V}_\theta \hat{\gamma}_1 \leq \mathbb{V}_\theta \hat{\gamma}_2 \quad \text{for all } \theta \in \Theta.$$

An unbiased estimator $\hat{\gamma}_*$ is called the **uniformly best estimator** if it is uniformly better than all other unbiased estimators for $\gamma(\theta)$.

Example 8.5.40. We observe values that, for some $b > 0$, are uniformly distributed on $[0, b]$. But the number $b > 0$ is unknown. In order to estimate it, one executes n independent trials and obtains as sample $x = (x_1, \dots, x_n)$. As point estimators for $b > 0$ one may either choose

$$\hat{b}_1(x) := \frac{n+1}{n} \max_{1 \leq i \leq n} x_i \quad \text{or} \quad \hat{b}_2(x) := \frac{2}{n} \sum_{i=1}^n x_i.$$

According to Problem 8.4, the estimators \hat{b}_1 and \hat{b}_2 are both unbiased. Furthermore, not too difficult calculations show that

$$\mathbb{V}_b \hat{b}_1 = \frac{b^2}{n(n+2)} \quad \text{and} \quad \mathbb{V}_b \hat{b}_2 = \frac{b^2}{3n^2}.$$

Therefore, $\mathbb{V}_b \hat{b}_1 \leq \mathbb{V}_b \hat{b}_2$ for all $b > 0$. This tells us that \hat{b}_1 is uniformly better than \hat{b}_2 .

Remark 8.5.41. A very natural question is whether there exists a lower bound for the precision of an estimator. In other words, are there estimators for which the risk function becomes arbitrarily small? The answer depends heavily on the inherent information in the statistical model. To explain this, let us come back once more to Example 8.5.4.

Suppose we had $\mathbb{P}_0(\{a\}) = 1$ and $\mathbb{P}_1(\{b\}) = 1$. Then the occurrence of “ a ” would tell us with 100 % confidence that $\theta = 0$ is the correct parameter. The risk for the corresponding estimator is then zero. On the contrary, if $\mathbb{P}_0(\{a\}) = \mathbb{P}_1(\{b\}) = 1/2$, then the occurrence of “ a ” or “ b ” does us tell nothing about the correct parameter.

To make the previous observation more precise, we have to introduce some quantity that measures the information contained in a statistical model.

Definition 8.5.42. Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a statistical model with log-likelihood function L introduced in Definition 8.5.10. For simplicity, assume $\Theta \subseteq \mathbb{R}$. Then the function $I : \Theta \rightarrow \mathbb{R}$ defined by

$$I(\theta) := \mathbb{E}_\theta \left(\frac{\partial L}{\partial \theta} \right)^2$$

is called the **Fisher information** of the model. Of course, we have to suppose that the derivatives and the expected value exist.

Example 8.5.43. Let us investigate the Fisher information for the model treated in Example 8.5.14. There we had

$$L(\mu, \sigma^2, x) = L(\mu, \sigma^2, x_1, \dots, x_n) = c - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2.$$

Fix σ^2 and take the derivative with respect to μ . This leads to

$$\frac{\partial L}{\partial \mu} = \frac{n\bar{x} - n\mu}{\sigma^2},$$

hence

$$\left(\frac{\partial L}{\partial \mu} \right)^2 = \frac{n^2}{\sigma^4} |\bar{x} - \mu|^2.$$

Recall that \bar{x} is $\mathcal{N}(\mu, \sigma^2/n)$ -distributed, hence the expected value of $|\bar{x} - \mu|^2$ is nothing else than the variance of \bar{x} , that is, it is σ^2/n . Consequently,

$$I(\mu) = \mathbb{E}_{\mu, \sigma^2} \left(\frac{\partial L}{\partial \mu} \right)^2 = \frac{n^2}{\sigma^4} \frac{\sigma^2}{n} = \frac{n}{\sigma^2}.$$

The following result answers the above question: how precise can an estimator become at the most?

Proposition 8.5.44 (Rao–Cramér–Frechet). *Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ be a parametric model for which the Fisher information $I : \Theta \rightarrow \mathbb{R}$ exists. If $\hat{\theta}$ is an unbiased estimator for θ , then*

$$\mathbb{V}_\theta \hat{\theta} \geq \frac{1}{I(\theta)}, \quad \theta \in \Theta. \quad (8.48)$$

Remark 8.5.45. Estimators $\hat{\theta}$ that attain the lower bound in estimate (8.48) are said to be **efficient**. That is, for those estimators $\mathbb{V}_\theta \hat{\theta} = 1/I(\theta)$ for all $\theta \in \Theta$. In other words, efficient estimators possess the best possible accuracy.

In view of Examples 8.5.37 and 8.5.43, for normally distributed populations the estimator $\hat{\mu}(x) = \bar{x}$ is an efficient estimator for μ . Other efficient estimators are those investigated in Examples 8.5.27 and 8.5.13. On the other hand, the estimator for σ^2 in Exam-

ple 8.5.28 is not efficient. But it can be shown that s_x^2 is a uniformly best estimator for σ^2 , that is, there do not exist efficient estimators in this case.

Summary: Let $\hat{y} : \mathcal{X} \rightarrow \mathbb{R}$ be a point estimator for $y(\theta)$. Using \hat{y} as estimator for $y(\theta)$, the mean quadratic error is measured by the risk function defined by

$$R(\theta, \hat{y}) := \mathbb{E}_\theta |y(\theta) - \hat{y}|^2, \quad \theta \in \Theta.$$

The smaller the risk function, the better the estimator \hat{y} . It holds that

$$R(\theta, \hat{y}) = |y(\theta) - \mathbb{E}_\theta \hat{y}|^2 + \mathbb{V}_\theta \hat{y}.$$

The first term is the systematic error which vanishes in the case of unbiased estimators, hence then $R(\theta, \hat{y}) = \mathbb{V}_\theta \hat{y}$. It depends on inner properties of the model $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$ how small $\mathbb{V}_\theta \hat{y}$ can be chosen at most. Estimators attaining this lower bound are said to be efficient.

8.6 Confidence regions and intervals

8.6.1 Construction of confidence regions

Point estimations provide us with a single value $\theta \in \Theta$. Further work or necessary decisions are then based on this estimated parameter. The disadvantage of this approach is that we have no knowledge about the precision of the obtained value. Is the estimated parameter far away from the true one or maybe very near? To explain the problem, let us come back to the situation described in Example 8.5.17. If the retailer observes 4 defective TV sets among 15 tested, then he estimates that there are 26 defective sets in the delivery of 100. But he does not know how precise his estimate of 26 is. Maybe there are much more defective sets in the delivery, or maybe less than 26. The only information he has is that the estimates are correct on average. But this does not say anything about the accuracy of a single estimate.

This disadvantage of point estimators is avoided when estimating a certain set of parameters, not only a single point. Then the true parameter is contained with great probability in this randomly chosen region. In most cases, these regions will be intervals of real or natural numbers.

Definition 8.6.1. Suppose the parametric statistical model is $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$. A mapping $C : \mathcal{X} \rightarrow \mathcal{P}(\Theta)$ is called an **interval estimator**, provided for fixed $\theta \in \Theta$,

$$\{x \in \mathcal{X} : \theta \in C(x)\} \in \mathcal{F}. \quad (8.49)$$

Remark 8.6.2. A better notation for the mapping C would be region or set estimator because $C(x) \subseteq \Theta$ may be an arbitrary subset, not necessarily an interval, but “interval estimator” is commonly accepted, therefore, we use it here also.

Remark 8.6.3. Condition (8.49) is quite technical and will play no role later on. But it is necessary because otherwise the next definition does not make sense.

Definition 8.6.4. Let α be a real number in $(0, 1)$. Suppose an interval estimator $C : \mathcal{X} \rightarrow \mathcal{P}(\Theta)$ satisfies, for each $\theta \in \Theta$, the condition

$$\mathbb{P}_\theta\{x \in \mathcal{X} : \theta \in C(x)\} \geq 1 - \alpha. \quad (8.50)$$

Then C is said to be a **$1 - \alpha$ interval estimator** (also $1 - \alpha$ estimator). The sets $C(x) \subseteq \Theta$ with $x \in \mathcal{X}$ are called **$1 - \alpha$ confidence regions** or **confidence intervals** (sometimes also called $100(1 - \alpha)\%$ confidence regions or intervals).

How does an interval estimator apply? Suppose $\theta \in \Theta$ is the “true” parameter. In a statistical experiment, one obtains some sample $x \in \mathcal{X}$ distributed according to \mathbb{P}_θ . Depending on the observed sample x , we choose a set $C(x)$ of parameters. Then with probability greater than or equal to $1 - \alpha$, the observed $x \in \mathcal{X}$ leads to a region $C(x)$ of parameters which contains the true parameter θ .

Remark 8.6.5. It is important to say that the set $C(x)$ is *random, not the unknown parameter* $\theta \in \Theta$. Metaphorically speaking, a fish (the true parameter θ) is in a pond at some fixed but unknown spot. We execute a certain statistical experiment to get some information about the place where the fish is situated. Depending on the result of the experiment, we throw a net into the pond. Doing so, we know that with probability greater than or equal to $1 - \alpha$, the result of the experiment leads to a net that catches the fish. In other words, the position of the fish is not random, it is the observed sample, hence also the thrown net.

Remark 8.6.6. It is quite self-evident that one should try to choose the confidence sets as small as possible, without violating condition (8.50). If we are not interested in “small” confidence sets, then we could always choose $C(x) = \Theta$. This is not forbidden, but completely useless because we do not get any information about the true value θ .

Construction of confidence regions via significance tests: For a better understanding of the subsequent construction, let us shortly recall the main assertions about hypothesis tests from a slightly different point of view.

Let $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\vartheta)_{\vartheta \in \Theta}$ be a statistical model. We choose a fixed, but arbitrary, $\theta \in \Theta$. With this chosen θ , we formulate the null hypothesis as $\mathbb{H}_0 : \vartheta = \theta$. The alternative hypothesis is then $\mathbb{H}_1 : \vartheta \neq \theta$. Let $\mathbf{T} = (\mathcal{X}_0, \mathcal{X}_1)$ be an α -significance test for \mathbb{H}_0 against \mathbb{H}_1 . Because the hypothesis, hence also the test, depends on the chosen $\theta \in \Theta$, we denote the null hypothesis by $\mathbb{H}_0(\theta)$ and write $\mathbf{T}(\theta) = (\mathcal{X}_0(\theta), \mathcal{X}_1(\theta))$ for the test. That is, $\mathbb{H}_0(\theta) : \vartheta = \theta$ and $\mathbf{T}(\theta)$ is an α -significance test for $\mathbb{H}_0(\theta)$. With this notation set (compare Figure 8.9)

$$C(x) := \{\theta \in \Theta : x \in \mathcal{X}_0(\theta)\}. \quad (8.51)$$

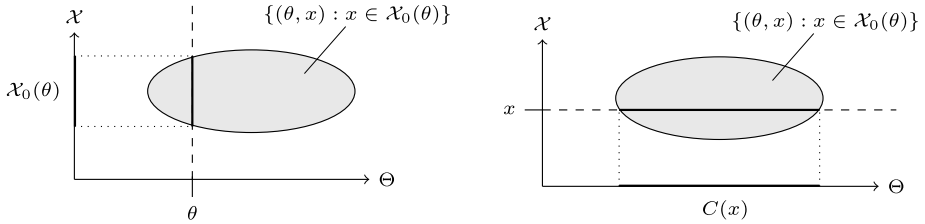


Figure 8.9: The equivalence between $x \in \mathcal{X}_0(\theta)$ and $\theta \in C(x)$.

Remark 8.6.7. Verbally said, an index $\theta \in \Theta$ belongs to $C(x)$ provided the observation of $x \in \mathcal{X}$ supports the hypothesis that θ is the correct parameter. For example, assume that there are only finitely many indices $\theta_1, \dots, \theta_n$. Then one may formulate n different hypothesis, namely $\mathbb{H}_0(\theta_1) : \vartheta = \theta_1$ up to $\mathbb{H}_0(\theta_n) : \vartheta = \theta_n$. Applying to each of these hypotheses an α -significance test, we obtain n regions $\mathcal{X}_0(\theta_1)$ up to $\mathcal{X}_0(\theta_n)$ of acceptance. Now, observing $x \in \mathcal{X}$, we choose those θ_j for which $x \in \mathcal{X}_0(\theta_j)$. That is, observing $x \in \mathcal{X}$, for those θ_j this would not lead to a rejection of $\mathbb{H}_0(\theta_j)$. So we set

$$C(x) = \{\theta_j : x \in \mathcal{X}_0(\theta_j)\}.$$

Example 8.6.8. Choose the hypothesis and the test as in Proposition 8.4.15. The statistical model is then given by $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(v, \sigma_0^2)^{\otimes n})_{v \in \mathbb{R}}$, where this time we denote the unknown expected value by v . For some fixed, but arbitrary, $\mu \in \mathbb{R}$ let

$$\mathbb{H}_0(\mu) : v = \mu \quad \text{and} \quad \mathbb{H}_1(\mu) : v \neq \mu.$$

The α -significance test $\mathbf{T}(\mu)$ constructed in Proposition 8.4.15 possesses the region of acceptance

$$\mathcal{X}_0(\mu) = \left\{ x \in \mathbb{R}^n : \sqrt{n} \left| \frac{\bar{x} - \mu}{\sigma_0} \right| \leq z_{1-\alpha/2} \right\}.$$

Thus, in this case, the set $C(x)$ in eq. (8.51) consists of those $\mu \in \mathbb{R}$ that satisfy the estimate $\sqrt{n} \left| \frac{\bar{x} - \mu}{\sigma_0} \right| \leq z_{1-\alpha/2}$. That is, given $x \in \mathbb{R}^n$, then $C(x)$ is the interval

$$C(x) = \left[\bar{x} - \frac{\sigma_0}{\sqrt{n}} z_{1-\alpha/2}, \bar{x} + \frac{\sigma_0}{\sqrt{n}} z_{1-\alpha/2} \right].$$

Let us come back to the general situation. The statistical model is $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\vartheta)_{\vartheta \in \Theta}$. Given $\theta \in \Theta$, let $\mathbf{T}(\theta)$ be an α -significance test for $\mathbb{H}_0(\theta)$ against $\mathbb{H}_1(\theta)$ where $\mathbb{H}_0(\theta)$ is the hypothesis $\mathbb{H}_0(\theta) : \vartheta = \theta$. Given $x \in \mathcal{X}$, define $C(x) \subseteq \Theta$ by eq. (8.51). Then the following is valid.

Proposition 8.6.9. *Let $\mathbf{T}(\theta)$ be as above an α -significance test for $\mathbb{H}_0(\theta)$ against $\mathbb{H}_1(\theta)$. Define $C(x)$ by eq. (8.51) where $\mathcal{X}_0(\theta)$ denotes the region of acceptance of $\mathbf{T}(\theta)$. Then the*

mapping $x \mapsto C(x)$ from \mathcal{X} into $\mathcal{P}(\Theta)$ is a $1 - \alpha$ interval estimator. Hence, $\{C(x) : x \in \mathcal{X}\}$ is a collection of $1 - \alpha$ confidence regions.

Proof. By assumption, $\mathbf{T}(\theta)$ is an α -significance test for $\mathbb{H}_0(\theta)$. The definition of those tests tells us that

$$\mathbb{P}_\theta(\mathcal{X}_1(\theta)) \leq \alpha, \quad \text{hence } \mathbb{P}_\theta(\mathcal{X}_0(\theta)) \geq 1 - \alpha.$$

Given $\theta \in \Theta$ and $x \in \mathcal{X}$, by the construction of $C(x)$, one has $\theta \in C(x)$ if and only if $x \in \mathcal{X}_0(\theta)$. Combining these two observations, given $\theta \in \Theta$, then it follows that

$$\mathbb{P}_\theta\{x \in \mathcal{X} : \theta \in C(x)\} = \mathbb{P}_\theta\{x \in \mathcal{X} : x \in \mathcal{X}_0(\theta)\} = \mathbb{P}_\theta(\mathcal{X}_0(\theta)) \geq 1 - \alpha.$$

This completes the proof. □

Summary: Suppose the parametric statistical model is $(\mathcal{X}, \mathcal{F}, \mathbb{P}_\theta)_{\theta \in \Theta}$. To each $x \in \mathcal{X}$, we assign a subset $C(x) \subseteq \Theta$ such that for a given $\alpha > 0$ the following holds: if \mathbb{P}_θ is the true probability measure, then with probability greater than $1 - \alpha$ we observe an $x \in \mathcal{X}$ such that $\theta \in C(x)$. That is, for all $\theta \in \Theta$ it follows that

$$\mathbb{P}_\theta\{x \in \mathcal{X} : \theta \in C(x)\} \geq 1 - \alpha.$$

The (random) sets $C(x)$ are called $1 - \alpha$ confidence sets.

There exists a tight relation between the construction of the confidence sets $C(x)$ and hypothesis tests. If $\mathbf{T}(\theta) = (\mathcal{X}_0(\theta), \mathcal{X}_1(\theta))$ is an α -test for the hypothesis $\mathbb{H}_0 : \text{“}\theta \text{ is the true parameter,“}$ then

$$C(x) = \{\theta \in \Theta : x \in \mathcal{X}_0(\theta)\}$$

are $1 - \alpha$ confidence sets. Verbally said, a parameter θ belongs to $C(x)$ if the occurrence of x does not contradict the hypotheses that θ is the correct parameter.

