

Lecture Notes for Math 447 - Probability

Student edition with proofs

Michael Fochler
Department of Mathematics
Binghamton University

Last update: May 6, 2024

Contents

| | | |
|----------|--|------------|
| 1 | Some Preliminaries | 4 |
| 1.1 | About This Document | 4 |
| 1.2 | A First Look at Probability | 4 |
| 2 | Sets, Numbers, Sequences and Functions | 15 |
| 2.1 | Sets – The Basics | 15 |
| 2.2 | The Proper Use of Language in Mathematics: Any vs All, etc | 22 |
| 2.3 | Numbers | 23 |
| 2.4 | Functions and Sequences | 26 |
| 2.5 | Cartesian Products | 34 |
| 2.6 | A Few Things You May Not Have Learned in Calculus | 37 |
| 2.7 | Exercises for Ch.2 | 38 |
| 2.7.1 | Exercises for Sets | 38 |
| 2.8 | Addenda to Ch.2 – EMPTY | 39 |
| 2.9 | Blank Page after Ch.2 | 40 |
| 3 | The Probability Model | 41 |
| 3.1 | Probability Spaces | 41 |
| 3.2 | Conditional Probability and Independent Events | 55 |
| 3.3 | Preimages and Indicator Functions | 60 |
| 3.4 | Random Variables and their Probability Distributions | 65 |
| 4 | Combinatorial Analysis | 71 |
| 4.1 | The Multiplication Rule | 71 |
| 4.2 | Permutations | 73 |
| 4.3 | Combinations, Binomial and Multinomial Coefficients | 74 |
| 5 | More on Probability | 84 |
| 5.1 | Total Probability and Bayes Formula | 84 |
| 5.2 | Random Sampling and Urn Models With and Without Replacement | 86 |
| 6 | Discrete Random Variables and Random Elements | 91 |
| 6.1 | Probability Mass Function and Expectation | 91 |
| 6.2 | Bernoulli Variables and the Binomial Distribution | 100 |
| 6.3 | Geometric + Negative Binomial + Hypergeometric Distributions | 102 |
| 6.4 | The Poisson Distribution | 108 |
| 6.5 | Moments, Central Moments and Moment Generating Functions | 110 |
| 6.6 | Exercises for Ch.6 | 113 |
| 7 | Continuous Random Variables | 115 |
| 7.1 | Cumulative Distribution Function of a Random Variable | 115 |
| 7.2 | Continuous Random Variables and Probability Density Functions | 116 |
| 7.3 | Expected Value, Variance and MGF of a Continuous Random Variable | 121 |
| 7.4 | The Uniform Probability Distribution | 128 |
| 7.5 | The Normal Probability Distribution | 132 |
| 7.6 | The Gamma Distribution | 135 |

| | | |
|-----------|--|------------|
| 7.7 | The Beta Distribution | 138 |
| 7.8 | Inequalities for Probabilities | 139 |
| 8 | Multivariate Probability Distributions | 143 |
| 8.1 | Multivariate CDFs, PMFs and PDFs | 143 |
| 8.2 | Marginal and Conditional Probability Distributions | 146 |
| 8.3 | Independence of Random Variables and Discrete Random Elements | 148 |
| 8.4 | The Multivariate Uniform Distribution | 154 |
| 8.5 | The Expected Value of a Function of Several Random Variables | 155 |
| 8.6 | The Covariance of Two Random Variables | 158 |
| 8.7 | Conditional Expectations and Conditional Variance | 165 |
| 8.7.1 | The Conditional Expectation With Respect to an Event ★ | 165 |
| 8.7.2 | The Conditional Expectation w.r.t a Random Variable or Random Element | 168 |
| 8.7.3 | Conditional Expectations as Optimal Mean Squared Distance Approximations | 171 |
| 8.8 | The Multinomial Probability Distribution | 174 |
| 8.9 | Order Statistics | 177 |
| 8.10 | The Bivariate Normal Distribution (Optional) | 188 |
| 9 | Functions of Random Variables and their Distribution | 190 |
| 9.1 | The Method of Distribution Functions | 190 |
| 9.2 | The Method of Transformations in One Dimension | 196 |
| 9.3 | The Method of Transformations in Multiple Dimension | 200 |
| 9.4 | The Method of moment-generating Functions | 206 |
| 10 | Limit Theorems | 213 |
| 10.1 | Four Kinds of Limits for Sequences of Random Variables | 214 |
| 10.2 | Two Laws of Large Numbers | 220 |
| 10.3 | Sampling Distributions | 222 |
| 10.4 | The Central Limit Theorem | 230 |
| 11 | Sample Problems for Exams | 237 |
| 11.1 | Practice Midterm 1 for Math 447 - Chris Haines | 237 |
| 12 | Other Appendices | 239 |
| 12.1 | Greek Letters | 239 |
| 12.2 | Notation | 239 |
| | References | 240 |
| | List of Symbols | 241 |
| | Index | 243 |

History of Updates:


| Date | Topic |
|------------|----------|
| 2020-12-23 | Created. |

1 Some Preliminaries

1.1 About This Document

These lecture notes are supporting material to the required text of this Math 447 course on probability theory. This text is [11] Wackerly, D. and Mendenhall, W. and Scheaffer, R.L.: Mathematical Statistics with Applications, 7th edition.

At this point in time (December, 2023) it focuses on some of the foundations of probability theory which cannot be found at a sufficient level of generality in that text. Examples are preimages and σ -algebras. It has not been determined at this point in time what further topics will be included at some future time.

Note the uses of the symbol  for material that will not appear on exams, quizzes and other graded assignments. Unless you see this symbol in a footnote, please note that I will utilize such material and build on it in my lectures. Thus, you should understand this material well enough to follow my lectures, even though you will not be directly tested on it.

Also we use colored boxes according to the following. Generally speaking,

These boxes contain important definitions or parts thereof.

These boxes contain important theorems and propositions or parts thereof.

These boxes contain other kinds of important items that are worth while to know.

1.2 A First Look at Probability

“All models are wrong, but some are useful”.

Attributed to the statistician George E. P. Box
(1919–2013)



This quote certainly applies to the probabilistic models and the role they play in answering statisti-

cal questions such as

- How do I collect data to predict next month's average unemployment rate?
- What is the risk that this prediction will be off by more than 0.5 percent?

You probably agree that we also could have formulated the second question as follows.

- What is the probability that this prediction will be off by more than 0.5 percent?

It is not easy to find a satisfactory answer to that question and it will depend on the assumptions that go into your model. We will consider probability in much simpler settings.

Example 1.1 (Empirical probability). The concept of probability serves as a model for quantifying how likely an event will happen that depends on chance. When we say that the probability of obtaining an even number when rolling a die equals 0.5, then we mean the following.

Assume that

- X_1 denotes the action of rolling that die for the first time.
- X_2 denotes the action of rolling that die for the second time.
- $\dots X_k$ denotes the action of rolling that die for the k th time.

We expect in the long run, i.e., for large k , close to half of X_1, X_2, \dots, X_k result in an even outcome.

We formulate this in the language of mathematics as follows:

- We write P for probability.
- We write $\{2, 4, 6\}$ for the event that rolling the die results in a 2 or a 4 or a 6, i.e., in an even outcome.
- We write n_k for the number of outcomes during those k rolls that result in a 2 or a 4 or a 6.
- We define $P\{2, 4, 6\} = \lim_{k \rightarrow \infty} \frac{n_k}{k}$ and call this limit the probability of the event $\{2, 4, 6\}$.

We expect this particular limit to be 0.5.

Any event associated with the roll of a die can be expressed as a list of integers $1 \leq i_1 < i_2 < \dots < i_m \leq 6$ if we interpret it in the following sense: The roll results in one of the numbers in that list.

Let us write Ω (this symbol denotes the Greek capital letter Omega) ¹ for the set of all potential **outcomes**. It is customary to drop the word "potential" and refer to the elements of Ω simply as outcomes. In this example, $\Omega = \{1, 2, 3, 4, 5, 6\}$ and the outcomes are $1, 2, \dots, 6$.

We stated above that events associated with the roll of a die can be expressed as a list $1 \leq i_1 < i_2 < \dots < i_m \leq 6$. This list corresponds to the set $A = \{i_1, i_2, \dots, i_m\}$. Observe that $A \subseteq \Omega$ i.e., each element of A also belongs to Ω . ² We call all subsets of Ω **events**.

We can apply the steps we used to determine $P\{2, 4, 6\}$ to ANY event $A \subseteq \Omega$. Now, n_k denotes the number of outcomes during the first k rolls that result in a number that is listed in A . We define

$$(1.1) \quad P(A) = \lim_{k \rightarrow \infty} \frac{n_k}{k}.$$

To be precise, this formula denotes the **empirical probability** of the event A .

Observe that the assignment $A \mapsto P(A)$ of (1.1) satisfies the following for all subsets A of Ω :

¹For a list of all Greek letters see Section 12.1 (Greek Letters) on page 239.

²See Definition 2.3 (Subsets and supersets) on p.16 on page 239.

- $0 \leq P(A) \leq 1$.
- $P(\emptyset) = 0$, since $n_k = 0$ for all k . (\emptyset is empty set which contains no elements.)
- $P(\Omega) = 1$, since $n_k = k$ for all k .
- If the subsets A, B of Ω have no elements in common (we speak of **mutually disjoint** sets), then the union $P(A \cup B)$ satisfies

$$(\star) \quad P(A \cup B) = P(A) + P(B).$$

To see the validity of (\star) , let $n_k(A)$ be the number of times an outcome in A is observed during k trials, and let $n_k(B)$ be defined likewise for B . Since an outcome ω is in $A \cup B$ if and only if ω either belongs to A or to B , we have $n_k(A \cup B) = n_k(A) + n_k(B)$, hence,

$$P(A \cup B) = \lim_{k \rightarrow \infty} \frac{n_k(A \cup B)}{k} = \lim_{k \rightarrow \infty} \frac{n_k(A)}{k} + \lim_{k \rightarrow \infty} \frac{n_k(B)}{k} = P(A) + P(B).$$

Note the following about the nature of the formula $P(A) = \lim_{k \rightarrow \infty} \frac{n_k}{k}$ for subsets A of Ω .

- It is a function $A \mapsto P(A) = \lim_{k \rightarrow \infty} \frac{n_k}{k}$ the same way $x \mapsto f(x) = x^2 + 4$ is a function.
- We are familiar with the latter: It assigns to each argument x (which happens to be a real number) the function value $f(x)$, also a real number. For example, $f(3) = 3^2 + 4 = 13$.
- The function $A \mapsto P(A)$ is harder to deal with only because its arguments A are not numbers or vectors of such numbers. Rather, those arguments are events, i.e., sets. □

You are strongly encouraged to take a first look at Section 2.4 (Functions and Sequences). It should be thorough enough to understand the following:

- The assignment $A \mapsto P(A)$ discussed at the end of Example 1.1 constitutes a function

$$P : \{ \text{all subsets of } \Omega \} \rightarrow [0, 1]$$

in the sense of Definition 2.14 on p.27.