Test of a hypothesis $\mathbb{H}_0 \Rightarrow$ fixed region $\mathcal{X}_0 \subseteq \mathcal{X}$ of acceptance,

Confidence regions \Rightarrow random region $C(x) \subseteq \Theta$ of probable parameters.

8.6.2 Normally distributed samples

The aim of this section is to apply Proposition 8.6.9 to transform results about two-sided significance tests for normally distributed samples into assertions about confidence intervals.

We start with the application of the two-sided Z-test for $\mathcal{N}(\mu, \sigma_0^2)$ -distributed samples with known variance $\sigma_0^2 > 0$.

Proposition 8.6.10. *Let $x = (x_1, \dots, x_n)$ be a sample of n independent $\mathcal{N}(\mu, \sigma_0^2)$ distributed numbers. Here $\mu \in \mathbb{R}$ is unknown while σ_0^2 is known. Then with probability greater than*

$1 - \alpha$ the observed sample $x \in \mathbb{R}^n$ leads to¹¹

$$\bar{x} - \frac{\sigma_0}{\sqrt{n}} z_{1-\alpha/2} \leq \mu \leq \bar{x} + \frac{\sigma_0}{\sqrt{n}} z_{1-\alpha/2}.$$

Recall that z_β denotes the β -quantile of the standard normal distribution.

Proof. This is a direct consequence of applying Proposition 8.6.9 to the result presented in Example 8.6.8. \square

Example 8.6.11. Choose $\alpha = 0.05$ and suppose we observed the nine values

$$10.1, \quad 9.2, \quad 10.2, \quad 10.3, \quad 10.1, \quad 9.9, \quad 10.0, \quad 9.7, \quad 9.8,$$

then $\bar{x} = 9.9222$. The variance σ_0 is known to be $\sigma_0^2 = 0.330824$. That is, we assume that s_x^2 is the correct variance. Because of $z_{1-\alpha/2} = z_{0.975} = 1.95996$, with a confidence of 95 % our sample of nine numbers leads to

$$9.7061 \leq \mu \leq 10.1384.$$

In the next result we describe the confidence intervals generated by the t-test treated in Proposition 8.4.15.

Proposition 8.6.12. Let $x = (x_1, \dots, x_n)$ be a sample of n independent $\mathcal{N}(\mu, \sigma^2)$ distributed numbers. Here $\mu \in \mathbb{R}$ and σ^2 are both unknown. Then with probability greater than $1 - \alpha$ the observed sample $x \in \mathbb{R}^n$ leads to

$$\bar{x} - \frac{s_x}{\sqrt{n}} t_{n-1;1-\alpha/2} \leq \mu \leq \bar{x} + \frac{s_x}{\sqrt{n}} t_{n-1;1-\alpha/2}.$$

Recall that $t_{n-1;\beta}$ denotes the β -quantile of the Student t_{n-1} -distribution.

Proof. This is a direct consequence of

$$\mathcal{X}_0(\mu) = \left\{ x \in \mathbb{R}^n : \sqrt{n} \left| \frac{\bar{x} - \mu}{s_x} \right| \leq t_{n-1;1-\alpha/2} \right\},$$

hence Proposition 8.6.9 implies that, given $x \in \mathbb{R}^n$, then

$$\begin{aligned} C(x) &= \left\{ \mu \in \mathbb{R} : \sqrt{n} \left| \frac{\bar{x} - \mu}{s_x} \right| \leq t_{n-1;1-\alpha/2} \right\} \\ &= \left[\bar{x} - \frac{s_x}{\sqrt{n}} t_{n-1;1-\alpha/2}, \bar{x} + \frac{s_x}{\sqrt{n}} t_{n-1;1-\alpha/2} \right]. \end{aligned} \quad \square$$

¹¹ We want to point this out again: not μ is random but the chosen interval is such. And $1 - \alpha$ is the probability to observe an x for which the random interval contains the correct μ .

Example 8.6.13. Let us explain the previous result by the concrete sample investigated in Example 8.4.20. There we had $\bar{x} = 22.072$, $s_x = 0.07554248$, and $n = 10$. If $\alpha = 0.05$, the quantile of t_9 equals $t_{9;0.975} = 2.26$. From this, we derive $[22.016, 22.126]$ as the 95 % confidence interval.

Verbally this says that with a confidence of 95 % we observed those x_1, \dots, x_{10} for which $\mu \in C(x) = [22.016, 22.126]$.

Finally, let us construct $1 - \alpha$ confidence intervals for the unknown variance of a normal sample.

Proposition 8.6.14. *Let $x = (x_1, \dots, x_n)$ be a sample of n independent $\mathcal{N}(\mu, \sigma^2)$ distributed numbers. If μ is known, then with probability greater than $1 - \alpha$ the observed sample $x \in \mathbb{R}^n$ leads to*

$$\frac{n\sigma_x^2}{\chi_{n;1-\alpha/2}^2} \leq \sigma^2 \leq \frac{n\sigma_x^2}{\chi_{n;\alpha/2}^2}$$

where, in contrast to eq. (8.14),

$$\sigma_x^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \mu)^2. \quad (8.52)$$

If, on the contrary, μ is unknown, then with probability greater than $1 - \alpha$ the observed sample $x \in \mathbb{R}^n$ gives

$$\frac{(n-1)s_x^2}{\chi_{n-1;1-\alpha/2}^2} \leq \sigma^2 \leq \frac{(n-1)s_x^2}{\chi_{n-1;\alpha/2}^2}.$$

Here s_x^2 is as in Definition 8.4.1 and $\chi_{n;\beta}^2$ denotes the β -quantile of a χ_n^2 distribution.

Proof. Both assertions easily follow from eqs. (8.27) and (8.28). Recall that they imply for known mean value μ that

$$\mathcal{X}_0(\sigma^2) = \left\{ x \in \mathbb{R}^n : \chi_{n;\alpha/2}^2 \leq \sum_{j=1}^n \frac{(x_j - \mu)^2}{\sigma^2} \leq \chi_{n;1-\alpha/2}^2 \right\}$$

while in the case of unknown μ one has

$$\mathcal{X}_0(\sigma^2) = \left\{ x \in \mathbb{R}^n : \chi_{n-1;\alpha/2}^2 \leq (n-1) \frac{s_x^2}{\sigma^2} \leq \chi_{n-1;1-\alpha/2}^2 \right\}.$$

Letting in both cases $C(x) = \{\sigma^2 > 0 : x \in \mathcal{X}_0(\sigma^2)\}$ completes the proof. \square

Example 8.6.15. A measuring instrument possesses an unknown precision. In order to get some information about it, we measure a certain item 9 times. The obtained re-

sults are

10.1, 10.3, 10.2, 10.7, 9.9, 10.0, 10.9, 8.9 and 11.0.

Let $x = (x_1, \dots, x_9)$ be the collection of these measurements. Then the mean value equals $\bar{x} = 10.222$, hence the unbiased variance of the observed sample is calculated by

$$s_x^2 = \frac{1}{8} \sum_{j=1}^9 (x_j - \bar{x})^2 = 0.4019.$$

If we choose $\alpha = 0.1$ as a significance level, then, in order to determine a 90 % confidence interval, we need the 0.05 and 0.95 quantiles of a χ_8^2 distribution. They are

$$\chi_{8,0.05}^2 = 2.73264 \quad \text{and} \quad \chi_{8,0.95}^2 = 15.5073.$$

So we finally obtain that there is chance of 90 % that the unknown variance σ^2 of the measuring instrument satisfies

$$\frac{8 \cdot 0.4019}{15.5073} = 0.207334 \leq \sigma^2 \leq \frac{8 \cdot 0.4019}{2.73264} = 1.17659.$$

Note that another n measurements by this instrument, of this or of a different item, will surely lead to different bounds for σ^2 .

To conclude this section, let us summarize the most important confidence intervals for normally distributed samples. Here σ_x^2 and s_x^2 are defined by eqs. (8.52) and (8.14), respectively.

Name	Parameters	$1 - \alpha$ Confidence Intervals
Confidence intervals for the mean value	$\sigma^2 > 0$ known	$\left[\bar{x} - \frac{\sigma}{\sqrt{n}} z_{1-\alpha/2}, \bar{x} + \frac{\sigma}{\sqrt{n}} z_{1-\alpha/2} \right]$
Confidence intervals for the mean value	$\sigma^2 > 0$ unknown	$\left[\bar{x} - \frac{s_x}{\sqrt{n}} t_{n-1;1-\alpha/2}, \bar{x} + \frac{s_x}{\sqrt{n}} t_{n-1;1-\alpha/2} \right]$
Confidence intervals for the variance	$\mu \in \mathbb{R}$ known	$\left[\frac{n \sigma_x^2}{\chi_{n;1-\alpha/2}^2}, \frac{n \sigma_x^2}{\chi_{n;\alpha/2}^2} \right]$
Confidence intervals for the variance	$\mu \in \mathbb{R}$ unknown	$\left[\frac{(n-1) s_x^2}{\chi_{n-1;1-\alpha/2}^2}, \frac{(n-1) s_x^2}{\chi_{n-1;\alpha/2}^2} \right]$

8.6.3 Binomial distributed populations

The aim of this section is to show how Proposition 8.6.9 applies in the case of binomial distributed populations. Thus, we execute n independent trials where every time occurs either success or failure. Hence, the number of successes is $B_{n,\theta}$ -distributed for a certain

$\theta \in [0, 1]$. But, in contrast to the investigations in Section 1.4.3, now the parameter θ is unknown.

Say we observed $0 \leq k \leq n$ times success. Then we look for a confidence interval $C(k) \subseteq [0, 1]$ such that very likely $\theta \in C(k)$. More precisely, given $0 < \alpha < 1$, we want that for any parameter $\theta \in [0, 1]$,

$$B_{n,\theta}\{k \leq n : \theta \in C(k)\} \geq 1 - \alpha.$$

That is, with probability greater than $1 - \alpha$ the observation of k successes leads to an interval $C(k)$ containing the correct parameter θ .

Proposition 8.6.16. *The statistical model is $(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{0 \leq \theta \leq 1}$ where the sample space is $\mathcal{X} = \{0, \dots, n\}$. Given $\alpha > 0$ and $k = 0, \dots, n$, define sets $C(k) \subseteq [0, 1]$ as follows:*

$$C(k) = \{\theta : B_{n,\theta}(\{0, \dots, k\}) > \alpha/2\} \cap \{\theta : B_{n,\theta}(\{k, \dots, n\}) > \alpha/2\}. \quad (8.53)$$

Then for each $k \leq n$, the set $C(k)$ is an $1 - \alpha$ confidence interval for $\theta \in [0, 1]$.

Proof. In order to get these confidence regions, we use Proposition 8.6.9. As shown in Proposition 8.3.1, the region of acceptance $\mathcal{X}_0(\theta)$ of an α -significance test $\mathbf{T}(\theta)$, where $\mathbb{H}_0 : \vartheta = \theta$, is given by

$$\mathcal{X}_0(\theta) = \{n_0(\theta), \dots, n_1(\theta)\}.$$

Here, the numbers $n_0(\theta)$ and $n_1(\theta)$ were defined by

$$n_0(\theta) := \min \left\{ k \leq n : \sum_{j=0}^k \binom{n}{j} \theta^j (1-\theta)^{n-j} > \alpha/2 \right\}$$

and

$$n_1(\theta) := \max \left\{ k \leq n : \sum_{j=k}^n \binom{n}{j} \theta^j (1-\theta)^{n-j} > \alpha/2 \right\}.$$

Applying Proposition 8.6.9, the sets

$$C(k) := \{\theta \in [0, 1] : k \in \mathcal{X}_0(\theta)\} = \{\theta \in [0, 1] : n_0(\theta) \leq k \leq n_1(\theta)\}, \quad k = 0, \dots, n,$$

are $1 - \alpha$ confidence regions. By the definition of $n_0(\theta)$ and $n_1(\theta)$, given $k \leq n$, then a number $\theta \in [0, 1]$ satisfies $n_0(\theta) \leq k \leq n_1(\theta)$ if and only if at the same time

$$B_{n,\theta}(\{0, \dots, k\}) = \sum_{j=0}^k \binom{n}{j} \theta^j (1-\theta)^{n-j} > \alpha/2 \quad \text{and}$$

$$B_{n,\theta}(\{k, \dots, n\}) = \sum_{j=k}^n \binom{n}{j} \theta^j (1-\theta)^{n-j} > \alpha/2.$$

Thus, as claimed,

$$C(k) = \{\theta : B_{n,\theta}(\{0, \dots, k\}) > \alpha/2\} \cap \{\theta : B_{n,\theta}(\{k, \dots, n\}) > \alpha/2\}$$

are $1 - \alpha$ confidence sets.

It remains to prove that these are indeed intervals, which by the definition of the sets $C(k)$ is not so obvious. We know by Lemma 8.3.6 that the function

$$f_{\{\geq k\}} : \theta \mapsto \sum_{j=k}^n \binom{n}{j} \theta^j (1-\theta)^{n-j} \quad (8.54)$$

is nondecreasing for any $k \geq 0$. Hence,

$$f_{\{\leq k\}} : \theta \mapsto \sum_{j=0}^k \binom{n}{j} \theta^j (1-\theta)^{n-j} = 1 - \sum_{j=k+1}^n \binom{n}{j} \theta^j (1-\theta)^{n-j} \quad (8.55)$$

is nonincreasing. Letting

$$\theta_k^- = \inf \left\{ \theta : f_{\{\geq k\}}(\theta) > \frac{\alpha}{2} \right\} \quad \text{and} \quad \theta_k^+ = \sup \left\{ \theta : f_{\{\leq k\}}(\theta) > \frac{\alpha}{2} \right\},$$

it follows that $C(k) = (\theta_k^-, \theta_k^+)$. This completes the proof. \square

Remark 8.6.17. The intervals $C(k) = (\theta_k^-, \theta_k^+)$ are usually called **100(1 - α)% Clopper–Pearson intervals** or also **exact confidence intervals** for the binomial distribution.

Since the functions $f_{\{\leq k\}}$ and $f_{\{\geq k\}}$ are continuous on $[0, 1]$, in the case $1 < k < n$, the numbers θ_k^- and θ_k^+ are also characterized by

$$f_{\{\geq k\}}(\theta_k^-) = \frac{\alpha}{2} \quad \text{and} \quad f_{\{\leq k\}}(\theta_k^+) = \frac{\alpha}{2}.$$

In other words, if $1 < k < n$, then the endpoints of the Clopper–Pearson intervals are the unique solution $\theta \in (0, 1)$ of

$$\sum_{j=k}^n \binom{n}{j} \theta^j (1-\theta)^{n-j} = \frac{\alpha}{2} \quad \text{or} \quad \sum_{j=0}^k \binom{n}{j} \theta^j (1-\theta)^{n-j} = \frac{\alpha}{2},$$

respectively. For the cases $k = 0$ and $k = n$, we refer to Problem 8.6.

Example 8.6.18. Suppose we execute $n = 20$ trials and observe $k = 5$ successes. What can be said about the underlying success probability θ ? The functions in eqs. (8.54) and (8.55) are in this case given by

$$f_{\{\geq 5\}}(\theta) = \sum_{j=5}^{20} \binom{20}{j} \theta^j (1-\theta)^{20-j} \quad \text{and} \quad f_{\{\leq 5\}}(\theta) = \sum_{j=0}^5 \binom{20}{j} \theta^j (1-\theta)^{20-j}.$$

If $\alpha = 0.1$, then numerical calculations lead to (see also Fig. 8.10)

$$\theta_5^- = \inf\{\theta : f_{\{\geq 5\}}(\theta) > 0.05\} \approx 0.1041 \quad \text{and}$$

$$\theta_5^+ = \sup\{\theta : f_{\{\leq 5\}}(\theta) > 0.05\} \approx 0.4566.$$

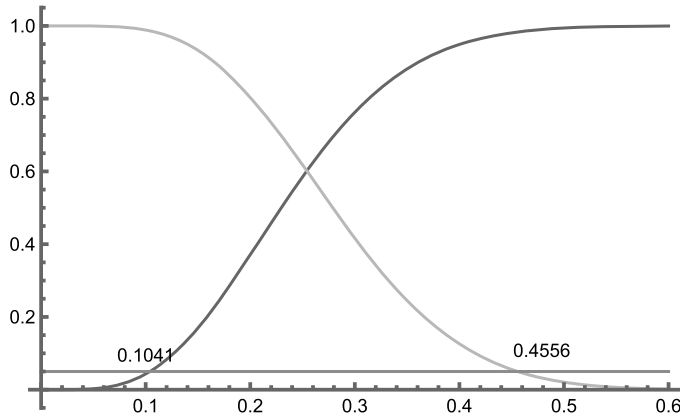


Figure 8.10: The increasing function $f_{\{\geq 5\}}$ and the decreasing one $f_{\{\leq 5\}}$, both taken with for $n = 20$. The horizontal line marks the significance level $\alpha/2 = 0.05$.

Consequently, with the probability of 90 % our observation of five successes implies that the underlying success probability θ satisfies

$$0.1041 < \theta < 0.4566. \quad (8.56)$$

The estimates for θ obtained in (8.56) are quite rough. This is mainly due to the fact that the number $n = 20$ of trials is pretty small. Also different numbers of success do not yield significantly tighter bounds. So, for example, if one observes 10 successes, then as the 90 % confidence interval one gets

$$0.30196 < \theta < 0.69804.$$

The next example provides sharper bounds for larger $n \geq 1$.

Example 8.6.19. In an urn there are white and black balls with an unknown proportion θ of white balls. In order to get some information about θ , we choose randomly 500 balls with replacement. Say 220 of the chosen balls are white. What is the 90 % confidence interval for θ based on this observation?

Answer: We have $n = 500$ and observed $k = 220$ white balls. Consequently, a 90 % confidence interval $C(220)$ consists of those $\theta \in [0, 1]$ for which at the same time

$$f_{\{\geq 220\}}(\theta) > \alpha/2 = 0.05 \quad \text{and} \quad f_{\{\leq 220\}}(\theta) > \alpha/2 = 0.05.$$

Note that in this case

$$f_{\{\leq 220\}}(\theta) = \sum_{j=0}^{220} \binom{500}{j} \theta^j (1-\theta)^{500-j} \quad \text{and}$$

$$f_{\{\geq 220\}}(\theta) = \sum_{j=220}^{500} \binom{500}{j} \theta^j (1-\theta)^{500-j}.$$

Numerical calculations tell us that

$$\theta_{220}^- \approx 0.4028 \quad \text{and} \quad \theta_{220}^+ \approx 0.4777.$$

Therefore, a 90 % confidence interval $C(220)$ is given by

$$C(220) = (0.4028, 0.4777).$$

For $n = 1000$ and 440 observed white balls, the calculations lead to the smaller, hence more significant, interval $C(440) = (0.4139, 0.4664)$.

Remark 8.6.20. The previous example already indicates that the determination of the Clopper–Pearson intervals becomes quite complicated for large n . Therefore, one looks for “approximate” intervals. Background for the construction is the central limit theorem in the form presented in Proposition 7.2.19. For S_n s distributed according to $B_{n,\theta}$, it implies

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \left| \frac{S_n - n\theta}{\sqrt{n\theta(1-\theta)}} \right| \leq z_{1-\alpha/2} \right\} = 1 - \alpha,$$

or, equivalently,

$$\lim_{n \rightarrow \infty} B_{n,\theta} \left\{ k \leq n : \left| \frac{k - n\theta}{\sqrt{n\theta(1-\theta)}} \right| \leq z_{1-\alpha/2} \right\} = 1 - \alpha.$$

Here $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile introduced in Definition 8.4.8. Thus, an “approximate” region of acceptance, testing the hypothesis “the unknown parameter is θ ,” is given by

$$\mathcal{X}_0(\theta) = \left\{ k \leq n : \left| \frac{k}{n} - \theta \right| \leq z_{1-\alpha/2} \sqrt{\frac{\theta(1-\theta)}{n}} \right\}. \quad (8.57)$$

An application of Proposition 8.6.9 leads to certain confidence regions, but mostly they cannot be described explicitly. Due to the term $\sqrt{\theta(1-\theta)}$ on the right-hand side of eq. (8.57), it is not possible, for a given $k \leq n$, to determine those θ s for which $k \in \mathcal{X}_0(\theta)$. To overcome this difficulty, we change $\mathcal{X}_0(\theta)$ yet again by replacing θ on the right-hand side by its MLE $\hat{\theta}(k) = \frac{k}{n}$. That is, we replace eq. (8.57) by

$$\tilde{\mathcal{X}}_0(\theta) = \left\{ k \leq n : \left| \frac{k}{n} - \theta \right| \leq z_{1-\alpha/2} \sqrt{\frac{\frac{k}{n}(1-\frac{k}{n})}{n}} \right\}.$$

Doing so, an application of Proposition 8.6.9 leads to the “approximate” confidence intervals $\tilde{C}(k)$, $k = 0, \dots, n$, defined as

$$\tilde{C}(k) = \left[\frac{k}{n} - z_{1-\alpha/2} \sqrt{\frac{\frac{k}{n}(1-\frac{k}{n})}{n}}, \frac{k}{n} + z_{1-\alpha/2} \sqrt{\frac{\frac{k}{n}(1-\frac{k}{n})}{n}} \right]. \quad (8.58)$$

Example 8.6.21. We investigate once more Example 8.6.19. Among 500 chosen balls we observed 220 white. This observation led to the “exact” 90 % confidence interval $C(220) = (0.4028, 0.4777)$.

Let us compare this result with the interval we get by using the approximative approach. Since the quantile $z_{1-\alpha/2}$ for $\alpha = 0.1$ equals $z_{0.95} = 1.64485$, the left and the right endpoints of the interval (8.58) with $k = 220$ are evaluated by

$$\begin{aligned} \frac{220}{500} - 1.64485 \cdot \sqrt{\frac{220 \cdot 280}{500^3}} &= 0.4035 \quad \text{and} \\ \frac{220}{500} + 1.64485 \cdot \sqrt{\frac{220 \cdot 280}{500^3}} &= 0.4765. \end{aligned}$$

Thus, the “approximate” 90 % confidence interval is $\tilde{C}(220) = (0.4035, 0.4765)$, which does not differ too much from $C(220) = (0.4028, 0.4777)$.

In the case of 1000 trials and 440 white balls, the endpoints of a confidence interval are evaluated by

$$\begin{aligned} \frac{440}{1000} - 1.64485 \cdot \sqrt{\frac{440 \cdot 560}{1000^3}} &= 0.414181 \quad \text{and} \\ \frac{440}{1000} + 1.64485 \cdot \sqrt{\frac{440 \cdot 560}{1000^3}} &= 0.4645819. \end{aligned}$$

That is, $\tilde{C}(440) = (0.4142, 0.4659)$ compared with $C(440) = (0.4139, 0.4664)$.

Example 8.6.22. A few days before an election, 1000 randomly chosen people are questioned whom they will vote for next week, either candidate *A* or candidate *B*. Suppose 540 of the interviewed people answered that they would vote for candidate *A*, the remaining 460 favored candidate *B*. Find a 90 % confidence interval for the expected result of candidate *A* in the election.

Solution: We have $n = 1000$, $k = 540$, and $\alpha = 0.1$. The quantile of level 0.95 of the standard normal distribution equals $z_{0.95} = 1.64485$ (see Example 8.6.21). This leads to $[0.514, 0.566]$ as “approximate” 90 % confidence interval for the expected result of candidate *A*.

If one questions another 1000 randomly chosen people, another confidence interval will occur. But, on average, in 9 of 10 cases questioning 1000 people will lead to an interval containing the correct value.

Summary: The statistical model is $(\mathcal{X}, \mathcal{P}(\mathcal{X}), B_{n,\theta})_{0 \leq \theta \leq 1}$ with $\mathcal{X} = \{0, \dots, n\}$. Given $\alpha > 0$ and $k = 0, \dots, n$, define numbers θ_k^- and θ_k^+ by

$$\theta_k^- = \inf \left\{ \theta \in [0, 1] : \sum_{j=k}^n \binom{n}{j} \theta^j (1-\theta)^{n-j} > \frac{\alpha}{2} \right\},$$

$$\theta_k^+ = \sup \left\{ \theta \in [0, 1] : \sum_{j=0}^k \binom{n}{j} \theta^j (1-\theta)^{n-j} > \frac{\alpha}{2} \right\}.$$

Letting $C(k) = [\theta_k^-, \theta_k^+]$, the $C(k)$ s, $0 \leq k \leq n$, are $1 - \alpha$ confidence intervals for the unknown $\theta \in [0, 1]$. That is, for all $\theta \in [0, 1]$,

$$B_{n,\theta} \{k \leq n : \theta \in C(k)\} \geq 1 - \alpha.$$

For large $n \geq 1$, the central limit theorem applies and leads to the approximate $1 - \alpha$ confidence intervals

$$\tilde{C}(k) = \left[\frac{k}{n} - z_{1-\alpha/2} \sqrt{\frac{\frac{k}{n}(1-\frac{k}{n})}{n}}, \frac{k}{n} + z_{1-\alpha/2} \sqrt{\frac{\frac{k}{n}(1-\frac{k}{n})}{n}} \right].$$

8.6.4 Hypergeometric distributed populations

Finally, we construct confidence intervals for *hypergeometric distributed populations*. Since the technique is quite similar to that used in the case of binomial distribution, we will not state all details. We refer everybody interested in them to Problem 8.7, to add missing details in the construction.

So assume that in an urn there are N balls, colored white and black. Among them there are $M \leq N$ white balls, hence $N - M$ are black. Hereby, the number M is unknown. We choose now at random $n \leq N$ balls. Denote by $m \leq n$ the number of chosen white balls. The aim is to determine for each $m \leq n$ a set $C(m) \subseteq \{0, \dots, N\}$ so that

$$H_{N,M,n} \{m \leq n : M \in C(m)\} > 1 - \alpha, \quad 0 \leq M \leq N.$$

Of course, thereby the confidence sets should be chosen as small as possible.

Before we can introduce confidence intervals for hypergeometric distributed samples, we first have to extend Proposition 8.2.19 to the case of two-sided tests. Since the proof of this two-sided case is very similar to that for binomial distributed samples in Proposition 8.3.1, we omit it. Furthermore, we formulate the two-sided tests already in a way appropriate for later use.

Proposition 8.6.23. *Let the statistical model be $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{N,K,n})_{K=0,\dots,N}$ with the sample space $\mathcal{X} = \{0, \dots, n\}$. Given an arbitrary $M \leq N$, an α -test for testing $K = M$ against $K \neq M$ is given by $\mathbf{T}(M) = (\mathcal{X}_0(M), \mathcal{X}_1(M))$ where the region of acceptance equals $\mathcal{X}_0(M) = \{m_0(M), \dots, m_1(M)\}$ with $m_0(M)$ and $m_1(M)$ defined by*

$$m_0(M) = \min \left\{ k \leq n : \sum_{m=0}^k \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}} > \alpha/2 \right\} \quad \text{and}$$

$$m_1(M) := \max \left\{ k \leq n : \sum_{m=k}^n \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}} > \alpha/2 \right\}.$$

Before proceeding further, let us introduce two functions similar to those in eqs. (8.54) and (8.55). For each $k = 0, \dots, n$ and $M \leq N$, set

$$f_{\{\geq k\}}(M) = \sum_{m=k}^n \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}} = H_{N,M,n}(\{k, \dots, n\})$$

and

$$f_{\{\leq k\}}(M) = \sum_{m=0}^k \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}} = H_{N,M,n}(\{0, \dots, k\}).$$

Lemma 8.2.20 implies that $f_{\{\geq k\}}$ is nondecreasing, hence $f_{\{\leq k\}}$ is nonincreasing. Verbally said, if there are M white balls in the urn, then $f_{\{\geq k\}}(M)$ is the probability to observe at least k white in the sample of size n , while $f_{\{\leq k\}}(M)$ tells us how likely it is for us to get k or less white balls.

Given $k = 0, \dots, n$, define two numbers M_k^- and M_k^+ (depending on α) by

$$M_k^- = \min \left\{ M : f_{\{\geq k\}}(M) > \frac{\alpha}{2} \right\} \quad \text{and} \quad M_k^+ = \max \left\{ M : f_{\{\leq k\}}(M) > \frac{\alpha}{2} \right\}.$$

Compare Figure 8.11 for an example with $N = 60$, $n = 15$, and $\alpha = 0.1$.

Proposition 8.6.24. *Let the statistical model be $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{N,M,n})_{M=0,\dots,N}$ with the sample space $\mathcal{X} = \{0, \dots, n\}$. Given an arbitrary $\alpha > 0$, for each $0 \leq k \leq n$, the sets*

$$C(k) = [M_k^-, M_k^+] = \left\{ M \leq N : f_{\{\geq k\}}(M) > \frac{\alpha}{2} \right\} \cap \left\{ M \leq N : f_{\{\leq k\}}(M) > \frac{\alpha}{2} \right\}$$

are $1 - \alpha$ confidential intervals for the unknown parameter $M \leq N$.

Proof. Applying Propositions 8.6.9 and 8.6.23, for each $k \leq n$, the sets

$$\{M \leq N : k \in \mathcal{X}_0(M)\} = \{M \leq N : m_0(M) \leq k \leq m_1(M)\}, \quad k = 0, \dots, n,$$

are $1 - \alpha$ confidential sets. But by the definition of $m_0(M)$ and $m_1(M)$, we have $m_0(M) \leq k \leq m_1(M)$ if and only if

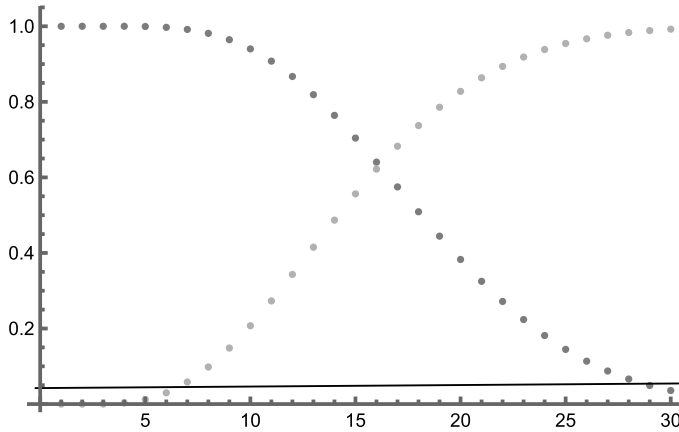


Figure 8.11: The increasing function $f_{\{\geq 4\}}$ and the decreasing $f_{\{\leq 4\}}$ in the case $N = 60$ and sample size $n = 15$. The horizontal line marks the significance level $\alpha/2 = 0.05$.

$$f_{\{\geq k\}}(M) > \frac{\alpha}{2} \quad \text{and} \quad f_{\{\leq k\}}(M) > \frac{\alpha}{2}.$$

This proves that

$$C(k) := \left\{ M \leq N : f_{\{\geq k\}}(M) > \frac{\alpha}{2} \right\} \cap \left\{ M \leq N : f_{\{\leq k\}}(M) > \frac{\alpha}{2} \right\}$$

are $1 - \alpha$ confidence sets. Finally, since $f_{\{\geq k\}}$ and $f_{\{\leq k\}}$ are monotone, we observe that $M_k^- \leq M \leq M_k^+$ if and only if $f_{\{\geq k\}}(M) > \frac{\alpha}{2}$ and $f_{\{\leq k\}}(M) > \frac{\alpha}{2}$. Hence, $C(k) = [M_k^-, M_k^+]$, which completes the proof. \square

Example 8.6.25. In an urn there are 200 balls. We choose randomly a sample of 60 balls. Among them 25 are white. What can be said about the number of white balls in the urn?

Answer: We are asking for a 90 % confidence set for the unknown number M of white balls in the urn. In this case the functions of interest are

$$f_{\{\leq 25\}}(M) = \sum_{m=0}^{25} \frac{\binom{M}{m} \binom{200-M}{60-m}}{\binom{200}{60}} \quad \text{and} \quad f_{\{\geq 25\}}(M) = \sum_{m=25}^{60} \frac{\binom{M}{m} \binom{200-M}{60-m}}{\binom{200}{60}}.$$

Since

$$f_{\{\geq 25\}}(65) = 0.0508, \quad f_{\{\geq 25\}}(64) = 0.0408, \quad \text{and} \\ f_{\{\leq 25\}}(103) = 0.048, \quad f_{\{\leq 25\}}(102) = 0.0507,$$

we conclude that $M_k^- = 65$ and $M_k^+ = 102$, hence $C(25) = [65, 102]$ is a 90 % confidence interval for the unknown number M of white balls. If one asks for sharper bounds, one has to relax the significance level. So, for example, as 80 % confidence interval one gets $[68, 98]$.

Note that Proposition 8.5.16 gives in this case $\hat{M}(25) = 83$ as a point estimator which is in the middle of both confidence intervals.

If there are only 10 white balls among the chosen 60, the 90 % confidence region is $[20, 50]$.

Remark 8.6.26. Proposition 1.4.39 shows the tight connection between the binomial and hypergeometric distributions. Hence, it could be of interest what happens if we replace in the preceding example the hypergeometric distribution by the binomial. To do so, the unknown success probability equals $M/200$. In other words, instead taking 60 balls without replacement we choose now 60 balls and replace every time the chosen ball. Applying eq. (8.58) with $n = 60$ and $k = 25$, our observation of 25 white balls leads to the 90 %-sure estimates

$$0.311977 < \frac{M}{200} < 0.521356.$$

Compare this with the hypergeometric case where we got the 90 %-sure estimates

$$0.325 = \frac{65}{200} < \frac{M}{200} < \frac{102}{200} = 0.51.$$

Thus, in this case nonreplacing of chosen balls leads to sharper bounds for the unknown number of white balls than the bounds one gets in the case of replacing. Of course, one has to assume that in both cases the number of chosen white balls coincides.

Summary: The statistical model is $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{N,M,n})_{M=0,\dots,N}$ with $\mathcal{X} = \{0, \dots, n\}$. If

$$M_k^- = \min \left\{ 0 \leq M \leq N : \sum_{m=k}^n \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}} > \frac{\alpha}{2} \right\},$$

$$M_k^+ = \max \left\{ 0 \leq M \leq N : \sum_{m=0}^k \frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}} > \frac{\alpha}{2} \right\},$$

then for $k = 0, \dots, n$ the sets $C(k) = [M_k^-, M_k^+]$ are $1 - \alpha$ confidence intervals for the unknown parameter $M \leq N$.