Remark 1.1. There are some issues with (1.1) as a definition of $P(A)$.

What if the limit $\lim_{k \rightarrow \infty} n_k/k$ does not exist? For example, the following is very unlikely but not impossible.

Let ω_k denote the outcome of the k th roll of the die. Assume that we obtain the following sequence of outcomes:

- $\omega_1 = 1$.
- From now on, only the number 6 appears until $n_k/k > 5$. We write $K(1)$ for that index k .
- From now on, only the number 1 appears until $n_k/k < 2$. We write $K(2)$ for that index k .
- From now on, only the number 6 appears until $n_k/k > 5$. We write $K(3)$ for that index k .
- From now on, only the number 1 appears until $n_k/k < 2$. We write $K(4)$ for that index k .
- and so on

The resulting sequence $K(1) < K(2) < K(3) < \dots$ satisfies the following:³

³A strict proof can be obtained by using the fact that the limit of a sequence does not depend its first k members, no matter how big k may be chosen.

- There are infinitely many indices $k = K(1), K(3), K(5), \dots$ such that $\frac{n_k}{k} > 5$.
- There are infinitely many indices $k = K(2), K(4), K(6), \dots$ such that $\frac{n_k}{k} < 2$.

Accordingly, $\lim_{k \rightarrow \infty} \frac{n_k}{k}$ does not exist and we were not able to determine $P(A)$.

But there are issues even if that limit exists. Consider again the event $A = \{2, 4, 6\}$. Let us assume that, by some freak of nature, all outcomes ω_k are 4.⁴ Accordingly, we declare that $P(A) = 1$. The teamleader has doubts about this result and asks for a repetition of the experiment. This time all outcomes ω_k are either 3 or 5.

What to do? Should we decide that $P(A) = 0$? Should the experiment be repeated once more? How about settling on the average $(1 + 0)/2 = 1/2$?

You may decide that this is a completely fictitious example without any bearing on reality and the author will agree. However, you should consider the following:

- The infinite repetition of an action such as rolling a die is in itself fictitious and so is the concept of the limit of a (infinite) sequence.
- In the real world the determination of probabilities $P(A)$ often is based on (1.1) as follows: It is decided to conduct an experiment of k trials. The larger this number k is chosen, the more confidence we will have that $P(A)$ is a good enough APPROXIMATION of the likelihood that the event A happens.

However, there are other factors to consider that will limit the size of k .

- The more repetitions, the longer it will take to obtain the result. If A is the event that the Old Faithful geyser in Yellowstone National Park erupts to a height of at least 150 feet and it is not possible for some reason to use the previously obtained records, then we must base the determination of $P(A)$ on a very small number of observations.
 - Money is another limiting factor. The more repetitions, the more it will cost to obtain the result.
-

Example 1.2 (Single roll of a die). To avoid the issues concerning the use of formula (1.1) (empirical probability) on p.5, we also could have used the concept of a fair die instead, i.e., a die for which each of the outcomes $1, 2, \dots, 6$ is equally likely, so each outcome must have the same likelihood (probability) of $1/6$. Since the even outcomes are 2, 4, 6, we obtain

$$(1.2) \quad P\{\text{even outcome}\} = P\{2, 4, 6\} = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = 0.5.$$

Note that fair dice do not exist in the real world. Matter of fact, if we had a sample of 1,000 dice and we were able to determine with infinite precision the probability that a throw of die $\#_k$ comes up even, chances are that we would obtain several different answers, due to imperfections in the manufacturing process.

But let us assume, for arguments sake, there is such a thing as a fair die. We model the **random action** of rolling such a fair die just once as follows.

⁴You will learn the following: If each j_1, j_2, \dots is a given potential outcome (an integer between 1 and 6), then $P\{\omega_1\} = j_1, P\{\omega_2\} = j_2, \dots, P\{\omega_k\} = j_k\} = (1/6)^k$. That number becomes very small for large k , since the sequence $(1/6)^k$ converges to zero. Nevertheless, $(1/6)^k > 0$ for each fixed k , so it is not impossible to obtain $\omega_k = 4$ for all k . (This is the case where $j_1 = j_2 = \dots = j_k = 4$ for all k).

- As in Example 1.1 (Empirical probability), on p.5, the set Ω of all (potential) outcomes is $\{1, 2, 3, 4, 5, 6\}$.
- We associate with each outcome $\omega \in \Omega$ the probability $P(\{\omega\}) = 1/6$.
- Let $A \subseteq \Omega$, i.e., A is a subset of Ω , i.e., each element of A also belongs to Ω .
- For each outcome $\omega \in \Omega$ there is a corresponding event $\{\omega\} \subseteq \Omega$.⁵ Such “atomar” events also are referred to as outcomes.
- We generalize (1.2) and associate with each event $A \subseteq \Omega$ the probability

$$(1.3) \quad P(A) = \sum_{\omega \in A} P(\{\omega\}).$$

Here, $\sum_{\omega \in A} P(\{\omega\})$ means that we sum up all those expressions $P(\{\omega\})$ that satisfy $\omega \in A$.

- For example, let $A = \{2, 4, 6\}$ and $B = \{\omega \in \Omega : \omega > 4\}$. Thus, A is the event of rolling an even outcome and B is that of rolling a 5 or 6. Then,

$$P(A) = \sum_{\omega \in A} P(\{\omega\}) = P(\{2\}) + P(\{4\}) + P(\{6\}) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2},$$

$$P(B) = P(\{5, 6\}) = P(\{5\}) + P(\{6\}) = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}.$$

It is customary to write $P\{...\}$ for $P(\{...\})$. Thus, the last equation can also be written as

$$P(B) = P\{5, 6\} = P\{5\} + P\{6\} = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}.$$

The assignment $A \mapsto P(A)$ satisfies the following for all subsets A of Ω :

- $0 \leq P(A) \leq 1$.
- $P(\emptyset) = 0$, since $n_k = 0$ for all k . (\emptyset is empty set which contains no elements.)
- $P(\Omega) = 1$, since $n_k = k$ for all k .
- If the subsets A, B of Ω have no elements in common (we speak of **mutually disjoint** sets), then the union $P(A \cup B)$ satisfies

$$(\star) \quad P(A \cup B) = P(A) + P(B).$$

Note that this was also the case for the assignment $A \mapsto P(A)$ in Example 1.1 (Empirical probability). \square

Example 1.3 (Two rolls of a die). Consider what happens if two fair dice are rolled. The set of outcomes is

$$\Omega = \{1, 2, \dots, 6\}^2 = \{1, 2, \dots, 6\} \times \{1, 2, \dots, 6\} = \{\omega : \omega = (i, j) \text{ and } i, j = 1, 2, \dots, 6\}.$$

⁵Such sets of size 1 are often called **singleton sets** or simply **singletons**.

- We make a willful decision to consider the outcomes (i, j) and (j, i) different for $i \neq j$. For example, if die #1 is red and #2 is white, we distinguish between the outcome of a red 2 and a white 5 and that of a red 5 and a white 2. Then Ω consists of 36 outcomes

$$(1, 1), (1, 2), \dots, (1, 6), (2, 1), (2, 2), \dots, (2, 6), \dots, (6, 1), (6, 2), \dots, (6, 6)$$

and each outcome $\omega \in \Omega$ has probability $P\{\omega\} = 1/36$.

- Just as in Example 1.2 those probabilities of the outcomes determine the probability of any event $A \in \Omega$ as was the case in formula (1.3) by

$$(1.4) \quad P(A) = \sum_{\omega \in A} P(\{\omega\}).$$

- For example, if $A = \{\text{die \#1 shows a 4}\} = \{(4, j) : j = 1, 2, \dots, 6\}$ then

$$\begin{aligned} P(A) &= \sum_{\omega \in A} P(\{\omega\}) = \sum_{(i,j) \in A} P(\{(i, j)\}) \\ &= P\{(4, 1)\} + P\{(4, 2)\} + \dots + P\{(4, 6)\} = 6 \left(\frac{1}{36} \right) = \frac{1}{6}. \end{aligned}$$

- As in examples 1.1 (Empirical probability) and 1.2 (Single roll of a die), we have a function $A \mapsto P(A)$ which assigns the events $A \subseteq \Omega$ to their probabilities $P(A)$ □

The assignment $A \mapsto P(A)$ in (1.4) satisfies the following for all subsets A of Ω :

- $0 \leq P(A) \leq 1$.
- $P(\emptyset) = 0$, since $n_k = 0$ for all k . (\emptyset is empty set which contains no elements.)
- $P(\Omega) = 1$, since $n_k = k$ for all k .
- If the subsets A, B of Ω have no elements in common (we speak of **mutually disjoint** sets), then the union $P(A \cup B)$ satisfies

$$(\star) \quad P(A \cup B) = P(A) + P(B).$$

Note that this was also the case for the assignment $A \mapsto P(A)$ in Example 1.1 (Empirical probability) and in Example 1.2. □

Example 1.4 (Sum of two die rolls). Consider what happens if two fair dice are rolled and we are interested in the sum of points obtained that way. For example,

- the outcome 8 is obtained when either of the following are rolled:
 - a 2 and a 6 □ a 3 and a 5 □ a 4 and a 4 □ a 5 and a 3 □ a 6 and a 2.
- the outcome 5 is obtained when either of the following are rolled:
 - a 1 and a 4 □ a 2 and a 3 □ a 3 and a 2. □ a 4 and a 1.

- Now, the set of outcomes is

$$\Omega = \{2, 3, \dots, 11, 12\}.$$

Since a roll of two dice has 36 outcomes $(1, 1), \dots, (6, 6)$ and each of those has probability $1/36$ (see Example 1.3), it follows for the outcomes 8 and 5 that

- $P(\{8\}) = \frac{5}{36}; \quad P(\{5\}) = \frac{4}{36}.$

Here is the complete list of outcome probabilities $P(\{\omega\})$:

$$(1.5) \quad \begin{aligned} P(\{2\}) = P(\{12\}) &= \frac{1}{36}; & P(\{3\}) = P(\{11\}) &= \frac{2}{36}; & P(\{4\}) = P(\{10\}) &= \frac{3}{36}; \\ P(\{5\}) = P(\{9\}) &= \frac{4}{36}; & P(\{6\}) = P(\{8\}) &= \frac{5}{36}; & P(\{7\}) &= \frac{6}{36}. \end{aligned}$$

- In the previous two examples each outcome had the same probability. We see that this is not the case for the sum of points obtained when rolling two dice.
- As in the previous examples, the probability of any event $A \in \Omega$ is obtained as the sum $P(\{\omega\})$ over all outcomes ω of the event:

$$(1.6) \quad P(A) = \sum_{\omega \in A} P(\{\omega\}).$$

- For example, if $A = \{\text{the sum is between 8 and 11}\}$, then

$$\begin{aligned} P(A) &= \sum_{\omega \in A} P(\{\omega\}) = \sum_{\omega=8}^{11} P(\{\omega\}) \\ &= P\{8\} + P\{9\} + P\{10\} + P\{11\} = (5 + 4 + 3 + 2) \left(\frac{1}{36}\right) = \frac{7}{18}. \end{aligned}$$

- As in examples 1.1, 1.2 and 1.3, we have a function $A \mapsto P(A)$ which assigns the events $A \subseteq \Omega$ to their probabilities $P(A)$ □

Let us examine what the examples we have studied so far have in common.