8.7 Problems

Problem 8.1. For some $b > 0$, let \mathbb{P}_b be the uniform distribution on $[0, b]$. The precise value of $b > 0$ is unknown. We claim that $b \leq b_0$ for a certain $b_0 > 0$. Thus, the hypotheses are

$$\mathbb{H}_0 : b \leq b_0 \quad \text{and} \quad \mathbb{H}_1 : b > b_0.$$

To test \mathbb{H}_0 , we chose randomly n numbers x_1, \dots, x_n distributed according to \mathbb{P}_b . Suppose the region of acceptance \mathcal{X}_0 of a hypothesis test \mathbf{T}_c is given by

$$\mathcal{X}_0 := \{(x_1, \dots, x_n) : \max_{1 \leq i \leq n} x_i \leq c\}$$

for some $c > 0$.

1. Determine those $c > 0$ for which \mathbf{T}_c is an α -significance test of level $\alpha < 1$.
2. Suppose \mathbf{T}_c is an α -significance test. For which of those $c > 0$ does the probability for a type II error become minimal?
3. Determine the power function of the α -test \mathbf{T}_c that minimizes the probability of the occurrence of a type II error.

Problem 8.2. For $\theta > 0$, let \mathbb{P}_θ be the probability measure with density p_θ defined by

$$p_\theta(s) = \begin{cases} \theta s^{\theta-1} & \text{if } s \in (0, 1], \\ 0 & \text{otherwise.} \end{cases}$$

1. Check whether the p_θ s are probability density functions.
2. In order to get information about θ , we execute n independent trials according to \mathbb{P}_θ . Which statistical model describes this experiment?
3. Find the maximum likelihood estimator for θ .

Problem 8.3. The lifetime of light bulbs is exponentially distributed with unknown parameter $\lambda > 0$. In order to determine λ , we switch on n light bulbs and record the number of light bulbs that burn out until a certain time $T > 0$. Determine a statistical model that describes this experiment. Find the MLE for λ .

Problem 8.4. Consider the statistical model in Example 8.5.40, that is, we have $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathbb{P}_b^{\otimes n})_{b>0}$ with uniform distribution \mathbb{P}_b on $[0, b]$. There are two natural estimators for $b > 0$, namely \hat{b}_1 and \hat{b}_2 , defined by

$$\hat{b}_1(x) := \frac{n+1}{n} \max_{1 \leq i \leq n} x_i \quad \text{and} \quad \hat{b}_2(x) := \frac{2}{n} \sum_{i=1}^n x_i, \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n.$$

Prove that \hat{b}_1 and \hat{b}_2 possess the following properties:

1. The estimators \hat{b}_1 and \hat{b}_2 are unbiased.
2. One has

$$\mathbb{V}_b \hat{b}_1 = \frac{b^2}{n(n+2)} \quad \text{and} \quad \mathbb{V}_b \hat{b}_2 = \frac{b^2}{3n^2}.$$

Problem 8.5. In a questionnaire, out of 2000 randomly chosen people 1420 answered that they regularly use the Internet. Find an “approximate” 90 % confidence interval for the proportion of people using the Internet regularly. Determine the inequalities that describe the exact intervals in eq. (8.53).

Problem 8.6. How do the confidence intervals $C(k)$ in eq. (8.53) look like for $k = 0$ and $k = n$? Is it possible that for some $k = 0, \dots, n$ and $n \geq 1$ it follows that $0 \in C(k)$ or $1 \in C(k)$? If so, when does this happen?

Problem 8.7. Suppose the statistical model is $(\mathcal{X}, \mathcal{P}(\mathcal{X}), H_{N,M,n})_{M \leq N}$ with the sample space $\mathcal{X} = \{0, \dots, n\}$ and with the hypergeometric distributions $H_{N,M,n}$ introduced in Definition 1.4.32.

1. For some $M_0 \leq M$, the hypotheses are $\mathbb{H}_0 : M = M_0$ against $\mathbb{H}_1 : M \neq M_0$. Find (optimal) numbers $0 \leq m_0 \leq m_1 \leq n$ such that $\mathcal{X}_0 = \{m_0, \dots, m_1\}$ is the region of acceptance of an α -significance test \mathbf{T} for \mathbb{H}_0 against \mathbb{H}_1 .

Hint: Modify the methods developed in Proposition 8.2.19 and compare the construction of two-sided tests for a binomial distributed population.

2. Use Proposition 8.6.9 to derive from \mathcal{X}_0 confidence intervals $C(k)$, $0 \leq k \leq n$, of level $1 - \alpha$ for the unknown parameter M .

Hint: Follow the methods in Proposition 8.6.16 for the binomial distribution.

A Appendix

A.1 Notations

Throughout the book, we use the following standard notations:

1. The **natural numbers** starting at 1 are always denoted by \mathbb{N} . In the case 0 is included, we write \mathbb{N}_0 .
2. As usual, the **integers** \mathbb{Z} are given by $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$.
3. By \mathbb{R} we denote the field of **real numbers** endowed with the usual algebraic operations and its natural order. The subset $\mathbb{Q} \subset \mathbb{R}$ is the union of all **rational numbers**, that is, of numbers m/n where $m, n \in \mathbb{Z}$ and $n \neq 0$.
4. Given $n \geq 1$, let \mathbb{R}^n be the **n -dimensional Euclidean vector space**, that is,

$$\mathbb{R}^n = \{x = (x_1, \dots, x_n) : x_j \in \mathbb{R}\}.$$

Addition and scalar multiplication in \mathbb{R}^n are carried out coordinate-wise,

$$x + y = (x_1, \dots, x_n) + (y_1, \dots, y_n) = (x_1 + y_1, \dots, x_n + y_n)$$

and if $\alpha \in \mathbb{R}$, then

$$\alpha x = (\alpha x_1, \dots, \alpha x_n).$$

A.2 Elements of set theory

A.2.1 Set operations

Given a set M , its **powerset** $\mathcal{P}(M)$ consists of all subsets of M . In the case that M is finite, we have $|\mathcal{P}(M)| = 2^{|M|}$, where $|A|$ denotes the **cardinality** (number of elements) of a finite set A .

If A and B are subsets of M , written as $A, B \subseteq M$ or also as $A, B \in \mathcal{P}(M)$, their **union** and their **intersection** are, as usual, defined by (compare Figure A.1)

$$A \cup B = \{x \in M : x \in A \text{ or } x \in B\} \quad \text{and} \quad A \cap B = \{x \in M : x \in A \text{ and } x \in B\}.$$

Of course, it always holds that

$$A \cap B \subseteq A \subseteq A \cup B \quad \text{and} \quad A \cap B \subseteq B \subseteq A \cup B.$$

In the same way, given subsets A_1, A_2, \dots of M their union $\bigcup_{j=1}^{\infty} A_j$ and intersection $\bigcap_{j=1}^{\infty} A_j$ is the set of those $x \in M$ that belong to at least one of the A_j or that belong to all A_j , respectively.

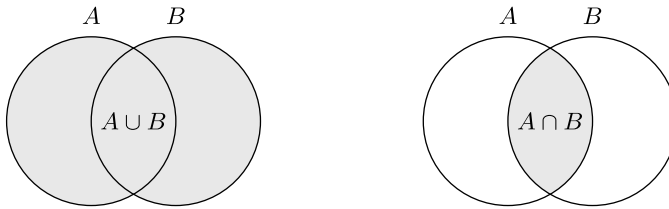


Figure A.1: The Venn diagrams of the union $A \cup B$ and the intersection $A \cap B$.

Quite often we use the **distributive law** for intersection and union. This asserts

$$A \cap \left(\bigcup_{j=1}^{\infty} B_j \right) = \bigcup_{j=1}^{\infty} (A \cap B_j).$$

Two sets A and B are said to be **disjoint**¹ provided that $A \cap B = \emptyset$. A sequence of sets A_1, A_2, \dots is called disjoint² whenever $A_i \cap A_j = \emptyset$ if $i \neq j$.

An element $x \in M$ belongs to the **set difference** $A \setminus B$ provided that $x \in A$ but $x \notin B$. Using the notion of the **complementary set** $B^c := \{x \in M : x \notin B\}$, the set difference may also be written as (compare Figure A.2)

$$A \setminus B = A \cap B^c.$$

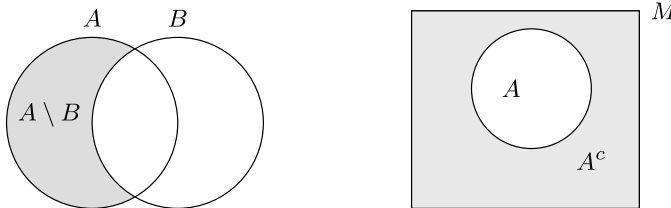


Figure A.2: The Venn diagrams of the set difference $A \setminus B$ and of the complement A^c of A with respect to a superset M .

Another useful identity is

$$A \setminus B = A \setminus (A \cap B).$$

Conversely, the complementary set may be represented as the set difference $B^c = M \setminus B$. We still mention the obvious $(B^c)^c = B$.

¹ Sometimes called “mutually exclusive.”

² More precisely, one should say “pairwise disjoint.”

Finally, we introduce the **symmetric difference** $A\Delta B$ of two sets A and B as (see Figure A.3)

$$A\Delta B := (A \setminus B) \cup (B \setminus A) = (A \cap B^c) \cup (B \cap A^c) = (A \cup B) \setminus (A \cap B). \quad (\text{A.1})$$

Note that an element $x \in M$ belongs to $A\Delta B$ if and only if x belongs exactly to one of the sets A or B .

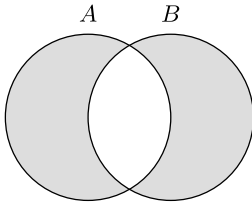


Figure A.3: The symmetric difference $A\Delta B$.

De Morgan's rules are very important and assert the following:

$$\left(\bigcup_{j=1}^{\infty} A_j \right)^c = \bigcap_{j=1}^{\infty} A_j^c \quad \text{and} \quad \left(\bigcap_{j=1}^{\infty} A_j \right)^c = \bigcup_{j=1}^{\infty} A_j^c.$$

Given sets A_1, \dots, A_n , their **Cartesian product** $A_1 \times \dots \times A_n$ is defined by (see Figure A.4 for an example of the Cartesian product)

$$A_1 \times \dots \times A_n := \{(a_1, \dots, a_n) : a_j \in A_j\}$$

with

$$(a_1, \dots, a_n) = (b_1, \dots, b_n) \Leftrightarrow a_1 = b_1, \dots, a_n = b_n.$$

Note that $|A_1 \times \dots \times A_n| = |A_1| \cdots |A_n|$.

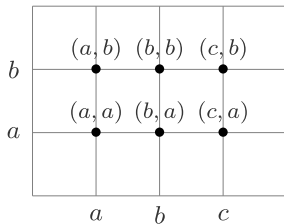


Figure A.4: The Cartesian product $\{a, b, c\} \times \{a, b\}$.

A.2.2 Preimages of sets

Let S be another set, for example, $S = \mathbb{R}$, and let $f : M \rightarrow S$ be some mapping from M to S . Given a subset $B \subseteq S$, we denote the **preimage** of B with respect to f by

$$f^{-1}(B) := \{x \in M : f(x) \in B\}. \quad (\text{A.2})$$

In other words, an element $x \in M$ belongs to $f^{-1}(B)$ if and only if its image with respect to f is an element of B (compare Figure A.5).

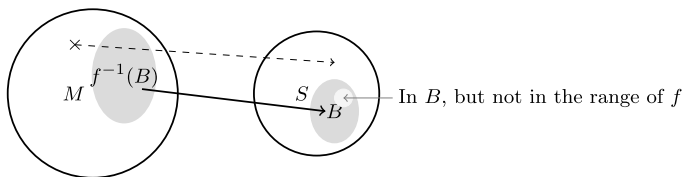


Figure A.5: Only elements from $f^{-1}(B)$ are mapped to B . All other elements in M have to have their image outside B . But not every element in B needs to be the image of an element in M .

We summarize some crucial properties of the preimage in a proposition.

Proposition A.2.1. *Let $f : M \rightarrow S$ be a mapping from M into another set S .*

- (1) $f^{-1}(\emptyset) = \emptyset$ and $f^{-1}(S) = M$.
- (2) For any subsets $B_j \subseteq S$, the following equalities are valid:

$$f^{-1}\left(\bigcup_{j \geq 1} B_j\right) = \bigcup_{j \geq 1} f^{-1}(B_j) \quad \text{and} \quad f^{-1}\left(\bigcap_{j \geq 1} B_j\right) = \bigcap_{j \geq 1} f^{-1}(B_j). \quad (\text{A.3})$$

Proof. We only prove the left-hand equality in eq. (A.3). The right-hand one is proved by the same methods. Furthermore, assertion (1) follows immediately.

Take $x \in f^{-1}(\bigcup_{j \geq 1} B_j)$. This happens if and only if

$$f(x) \in \bigcup_{j \geq 1} B_j \quad (\text{A.4})$$

is satisfied. But this is equivalent to the existence of a certain $j_0 \geq 1$ with $f(x) \in B_{j_0}$. By definition of the preimage, the last statement may be reformulated as follows: there exists a $j_0 \geq 1$ such that $x \in f^{-1}(B_{j_0})$. But this implies

$$x \in \bigcup_{j \geq 1} f^{-1}(B_j). \quad (\text{A.5})$$

Consequently, an element $x \in M$ satisfies condition (A.4) if and only if property (A.5) holds. This proves the left-hand identity in formulas (A.3). \square

A.2.3 Problems

Problem A.1. For any three sets A, B , and C , prove the following assertions:

$$A \setminus (B \cup C) = (A \setminus B) \cap (A \setminus C) \quad \text{and} \quad (A \cap B) \cup C = (A \cup C) \cap (B \cup C).$$

Problem A.2. For any sets $A, B \subseteq M$, prove that

$$(B \setminus A)^c \cap B = A \cap B \quad \text{and} \quad (A \cup B)^c \cap B = \emptyset.$$

Problem A.3. Let A and B two finite sets. Show that

$$|A \cup B| = |A| + |B| - |A \cap B|.$$

Let C be another finite set. Find a similar formula for $|A \cup B \cup C|$ and prove it.

Problem A.4. Prove that all three expressions in eq. (A.1) define the symmetric difference of two sets. That is, show that

$$(A \setminus B) \cup (B \setminus A) = (A \cap B^c) \cup (B \cap A^c) = (A \cup B) \setminus (A \cap B)$$

for all sets A and B .

Problem A.5. Let A, B , and C be three sets. Show that an element x belongs to $A \Delta B \Delta C$ if and only if x belongs either to all three sets or to exactly one of those. In other words, $x \notin A \Delta B \Delta C$ if and only if x is either in none of the three sets or exactly in two of them.

Problem A.6. Let A and B be two subsets of a set M . Which of the following equations are valid? Prove the correct identities, give counterexamples for the false ones:

$$\begin{aligned} (A \times B)^c &= A^c \times B^c, \\ (A \times B)^c &= (A^c \times B^c) \cup (A^c \times M) \cup (M \times B^c), \\ (A \times B)^c &= (A^c \times B) \cup (A \times B^c), \\ (A \times B)^c &= (A^c \times M) \cup (M \times B^c). \end{aligned}$$

Problem A.7. Define f from \mathbb{N} to \mathbb{Z} by

$$f(n) = \begin{cases} 0 & \text{if } n \text{ is even,} \\ 1 & \text{if } n \text{ is odd.} \end{cases}$$

Describe $f^{-1}(B)$ for all $B \subseteq \mathbb{Z}$.

Problem A.8. Determine

$$f^{-1}([0, 5]) \quad \text{and} \quad f^{-1}([0, \infty))$$

where $f : \mathbb{R} \rightarrow \mathbb{R}$ denotes the floor function. That is,

$$f(x) = \lfloor x \rfloor, \quad x \in \mathbb{R}.$$

A.3 Combinatorics

A.3.1 Binomial coefficients

A one-to-one mapping π from $\{1, \dots, n\}$ to $\{1, \dots, n\}$ is called a **permutation** (of order n). Any permutation reorders the numbers from 1 to n as $\pi(1), \pi(2), \dots, \pi(n)$ and, vice versa, each reordering of these numbers generates a permutation. One way to write a permutations is

$$\pi = \begin{pmatrix} 1 & 2 & \dots & n \\ \pi(1) & \pi(2) & \dots & \pi(n) \end{pmatrix}.$$

For example, if $n = 3$, then $\pi = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 3 & 1 \end{pmatrix}$ is equivalent to the order 2, 3, 1 or to $\pi(1) = 2$, $\pi(2) = 3$, and $\pi(3) = 1$.

Let S_n be the set of all permutations of order n . Then one may ask for $|S_n|$ or, equivalently, for the number of possible orderings of the numbers $\{1, \dots, n\}$.

To treat this problem, we need the following definition.

Definition A.3.1. For $n \in \mathbb{N}$, we define **n -factorial** by setting

$$n! = 1 \cdot 2 \cdot \dots \cdot (n-1) \cdot n.$$

Furthermore, let $0! = 1$.

Now we may answer the question about the cardinality of S_n .

Proposition A.3.2. *We have*

$$|S_n| = n! \tag{A.6}$$

or, equivalently, there are $n!$ different ways to order n distinguishable objects.

Proof. The proof is done by induction over n . If $n = 1$ then $|S_1| = 1 = 1!$ and eq. (A.6) is valid.

Now suppose that eq. (A.6) is true for n . In order to prove eq. (A.6) for $n + 1$, we split S_{n+1} as follows:

$$S_{n+1} = \bigcup_{k=1}^{n+1} A_k,$$

where

$$A_k = \{\pi \in S_{n+1} : \pi(n+1) = k\}, \quad k = 1, \dots, n+1.$$

Each $\pi \in A_k$ generates a one-to-one mapping $\tilde{\pi}$ from $\{1, \dots, n\}$ onto the set $\{1, \dots, k-1, k+1, \dots, n+1\}$ by letting $\tilde{\pi}(j) = \pi(j)$, $1 \leq j \leq n$. Vice versa, each such $\tilde{\pi}$ defines a permutation $\pi \in A_k$ by setting $\pi(j) = \tilde{\pi}(j)$, $j \leq n$, and $\pi(n+1) = k$. Consequently, since eq. (A.6) holds for n , we get $|A_k| = n!$. Furthermore, the A_k s are disjoint, and

$$|S_{n+1}| = \sum_{k=1}^{n+1} |A_k| = (n+1) \cdot n! = (n+1)!,$$

hence eq. (A.6) also holds for $n+1$. This completes the proof. \square

Next we treat a tightly related problem. Say we have n different objects and we want to distribute them into two disjoint groups, one having k elements, the other $n-k$. Hereby it is of no interest in which order the elements are distributed, only the composition of the two sets matters.

Example A.3.3. There are 52 cards in a deck that are distributed to two players, so that each of them gets 26 cards. For this game, it is only important which cards each player has, not in which order the cards were received. Here $n = 52$ and $k = n - k = 26$.

The main question is: how many ways can n elements be distributed, say the numbers from 1 to n , into one group of k elements and into another of $n - k$ elements? In the above example, that is how many ways can 52 cards be distributed into two groups of 26.

To answer this question, we use the following auxiliary model. Let us take any permutation $\pi \in S_n$. We place the numbers $\pi(1), \dots, \pi(k)$ into group 1 and the remaining $\pi(k+1), \dots, \pi(n)$ into group 2. In this way, we obtain all possible distributions but many of them appear several times. Say that two permutations π_1 and π_2 are equivalent if (as sets)

$$\{\pi_1(1), \dots, \pi_1(k)\} = \{\pi_2(1), \dots, \pi_2(k)\}.$$

Of course, this also implies

$$\{\pi_1(k+1), \dots, \pi_1(n)\} = \{\pi_2(k+1), \dots, \pi_2(n)\},$$

and two permutations generate the same partition if and only if they are equivalent. Equivalent permutations are achieved by taking one fixed permutation π , then permuting $\{\pi(1), \dots, \pi(k)\}$ and also $\{\pi(k+1), \dots, \pi(n)\}$. Consequently, there are exactly $k!(n-k)!$ permutations that are equivalent to a given one. Summing up, we get that there are $\frac{n!}{k!(n-k)!}$ different classes of equivalent permutations. Let

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}.$$

! There are $\binom{n}{k}$ different ways to distribute n objects into one group of k and into another one of $n-k$ elements. For any $n \geq 0$, we set $\binom{n}{0} = 1$ and $\binom{n}{k} = 0$ in case of $k > n$ or $k < 0$.

Definition A.3.4. The numbers

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}, \quad n = 0, 1, \dots \text{ and } k = 0, \dots, n,$$

are called **binomial coefficients**, read “ n choose k .”

Example A.3.5. A digital word of length n consists of n zeroes or ones. Since at every position we may have either “0” or “1”, there are 2^n different words of length n . How many of these words possess exactly k ones or, equivalently, $n-k$ zeroes? To answer this, put all positions where there is a “1” into a first group and those where there is a “0” into a second one. In this way, the numbers from 1 to n are divided into two different groups of size k and $n-k$, respectively. But we already know how many such partitions exist, namely $\binom{n}{k}$. As a consequence, we get

! There are $\binom{n}{k}$ words of length n possessing exactly k ones and $n-k$ zeroes.

The next proposition summarizes some crucial properties of binomial coefficients.

Proposition A.3.6. Let n be a natural number; $k = 0, \dots, n$, and let $r \geq 1$ be an integer. Then the following equations hold:

$$\binom{n}{k} = \binom{n}{n-k}, \tag{A.7}$$

$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}, \quad \text{and} \tag{A.8}$$

$$\binom{n+r}{n} = \sum_{j=0}^n \binom{n+r-j-1}{n-j} = \sum_{j=0}^n \binom{r+j-1}{j}. \tag{A.9}$$

Proof. Equations (A.7) and (A.8) follow immediately by the definition of the binomial coefficients. Note that eq. (A.8) also holds if $k = n$ because we agreed that $\binom{n-1}{n} = 0$.

If $k < n$, then an iteration of eq. (A.8) leads to

$$\binom{n}{k} = \sum_{j=0}^k \binom{n-j-1}{k-j}.$$

Replacing in the last equation n by $n+r$, as well as k by n (note that $n+r > n$), we obtain the left-hand identity (A.9). The right-hand equation follows by inverting the summation,

that is, one replaces j by $n - j$. Observe that (A.9) becomes wrong in the case $r = 0$. Then the left-hand side is 1 while the right-hand one equals 0. \square

Remark A.3.7. Equation (A.8) allows a graphical interpretation by **Pascal's triangle**. The coefficient $\binom{n}{k}$ in the n th row follows by summing the two values $\binom{n-1}{k-1}$ and $\binom{n-1}{k}$ above $\binom{n}{k}$ in the $(n - 1)$ th row. Look at Figure A.6 for a visualization of the triangle.

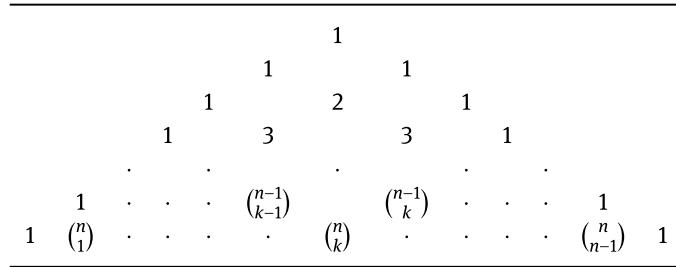


Figure A.6: Pascal's triangle.

Next we state and prove an important binomial theorem.

Proposition A.3.8 (Binomial theorem). *For real numbers a, b , and any $n \in \mathbb{N}_0$,*

$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}. \tag{A.10}$$

Proof. The binomial theorem is proved by induction over n . If $n = 0$, then eq. (A.10) holds trivially.

Suppose now that eq. (A.10) has been proven for $n - 1$. Our aim is to verify that it is also true for n . Using that the expansion holds for $n - 1$, we get

$$\begin{aligned} (a + b)^n &= (a + b)(a + b)^{n-1} = (a + b) \sum_{k=0}^{n-1} \binom{n-1}{k} a^k b^{n-1-k} \\ &= \sum_{k=0}^{n-1} \binom{n-1}{k} a^{k+1} b^{n-1-k} + \sum_{k=0}^{n-1} \binom{n-1}{k} a^k b^{n-k} \\ &= a^n + \sum_{k=0}^{n-2} \binom{n-1}{k} a^{k+1} b^{n-1-k} + b^n + \sum_{k=1}^{n-1} \binom{n-1}{k} a^k b^{n-k} \\ &= a^n + b^n + \sum_{k=1}^{n-1} \left[\binom{n-1}{k-1} + \binom{n-1}{k} \right] a^k b^{n-k} = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}, \end{aligned}$$

where we used eq. (A.8) in the last step. \square

The following property of binomial coefficients plays an important role when introducing the hypergeometric distribution (see Proposition 1.4.31). It is also used during the investigation of sums of independent binomial distributed random variables (see Proposition 4.6.1).

Proposition A.3.9 (Vandermonde's identity). *If $k, m,$ and n are all in \mathbb{N}_0 , then*

$$\sum_{j=0}^k \binom{n}{j} \binom{m}{k-j} = \binom{n+m}{k}. \quad (\text{A.11})$$

Proof. An application of the binomial theorem leads to

$$(1+x)^{n+m} = \sum_{k=0}^{n+m} \binom{n+m}{k} x^k, \quad x \in \mathbb{R}. \quad (\text{A.12})$$

On the other hand, another use of Proposition A.3.8 implies³

$$\begin{aligned} (1+x)^{n+m} &= (1+x)^n (1+x)^m \\ &= \left[\sum_{j=0}^n \binom{n}{j} x^j \right] \left[\sum_{i=0}^m \binom{m}{i} x^i \right] = \sum_{j=0}^n \sum_{i=0}^m \binom{n}{j} \binom{m}{i} x^{i+j} \\ &= \sum_{k=0}^{n+m} \left[\sum_{i+j=k} \binom{n}{j} \binom{m}{i} \right] x^k = \sum_{k=0}^{n+m} \left[\sum_{j=0}^k \binom{n}{j} \binom{m}{k-j} \right] x^k. \end{aligned} \quad (\text{A.13})$$

The coefficients in an expansion of a polynomial are unique. Hence, in view of eqs. (A.12) and (A.13), we get for all $k \leq m+n$ the identity

$$\binom{n+m}{k} = \sum_{j=0}^k \binom{n}{j} \binom{m}{k-j}.$$

Hereby note that both sides of eq. (A.11) become zero whenever $k > n+m$. This completes the proof. \square

Remark A.3.10. Another, more heuristic, way to prove Vandermonde's identity is as follows. Suppose one has $n+m$ fruits, n apples and m oranges. There are $\binom{n+m}{k}$ ways to choose k fruits out of the $n+m$ ones. These possibilities split into the following $k+1$ disjoint events: among the chosen k fruits there are zero apples and k oranges, or one apple and $k-1$ oranges up to k apples and zero oranges. If the number of apples in the sample is j , then there are $\binom{n}{j}$ ways to choose the apples and $\binom{m}{k-j}$ ways to choose the

³ When passing from line 2 to line 3, the order of summation is changed. One no longer sums over the rectangle $[0, m] \times [0, n]$. Instead, one sums along the diagonals, where $i+j=k$.

oranges, respectively. Summing over all possibilities $j = 0, \dots, k$ proves Vandermonde's identity.

Our next objective is to generalize the binomial coefficients. In view of

$$\binom{n}{k} = \frac{n(n-1)\cdots(n-k+1)}{k!}$$

for $k \geq 1$ and $n \in \mathbb{N}$, the **generalized binomial coefficient** is introduced as

$$\binom{-n}{k} := \frac{-n(-n-1)\cdots(-n-k+1)}{k!}. \quad (\text{A.14})$$

The next lemma shows the tight relation between generalized and “ordinary” binomial coefficients.

Lemma A.3.11. For $k \geq 1$ and $n \in \mathbb{N}$,

$$\binom{-n}{k} = (-1)^k \binom{n+k-1}{k}.$$

Proof. By definition of the generalized binomial coefficient, we obtain

$$\begin{aligned} \binom{-n}{k} &= \frac{(-n)(-n-1)\cdots(-n-k+1)}{k!} \\ &= (-1)^k \frac{(n+k-1)(n+k-2)\cdots(n+1)n}{k!} = (-1)^k \binom{n+k-1}{k}. \end{aligned}$$

This completes the proof. □

For example, Lemma A.3.11 implies $\binom{-1}{k} = (-1)^k$ and $\binom{-n}{1} = -n$.

A.3.2 Drawing balls out of an urn

Assume that there are n balls labeled from 1 to n in an urn. We draw k balls out of the urn, thus observing a sequence of length k with entries from $\{1, \dots, n\}$. How many different results (sequences) may be observed? To answer this question, we have to decide on the arrangement of drawing. Do we or do we not replace the chosen ball? Is it important in which order the balls were chosen or is it only of importance which balls were chosen at all? Thus, we see that there are four different ways to answer this question (replacement or nonreplacement, recording the order or nonrecording).

Example A.3.12. Let us regard the drawing of two balls out of four, that is, $n = 4$ and $k = 2$. Depending on the different arrangements, the following results may be observed. Note, for example, that in the two latter cases (nonrecording of the order) the pair $(3, 2)$ does not appear because it is identical to $(2, 3)$.

Replacement and the order is important

(1, 1)	(1, 2)	(1, 3)	(1, 4)
(2, 1)	(2, 2)	(2, 3)	(2, 4)
(3, 1)	(3, 2)	(3, 3)	(3, 4)
(4, 1)	(4, 2)	(4, 3)	(4, 4)

16 different results

Nonreplacement and the order is important

·	(1, 2)	(1, 3)	(1, 4)
(2, 1)	·	(2, 3)	(2, 4)
(3, 1)	(3, 2)	·	(3, 4)
(4, 1)	(4, 2)	(4, 3)	·

12 different results

Replacement and the order is not important

(1, 1)	(1, 2)	(1, 3)	(1, 4)
·	(2, 2)	(2, 3)	(2, 4)
·	·	(3, 3)	(3, 4)
·	·	·	(4, 4)

10 different results

Nonreplacement and the order is not important

·	(1, 2)	(1, 3)	(1, 4)
·	·	(2, 3)	(2, 4)
·	·	·	(3, 4)
·	·	·	·

6 different results

Let us come back now to the general situation of n different balls from which we choose k at random:

Case 1. Drawing with replacement and taking the order into account. We have n different possibilities for the choice of the first ball and, since the chosen ball is placed back, there are also n possibilities for the second one, and so on. Thus, there are n possibilities for each of the k balls, leading to the following result.

! The number of different results in this case is n^k .

Example A.3.13. Letters in Braille, a scripture for blind people, are generated by dots or nondots at six different positions. How many letters may be generated in that way?

Answer: It holds that $n = 2$ (dot or no dot) at $k = 6$ different positions. Hence, the number of possible representable letters is $2^6 = 64$. In fact, there are only 63 possibilities because we have to rule out the case of no dots at all 6 positions.

Case 2. Drawing without replacement and taking the order into account. This case only makes sense if $k \leq n$. There are n possibilities to choose the first ball. After that there are still $n - 1$ balls in the urn. Hence there are only $n - 1$ possibilities for the second choice, $n - 2$ for the third, and so on. Summing up, we get the following.

! The number of possible results in this case equals $n(n - 1) \cdots (n - k + 1) = \frac{n!}{(n - k)!}$.

Example A.3.14. In a lottery, 6 numbers are chosen out of 49. Of course, the chosen numbers are not replaced. If we record the numbers as they appear (not putting them in order), how many different sequences of six numbers exist?

Answer: Here we have $n = 49$ and $k = 6$. Hence the wanted number equals

$$\frac{49!}{43!} = 49 \cdots 44 = 10,068,347,520.$$

Case 3. Drawing with replacement not taking the order into account. This case is more complicated and requires a different point of view. We count how often each of the n balls was chosen during the k trials. Let $k_1 \geq 0$ be the frequency of the first ball, $k_2 \geq 0$ that of the second one, and so on. In this way we obtain n nonnegative integers k_1, \dots, k_n satisfying

$$k_1 + \cdots + k_n = k.$$

Indeed, since we choose k balls, the frequencies have to sum to k . Consequently, the number of possible results when drawing k of n balls with replacement and not taking the order into account coincides with

$$|\{(k_1, \dots, k_n), k_j \in \mathbb{N}_0, k_1 + \cdots + k_n = k\}|. \tag{A.15}$$

In order to determine the cardinality (A.15), we use the following auxiliary model:

Let B_1, \dots, B_n be n boxes. Given n nonnegative integers k_1, \dots, k_n , summing to k , we place exactly k_1 dots into B_1 , k_2 dots into B_2 , and so on. At the end we distributed k indistinguishable dots into n different boxes. Thus, we see that the value of (A.15) coincides with the number of different possibilities to distribute k indistinguishable dots into n boxes. Now assume that the boxes are glued together; on the very left we put box B_1 , on its right we put box B_2 , and continue in this way up to box B_n on the very right. In this way, we obtain $n + 1$ dividing walls, two outer and $n - 1$ inner ones. Now we get all possible distributions of k dots into n boxes by shuffling the k dots and the $n - 1$ inner dividing walls. For example, if we get the order $d, d, d, w, w, d, w, \dots$ then this means that there are three dots in B_1 , none in B_2 , and one in B_3 , and so on (compare Figure A.7 for a slightly more general example).

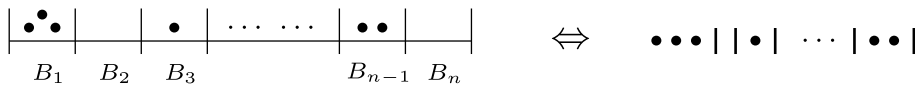


Figure A.7: The case $k_1 = 3, k_2 = 0, k_3 = 1, \dots, k_{n-1} = 2, k_n = 0$: k dots and $n - 1$ inner walls.

Summing up, we have $N = n + k - 1$ objects, k of them are dots and $n - 1$ are walls. As we know, there are $\binom{N}{k}$ different ways to order these N objects. Hence we arrived at the following result.



The number of possibilities to distribute k anonymous dots into n boxes equals

$$\binom{n+k-1}{k} = \binom{n+k-1}{n-1}.$$

It coincides with $|\{(k_1, \dots, k_n), k_j \in \mathbb{N}_0, k_1 + \dots + k_n = k\}|$, as well as with the number of different results when choosing k balls out of n with replacement and not taking order into account.

Example A.3.15. Dominoes are marked on each half either with no dots, one dot, or up to six dots. Hereby the dominoes are symmetric, that is, a tile with three dots on the left-hand side and two ones on the right-hand side is identical with that having two dots on the left-hand side and three dots on the right-hand side. How many different dominoes exist?

Answer: It holds⁴ $n = 7$ and $k = 2$, hence the number of different dominoes equals

$$\binom{7+2-1}{2} = \binom{8}{2} = 28.$$

Case 4. Drawing without replacement not taking the order into account. Here we also have to assume $k \leq n$. We already investigated this case when we introduced the binomial coefficients. The k chosen numbers are put in group 1, the remaining $n - k$ balls in group 2. As we know, there are $\binom{n}{k}$ ways to split the n numbers into such two groups. Hence we obtained the following.