Remark 1.2. In the examples given so far a probability $P(A)$ was assigned to each event $A \subseteq \Omega$. In each case this assignment $A \mapsto P(A)$ satisfies the following.

$$(1.7) \quad 0 \leq P(A) \leq 1.$$

$$(1.8) \quad P(\emptyset) = 0. \quad \text{Here } \emptyset \text{ denotes the empty set which contains no elements.}$$

$$(1.9) \quad P(\Omega) = 1.$$

If the subsets A, B of Ω have no elements in common (we speak of **mutually disjoint** sets), then the union $P(A \cup B)$ satisfies

$$(1.10) \quad P(A \cup B) = P(A) + P(B).$$

- The probabilist likes to speak of the **probability space** Ω , since it comes with a **probability measure** (WMS: probability function), $A \mapsto P(A)$, which assigns to the events A of Ω , the probability $P(A)$ that this event might “occur” or “happen”.
- Statisticians tend to speak of the **sample space** S (that’s S as in **S**ample). An element s of S still is referred to as an **outcome**, but some, like WMS, also call the elements of S the **sample points** of S .

We translate some of the examples already encountered into the language of sample spaces and sample points.

- In example 1.2 (Single roll of a die) on p.7, the elements 1, 2, 3, 4, 5, 6 of the the sample space $S = \{1, 2, \dots, 6\}$ (the outcomes or sample points of S represent all possible “samples” that can be obtained from the single roll of a die. Note that each of those six potential samples has size $n = 1$.
- In example 1.3 (Two rolls of a die) on p.8, the sample points $(1, 1), (1, 2), \dots, (6, 5), (6, 6)$ of the the sample space $S = \{1, 2, \dots, 6\}^2$ represent all possible samples that can be obtained from two rolls of a die. Note that each of those 36 potential samples has size $n = 2$.
- In example 1.4 (Sum of two die rolls) on p.9, the sample points 2, 3, \dots , 12 of the the sample space $S = \{2, 3, \dots, 12\}$ represent all possible “samples” that can be obtained from adding the points of two roll of a die. Note that each of those 11 potential samples has size $n = 1$. \square

Example 1.5. This example needs more computational skills than the ones we have encountered so far.

- To understand whether a traffic light works as expected, the following experiment is conducted. 200 cars are observed and a record is made for each one of those cars whether it reached the intersection on red, green or yellow.
- This “**sampling action**” of observing those 200 cars results in ONE sample point of size 200. Its actual outcome depends on chance
- Once the experiment is completed, the sampling action has resulted in a **realization** of the sampling action (the SPECIFIC sample point that was obtained). If we write r for red, g for green, y for yellow, this realization might be, e.g., $\{r, r, y, g, g, g, r, y, \dots, r\}$.
- Once that realization has been obtained, the sampling action has lost its random character.
- The sample space S of all (potential) sample points for this experiment is huge: It contains 3^{200} sample points. This will be discussed in Chapter 4 (Combinatorial Analysis)
- Each event $A \subseteq S$ comes with a probability $P(A)$ and one can show that the assignment $A \mapsto P(A)$ satisfies the formulas (1.7) – (1.10) of Remark 1.2 on p.10. \square

Here is a formal definition of probability. It is based on the formulas (1.7) – (1.10) of Remark 1.2 on p.10. **This definition is PRELIMINARY and will be amended!** Since sampling does not play a role, we will talk about a probability space Ω rather than a sample space S .

Definition 1.1 (Probability measure - Preliminary Definition). A **probability measure** P on a set Ω is a function ⁶ which assigns to each subset A of Ω a real number $P(A)$ between 0 and 1 as follows.

⁶we’ll review functions briefly in Section 2.1 (Sets, Numbers, Sequences and Functions) on page 15.

- (a) $P(\emptyset) = 0$ and $P(\Omega) = 1$. Here \emptyset denotes the empty set which contains no elements.
 (b) If the subsets A, B of Ω have no elements in common, then probability is **additive**:

$$(1.11) \quad P(A \cup B) = P(A) + P(B). \quad \square$$

Remark 1.3.

- Note that Definition 1.1 makes no mention about how one should interpret the number $P(A)$!

For example, one could take a fair coin and define $P(H) := 0.1$, $P(T) := 0.9$. Here, H denotes Heads and T denotes Tails. This defines a probability $A \mapsto P(A)$ on the sample space $S := \{H, T\}$.
⁷

If this example strikes you as nonsensical, here is a model used by Wall Street that uses a probability measure in which the probability of an event is different from the chance that this event will happen.

The so called binomial asset model is a probabilistic model to determine today's price of a stock option which will be exercised at some future point in time.⁸ In this model, trading of a specific stock (e.g., IBM or Amazon), happens at times $0, 1, 2, \dots$. There are only two possible ways that stock price can change and there are two "real world" probabilities, one for each possibility:

- $p_u := P\{\text{the price of a share of stock changes by the factor } u\}$.
- $p_d := P\{\text{the price of a share of stock changes by the factor } d < u = 1 - p_u\}$.

These two numbers p_u and p_d are sufficient to determine a probability space Ω and probability measure P for trading in that stock.

Strangely enough, p_u and p_d are replaced by the so-called risk-neutral probability \tilde{p}_u and \tilde{p}_d which are sufficient to determine an alternate probability measure \tilde{P} on that same probability space Ω .

Even stranger, the real world probability measure P has no bearing on the determination of \tilde{P} , i.e., of \tilde{p}_u and \tilde{p}_d .⁹ And yet, even though \tilde{p}_u and \tilde{p}_d do not reflect the actual probabilities that govern the stock price, they are used to set today's price of an option on that stock that can be redeemed only, say, 90 days from today. \square

In the next example we combine Example 1.3 and Example 1.4.

Example 1.6. When computing the outcome probabilities of the sum of points obtained by rolling two dice, we argued with a result obtained in Example 1.3: There the probability of an outcome (i, j) was $1/36$ for all $i, j = 1, 2, \dots, 6$. It should not be surprising that there is a connection between the probability models of those examples. Each one of those examples had a set of outcomes which we denoted Ω and a function $P : A \mapsto P(A)$ which associated a probability $P(A)$ with each event $A \subseteq \Omega$. Since this example deals with both outcome sets and both probability assignments, we must change our notation. We proceed as follows.

⁷To complete the definition of P , we define $P(S) := P(H) + P(T) = 1$ and $P(\emptyset) := 0$.

⁸Since this is not a course on probabilistic finance, we must refer you to the literature for details. Some references are [10] Shreve, Steve: Stochastic Calculus for Finance I: The Binomial Asset Pricing Model, [2] Björk, Thomas: Arbitrage Theory in Continuous Time and this author's [Math 454 lecture notes](#) (Spring 2023).

⁹Rather, the interest earned by depositing money in a bank plays a major role.

- We keep the symbols used in Example 1.3 and define

$$\begin{aligned}\Omega &:= \{1, 2, \dots, 6\} \times \{1, 2, \dots, 6\} = \{\omega : \omega = (i, j) \text{ and } i, j = 1, 2, \dots, 6\}, \\ P\{(i, j)\} &:= \frac{1}{36} \quad \text{for } i, j = 1, 2, \dots, 6.\end{aligned}$$

- For the outcome set and probability assignment of Example 1.4, we write

$$\begin{aligned}\Omega' &:= \{2, 3, \dots, 11, 12\}, \\ P'\{2\} &:= P'\{12\} := \frac{1}{36}, \quad P'\{3\} := P'\{11\} := \frac{2}{36}, \dots \quad \text{See (1.5) on p.10.}\end{aligned}$$

Note that $P'\{k\}$ equals the probability that the sum of the two die rolls equals k , since both probabilities are defined by the same formula, (1.5).

We establish a relationship between the “probability spaces” (Ω, P) and (Ω', P') as follows.

Let $(i, j) \in \Omega$, i.e., i is the outcome of rolling die #1 and j is that of rolling die #2. The assignment

$$(1.12) \quad (i, j) \mapsto Y(i, j) := i + j$$

then associates with this outcome an integer between 2 and 12, i.e., an outcome in Ω' .¹⁰

For $k \in \Omega'$, let

$$(1.13) \quad A_k := \{(i, j) \in \Omega : Y(i, j) = k\}.$$

We claim that

$$(1.14) \quad P'\{k\} = P(A_k).$$

That claim is proved by the following chain of equations:

$$\begin{aligned}P'\{k\} &= \text{probability that total points of both dice rolled is } k \\ &= \frac{1}{36} \cdot \text{the number of elements in } A_k. \\ &= \sum_{\omega \in A_k} \frac{1}{36} = \sum_{(i, j) \in A_k} P\{(i, j)\} = P(A_k).\end{aligned}$$

Equation #1 is merely the definition of $P'\{k\}$, #2 follows from (1.5) on p.10 and the definition of A_k , #3 is true by the definition of $\sum \dots$, #4 is just the definition of $P\{(i, j)\}$ and the last equation follows from (1.4) on p.9. We can rewrite (1.14) as follows.

- If $B = \{k\}$ and $A = \{(i, j) \in \Omega : Y(i, j) \in B\}$, then $P'(B) = P(A)$

Next, let $k \neq \ell$. Note that A_k and A_ℓ have no outcomes in common. Both P and P' are probabilities and satisfy the additivity formula (1.11) on p.12. Thus,

$$(1.15) \quad \begin{aligned}P'\{k, \ell\} &= P(A_k \uplus A_\ell) = P\{(i, j) \in \Omega : Y(i, j) = k \text{ or } Y(i, j) = \ell\} \\ &= P\{(i, j) \in \Omega : Y(i, j) \in \{k, \ell\}\}\end{aligned}$$

¹⁰Hence, this assignment is a function $Y : \Omega \rightarrow \Omega'$ in the sense of Definition 2.14 on p.27.

We can rewrite (1.15) as follows.

- If $B = \{k, \ell\}$ and $A = \{(i, j) \in \Omega : Y(i, j) \in B\}$, then $P'(B) = P(A)$

Generalization of (1.14) and (1.15): We will see the following in Section 3.4 (Random Variables and their Probability Distributions)

$$(1.16) \quad \text{If } B \subseteq \Omega' \text{ and } A = \{(i, j) \in \Omega : Y(i, j) \in B\}, \text{ then } P'(B) = P(A).$$

We will then call a function Y that assigns elements of Ω to elements of Ω' a random element and speak of the probability measure P' on Ω' , which is uniquely specified by the probability measure P on Ω and the function Y , as the distribution of Y . \square

2 Sets, Numbers, Sequences and Functions

Introduction 2.1. \square

The student should read this chapter carefully, with the expectation that it contains material that they are not familiar with, as much of it will be used in lecture without comment. Very likely candidates are power sets, a function $f : X \rightarrow Y$ where domain X and codomain Y are part of the definition.