The number of different results in this case is $\binom{n}{k}$.

Example A.3.16. If the order of the six numbers is not taken into account in Example A.3.14, that is, we ignore which number was chosen first, which second, and so on, the number of possible results equals

$$\binom{49}{6} = \frac{49 \cdot \dots \cdot 43}{6!} = 13,983,816.$$

Let us summarize the four different cases in a table. Here **O** and **NO** stand for recording or nonrecording of the order while **R** and **NR** represent replacement or nonreplacement, and one chooses k balls out of n possible.

⁴ There are 7 boxes B_0 up to B_6 and two particles distributed into these 7 boxes. The number of the box containing a particle corresponds to the number of dots on the tile. For example, if one particle is in box B_2 and the other in B_4 , then the corresponding tile is that with 2 and 4 dots on it. Or both particles in B_6 means that our tile has 6 dots at each side. Another possibility to describe these tiles is to represent 2 as $2 = 0 + 0 + 1 + 0 + 1 + 0 + 0$ or $2 = 0 + 0 + 0 + 0 + 0 + 0 + 2$, respectively.


	R	NR
O	n^k	$\frac{n!}{(n-k)!}$
NO	$\binom{n+k-1}{k}$	$\binom{n}{k}$

A.3.3 Multinomial coefficients

The binomial coefficient $\binom{n}{k}$ describes the number of possibilities to distribute n objects into two groups of k and $n - k$ elements. What happens if we have not only two groups but $m \geq 2$? Say the first group has k_1 elements, the second has k_2 elements, and so on, up to the m th group that has k_m elements. Of course, if we distribute n elements the k_j s have to satisfy

$$k_1 + \cdots + k_m = n.$$

Using exactly the same arguments as in the case where $m = 2$, we get the following.

There exist exactly $\frac{n!}{k_1! \cdots k_m!}$ different ways to distribute n elements into m groups of sizes k_1, k_2, \dots, k_m where $k_1 + \cdots + k_m = n$. 

In accordance with the binomial coefficient, we write

$$\binom{n}{k_1, \dots, k_m} := \frac{n!}{k_1! \cdots k_m!}, \quad k_1 + \cdots + k_m = n, \quad (\text{A.16})$$

and call $\binom{n}{k_1, \dots, k_m}$ a **multinomial coefficient**, read “ n chose k_1 up to k_m .”

Remark A.3.17. If $m = 2$, then $k_1 + k_2 = n$, and

$$\binom{n}{k_1, k_2} = \binom{n}{k_1, n - k_1} = \binom{n}{k_1} = \binom{n}{k_2}.$$

Example A.3.18. A deck of cards for playing skat consists of 32 cards. Three players each gets 10 cards; the remaining two cards (called “skat”) are placed on the table. How many different distributions of the cards exist?

Answer: Let us first define what it means for two distribution of cards to be identical. Say this happens if each of the three players has exactly the same cards as in the previous game. Therefore, the remaining two cards on the table are also identical. Hence we distribute 32 cards into 4 groups possessing 10, 10, 10, and 2 elements. Consequently, the number of different distributions equals⁵

⁵ The huge size of this number explains why playing skat never becomes boring.

$$\binom{32}{10, 10, 10, 2} = \frac{32!}{(10!)^3 2!} = 2.753294409 \times 10^{15}.$$

Remark A.3.19. One may also look at multinomial coefficients from a different point of view. Suppose we are given n balls of m different colors. Say there are k_1 balls of the first color, k_2 balls of the second color, up to k_m balls of color m where, of course, we have $k_1 + \cdots + k_m = n$. Then there exist

$$\binom{n}{k_1, \dots, k_m}$$

different ways to order these n balls. This is followed by the same arguments as we used in Example A.3.5 for $m = 2$.

For instance, given 3 blue, 4 red, and 2 white balls, there are

$$\binom{9}{3, 4, 2} = \frac{9!}{3! 4! 2!} = 1260$$

different ways to order them.

Finally, let us still mention that in the literature one sometimes finds another (equivalent) way of introducing the multinomial coefficients. Given nonnegative integers k_1, \dots, k_m with $k_1 + \cdots + k_m = n$, it follows that

$$\binom{n}{k_1, \dots, k_m} = \binom{n}{k_1} \binom{n-k_1}{k_2} \binom{n-k_1-k_2}{k_3} \cdots \binom{n-k_1-\cdots-k_{m-1}}{k_m}. \quad (\text{A.17})$$

A direct proof of this fact is easy and left as an exercise.

There is a combinatorial interpretation of the expression on the right-hand side of eq. (A.17). To reorder n balls of m different colors, one chooses first the k_1 positions for balls of color 1. There are $\binom{n}{k_1}$ ways to do this. Thus, there remain $n - k_1$ possible positions for balls of color 2, and there are $\binom{n-k_1}{k_2}$ possible choices for this, and so on. Note that at the end there remain k_m positions for balls of color m ; hence, the last term on the right-hand side of eq. (A.17) equals 1.

Let us come now to the announced generalization of Proposition A.3.8.

Proposition A.3.20 (Multinomial theorem). *Let $n \geq 0$. Then for any $m \geq 1$ and real numbers x_1, \dots, x_m ,*

$$(x_1 + \cdots + x_m)^n = \sum_{\substack{k_1 + \cdots + k_m = n \\ k_i \geq 0}} \binom{n}{k_1, \dots, k_m} x_1^{k_1} \cdots x_m^{k_m}. \quad (\text{A.18})$$

Proof. Equality (A.18) is proved by induction. In contrast to the proof of the binomial theorem, now induction is done over m , the number of summands.

If $m = 1$, the assertion is valid due to trivial reasons.

Suppose now eq. (A.18) holds for m , all $n \geq 1$, and all real numbers x_1, \dots, x_m . We have to show the validity of eq. (A.18) for $m + 1$ and all $n \geq 1$. Given real numbers x_1, \dots, x_{m+1} and $n \geq 1$ set $y := x_1 + \dots + x_m$. Using Proposition A.3.8, by the validity of eq. (A.18) for m and all $n - j, 0 \leq j \leq n$, we obtain

$$\begin{aligned} (x_1 + \dots + x_{m+1})^n &= (y + x_{m+1})^n = \sum_{j=1}^n \frac{n!}{j!(n-j)!} x_{m+1}^j y^{n-j} \\ &= \sum_{j=1}^n \frac{n!}{j!(n-j)!} \sum_{\substack{k_1 + \dots + k_m = n-j \\ k_i \geq 0}} \frac{(n-j)!}{k_1! \dots k_m!} x_1^{k_1} \dots x_m^{k_m} x_{m+1}^j. \end{aligned}$$

Replacing j by k_{m+1} and combining both sums leads to

$$(x_1 + \dots + x_{m+1})^n = \sum_{\substack{k_1 + \dots + k_{m+1} = n \\ k_i \geq 0}} \frac{n!}{k_1! \dots k_{m+1}!} x_1^{k_1} \dots x_{m+1}^{k_{m+1}},$$

hence eq. (A.18) is also valid for $m + 1$. This completes the proof. \square

Remark A.3.21. The number of summands in eq. (A.18) equals⁶ $\binom{n+m-1}{n}$.

Example A.3.22. Let w, x, y , and z be four real numbers. Then we get

$$(w + x + y + z)^5 = \sum_{k_1 + \dots + k_4 = 5} \binom{5}{k_1, k_2, k_3, k_4} w^{k_1} x^{k_2} y^{k_3} z^{k_4}.$$

So, for example, the coefficient of w^2xyz is $\binom{5}{2,1,1,1} = 60$ or that of w^2x^2y equals $\binom{5}{2,2,1,0} = 30$.

A.3.4 Problems

Problem A.9. Given a set M of cardinality $n \geq 1$. Argue why there are exactly $\binom{n}{k}$ subsets $A \subseteq M$ with cardinality $k \leq n$.

Problem A.10. Determine, with proof, the number of ordered triples (A_1, A_2, A_3) of sets such that

$$A_1 \cup A_2 \cup A_3 = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \quad \text{and} \quad A_1 \cap A_2 \cap A_3 = \emptyset.$$

Problem A.11. Prove that

$$\sum_{k=0}^n \binom{n}{k}^2 = \binom{2n}{n}. \quad (\text{A.19})$$

⁶ Compare with Case 3 in Section A.3.2.

Problem A.12. Let n be a natural number. Prove that

$$\frac{2^{2n}}{2n+1} < \binom{2n}{n} < 2^{2n} \quad \text{and} \quad \frac{2^{2n+1}}{2n+2} < \binom{2n+1}{n} = \binom{2n+1}{n+1} < 2^{2n+1}.$$

Problem A.13. Let n be a natural number. Give proofs of the identities

$$\binom{n}{0} + \binom{n}{1} + \cdots + \binom{n}{n} = 2^n$$

and

$$\binom{n}{0} - \binom{n}{1} + \cdots + (-1)^n \binom{n}{n} = 0.$$

Problem A.14. Given $r \in \mathbb{N}$ and an integer $n \geq r$, show that

$$\binom{r}{r} + \binom{r+1}{r} + \cdots + \binom{n}{r} = \binom{n+1}{r+1}.$$

For example, if $n \geq 4$, then

$$\binom{4}{4} + \binom{5}{4} + \cdots + \binom{n}{4} = \binom{n+1}{5}.$$

Problem A.15. What is the coefficient of x^2 in $(3x^2 - 2x^{-1})^7$ where $x \neq 0$ is some variable?

Problem A.16. Given integers $n \geq 1$ and $k \geq n$, how many vectors (k_1, \dots, k_n) of integers exist for which $k_j \geq 1$ and $k_1 + \cdots + k_n = k$. How about if we ask for k_j s with $k_j \geq M$, $1 \leq j \leq n$, for some integer $M \geq 1$? Of course, this question makes only sense if $k \geq nM$.

Problem A.17. Give an algebraic proof of equation (A.17).

A.4 Vectors and matrices

The aim of this section is to summarize results and notations about vectors and matrices used throughout this book. For more detailed reading, we refer to any book about Linear Algebra, for example, [Ax15].

Given two vectors x and y in \mathbb{R}^n , their⁷ **scalar product** is defined as

$$\langle x, y \rangle := \sum_{j=1}^n x_j y_j, \quad x = (x_1, \dots, x_n), \quad y = (y_1, \dots, y_n).$$

If $x \in \mathbb{R}^n$, then

⁷ Sometimes also called “dot-product.”

$$|x| := \langle x, x \rangle^{1/2} = \left(\sum_{j=1}^n x_j^2 \right)^{1/2}$$

denotes the **Euclidean distance** from x to 0. Thus, $|x|$ may also be regarded as the **length** of the vector $x \in \mathbb{R}^n$. In particular, we have $|x| > 0$ for all nonzero $x \in \mathbb{R}^n$.

Any matrix $A = (a_{ij})_{i,j=1}^n$ of real numbers a_{ij} generates a **linear⁸ mapping** (also denoted by A) via

$$Ax = \left(\sum_{j=1}^n a_{1j}x_j, \dots, \sum_{j=1}^n a_{nj}x_j \right), \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n. \quad (\text{A.20})$$

Conversely, any linear mapping $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ defines a matrix $(a_{ij})_{i,j=1}^n$ by representing $Ae_j \in \mathbb{R}^n$ as

$$Ae_j = (a_{1j}, \dots, a_{nj}), \quad j = 1, \dots, n.$$

Here $e_j = (0, \dots, 0, \underbrace{1}_j, 0, \dots, 0)$ denotes the j th unit vector in \mathbb{R}^n . With this generated matrix $(a_{ij})_{i,j=1}^n$, the linear mapping A acts as stated in eq. (A.20). Consequently, we may always identify linear mappings in \mathbb{R}^n with $n \times n$ -matrices $(a_{ij})_{i,j=1}^n$.

Given two $n \times n$ matrices $A = (a_{ij})_{i,j=1}^n$ and $B = (b_{ij})_{i,j=1}^n$, their product $A \circ B$, or AB in short, is the matrix $C = (c_{ij})_{i,j=1}^n$ where

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}, \quad 1 \leq i, j \leq n.$$

Note that in general $AB \neq BA$. An important formula for the determinant of the product of two matrices is

$$\det(AB) = \det(A) \cdot \det(B). \quad (\text{A.21})$$

An $n \times n$ matrix A is said to be **regular⁹** if the generated linear mapping is one-to-one, that is, if $Ax = 0$ implies $x = 0$. This is equivalent to the fact that the determinant $\det(A)$ is nonzero.

Let $A = (a_{ij})_{i,j=1}^n$ be an $n \times n$ matrix. Then its **transposed** matrix is defined as $A^T := (a_{ji})_{i,j=1}^n$. With this notation, it follows for $x, y \in \mathbb{R}^n$ that

$$\langle Ax, y \rangle = \langle x, A^T y \rangle.$$

Moreover, we have $(AB)^T = B^T A^T$ for any two $n \times n$ matrices A and B , and, of course,

⁸ A mapping $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is said to be linear if $A(\alpha x + \beta y) = \alpha Ax + \beta Ay$ for all $\alpha, \beta \in \mathbb{R}$ and $x, y \in \mathbb{R}^n$.

⁹ Sometimes also called **nonsingular** or **invertible**.

$(A^T)^T = A$. The following property of the transposed matrix is crucial:

$$\det(A^T) = \det(A).$$

In particular, the matrix A is regular if and only if A^T is regular.

A matrix A with $A = A^T$ is said to be **symmetric**. Equivalently, A satisfies

$$\langle Ax, y \rangle = \langle x, Ay \rangle, \quad x, y \in \mathbb{R}^n.$$

An $n \times n$ matrix $R = (r_{ij})_{i,j=1}^n$ is **positive definite** (or **positive** in short) provided it is symmetric and

$$\langle Rx, x \rangle = \sum_{i,j=1}^n r_{ij}x_i x_j > 0, \quad x = (x_1, \dots, x_n) \neq 0.$$

We will write $R > 0$ in this case. In particular, each positive matrix R is regular and its determinant satisfies $\det(R) > 0$.

Let $A = (a_{ij})_{i,j=1}^n$ be an arbitrary regular $n \times n$ matrix. Set

$$R := AA^T, \tag{A.22}$$

that is, the entries r_{ij} of R are computed by

$$r_{ij} = \sum_{k=1}^n a_{ik}a_{jk}, \quad 1 \leq i, j \leq n.$$

Proposition A.4.1. *Suppose the matrix R is defined by eq. (A.22) for some regular A . Then it follows that $R > 0$.*

Proof. Because of

$$R^T = (AA^T)^T = (A^T)^T A^T = AA^T = R,$$

the matrix R is symmetric. Furthermore, for $x \in \mathbb{R}^n$ with $x \neq 0$, we obtain

$$\langle Rx, x \rangle = \langle AA^T x, x \rangle = \langle A^T x, A^T x \rangle = |A^T x|^2 > 0.$$

Hereby we used that for a regular A , the transposed matrix A^T is regular, too. Consequently, if $x \neq 0$, then $A^T x \neq 0$, and thus $|A^T x| > 0$. This completes the proof. \square

The **identity matrix** I_n is defined as the $n \times n$ matrix with entries δ_{ij} , $1 \leq i, j \leq n$, where

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases} \tag{A.23}$$

Of course, $I_n x = x$ for $x \in \mathbb{R}^n$ and, moreover, $\det(I_n) = 1$.

Given a regular $n \times n$ matrix A , there is a unique matrix B such that $AB = I_n$. The matrix B is called the **inverse matrix** of A and denoted by A^{-1} . Recall that also $A^{-1}A = I_n$ and, moreover, $(A^T)^{-1} = (A^{-1})^T$. Equation (A.21) lets us conclude that for any regular matrix A , it follows that

$$1 = \det(I_n) = \det(AA^{-1}) = \det(A) \cdot \det(A^{-1}) \quad \Rightarrow \quad \det(A^{-1}) = \frac{1}{\det(A)}.$$

An $n \times n$ matrix U is said to be **unitary** or **orthogonal** provided that

$$UU^T = U^T U = I_n \tag{A.24}$$

with identity matrix I_n . Another way to express this is either that $U^T = U^{-1}$ or, equivalently, that U satisfies

$$\langle Ux, Uy \rangle = \langle x, y \rangle, \quad x, y \in \mathbb{R}^n.$$

In particular, for each $x \in \mathbb{R}^n$ it follows that

$$|Ux|^2 = \langle Ux, Ux \rangle = \langle x, x \rangle = |x|^2,$$

that is, U preserves the length of vectors in \mathbb{R}^n . Indeed, this property characterizes unitary matrices. It is a nice task to prove this.

It is easy to see that an $n \times n$ matrix U is unitary if and only if its column vectors u_1, \dots, u_n form an orthonormal basis in \mathbb{R}^n . That is, $\langle u_i, u_j \rangle = \delta_{ij}$ with δ_{ij} s as in (A.23). This characterization of unitary matrices remains valid when we take the column vectors instead of those generated by the rows.

We saw in Proposition A.4.1 that each matrix R of the form (A.22) is positive. Next we prove that conversely, each $R > 0$ may be represented in this way.

Proposition A.4.2. *Let R be an arbitrary positive $n \times n$ matrix. Then there exists a regular matrix A such that $R = AA^T$.*

Proof. Since R is symmetric, we may apply the principal axis transformation for symmetric matrices. It asserts that there exists a diagonal matrix¹⁰ D and a unitary matrix U such that

$$R = UDU^T.$$

Let $\delta_1, \dots, \delta_n$ be the entries of D at its diagonal. From $R > 0$ we derive $\delta_j > 0, 1 \leq j \leq n$. To see this fix $j \leq n$ and set $x := Ue_j$ where as above e_j denotes the j th unit vector in \mathbb{R}^n . Then $U^T x = e_j$, hence

¹⁰ The entries d_{ij} of D satisfy $d_{ij} = 0$ if $i \neq j$.

$$0 < \langle Rx, x \rangle = \langle UDU^T x, x \rangle = \langle DU^T x, U^T x \rangle = \langle De_j, e_j \rangle = \delta_j.$$

Because of $\delta_j > 0$, we may define $D^{1/2}$ as diagonal matrix with entries $\delta_j^{1/2}$ on its diagonal. Setting $A := UD^{1/2}$ and because $(D^{1/2})^T = D^{1/2}$, it follows that

$$R = (UD^{1/2})(UD^{1/2})^T = AA^T.$$

Since

$$\det(A)^2 = \det(A)\det(A) = \det(A)\det(A^T) = \det(AA^T) = \det(R) > 0,$$

the matrix A is regular, and this completes the proof.

Another way to argue without using properties of determinants is as follows. Assume the matrix A is not regular. Then this is also true for A^T . Hence there is a nonzero vector $x \in \mathbb{R}^n$ such that $A^T x = 0$. But this implies $Rx = A(A^T x) = A(0) = 0$ which contradicts the regularity of R . Thus, A has to be regular. \square

Remark A.4.3. Note that representation (A.22) is *not* unique. Indeed, if $R = AA^T$, then we also have $R = (AV)(AV)^T$ for any unitary matrix V .

A.5 Some analytic tools

The aim of this section is to present some special results of Calculus that play an important role in the book. Hereby we restrict ourselves to those topics that are maybe less known and that are not necessarily taught in a basic Calculus course. For a general introduction to Calculus, including those topics as convergence of power series, fundamental theorem of Calculus, mean-value theorem, and so on, we refer to the books [Spi08] and [Ste15].

We start with a result that is used in the proof of Poisson's limit theorem (Theorem 1.4.22). From Calculus, it is well known that for $x \in \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n = e^x. \quad (\text{A.25})$$

The probably easiest proof of this fact is via the approach presented in [Spi08]. There the logarithm function is defined by $\ln t = \int_1^t \frac{1}{s} ds$, $t > 0$. Hence, l'Hôpital's rule implies

$$\lim_{t \rightarrow \infty} t \ln\left(1 + \frac{x}{t}\right) = x, \quad x \in \mathbb{R}.$$

From this, eq. (A.25) easily follows by the continuity of the exponential function.

The next proposition may be viewed as a slight generalization of eq. (A.25).

Proposition A.5.1. Let $(x_n)_{n \geq 1}$ be a sequence of real numbers with $\lim_{n \rightarrow \infty} x_n = x$ for some $x \in \mathbb{R}$. Then we get

$$\lim_{n \rightarrow \infty} \left(1 + \frac{x_n}{n}\right)^n = e^x.$$

Proof. Because of eq. (A.25), it suffices to verify that

$$\lim_{n \rightarrow \infty} \left| \left(1 + \frac{x_n}{n}\right)^n - \left(1 + \frac{x}{n}\right)^n \right| = 0. \quad (\text{A.26})$$

Since the sequence $(x_n)_{n \geq 1}$ is convergent, it is bounded. Consequently, there is a $c > 0$ such that for all $n \geq 1$, we have $|x_n| \leq c$. Of course, we may also assume $|x| \leq c$. Fix for a moment $n \geq 1$ and set

$$a := 1 + \frac{x_n}{n} \quad \text{and} \quad b := 1 + \frac{x}{n}.$$

The choice of $c > 0$ yields $|a| \leq 1 + c/n$, as well as $|b| \leq 1 + c/n$. Hence it follows

$$\begin{aligned} |a^n - b^n| &= |a - b| |a^{n-1} + a^{n-2}b + \dots + ab^{n-2} + b^{n-1}| \\ &\leq |a - b| (|a|^{n-1} + |a|^{n-2}|b| + \dots + |a||b|^{n-2} + |b|^{n-1}) \\ &\leq |a - b| n \left(1 + \frac{c}{n}\right)^{n-1} \leq C n |a - b|. \end{aligned}$$

Here $C > 0$ is some constant that exists since $(1+c/n)^{n-1}$ converges to e^c . By the definition of a and b ,

$$\left| \left(1 + \frac{x_n}{n}\right)^n - \left(1 + \frac{x}{n}\right)^n \right| \leq C n \frac{|x_n - x|}{n} = C |x - x_n|.$$

Since $x_n \rightarrow x$, this immediately implies eq. (A.26) and proves the proposition. \square

Our next objective is to present some properties of power series and functions generated by them. Hereby we restrict ourselves to such assertions that we will use in this book. For further reading, we refer to Part IV in [Spi08].

Let $(a_k)_{k \geq 0}$ be a sequence of real numbers. Then its **radius of convergence** $r \in [0, \infty]$ is defined by

$$r := \frac{1}{\limsup_{k \rightarrow \infty} |a_k|^{1/k}}.$$

Hereby we let $1/0 := \infty$ and $1/\infty := 0$. If $0 < r \leq \infty$ and $|x| < r$, then the infinite series

$$f(x) := \sum_{k=0}^{\infty} a_k x^k \quad (\text{A.27})$$

converges (even absolutely). Hence the function f generated by eq. (A.27) is well defined on its **region of convergence** $\{x \in \mathbb{R} : |x| < r\}$. We say that f is represented as a **power series** on $\{x \in \mathbb{R} : |x| < r\}$.

The function f defined by eq. (A.27) is infinitely often differentiable on its region of convergence and (compare with [Spi08, Chapter 27, Theorem 6])

$$\begin{aligned} f^{(n)}(x) &= \sum_{k=n}^{\infty} k(k-1)\cdots(k-n+1) a_k x^{k-n} \\ &= \sum_{k=0}^{\infty} (k+n)(k+n-1)\cdots(k+1) a_{k+n} x^k \\ &= n! \sum_{k=0}^{\infty} \binom{n+k}{k} a_{k+n} x^k. \end{aligned} \quad (\text{A.28})$$

The coefficients $n! \binom{n+k}{k} a_{k+n}$ in the series representation of the n th derivative $f^{(n)}$ possess the same radius of convergence as the original sequence $(a_k)_{k \geq 0}$. This is easy to see for $n = 1$. The general case then follows by induction.

Furthermore, eq. (A.28) implies $a_n = f^{(n)}(0)/n!$, which, in particular, tells us that given f , the coefficients $(a_k)_{k \geq 0}$ in representation (A.27) are unique.

Proposition A.5.2. *If $n \geq 1$ and $|x| < 1$ then*

$$\frac{1}{(1+x)^n} = \sum_{k=0}^{\infty} \binom{-n}{k} x^k. \quad (\text{A.29})$$

Proof. Using the formula to add a geometric series and applying $\binom{-1}{k} = (-1)^k$ yields for $|x| < 1$ that

$$\frac{1}{1+x} = \sum_{k=0}^{\infty} (-1)^k x^k = \sum_{k=0}^{\infty} \binom{-1}{k} x^k.$$

Consequently Proposition A.5.2 holds for $n = 1$.

Assume now we have proven the proposition for $n - 1$, that is, if $|x| < 1$, then

$$\frac{1}{(1+x)^{n-1}} = \sum_{k=0}^{\infty} \binom{-n+1}{k} x^k.$$

Differentiating this equality in the region $\{x : |x| < 1\}$ implies

$$-\frac{n-1}{(1+x)^n} = \sum_{k=1}^{\infty} \binom{-n+1}{k} k x^{k-1} = \sum_{k=0}^{\infty} \binom{-n+1}{k+1} (k+1) x^k. \quad (\text{A.30})$$

Direct calculations give

$$\begin{aligned} -\frac{k+1}{n-1} \binom{-n+1}{k+1} &= -\frac{k+1}{n-1} \cdot \frac{(-n+1)(-n) \cdots (-n+1-(k+1)+1)}{(k+1)!} \\ &= \frac{(-n)(-n-1) \cdots (-n-k+1)}{k!} = \binom{-n}{k}, \end{aligned}$$

which, together with eq. (A.30), leads to

$$\frac{1}{(1+x)^n} = \sum_{k=0}^{\infty} \binom{-n}{k} x^k.$$

This completes the proof of Proposition A.5.2. \square

The next proposition may be viewed as a counterpart to eq. (A.11) in the case of generalized binomial coefficients.

Proposition A.5.3. For $k \geq 0$ and $m, n \in \mathbb{N}$,

$$\sum_{j=0}^k \binom{-n}{j} \binom{-m}{k-j} = \binom{-n-m}{k}.$$

Proof. The proof is similar to that of Proposition A.3.9. Using Proposition A.5.2, we represent the function $(1+x)^{-n-m}$ as a power series in two different ways. On the one hand, for $|x| < 1$, we have the representation

$$\frac{1}{(1+x)^{n+m}} = \sum_{k=0}^{\infty} \binom{-n-m}{k} x^k \quad (\text{A.31})$$

and, on the other hand,

$$\begin{aligned} \frac{1}{(1+x)^{n+m}} &= \left[\sum_{j=0}^{\infty} \binom{-n}{j} x^j \right] \left[\sum_{l=0}^{\infty} \binom{-m}{l} x^l \right] \\ &= \sum_{k=0}^{\infty} \left[\sum_{j+l=k} \binom{-n}{j} \binom{-m}{l} \right] x^k = \sum_{k=0}^{\infty} \left[\sum_{j=0}^k \binom{-n}{j} \binom{-m}{k-j} \right] x^k. \end{aligned} \quad (\text{A.32})$$

As observed above, the coefficients in a power series are uniquely determined. Thus, the coefficients in eqs. (A.31) and (A.32) have to coincide, which implies

$$\sum_{j=0}^k \binom{-n}{j} \binom{-m}{k-j} = \binom{-n-m}{k},$$

as asserted. \square

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function. How does one define the integral $\int_{\mathbb{R}^n} f(x) \, dx$? To simplify the notation, let us restrict ourselves to the case $n = 2$. The main problems already become clear in this case and the obtained results easily extend to higher dimensions.

The easiest way to introduce the integral of a function of two variables is as follows:

$$\int_{\mathbb{R}^2} f(x) \, dx := \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f(x_1, x_2) \, dx_2 \right] dx_1.$$

In order for this double integral to be well defined, we have to assume the existence of the inner integral for each fixed $x_1 \in \mathbb{R}$ and then the existence of the integral of the function

$$x_1 \mapsto \int_{-\infty}^{\infty} f(x_1, x_2) \, dx_2.$$

In doing so, the following question arises immediately: why do we not define the integral in reversed order, that is, first integrating with respect to x_1 and then with respect to x_2 ?

To see the difficulties that may appear, let us consider the following example.

Example A.5.4. The function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is defined as follows (see Fig. A.8): If either $x_1 < 0$ or $x_2 < 0$ set $f(x_1, x_2) = 0$. If $x_1, x_2 \geq 0$ define f by

$$f(x_1, x_2) := \begin{cases} +1 & \text{if } x_1 \leq x_2 < x_1 + 1, \\ -1 & \text{if } x_1 + 1 \leq x_2 \leq x_1 + 2, \\ 0 & \text{otherwise.} \end{cases}$$

We immediately see that

$$\int_0^{\infty} f(x_1, x_2) \, dx_2 = 0 \quad \text{for all } x_1 \in \mathbb{R}, \quad \text{hence} \quad \int_0^{\infty} \left[\int_0^{\infty} f(x_1, x_2) \, dx_2 \right] dx_1 = 0.$$

On the other hand, it follows

$$\int_0^{\infty} f(x_1, x_2) \, dx_1 = \begin{cases} \int_0^{x_2} (+1) \, dx_1 = x_2 & \text{if } 0 \leq x_2 < 1, \\ \int_0^{x_2-1} (-1) \, dx_1 + \int_{x_2-1}^{x_2} (+1) \, dx_1 = 2 - x_2 & \text{if } 1 \leq x_2 \leq 2, \\ \int_{x_2-2}^{x_2} f(x_1, x_2) \, dx_1 = 0 & \text{if } 2 < x_2 < \infty, \end{cases}$$

leading to

$$\int_0^{\infty} \left[\int_0^{\infty} f(x_1, x_2) \, dx_1 \right] dx_2 = \int_0^1 x_2 \, dx_2 + \int_1^2 (2 - x_2) \, dx_2 = 1.$$

Thus, in this case

$$\int_0^{\infty} \left[\int_0^{\infty} f(x_1, x_2) dx_1 \right] dx_2 \neq \int_0^{\infty} \left[\int_0^{\infty} f(x_1, x_2) dx_2 \right] dx_1 .$$

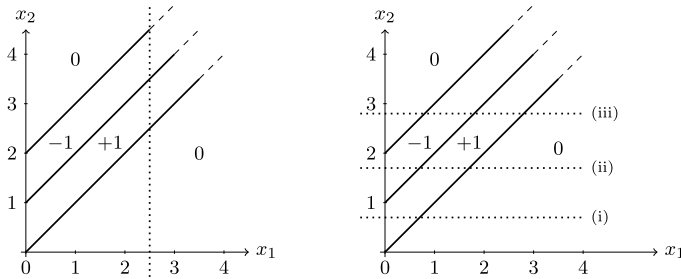


Figure A.8: On the left-hand side, one first integrates f over x_2 with x_1 fixed. On the right-hand side, the integration of f is done over x_1 with x_2 fixed. In this case three different regions for the choice of x_2 have to be considered.

Example A.5.4 shows that neither the definition of the integral of functions of several variables nor the interchange of integrals are unproblematic. Fortunately, we have the following positive result (see [Dur19, Section 1.7], for more information).

Proposition A.5.5 (Fubini’s theorem). *If $f(x_1, x_2) \geq 0$ for all $(x_1, x_2) \in \mathbb{R}^2$, then one may interchange the order of integration. In other words,*

$$\int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f(x_1, x_2) dx_1 \right] dx_2 = \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f(x_1, x_2) dx_2 \right] dx_1 . \tag{A.33}$$

Hereby we do not exclude that one of the two, hence also the other, iterated integral is infinite.

Furthermore, in the general case (the function f may attain also negative values) equality (A.33) holds provided that one of the iterated integrals, for example,

$$\int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} |f(x_1, x_2)| dx_1 \right] dx_2$$

is finite. Due to the first part, then we also have

$$\int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} |f(x_1, x_2)| dx_2 \right] dx_1 < \infty .$$

Whenever a function f on \mathbb{R}^2 satisfies one of the two assumptions in Proposition A.5.5, by

$$\int_{\mathbb{R}^2} f(x) \, dx := \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f(x_1, x_2) \, dx_1 \right] dx_2 = \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f(x_1, x_2) \, dx_2 \right] dx_1$$

the integral of f is well defined. Given a subset $B \subseteq \mathbb{R}^2$, we set

$$\int_B f(x) \, dx := \int_{\mathbb{R}^2} f(x) \mathbb{1}_B(x) \, dx,$$

provided the integral exists. Recall that $\mathbb{1}_B$ denotes the indicator function of B introduced in eq. (3.21).