2.1 Sets – The Basics

An entire book can be filled with a mathematically precise theory of sets. For our purposes the following “naive” definition suffices:

Definition 2.1 (Sets).

- A **set** is a collection of stuff called **members** or **elements** which satisfies the following rules: The order in which you write the elements does not matter and if you list an element two or more times then **it only counts once**.
- We write $x_1 \in X$ to denote that an item x_1 is an element of the set X and $x_2 \notin X$ to denote that an item x_2 is not an element of the set X .
- Occasionally we are less formal and write x_1 **in** X for $x_1 \in X$ and x_2 **not in** X for $x_2 \notin X$.

We write a set by enclosing within curly braces the elements of the set. This can be done by listing all those elements or giving instructions that describe those elements. For example, to denote by X the set of all integer numbers between 18 and 24 we can write either of the following:

$$X := \{18, 19, 20, 21, 22, 23, 24\} \quad \text{or} \quad X := \{n : n \text{ is an integer and } 18 \leq n \leq 24\}$$

Both formulas clearly define the same collection of all integers between 18 and 24. On the left the elements of X are given by a complete list, on the right **setbuilder notation**, i.e., instructions that specify what belongs to the set, is used instead.

For the above example we have $20 \in X$, $27 - 6 \in X$, $38 \notin X$, ‘Jimmy’ $\notin X$.

It is customary to denote sets by capital letters and their elements by small letters We try to adhere to this convention as much as possible. \square

Example 2.1. We looked in the introduction at the set $\Omega = \{1, 2, 3, 4, 5, 6\}$ of potential outcomes for the roll of a die. Then $3 \in \Omega$, $5 \in \Omega$, $-2 \notin \Omega$, $2.34 \notin \Omega$. \square

Example 2.2 (No duplicates in sets). The following collection of alphabetic letters is a set:

$$S_1 = \{a, e, i, o, u\}$$

and so is this one:

$$S_2 = \{a, e, e, i, i, i, o, o, o, o, u, u, u, u, u\}$$

Did you notice that those two sets are equal? \square

Remark 2.1. The symbol n in the definition of $X = \{n : n \text{ is an integer and } 18 \leq n \leq 24\}$ is a **dummy variable** in the sense that it does not matter what symbol you use. The following sets all are equal to X :

$$\begin{aligned} &\{x : x \text{ is an integer and } 18 \leq x \leq 24\}, \\ &\{\alpha : \alpha \text{ is an integer and } 18 \leq \alpha \leq 24\}, \\ &\{\mathfrak{J} : \mathfrak{J} \text{ is an integer and } 18 \leq \mathfrak{J} \leq 24\} \quad \square \end{aligned}$$

Definition 2.2 (empty set).

\emptyset denotes the **empty set**. It is the set that does not contain any elements. \square

Definition 2.3 (subsets and supersets).

- We say that a set A is a **subset** of the set B and we write $A \subseteq B$ if any element of A also belongs to B . Equivalently we say that B is a **superset** of the set A and we write $B \supseteq A$. We also say that B includes A or A is included by B . Note that $A \subseteq A$ and $\emptyset \subseteq A$ is true for any set A .
- If $A \subseteq B$ but $A \neq B$, i.e., there is at least one $x \in B$ such that $x \notin A$, then we say that A is a **strict subset** or a **proper subset** of B . We write " $A \subsetneq B$ ". Alternatively we say that B is a **strict superset** or a **proper superset** of A and we write " $B \supsetneq A$ ".

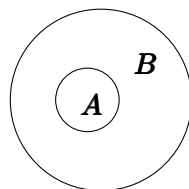


Figure 2.1: Set inclusion: $A \subseteq B$, $B \supseteq A$ \square

Remark 2.2. (a) We STRONGLY discourage the use of " $A \subset B$ " in place of " $A \subsetneq B$ " and of " $B \supset A$ " in place of " $A \supsetneq B$ ". These are outdated symbols for $A \subseteq B$ and $A \supseteq B$

(b) Two sets A and B are equal means that they both contain the same elements. In other words, since $U \subseteq V$ means that the set V contains all elements of the set U ,

$$(2.1) \quad A = B \Leftrightarrow [A \subseteq B \text{ and } B \subseteq A].$$

In the above, " \Leftrightarrow " denotes the phrase "if and only if": The expression to the left (" $A = B$ ") means the same as the expression to the right (" $A \subseteq B$ and $B \subseteq A$ "). The square brackets only serve to clarify that everything inbetween belongs to the scope of the right-hand side of " \Leftrightarrow ". \square

Definition 2.4 (unions, intersections and disjoint unions). Given are two arbitrary sets A and B . No assumption is made that either one is contained in the other or that either one is not empty!

- The **union** $A \cup B$ (pronounced "A union B") is defined as the set of all elements which belong to at least one of A, B .
- The **intersection** $A \cap B$ (pronounced "A intersection B") is defined as the set of all elements which belong to both A and B .
- We call A and B **disjoint**, also **mutually disjoint**, if $A \cap B = \emptyset$. We then often write $A \uplus B$ (pronounced "A disjoint union B") rather than $A \cup B$.

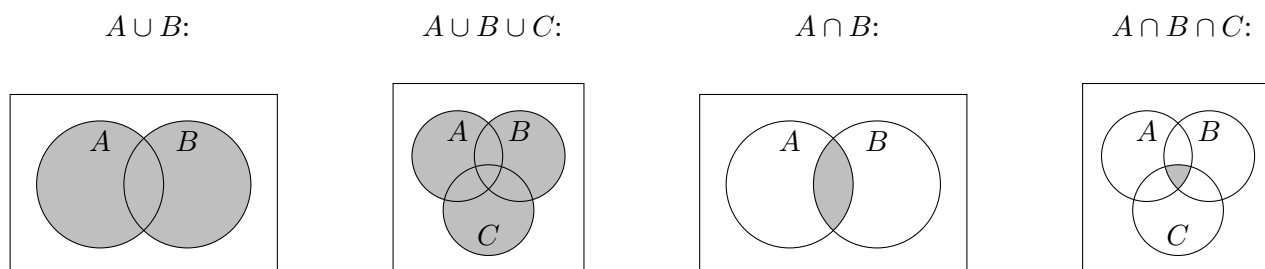


Figure 2.2: Union and intersection of sets

Since $A \cup B = B \cup A$ and $A \cap B = B \cap A$ and $A \uplus B = B \uplus A$, it is obvious how to specify those operations to any finite or infinite collection of sets. Let J be a nonempty, finite or infinite subset of the set $\mathbb{Z} = \{0, \pm 1, \pm 2, \pm 3, \dots\}$ of all integers. In particular, $J = \mathbb{Z}$ is allowed. Assume that each $j \in J$ is associated with a set A_j .¹¹ We say that

- The **union** $\bigcup_{j \in J} A_j$ is defined as the set of all elements which belong to at least one A_j , where $j \in J$.
- The **intersection** $\bigcap_{j \in J} A_j$ is defined as the set of all elements which belong to each A_j , where $j \in J$.
- We call this collection of sets **disjoint**, also **mutually disjoint**, if $A_i \cap A_j = \emptyset$ whenever $i, j \in J$ and $i \neq j$. We then often write $\biguplus_{j \in J} A_j$ rather than $\bigcup_{j \in J} A_j$. \square

Remark 2.3. If $J = \{k_*, k_* + 1, k_* + 2, \dots, k^* - 1, k^*\}$, we also write

$$\bigcup_{j=k_*}^{k^*} A_j, \quad \bigcap_{j=k_*}^{k^*} A_j, \quad \biguplus_{j=k_*}^{k^*} A_j, \quad \text{for} \quad \bigcup_{j \in J} A_j, \quad \bigcap_{j \in J} A_j, \quad \biguplus_{j \in J} A_j.$$

¹¹You might call this a **collection** of sets A_i which is **indexed by** the elements of J and write $(A_j)_{j \in J}$ for this **indexed collection**.

If $J = \{k_*, k_* + 1, k_* + 2, \dots\}$, we also write

$$\bigcup_{j=k_*}^{\infty} A_j, \quad \bigcap_{j=k_*}^{\infty} A_j, \quad \bigoplus_{j=k_*}^{\infty} A_j, \quad \text{for} \quad \bigcup_{j \in J} A_j, \quad \bigcap_{j \in J} A_j, \quad \bigoplus_{j \in J} A_j.$$

Examples: If $I = \{-1, 0, 1, 2\}$, then $\bigcap_{i \in I} A_i = \bigcap_{i=-1}^2 A_i = A_{-1} \cap A_0 \cap A_1 \cap A_2$.

If $U = \{5, 6, 7, \dots\}$, then $\bigcup_{j \in U} C_j = \bigcap_{j=5}^{\infty} C_j = C_5 \cup C_6 \cup C_7 \cup \dots$. \square

Remark 2.4. Convince yourself that for any sets A, B and C .

$$(2.2) \quad A \cap B \subseteq A \subseteq A \cup B,$$

$$(2.3) \quad A \subseteq B \Rightarrow A \cap B = A \text{ and } A \cup B = B,$$

$$(2.4) \quad A \subseteq B \Rightarrow A \cap C \subseteq B \cap C \text{ and } A \cup C \subseteq B \cup C.$$

The symbol \Rightarrow stands for “allows us to conclude that”. So $A \subseteq B \Rightarrow A \cap B = A$ means “From the truth of $A \subseteq B$ we can conclude that $A \cap B = A$ is true”. Shorter: “From $A \subseteq B$ we can conclude that $A \cap B = A$ ”. Shorter: “If $A \subseteq B$, then it follows that $A \cap B = A$ ”. Shorter: “If $A \subseteq B$, then $A \cap B = A$ ”. More technical: $A \subseteq B$ implies $A \cap B = A$. \square

Definition 2.5 (set differences and symmetric differences). Given are two arbitrary sets A and B . No assumption is made that either one is contained in the other or contains any elements!

- The **difference set** or **set difference** $A \setminus B$ (pronounced “A minus B”) is defined as the set of all elements which belong to A but not to B :

$$(2.5) \quad A \setminus B := \{x \in A : x \notin B\}$$

- The **symmetric difference** $A \Delta B$ (pronounced “A delta B”) is defined as the set of all elements which belong to either A or B but not to both A and B :

$$(2.6) \quad A \Delta B := (A \cup B) \setminus (A \cap B) \quad \square$$

Definition 2.6 (Universal set).

Usually there always is a big set Ω that contains everything we are interested in and we then deal with all kinds of subsets $A \subseteq \Omega$. Such a set is called a “**universal**” set. \square

Example 2.3.