For example, let K_1 be the unit circle in \mathbb{R}^2 , that is, $K_1 = \{(x_1, x_2) : x_1^2 + x_2^2 \leq 1\}$, then it follows

$$\int_{K_1} f(x) \, dx = \int_{-1}^1 \int_{-\sqrt{1-x_1^2}}^{\sqrt{1-x_1^2}} f(x_1, x_2) \, dx_2 \, dx_1.$$

Or, if $B = \{(x_1, x_2, x_3) \in \mathbb{R}^3 : x_1 \leq x_2 \leq x_3\}$, we have

$$\int_B f(x) \, dx = \int_{-\infty}^{\infty} \int_{-\infty}^{x_3} \int_{-\infty}^{x_2} f(x_1, x_2, x_3) \, dx_1 \, dx_2 \, dx_3.$$

Remark A.5.6. Proposition A.5.5 is also valid for infinite double series. Let a_{ij} be real numbers either satisfying $a_{ij} \geq 0$ or $\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} |a_{ij}| < \infty$, then this implies

$$\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} a_{ij} = \sum_{j=0}^{\infty} \sum_{i=0}^{\infty} a_{ij} = \sum_{ij=0}^{\infty} a_{ij}.$$

Even more generally, if the sets $I_k \subseteq \mathbb{N}_0^2$, $k \in \mathbb{N}_0$, form a disjoint partition of \mathbb{N}_0^2 , then

$$\sum_{ij=0}^{\infty} a_{ij} = \sum_{k=0}^{\infty} \sum_{(i,j) \in I_k} a_{ij}.$$

For example, if $I_k = \{(i, j) \in \mathbb{N}^2 : i + j = k\}$, then

$$\sum_{ij=0}^{\infty} a_{ij} = \sum_{k=0}^{\infty} \sum_{(i,j) \in I_k} a_{ij} = \sum_{k=0}^{\infty} \sum_{i=0}^k a_{ik-i}.$$

Bibliography

- [Art64] Emil Artin. *The Gamma Function*. Athena Series: Selected Topics in Mathematics, Holt, Rinehart and Winston, New York, Toronto, London, 1964.
- [Axl15] Sheldon Axler. *Linear Algebra Done Right*. Springer International Publishing, Cham Heidelberg, New York, Dordrecht, London, 3rd edition, 2015.
- [Bau96] Heinz Bauer. *Probability Theory*. De Gruyter Studies in Mathematics, Walter de Gruyter & Co., Berlin, 1996.
- [Bau01] Heinz Bauer. *Measure and Integration Theory*. De Gruyter Studies in Mathematics, Walter de Gruyter & Co., Berlin, 2001.
- [Bil12] Patrick Billingsley. *Probability and Measure*. John Wiley and Sons, Inc., Hoboken, 4th edition, 2012.
- [CB02] George Casella and Roger L. Berger. *Statistical Inference*. Duxburg Press, Pacific Grove, CA, 2nd edition, 2002.
- [CL23] Sebastian M. Cioabă and Werner Linde. *A Bridge to Advanced Mathematics: From Natural to Complex Numbers*. Pure and Applied Undergraduate Texts, **58**, American Mathematical Society, Providence, RI, 2023.
- [Coh13] Donald L. Cohn. *Measure Theory*. Birkhäuser Advanced Texts, Birkhäuser, Springer, New York, 2nd edition, 2013.
- [Dud02] Richard M. Dudley. *Real Analysis and Probability*. Cambridge University Press, Cambridge, 2002.
- [Dur19] Richard Durrett. *Probability: Theory and Examples*. Cambridge University Press, New York, 5th edition, 2019.
- [Eri73] Bruce K. Erickson. *The strong law of large numbers when the mean is undefined*. Trans. Amer. Math. Soc. **185** (1973), 371–381.
- [Fel68] William Feller. *An Introduction to Probability Theory and its Applications*, volume 1. John Wiley and Sons, New York, London, Sydney, 3rd edition, 1968.
- [Fis11] Hans Fischer. *A History of the Central Limit Theorem*. Ergebnisse der Mathematik und ihrer Grenzgebiete, Springer, New York, 2011.
- [Fsh71] Ronald Aymer Fisher. *The Design of Experiments*. Hafner Press, London, 9th edition, 1971.
- [Gha19] Saeed Ghahramani. *Fundamentals of Probability*. Pearson Education, Inc., Upper Saddle River, NJ, 4th edition, 2019.
- [GS01] Geoffrey R. Grimmett and David R. Stirzacker. *One Thousand Exercises in Probability*. Oxford University Press, Oxford, New York, 1st edition, 2001.
- [GS20] Geoffrey R. Grimmett and David R. Stirzacker. *Probability and Random Processes*. Oxford University Press, Oxford, New York, 4th edition, 2020.
- [Hal14] Paul R. Halmos. *Measure Theory*. Springer, New York, NY, 2014.
- [Kal21] Olav Kallenberg. *Foundations of Modern Probability*. Probability Theory and Stochastic Modeling, **99**, Springer, Cham, 3rd edition, 2021.
- [Kho07] Davar Khoshnevisan. *Probability*. Graduate Studies in Mathematics, **80**, American Mathematical Society, New York, 2007.
- [Kle20] Achim Klenke. *Probability Theory – A Comprehensive Course*. Universitext, Springer, Cham, 3rd edition, 2020.
- [Kol33] Andrey Nikolajewitsch Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*. Ergebnisse der Mathematik und ihrer Grenzgebiete, Julius Springer, Berlin, 1933.
- [Lag13] Jeffrey C. Lagarias. *Euler's constant: Euler's work and modern developments*. Bull. Amer. Math. Soc. (N. S.) **50** (2013), 527–628.
- [LG22] Jean-François Le Gall. *Measure Theory, Probability, and Stochastic Processes*. Graduate Texts in Mathematics, Springer, Cham, 2022.
- [Mor16] Samuel G. Moreno. *A Short and elementary proof of the Basel problem*. College Math. J. **47** (2016), 134–135.

- [Pao06] Marc S. Paoella. *Fundamental Probability. A Computational Approach*. John Wiley and Sons, Chichester, 2006.
- [Par05] Kalyanapuram Rangachari Parthasarathy. *Introduction to Probability and Measure*. Texts and Readings in Mathematics, **33**, Corrected reprint of the 1977 original, Hindustan Book Agency, New Delhi, 2005.
- [Rev13] Pál Révész. *Random Walk in Random and Non-random Environments*. World Scientific Publishing Co. Pte. Ltd., Hackensack, NJ, 3rd edition, 2013.
- [Ros06] Jeffrey S. Rosenthal. *A First Look at Rigorous Probability Theory*. World Scientific Publishing Co. Pte. Ltd., Hackensack, NJ, 2nd edition, 2006.
- [Rss14] Sheldon Ross. *A First Course in Probability*. Pearson Education Limited, Essex, 9th edition, 2014.
- [Sam04] Dov Samet, Iddo Samet and David Schmeidler. *One observation behind two-envelope puzzles*. *Amer. Math. Monthly* **111** (2004), 347–351.
- [Sch17] René L. Schilling. *Measures, Integrals and Martingales*. Cambridge University Press, Cambridge, 2nd edition, 2017.
- [Spi08] Michael Spivak. *Calculus*. Publish or Perish, Houston, TX, 4th edition, 2008.
- [Ste15] James Stewart. *Calculus*. Cengage Learning, Boston, 8th edition, 2015.
- [Sti03] David R. Stirzaker. *Elementary Probability*. Cambridge University Press, Cambridge, 2nd edition, 2003.

Index

- Absolute n th moment
 - of a random variable 257
- α -significance test 381
 - most powerful 382
- Arcsine distribution 68
- Arcsine law
 - for random walks 303
- $\mathcal{B}_{\alpha,\beta}$, beta distribution 66
- Banach's matchbox problem 46
- Basel problem 19
- Bayes' formula 118
- Bernoulli trial 211
- Bernstein polynomial 367
- Berry–Esséen theorem 366
- Bertrand paradox 109
- Beta distribution 66
- Beta function 65
- Bias
 - of an estimator 428
- Binary fraction 189
- Binomial coefficient 458
 - generalized 461
- Binomial distribution 28
- Binomial theorem 459
- $B_{n,p}$, binomial distribution 28
- $B_{n,p}^-$, negative binomial distribution 45
- Boole's inequality 11
- Borel σ -field
 - on \mathbb{R} 6
 - on \mathbb{R}^n 81
- Borel set
 - in \mathbb{R} 6
 - in \mathbb{R}^n 81
- Borel–Cantelli lemma 335
- Box
 - n -dimensional 81
- Boy or girl paradox 279
- Buffon's needle test 88
- Cantor set 54
- Cardinality of a set 451
- Cartesian product 453
- Cauchy distribution 69
 - general 109
- Cauchy–Schwarz inequality 276
- CDF 70
- Central limit theorem 350
 - for Γ -distributed random variables 363
 - for binomial random variables 359
 - for Poisson random variables 361
- Chain rule
 - for conditional probabilities 129
- Chebyshev's inequality 328
- χ^2 -distribution 64
- χ^2 -tests
 - known expected value 406
 - unknown expected value 407
- Clopper–Pearson intervals 440
- Complementary set 452
- Completely normal numbers 347
- Conditional distribution 115
- Conditional probability 115
- Confidence intervals 433
 - for binomial populations 438
 - approximative ones 443
 - exact ones 440
 - for hypergeometric populations 444
 - for normal populations 435
- Confidence regions 433
- Continuity correction
 - for normal approximation 353
- Continuity of a probability measure
 - from above 11
 - from below 11
- Convergence
 - almost surely 342
 - in distribution 349
 - in probability 341
- Convolution
 - of two functions 208
- Convolution formula
 - \mathbb{N}_0 -valued random variables 205
 - \mathbb{Z} -valued random variables 204
 - continuous random variables 208
- Coordinate mappings
 - of a random vector 146
- Correlated random variables 274
- Correlation coefficient 276
- Coupon collector's problem 255
- Covariance
 - of two random variables 271
 - properties 272

- Covariance matrix
 - of a normal vector 320
 - of a random vector 318
- Critical region 377
- Cumulative distribution function
 - of a probability measure 70
 - of a random variable 139
- De Morgan's rules 453
- Density
 - of a probability measure
 - multivariate 81
 - univariate 50
- Density function
 - of a probability measure
 - multivariate 81
 - univariate 49
 - of a random variable 136
 - of a random vector 154
- Dependence of events 120
- Dilemma
 - of hypothesis testing 380
- Dirac measure 21
- Discrete random variable 136
- Disjoint sets 452
- Distributing particles 23, 32
 - anonymous ones 24
 - distinguishable ones 24
- Distribution
 - of a random variable 135
 - of a random vector 147
- Distribution assumption 371
- Distribution density
 - of a random variable 136
 - of a random vector 154
- Distribution function
 - of a probability measure 70
 - of a random variable 139
- Distributive law
 - intersection and union 452
- Double factorial 58
- Drawing with replacement
 - no order 463
 - with order 462
- Drawing without replacement
 - no order 463
 - with order 462
- Dyadic rational number 189
- E_λ , exponential distribution 61
- $E_{\lambda,n}$, Erlang distribution 63
- Elementary event 2
- Envelope exchange paradox 287
- Erlang distribution 63
- Error
 - of the first kind 378
 - of the second kind 378
- Error function
 - Gaussian 72
- Estimator 415
 - efficient 431
 - maximum likelihood 418
 - unbiased 424
 - uniformly best 430
- Euclidean distance
 - in \mathbb{R}^n 468
- Euler's constant 256
- Event 2
 - certain 3
 - elementary 2
 - impossible 3
- Expected value
 - of continuous random variables 245
 - of discrete random variables 235
 - of nonnegative random variables
 - continuous case 243
 - discrete case 236
 - of random vectors 318
 - properties 250
- Exponential distribution 61
- F-distribution 229
- F-tests 412
- f. a. 334
- Factorial 456
- Finite additivity 8
- Fisher information 431
- Fisher–Snedecor distribution 229
- Fisher's
 - lemma 392
 - theorem 393
- Floor function 456
- Frequency
 - absolute 7
 - relative 7
- Fubini's theorem 477
- Function
 - absolutely integrable 244

- integrable 243
- measurable 180
- Gambler's ruin 293
- $\Gamma_{\alpha,\beta}$, gamma distribution 61
- Gamma function 58
- Gauss test
 - one-sided 400
 - two-sided 401
- Gaussian Φ -function 72
- Gaussian error function 72
- Generalized binomial coefficient 461
- Generated σ -field 4
- Generating function
 - of a nonnegative random variable 308
 - of an \mathbb{N}_0 -valued random variable 231
- Geometric distribution 41
- G_p , geometric distribution 41
- Histogram correction
 - for normal approximation 353
- Hitting time theorem 302
- $H_{N,M,n}$, hypergeometric distribution 39
- Hypergeometric distribution 39
- Hypothesis
 - alternative 376
 - null 376
- Hypothesis test 377
- Identically distributed 135
- Identity matrix 470
- Inclusion–exclusion formula 14, 104
- Independence
 - of infinitely many events 335
 - of n events 123
 - of two events 120
- Independent random variables 157
 - continuous case 165
 - discrete case 161
 - infinitely many 193
- Independent repetition
 - of an experiment 372
- Indicator function
 - of a set 164
- Inequality
 - Boole's 11
 - Cauchy–Schwarz 276
 - Chebyshev's 328
 - Lyapunov's 258
- Initial model
 - of a statistical experiment 372
- Intersection
 - of sets 451
- Interval estimator 432
- i. o. 334
- Joint density
 - of n random variables 154
- Joint distribution
 - of n random variables 147
- Laplace distribution 21
- Law
 - of iterated logarithm 357
 - of multiplication 112
 - of total probability 115
- Lemma
 - Borel–Cantelli 335
 - Fisher's 392
- Likelihood function
 - continuous case 417
 - discrete case 417
- Log-likelihood function 418
- Loss function
 - of an estimator 427
- Lower limit
 - of events 333
- Marginal distributions
 - of a random vector 148
 - continuous case 154
 - discrete case 150
- Marriage problem 283
- Matchbox problem 46
- Matrix
 - identity 470
 - inverse 471
 - invertible 469
 - nonsingular 469
 - orthogonal 471
 - positive definite 470
 - regular 469
 - symmetric 470
 - transposed 469
 - unitary 471
- Maximum likelihood estimator 418
- Measurable function 180

- Median
 - of a random variable 247
- MLE 418
- Moments
 - of a random variable 257
- Monkey at the cliff 299
- Monte Carlo method 344
- Monty Hall problem 106
- Multinomial
 - coefficient 465
 - random vector 152
 - theorem 466
- Multinomial distribution 31

- \mathbb{N} , natural numbers 451
- Needle test 88
- Negative binomial distribution 45
- Negatively correlated 278
- $\mathcal{N}(\mu, R)$, normal distribution
 - multivariate 316
- $\mathcal{N}(\mu, \sigma^2)$, normal distribution
 - univariate 57
- Normal distribution
 - multivariate 316
 - univariate 57
- Normal numbers 346
- \mathbb{N}_0 , natural numbers with zero 451

- Occurrence
 - of an event 3
- Occurrence of events
 - finally always 334
 - infinitely often 334
- Order statistics 171
 - density 174
 - distribution function 172

- Pairwise independence 122
- Paradox
 - boy or girl 279
 - envelope exchange 287
 - of Bertrand 109
 - of Chevalier de Méré 105
- Parameter set 374
- Pascal's triangle 459
- Permutation 456
- Point estimator 415
- Point measure 21
- Pois_λ , Poisson distribution 34

- Poisson distribution 34
- Poisson's limit theorem 35
- Positively correlated 278
- Power function
 - of a test 379
- Power series 473
- Powerset 451
- Preimage 454
- Principal axis transformation 471
- Probabilities
 - *a posteriori* 117
 - *a priori* 117
- Probability density function
 - multivariate 81
 - univariate 49
- Probability distribution 9
 - of a random variable 135
 - continuous case 139
 - discrete case 137
 - of a random vector 147
- Probability mass function 137
- Probability measure 9
 - continuous
 - multivariate 81
 - univariate 50
 - discrete 20
- Probability space 9
- Product σ -field 91
- Product measure 93
 - of continuous probabilities 95
 - of discrete probabilities 95
- Pseudoinverse
 - of a distribution function 199

- \mathbb{Q} , rational numbers 451
- Quantile
 - $F_{m,n}$ -distribution 399
 - χ_n^2 -distribution 398
 - t_n -distribution 399
 - general setting 395
 - standard normal distribution 396

- \mathbb{R} , real numbers 451
- Radius of convergence 473
- Raisins in dough 211
- Random experiment 1
- Random real number 132
- Random variable 132
 - continuous 136

- discrete 136
- real-valued 132
- singularly continuous 140
- vector valued 146
- Random variables
 - identically distributed 135
 - independent 157
- Random vector 146
 - $\mathcal{N}(\mu, R)$ -distributed 315
 - continuous 154
 - discrete 149
 - multinomial distributed 152
 - normally distributed 312
 - standard normally distributed 311
- Random walk
 - limit behavior 355
 - (next neighbor) on \mathbb{Z} 183
 - starting at an integer 293
 - symmetric 184
- Randomized test 377
- Reduction
 - of the sample space 113
- Region
 - of acceptance 377
 - of rejection 377
- Region of convergence
 - of a power series 474
- Risk
 - of the buyer 378
 - of the trader 378
- Risk function
 - of an estimator 427
- \mathbb{R}^n , n -dimensional Euclidean space 451
- Roulette
 - chance of winning 296
- Round-off errors 357

- Sample
 - random 371
- Sample mean 391
- Sample space 1
- Sample variance
 - biased 391
 - unbiased 391
- Scalar product
 - of two vectors 468
- Secretary problem 283
- Sequence
 - absolutely summable 236
 - summable 236
- Set difference 452
- σ -additivity 8
- σ -field 4
 - generated 5
- Significance level 381
- Significance test 381
 - for a binomial population
 - one-sided 389
 - two-sided 386
 - for a hypergeometric population
 - one-sided 382
 - two-sided 449
- Simulation
 - of a random variable
 - continuous case 198
 - discrete case 196
- Standard normal distribution
 - multivariate 102, 317
 - univariate 57
- Statistical model
 - nonparametric 371
 - parametric 374
- Stirling's formula
 - for n -factorial 60
 - for the Γ -function 59
- Strong law of large numbers 342
- Student's t -distribution 228
- Success probability 30
- Sultan's dowry problem 283
- Symmetric difference 453
- Systematic error
 - of an estimator 428

- t -distribution 228
- t -test
 - one-sided 405
 - two-sided 405
- Theorem
 - Berry–Esséen 366
 - binomial 459
 - De Moivre–Laplace 359
 - Fisher's 393
 - Fubini's 477
 - multinomial 466
 - Poisson's limit 35
 - Rao–Cramér–Frechet 431
- 37% rule 287
- Three sigma rule 187

- Tossing a coin
 - infinitely often 193
- Two-envelope paradox 287
- Two-number problem 292
- Two-sample t-tests 410
- Two-sample Z-tests 408
- Type I error 378
- Type II error 378

- Unbiased estimator 424
- Uncorrelated random variables 274
- Uniform distribution
 - on a finite set 21
 - on a set in \mathbb{R}^n 87
 - on an interval 52
- Uniformly best estimator 430
- Union
 - of sets 451

- Upper limit
 - of events 333

- Vandermonde's identity 460
- Variance
 - of a random variable 261
 - properties 262
- Volume
 - n -dimensional 83

- Weak law of large numbers 341

- \mathbb{Z} , integers 451
- Z-test
 - one-sided 400
 - two-sided 401