- (a) Often the context are the real numbers and their subsets. An appropriate universal set will then be \mathbb{R} .¹²
- (b) We will discuss at length why the set $\{1, 2, 3, 4, 5, 6\}$ can be considered a universal set in the context of rolling a die. See Section 1.2 (A First Look at Probability). \square

If there is a universal set, it makes perfect sense to talk about the complement of a set:

Definition 2.7 (Complement of a set). Let Ω be a universal set. The **complement** of a set $A \subseteq \Omega$ consists of all elements of Ω which do not belong to A . We write A^c . In other words:

$$(2.7) \quad A^c = \Omega \setminus A = \{\omega \in \Omega : \omega \notin A\} \quad \square$$

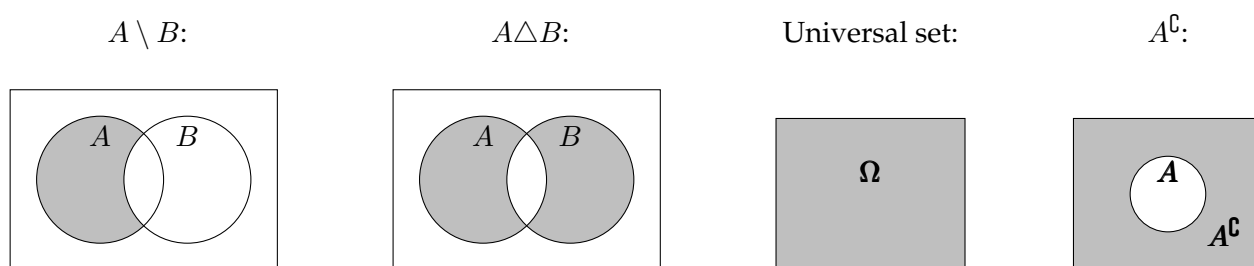


Figure 2.3: Difference, symmetric difference, universal set, complement

Remark 2.5. Note that for any kind of universal set Ω it is true that

$$(2.8) \quad \Omega^c = \emptyset, \quad \emptyset^c = \Omega. \quad \square$$

Example 2.4 (Complement of a set relative to the unit interval). Assume we are exclusively dealing with the unit interval, i.e., $\Omega = [0, 1] = \{x \in \mathbb{R} : 0 \leq x \leq 1\}$. Let $a \in [0, 1]$ and $\delta > 0$ and

$$(2.9) \quad A = \{x \in [0, 1] : a - \delta < x < a + \delta\}$$

the “ δ -neighborhood”¹³ of a (with respect to $[0, 1]$) because numbers outside the unit interval are not considered part of our universe). Then the complement of A is

$$A^c = \{x \in [0, 1] : x \leq a - \delta \text{ or } x \geq a + \delta\}. \quad \square$$

Draw some Venn diagrams to visualize the following formulas. It is very important that you understand each one of them rather than simply trying to memorize them.

¹² \mathbb{R} is the set of all real numbers, i.e., the kind of numbers that make up the x -axis and y -axis in a beginner’s calculus course (see Section 2.3 (Numbers) on p.23).

¹³Draw a picture: The δ -neighborhood of a is the set of all points (in the universal set $[0, 1]$) with distance less than δ from a .

Proposition 2.1. *Let A, B, X be subsets of a universal set Ω and assume $A \subseteq X$. Then*

$$\begin{aligned}
 (2.10a) \quad & A \cup \emptyset = A; & A \cap \emptyset = \emptyset \\
 (2.10b) \quad & A \cup \Omega = \Omega; & A \cap \Omega = A \\
 (2.10c) \quad & A \cup A^c = \Omega; & A \cap A^c = \emptyset \\
 (2.10d) \quad & A \Delta B = (A \setminus B) \cup (B \setminus A) \\
 (2.10e) \quad & A \setminus A = \emptyset \\
 (2.10f) \quad & A \Delta \emptyset = A; & A \Delta A = \emptyset \\
 (2.10g) \quad & X \Delta A = X \setminus A \\
 (2.10h) \quad & A \cup B = (A \Delta B) \cup (A \cap B) \\
 (2.10i) \quad & A \cap B = (A \cup B) \setminus (A \Delta B) \\
 (2.10j) \quad & A \Delta B = \emptyset \text{ if and only if } B = A
 \end{aligned}$$

PROOF: The proof is left as exercise 2.2. See p.38. ■

Next we give a very detailed and rigorous proof of a simple formula for sets. You definitely want to remember the formulas, but it's perfectly OK to skip the proof.

Proposition 2.2 (Distributivity of unions and intersections for two sets). *Let A, B, C be sets. Then*

$$(2.11) \quad (A \cup B) \cap C = (A \cap C) \cup (B \cap C),$$

$$(2.12) \quad (A \cap B) \cup C = (A \cup C) \cap (B \cup C).$$

PROOF: ★ We only prove (2.11). The proof of (2.12) is left as exercise 2.1.

PROOF of “ \subseteq ”: Let $x \in (A \cup B) \cap C$. It follows from (2.2) on p.18 that $x \in (A \cup B)$, i.e., $x \in A$ or $x \in B$ (or both). It also follows from (2.2) that $x \in C$. We must show that $x \in (A \cap C) \cup (B \cap C)$ regardless of whether $x \in A$ or $x \in B$.

Case 1: $x \in A$. Since also $x \in C$, we obtain $x \in A \cap C$, hence, again by (2.2), $x \in (A \cap C) \cup (B \cap C)$, which is what we wanted to prove.

Case 2: $x \in B$. We switch the roles of A and B . This allows us to apply the result of case 1, and we again obtain $x \in (A \cap C) \cup (B \cap C)$.

PROOF of “ \supseteq ”: Let $x \in (A \cap C) \cup (B \cap C)$, i.e., $x \in A \cap C$ or $x \in B \cap C$ (or both). We must show that $x \in (A \cup B) \cap C$ regardless of whether $x \in A \cap C$ or $x \in B \cap C$.

Case 1: $x \in A \cap C$. It follows from $A \subseteq A \cup B$ and (2.4) on p.18 that $x \in (A \cup B) \cap C$, and we are done in this case.

Case 2: $x \in B \cap C$. This time it follows from $A \subseteq A \cup B$ that $x \in (A \cup B) \cap C$. This finishes the proof of (2.11).

Epilogue: The proofs both of “ \subseteq ” and of “ \supseteq ” were **proofs by cases**, i.e., we divided the proof into several cases (to be exact, two for each of “ \subseteq ” and “ \supseteq ”), and we proved each case separately. For example we proved that $x \in (A \cup B) \cap C$ implies $x \in (A \cap C) \cup (B \cap C)$ separately for the cases $x \in A$ and $x \in B$. Since those two cases cover all possibilities for x the assertion “if $x \in (A \cup B) \cap C$ then $x \in (A \cap C) \cup (B \cap C)$ ” is proven. ■

Proposition 2.3 (De Morgan’s Law for two sets). *Let $A, B \subseteq \Omega$. Then the complement of the union is the intersection of the complements, and the complement of the intersection is the union of the complements:*

$$(2.13) \quad \text{a. } (A \cup B)^c = A^c \cap B^c \quad \text{b. } (A \cap B)^c = A^c \cup B^c$$

PROOF:

1) First we prove that $(A \cup B)^c \subseteq A^c \cap B^c$:

Assume that $x \in (A \cup B)^c$. Then $x \notin A \cup B$, which is the same as saying that x does not belong to at least one of A and B . That in turn means that x belongs to all complements, i.e., to both A^c and B^c and hence, also to the intersection $A^c \cap B^c$.

2) Now we prove that $(A \cup B)^c \supseteq A^c \cap B^c$:

Let $x \in A^c \cap B^c$. Then x belongs to each one of A^c, B^c , hence to none of A, B , hence $x \notin A \cup B$. Therefore x belong to the complement of $A \cup B$. This completes the proof of formula **a**.

PROOF of **b**: The proof is very similar to that of formula **a** and left as an exercise. ■

Definition 2.8 (Power set).

The power set

$$2^\Omega := \{A : A \subseteq \Omega\}$$

of a set Ω is the set of all its subsets. Note that many older texts also use the notation $\mathfrak{P}(\Omega)$ for the power set. □

Remark 2.6. Note that $\emptyset \in 2^\Omega$ for any set Ω , even if $\Omega = \emptyset$: $2^\emptyset = \{\emptyset\}$. It follows that the power set of the empty set is not empty. □

Definition 2.9 (Partition). Let Ω be a set and $\mathfrak{A} \subseteq 2^\Omega$, i.e., the elements of \mathfrak{A} are subsets of Ω .

We call \mathfrak{A} a **partition** or a **partitioning** of Ω if

- (a) If $A, B \in \mathfrak{A}$ such that $A \neq B$ then $A \cap B = \emptyset$. In other words, \mathfrak{A} consists of mutually disjoint subsets of Ω .
- (b) Each $x \in \Omega$ is an element of some $A \in \mathfrak{A}$. □

Remark 2.7. Let Ω be a set and $\mathfrak{A} \subseteq 2^\Omega$. Then \mathfrak{A} is a partition of Ω if and only if

For each $x \in \Omega$, there exists a UNIQUE $A \in \mathfrak{A}$ such that $x \in A$. □

Example 2.5.

- a. For $n \in \mathbb{Z}$ let $A_n := \{n\}$. Then $\mathfrak{A} := \{A_n : n \in \mathbb{Z}\}$ is a partition of \mathbb{Z} . \mathfrak{A} is not a partition of \mathbb{N} because not all its members are subsets of \mathbb{N} and it is not a partition of \mathbb{Q} or \mathbb{R} . The reason: $\frac{1}{2} \in \mathbb{Q}$ and hence $\frac{1}{2} \in \mathbb{R}$, but $\frac{1}{2} \notin A_n$ for any $n \in \mathbb{Z}$, hence condition **b** of def.2.9 is not satisfied.
- b. For $n \in \mathbb{N}$ let $B_n := [n^2, (n+1)^2[= \{x \in \mathbb{R} : n^2 \leq x < (n+1)^2\}$. Then $\mathfrak{B} := \{B_n : n \in \mathbb{N}\}$ is a partition of $[1, \infty[$. □

Definition 2.10 (Size of a set).

- a. Let X be a finite set, i.e., a set which only contains finitely many elements. We write $|X|$ for the number of its elements, and we call $|X|$ the **size** of the set X .
- b. For infinite, i.e., not finite sets Y , we define $|Y| := \infty$. \square

More will be said about sets later.

2.2 The Proper Use of Language in Mathematics: Any vs All, etc

Mathematics must be very precise in its formulations. Such precision is achieved not only by means of symbols and formulas, but also by its use of the English language. We will list some important points to consider early on in this document.

2.2.0.1 All vs. ANY

Assume for the following that X is a set of numbers. Do the following two statements mean the same?

- (1) It is true for ALL $x \in X$ that x is an integer.
- (2) It is true for ANY $x \in X$ that x is an integer.

You will hopefully agree that there is no difference and that one could rewrite them as follows:

- (3) ALL $x \in X$ are integers.
- (4) ANY $x \in X$ is an integer.
- (5) EVERY $x \in X$ is an integer.
- (6) EACH $x \in X$ is an integer.
- (7) IF $x \in X$ THEN x is an integer.

Is it then always true that ALL and ANY means the same? Consider

- (8a) It is NOT true for ALL $x \in X$ that x is an integer.
- (8b) It is NOT true for ANY $x \in X$ that x is an integer.

Completely different things have been said: Statement (8) asserts that as few as one item and as many as all items in X are not integers, whereas (9) states that no items, i.e., exactly zero items in X , are integers.

My suggestion: Express formulations like (8b) differently. You could have written instead

- (8c) There is no $x \in X$ such that x is an integer.

2.2.0.2 AND vs. IF ... THEN

Some people abuse the connective AND to also mean IF ... THEN. However, mathematicians use the phrase “p AND q” exclusively to mean that something applies to both p and q. Contrast the use of AND in the following statements:

- (9) “Jane is a student AND Joe likes baseball”. This phrase means that both are true: Jane is indeed a student and Joe indeed likes baseball.
- (10) “You hit me again AND you’ll be sorry”. **Never, ever use the word AND in this context!** A mathematician would express the above as “IF you hit me again THEN you’ll be sorry”.

2.2.0.3 OR vs. EITHER ... OR

The last topic we address is the proper use of “OR”. In mathematics the phrase

(11) “ p is true OR q is true”

is always to be understood as

(12) “ p is true OR q is true OR BOTH are true”, i.e., at least one of p, q is true.

This is in contrast to everyday language where “ p is true OR q is true” often means that exactly one of p and q is true, but not not both.

When referring to a collection of items then the use of “OR” also is inclusive. If the items a, b, c, \dots belong to a collection \mathcal{C} , e.g., if those items are elements of a set, then

(13) “ a OR b OR c OR ...” means that we refer to at least one of a, b, c, \dots .

Note that “OR” in mathematics always is an **inclusive or**, i.e., “A OR B” means “A OR B OR BOTH”. More generally, “A OR B OR ...” means “at least one of A, B, ...”. To rule out that more than one of the choices is true you must use a phrase like “EXACTLY ONE OF A, B, C, ...” or “EITHER A OR B OR C OR ...”. We refer to this as an **exclusive or**.

2.2.0.4 Some Convenient Shorthand Notation We have previously encountered the notation “ $P \Rightarrow Q$ ” for “if P then Q ”, i.e., if P is true, then Q is true, and “ $P \Leftrightarrow Q$ ” for “ P iff Q ”, i.e., “ P is true exactly when Q is true”. We list them here again with some additional convenient abbreviations.

- $\forall x \dots$ For all $x \dots$
- $\exists x \text{ s.t. } \dots$ There exists an x such that \dots
- $\exists! x \text{ s.t. } \dots$ There exists a **UNIQUE** x such that \dots
- $P \Rightarrow Q$ If P then Q
- $P \Leftrightarrow Q$ P iff Q , i.e., P if and only if Q

It is important that you are clear about the difference between \exists and $\exists!$.

$\exists x$: you can find at least one x but there might be more; potentially infinitely many!

$\exists! x$: you can find one and only one x ; not zero, not two, not 200, ... \square

2.3 Numbers

We start with an informal classification of numbers.

Definition 2.11 (Types of numbers). Here is a definition of the various kinds of numbers in a nutshell.

$\mathbb{N} := \{1, 2, 3, \dots\}$ denotes the set of **natural numbers**.
 $\mathbb{Z} := \{0, \pm 1, \pm 2, \pm 3, \dots\}$ denotes the set of all **integers**.
 $\mathbb{Q} := \{n/d : n \in \mathbb{Z}, d \in \mathbb{N}\}$ (fractions of integers) denotes the set of all **rational numbers**.
 $\mathbb{R} := \{\text{all integers or decimal numbers with finitely or infinitely many decimal digits}\}$ denotes the set of all **real numbers**.
 $\mathbb{R} \setminus \mathbb{Q} = \{\text{all real numbers which cannot be written as fractions of integers}\}$ denotes the set of all **irrational numbers**. There is no special symbol for irrational numbers. Example: $\sqrt{2}$ and π are irrational. \square

Here are some customary abbreviations of some often referenced sets of numbers:

$$\begin{aligned} \mathbb{N}_0 &:= \mathbb{Z}_+ := \mathbb{Z}_{\geq 0} := \{0, 1, 2, 3, \dots\} \text{ denotes the set of nonnegative integers,} \\ \mathbb{R}_+ &:= \mathbb{R}_{\geq 0} := \{x \in \mathbb{R} : x \geq 0\} \text{ denotes the set of all nonnegative real numbers,} \\ \mathbb{R}^+ &:= \mathbb{R}_{> 0} := \{x \in \mathbb{R} : x > 0\} \text{ denotes the set of all positive real numbers,} \\ \mathbb{R}_{\neq 0} &:= \{x \in \mathbb{R} : x \neq 0\}. \quad \square \end{aligned}$$

Examples of rational numbers are

$$\frac{3}{4}, -0.75, -\frac{1}{3}, \bar{3}, \frac{7}{1}, 16, \frac{13}{4}, -5, 2.99\bar{9}, -37\frac{2}{7}.$$

Note that a mathematician does not care whether a rational number is written as a fraction

$$\frac{\text{numerator}}{\text{denominator}}$$

or as a decimal numeral. The following all are representations of one third:

$$(2.14) \quad 0.\bar{3} = \bar{3} = 0.3333333333\dots = \frac{1}{3} = \frac{-1}{-3} = \frac{2}{6},$$

and here are several equivalent ways of expressing the number minus four:

$$(2.15) \quad -4 = -4.000 = -3.\bar{9} = -\frac{12}{3} = \frac{4}{-1} = \frac{-4}{1} = \frac{12}{-3} = -\frac{400}{100}.$$

Definition 2.12 (Intervals of Numbers). For $a, b \in \mathbb{R}$ we have the following intervals.

- $[a, b] := \{x \in \mathbb{R} : a \leq x \leq b\}$ is the **closed interval** with endpoints a and b .
- $]a, b[:= \{x \in \mathbb{R} : a < x < b\}$ is the **open interval** with endpoints a and b .
- $[a, b[:= \{x \in \mathbb{R} : a \leq x < b\}$ and $]a, b] := \{x \in \mathbb{R} : a < x \leq b\}$ are **half-open intervals** with endpoints a and b .

The symbol “ ∞ ” stands for an object which itself is not a number but is larger than any (real) number, and the symbol “ $-\infty$ ” stands for an object which itself is not a number but is smaller than any number. We thus have $-\infty < x < \infty$ for any number x . This allows us to define the following intervals of “infinite length”:

$$(2.16) \quad \begin{aligned}]-\infty, a] &:= \{x \in \mathbb{R} : x \leq a\}, \quad]-\infty, a[:= \{x \in \mathbb{R} : x < a\}, \\]a, \infty[&:= \{x \in \mathbb{R} : x > a\}, \quad [a, \infty[:= \{x \in \mathbb{R} : x \geq a\}, \quad]-\infty, \infty[:= \mathbb{R} \end{aligned}$$

You should always work with $a < b$. In case you don't, you get

- $[a, a] = \{a\}; [a, a[=]a, a[=]a, a] = \emptyset$
- $[a, b] = [a, b[=]a, b[=]a, b] = \emptyset$ for $a \geq b$ \square

Notation 2.1 (Notation Alert for intervals of integers or rational numbers).

It is at times convenient to also use the notation $[\dots],]\dots[, [\dots[,]\dots]$, for intervals of integers or rational numbers. We will subscript them with \mathbb{Z} or \mathbb{Q} . For example,

$$\begin{aligned} [3, n]_{\mathbb{Z}} &= [3, n] \cap \mathbb{Z} = \{k \in \mathbb{Z} : 3 \leq k \leq n\}, \\]-\infty, 7]_{\mathbb{Z}} &=]-\infty, 7] \cap \mathbb{Z} = \{k \in \mathbb{Z} : k \leq 7\} = \mathbb{Z}_{\leq 7}, \\]a, b[_{\mathbb{Q}} &=]a, b[\cap \mathbb{Q} = \{q \in \mathbb{Q} : a < q < b\}. \end{aligned}$$

An interval which is not subscripted always means an interval of real numbers, but we will occasionally write, e.g., $[a, b]_{\mathbb{R}}$ rather than $[a, b]$, if the focus is on integers or rational numbers and an explicit subscript helps to avoid confusion. \square

Definition 2.13 (Absolute value, positive and negative part). Let $x, y \in \mathbb{R}$. We define the following.

$$\begin{aligned} \text{absolute value: } |x| &= \begin{cases} x & \text{if } x \geq 0, \\ -x & \text{if } x < 0. \end{cases} \\ \text{maximum: } \max(x, y) &= \begin{cases} x & \text{if } x \geq y, \\ y & \text{if } x \leq y. \end{cases} \\ \text{minimum: } \min(x, y) &= \begin{cases} y & \text{if } x \geq y, \\ x & \text{if } x \leq y. \end{cases} \quad \square \end{aligned}$$

Assumption 2.1 (Square roots are always assumed nonnegative). Remember that for any number a it is true that

$$a \cdot a = (-a)(-a) = a^2, \quad \text{e.g., } 2^2 = (-2)^2 = 4,$$

or that, expressed in form of square roots, for any number $b \geq 0$

$$(+\sqrt{b})(+\sqrt{b}) = (-\sqrt{b})(-\sqrt{b}) = b.$$

We will always assume that “ \sqrt{b} ” is the **positive** value unless the opposite is explicitly stated.

Example: $\sqrt{9} = +3$, not -3 . \square

Remark 2.8. For any real number x we have

$$(2.17) \quad \sqrt{x^2} = |x|. \quad \square$$

Proposition 2.4 (The Triangle Inequality for real numbers). *The following inequality is used all the time in mathematical analysis to show that the size of a certain expression is limited from above:*

$$(2.18) \quad \text{Triangle Inequality: } |a_1 + a_2 + \dots + a_n| \leq |a_1| + |a_2| + \dots + |a_n|$$

This inequality is true for any list of real numbers a_1, a_2, \dots, a_n .

PROOF:

It is easy to prove this for $n = 2$: Just look separately at the three cases where both numbers are nonnegative, both are negative, or one of each is positive and negative. ■

2.4 Functions and Sequences

Introduction 2.2. You are familiar with functions from calculus. Examples are $f_1(x) = \sqrt{x}$ and $f_2(x, y) = \ln(x - y)$. Sometimes $f_1(x)$ means the entire graph, i.e., the entire collection of points (x, \sqrt{x}) in the plane and sometimes it just refers to the function value \sqrt{x} for a “fixed but arbitrary” number x . In case of the function $f_2(x, y)$: Sometimes $f_2(x, y)$ means the entire graph, i.e., the entire collection of points $((x, y), \ln(x - y))$ in threedimensional space. At other times this expression just refers to the function value $\ln(x - y)$ for a pair of “fixed but arbitrary” numbers (x, y) .

To obtain a usable definition of a function there are several things to consider. In the following $f_1(x)$ and $f_2(x, y)$ again denote the functions $f_1(x) = \sqrt{x}$ and $f_2(x, y) = \ln(x - y)$.

- a. The source of all allowable arguments (x -values in case of $f_1(x)$ and (x, y) -values in case of $f_2(x, y)$) will be called the **domain** of the function. The domain is explicitly specified as part of a function definition and it may be chosen for whatever reason to be only a subset of all arguments for which the function value is a valid expression. In case of the function $f_1(x)$ this means that the domain must be a subset of the interval $[0, \infty[$ because the square root of a negative number cannot be taken. In case of the function $f_2(x, y)$ this means that the domain must be a subset of

$$\{ (x, y) : x, y \in \mathbb{R} \text{ and } x - y > 0 \},$$

because logarithms are only defined for strictly positive numbers.

- b. The set to which all possible function values belong will be called the **codomain** of the function. As is the case for the domain, the codomain also is explicitly specified as part of a function definition. It may be chosen as any superset of the set of all function values for which the argument belongs to the domain of the function.

For the function $f_1(x)$ this means that we are OK if the codomain is a superset of the interval $[0, \infty[$. Such a set is big enough because square roots are never negative. It is OK to specify the interval $] - 3.5, \infty[$ or even the set \mathbb{R} of all real numbers as the codomain. In case of the function $f_2(x, y)$ this means that we are OK if the codomain contains \mathbb{R} . Not that it would make a lot of sense, but the set $\mathbb{R} \cup \{ \text{all inhabitants of Chicago} \}$ also is an acceptable choice for the codomain.

- c. A function $y = f(x)$ is not necessarily something that maps (assigns) numbers or pairs of numbers to numbers. Rather domain and codomain can be a very different kind of animal. The following example will be very relevant for the remainder of the course:

At the end of Section 1.2 (A First Look at Probability) We informally defined the probability associated with rolling a die as a function $A \mapsto P(A)$ which maps subsets A of $\Omega = \{1, 2, \dots, 6\}$ to a real number $0 \leq P(A) \leq 1$. Thus, the domain here is 2^Ω , the power set of Ω ; the codomain is $[0, 1]$ (or any superset of $[0, 1]$).

- d. Considering all that was said so far one can think of the graph of a function $f(x)$ with domain D and codomain C (see earlier in this note) as the set

$$\Gamma_f := \{(x, f(x)) : x \in D\}.$$

Alternatively one can characterize this function by the assignment rule which specifies how $f(x)$ depends on any given argument $x \in D$. We write “ $x \mapsto f(x)$ ” to indicate this. You can also write instead $f(x) =$ whatever the actual function value will be.

This is possible if one does not write about functions in general but about specific functions such as $f_1(x) = \sqrt{x}$ and $f_2(x, y) = \ln(x - y)$. We further write

$$f : D \longrightarrow C$$

as a short way of saying that the function $f(x)$ has domain D and codomain C .

In case of the function $f_1(x) = \sqrt{x}$ for which we might choose the interval $X := [2.5, 7]$ as the domain (small enough because $X \subseteq [0, \infty[$) and $Y :=]1, 3[$ as the codomain (big enough because $1 < \sqrt{x} < 3$ for any $x \in X$) we specify this function as

$$\text{either } f_1 : [2.5, 7] \rightarrow]1, 3[; \quad x \mapsto \sqrt{x} \quad \text{or } f_1 : [2.5, 7] \rightarrow]1, 3[; \quad f(x) = \sqrt{x}.$$

Let us choose $U := \{(x, y) : x, y \in \mathbb{R} \text{ and } 1 \leq x \leq 10 \text{ and } y < -2\}$ as the domain and $V := [0, \infty[$ as the codomain for $f_2(x, y) = \ln(x - y)$. These choices are OK because $x - y \geq 1$ for any $(x, y) \in U$ and hence $\ln(x - y) \geq 0$, i.e., $f_2(x, y) \in V$ for all $(x, y) \in U$. We specify this function as

$$\text{either } f_2 : U \rightarrow V, \quad (x, y) \mapsto \ln(x - y) \quad \text{or } f_2 : U \rightarrow V, \quad f(x, y) = \ln(x - y). \quad \square$$

We incorporate what we noted above into this definition of a function.

Definition 2.14 (Function).

A **function** f consists of two nonempty sets X and Y and an assignment rule $x \mapsto f(x)$ which assigns any $x \in X$ uniquely to some $y \in Y$. We write $f(x)$ for this assigned value and call it the **function value** of the **argument** x . X is called the **domain** and Y is called the **codomain** of f . We write

$$(2.19) \quad f : X \rightarrow Y, \quad x \mapsto f(x).$$

We read “ $a \mapsto b$ ” as “ a is assigned to b ” or “ a maps to b ” and refer to \mapsto as the **maps to operator** or **assignment operator**. The **graph** of such a function is the collection of pairs

$$(2.20) \quad \Gamma_f := \{(x, f(x)) : x \in X\},$$

and the subset $f(X) := \{f(x) : x \in X\}$ of Y is called the **range** of the function f . \square

Note that the codomain Y of f and its range $f(X)$ can be vastly different. For example, if $f : \mathbb{R} \rightarrow \mathbb{R}$

is given by the assignment $f(x) = \sin(x)$ then $f(\mathbb{R}) = [-1, 1]$ is a very small part of the codomain!

Remark 2.9. The name given to the argument variable is irrelevant. Let f_1, f_2, X, Y, U, V be as defined in **d** of the introduction to ch.2.4 (A First Look at Functions and Sequences). The function

$$g_1 : X \rightarrow Y, \quad p \mapsto \sqrt{p}$$

is identical to the function f_1 . The function

$$g_2 : U \rightarrow V, \quad (t, s) \mapsto \ln(t - s)$$

is identical to the function f_2 and so is the function

$$g_3 : U \rightarrow V, \quad (s, t) \mapsto \ln(s - t).$$

The last example illustrates the fact that you can swap function names as long as you do it consistently in all places. \square

We all know what it means that $f : \mathbb{R} \rightarrow]0, \infty]$; $x \mapsto e^x$ has $f^{-1}(x) = \ln(x)$ as its inverse function:

- The arguments of f^{-1} will be the function values of f and the function values of f^{-1} will be the arguments of f : $f(x) = e^x = y \Leftrightarrow g(y) = \ln(y) = x$.
- f and f^{-1} cancel each other, i.e.,

$$f^{-1}(f(y)) = y \quad \text{and} \quad f(f^{-1}(x)) = x.$$

- Not so obvious but very useful: We want both codomains to be so small that $f^{-1}(f(y)) = y$ is true for all y in the codomain of f and $f(f^{-1}(x)) = x$ is true for all x in the codomain of f^{-1} . One can show that this requires

$$\text{domain of } f = \text{codomain of } f^{-1} \quad \text{and} \quad \text{domain of } f^{-1} = \text{codomain of } f.$$

This leads to the following definition for the inverse of a function.

Definition 2.15 (Inverse function).

Given are two nonempty sets X and Y and a function $f : X \rightarrow Y$ with domain X and codomain Y . We say that f has an **inverse function** if it satisfies all of the following conditions which uniquely determine this inverse function, so that we are justified to give it the symbol f^{-1} :

- (a) $f^{-1} : Y \rightarrow X$, i.e., f^{-1} has domain Y and codomain X .
- (b) $f^{-1}(f(x)) = x$ for all $x \in X$, and $f(f^{-1}(y)) = y$ for all $y \in Y$. \square

Definition 2.16 (Surjective, injective and bijective functions).

Given are two nonempty sets X and Y and a function $f : X \rightarrow Y$ with domain X and codomain Y . We say that

- (a) f is “one-one” or **injective**, if for each $y \in Y$ there is at most one $x \in X$ such that $f(x) = y$.
- (b) f is “onto” or **surjective**, if for each $y \in Y$ there is at least one $x \in X$ such that $f(x) = y$.
- (c) f is **bijective**, f is both injective and surjective. \square

Remark 2.10. that One can show that a function f has an inverse f^{-1} if and only if f is bijective. \square

Remark 2.11. that If the inverse function f^{-1} exists and if $x \in X$ and $y \in Y$, then we have the relation

$$y = f(x) \Leftrightarrow x = f^{-1}(y).$$

Example 2.6. If h is a function, we write Dom_h and Cod_h for its domain and codomain. Be sure you understand the following:

- (a) $f : \mathbb{R} \rightarrow \mathbb{R}; x \rightarrow e^x$ does not have an inverse $f^{-1}(y) = \ln(y)$ since its domain $Dom_{f^{-1}}$ would have to be the codomain \mathbb{R} of f and $\ln(y)$ is not defined for $y \leq 0$.
- (b) $g : \mathbb{R} \rightarrow]0, \infty[; x \rightarrow e^x$ has the inverse $g^{-1} :]0, \infty[\rightarrow \mathbb{R}; g^{-1}(y) = \ln(y)$ since

$$\begin{aligned} Dom_{g^{-1}} = Cod_g =]0, \infty[, & \quad Cod_{g^{-1}} = Dom_g = \mathbb{R}, \\ e^{\ln(y)} = y \text{ for } 0 < y < \infty, & \quad \ln(e^x) = x \text{ for all } x \in \mathbb{R}. \quad \square \end{aligned}$$

Definition 2.17 (Restriction/Extension of a function). ★ Given are three nonempty sets A, X and Y such that $A \subseteq X$, and a function $f : X \rightarrow Y$ with domain X . We define the **restriction of f to A** as the function

$$(2.21) \quad f|_A : A \rightarrow Y \quad \text{defined as} \quad f|_A(x) := f(x) \text{ for all } x \in A.$$

Conversely let $f : A \rightarrow Y$ and $\varphi : X \rightarrow Y$ be functions such that $f = \varphi|_A$. We then call φ an **extension** of f to X . \square

We now briefly address sequences and subsequences.

Definition 2.18. Let n_* be an integer and assume that an item x_j associated

- **either** with each integer $j \geq n_*$, In other words, we have an item x_j assigned to each $j = n_*, n_* + 1, n_* + 2, \dots$
- **or** with each integer j such that $n_* \leq j \leq n^*$. In this case an item x_j is assigned to each $j = n_*, n_* + 1, \dots, n^*$.

Such items can be anything, but we usually deal with numbers or outcomes or sets of outcomes of an experiment.

- In the first case we usually write $x_{n_*}, x_{n_*+1}, x_{n_*+2}, \dots$ or $(x_n)_{n \geq n_*}$ for such a collection of items and we call it a **sequence with start index n_*** .
- In the second case we speak of a **finite sequence**, which starts at n_* and ends at n^* . We write $(x_n)_{n_* \leq n \leq n^*}$ or $x_{n_*}, x_{n_*+1}, \dots, x_{n^*}$ for such a finite collection of items.
- If we refer to a sequence $(x_n)_n$ without qualifying it as finite then we imply that we deal with an **infinite sequence**, $x_{n_*}, x_{n_*+1}, x_{n_*+2}, \dots$. \square

Example 2.7.

- (1) If $u_k = k^2$ for $k \in \mathbb{Z}$, then $(u_k)_{k \geq -2}$ is the sequence of integers 4, 1, 0, 1, 4, 9, 16, ...
- (2) If $A_j = [-1 - \frac{1}{j}, 1 + \frac{1}{j}] = \{x \in \mathbb{R} : -1 - \frac{1}{j} \leq x \leq 1 + \frac{1}{j}\}$, then $(A_j)_{j \geq 3}$ is the sequence of intervals of real numbers $[-\frac{4}{3}, \frac{4}{3}]$, $[-\frac{5}{4}, \frac{5}{4}]$, $[-\frac{6}{5}, \frac{6}{5}]$, ... This is a sequence of sets! \square

Remark 2.12 (Sequences are functions). that

- One can think of a sequence $(x_i)_{i \geq n_*}$ in terms of the assignment $i \mapsto x_i$. This sequence can then be interpreted as the function

$$x(\cdot) : [n_*, \infty[\mathbb{Z} \longrightarrow \text{suitable codomain}; \quad i \mapsto x(i) := x_i,$$

where that “suitable codomain” depends on the nature of the items x_i .

- In Example 2.7(1), we could chose \mathbb{Z} as that codomain. In Example 2.7(2) $2^{\mathbb{R}}$, the power set of \mathbb{R} would be an appropriate choice. \square

Definition 2.19.

- If $(x_n)_n$ is a finite or infinite sequence and one pares down the full set of indices to a subset $\{n_1, n_2, n_3, \dots\}$ such that $n_1 < n_2 < n_3 < \dots$, then we call the corresponding thinned out sequence $(x_{n_j})_{j \in \mathbb{N}}$ a **subsequence** of that sequence.
- If this subset of indices is finite, i.e., we have $n_1 < n_2 < \dots < n_K$ for some suitable $K \in \mathbb{N}$, then we call $(x_{n_j})_{j \leq K}$ a **finite subsequence** of the original sequence. \square

Note that subsequences of finite sequences are necessarily finite whereas subsequences of infinite sequences can be finite or infinite.

Remark 2.13. Does it matter whether we look at a sequence $(x_j)_{j \in J}$ or at the corresponding set $\{x_j : j \in J\}$? The answer: **THIS CAN MATTER GREATLY!** Consider the sequence

$$x_1 = -1, x_2 = 1, x_3 = -1, x_4 = 1, \dots; \quad \text{i.e., } x_n = (-1)^n \text{ for } n \in \mathbb{N}$$

- The sequence is infinite, since the index set \mathbb{N} is infinite
- Let $A := \{x_j : j \in \mathbb{N}\}$. Since **sets have no duplicates**, $A = \{-1, 1\}$ has only two elements.
- The ordering of the indices j is lost when considering the set: There is no difference between $\{-1, 1\}$ and $\{1, -1\}$!

Considering the last point, do not confuse the ordering of the indices j with a possible ordering of the x_j ! The order may be reversed (e.g., $x_j = 5 - j$), neither increasing nor decreasing ($x_j = \sin(j)$), or there is no ordering ($x_j = \text{eye color of person } j$). \square

Definition 2.20. We give some convenient definitions and notations for monotone sequences of numbers, functions and sets.

- (a) Let x_n be a sequence of extended real-valued numbers.
- We call x_n a **nondecreasing** or **increasing** sequence, if $j < n \Rightarrow x_j \leq x_n$.
 - We call x_n a **strictly increasing** sequence, if $j < n \Rightarrow x_j < x_n$.

- We call x_n a **nonincreasing** or **decreasing** sequence, if $j < n \Rightarrow x_j \geq x_n$.
- We call x_n a **strictly decreasing** sequence, if $j < n \Rightarrow x_j > x_n$.
- We write $x_n \uparrow$ for nondecreasing x_n , and $x_n \uparrow x$ to indicate that $\lim_{n \rightarrow \infty} x_n = x$,
- We write $x_n \downarrow$ for nonincreasing x_n , $x_n \downarrow x$ to indicate that $\lim_{n \rightarrow \infty} x_n = x$. \square

(b) Let A_n be a sequence of sets.

- We call A_n a **nondecreasing** or **increasing** sequence, if $j < n \Rightarrow A_j \subseteq A_n$.
- We call A_n a **strictly increasing** sequence, if $j < n \Rightarrow A_j \subsetneq A_n$.
- We call A_n a **nonincreasing** or **decreasing** sequence, if $j < n \Rightarrow A_j \supseteq A_n$.
- We call A_n a **strictly decreasing** sequence, if $j < n \Rightarrow A_j \supsetneq A_n$.
- We write $A_n \uparrow$ for nondecreasing A_n , and $A_n \uparrow A$ to indicate that $\bigcup_n A_n = A$,
- We write $A_n \downarrow$ for nonincreasing A_n , $A_n \downarrow A$ to indicate that $\bigcap_n A_n = A$. \square

Example 2.8.

- (a) The sequence $x_n = -\frac{1}{n}$ is strictly increasing.
- (b) The sequence $y_n = \frac{1}{n}$ is strictly decreasing.
- (c) The sequence $a_1 = 1, a_{n+1} = a_n$ for even n and $a_{n+1} = -\frac{1}{n}$ for odd n , is nonincreasing.
- (c) The sequence $b_1 = 1, b_{n+1} = b_n$ for even n and $b_{n+1} = \frac{1}{n}$ for odd n , is nondecreasing. \square

There are different degrees of infinity for the size of a set. Finite sets and many infinite sets are “small enough” to list all their elements in a finite or infinite sequence. Other infinite sets are too big for that.

Definition 2.21 (Countable and uncountable sets). Let X be a set.

- (a) We call X **countable** if its elements can be written as a finite sequence (those are the finite sets) $X = \{x_1, x_2, \dots, x_n\}$ or as an infinite sequences. $X = \{x_1, x_2, \dots\}$.
- (b) We call a nonempty set **uncountable** if it is not countable, i.e., its elements cannot be sequenced.
- (c) By convention the empty set, \emptyset , is countable. \square

Fact 2.1. One can prove the following important facts:

- (a) The integers are countable. (Easy: $\mathbb{Z} = \{0, -1, 1, -2, 2, -3, 3, \dots\}$) lists all elements of \mathbb{Z} in a sequence.
- (b) Subsets of countable sets are countable. (Easy: If $X = \{x_1, x_2, \dots\}$ and $A \subseteq X$, then remove all x_j that are not in A . That subsequence lists the elements of A .)
- (c) Countable unions of countable sets are countable: If A_1, A_2, \dots is a finite or infinite sequence of sets, then $A_1 \cup A_2 \cup \dots$ is countable.
- (d) The rational numbers \mathbb{Q} are countable. A proof is given below.
- (e) The real numbers \mathbb{R} are uncountable! \square

★ Here is a proof that \mathbb{Q} is countable. For fixed $d \in \mathbb{N}$, let $A_d := \{n/d : n \in \mathbb{Z}\}$ (“d” for denominator). Then is countable since it can be sequenced as follows.

$$A_d = \left\{0, -\frac{1}{d}, \frac{1}{d}, -\frac{2}{d}, \frac{2}{d}, \dots\right\}$$

The assertion follows from fact (c) and $\mathbb{Q} = \bigcup_{d=1}^{\infty} A_d$ (WHY?)

Example 2.9. ★ For $a, b, r \in \mathbb{R}$, let $A_{(a,b,r)} := \{(x, y) \in \mathbb{R}^2\}$ such that $(x - a)^2 + (y - b)^2 = r^2$, i.e., $A_{(a,b,r)}$ is the circle with radius $|r|$ around the point (a, b) in the plane. It is not possible to write the indexed collection

$$(A_{(a,b,r)})_{(a,b,r) \in \mathbb{R}^3}$$

as a sequence, since \mathbb{R}^3 is bigger than the uncountable set \mathbb{R} , hence cannot be sequenced. \square

There is a name for those “generalized sequences” $(x_i)_{i \in I}$ which have an index set that not necessarily consists of integers $n_*, n_* + 1, \dots, n^*$ or $n_*, n_* + 1, \dots$ or of a subset of such a set. The next definition is marked as optional and you not need remember it for quizzes or exams. But you must remember it well enough to understand problems and propositions which refer to families.

Definition 2.22 (Families). ★

Let I and X be nonempty sets such that each $i \in I$ is associated with some $x_i \in X$. Then

- a. $(x_i)_{i \in I}$ is called an **indexed family** or simply a **family** in X .
- b. I is called the **index set** of the family.
- c. For each $i \in I$, x_i is called a **member of the family** $(x_i)_{i \in I}$. \square

Remark 2.14 (Families are functions). that

We saw in example 2.12 on p.30 that sequences $(x_n)_n$ can be interpreted as functions with domain = index set and codomain = a set that contains all members x_n . This also holds true for families and is particularly easily understood if the family $(x_i)_{i \in I}$ in X is written in a way that each member explicitly tracks the index that it is associated with, i.e., we write $(i, x_i)_{i \in I}$. The set

$$\Gamma_f := \{(i, x_i) : i \in I\}$$

is the graph Γ_f of the function

$$f : I \longrightarrow X; \quad i \mapsto f(i) := x_i.$$

At the end of Definition 2.4 on p.17 we defined unions and intersections of any collection of sets $(A_i)_{i \in J}$ which is indexed by integers, i.e., $J \subseteq \mathbb{Z}$. We did so by saying that ¹⁴

$$\bigcup_{i \in J} A_i = \{x : \exists i_0 \in J \text{ s.t. } x \in A_{i_0}\} \quad \text{and} \quad \bigcap_{i \in J} A_i = \{x : \forall i \in J : x \in A_i\}.$$

This allows us to generalize unions and intersections of finite and infinite sequences of sets to collections of sets with an arbitrary index set. Note the following:

- The next definition is NOT marked as OPTIONAL
- It contains Definition 2.4 as a special case!

¹⁴See paragraph 2.2.0.4 (Some Convenient Shorthand Notation) on p.23 about \forall and \exists .

Definition 2.23 (Arbitrary unions and intersections). Let J be an arbitrary, nonempty set and $(A_j)_{j \in J}$ a family of sets with index set J . We define

- The **union** $\bigcup_{j \in J} A_j := \{x : \exists i_0 \in J \text{ s.t. } x \in A_{i_0}\}$.
- The **intersection** $\bigcap_{j \in J} A_j = \{x : \forall i \in J : x \in A_i\}$.
- If the sets A_i are disjoint, we often write $\bigsqcup_{j \in J} A_j$ rather than $\bigcup_{j \in J} A_j$.
- Let $(B_j)_{j \in J}$ be a family of subsets of a set X . We call this family a **partition** or a **partitioning** of X if the corresponding set of sets $\{B_i : i \in J\}$ is a partition of X :
 - (a) $i \neq j \Rightarrow B_i \cap B_j = \emptyset$ (b) $X = \bigsqcup_{j \in J} B_j$. See Definition 2.9 on p.21. \square

Remark 2.15. ★ For typographical reasons I sometimes use the following notation.

$$\bigcup [A_i; i \in I] := \bigcup_{i \in I} A_i.$$

Analogous notation exists for \bigcap , \bigsqcup and even summation. For example, assume that $g : \mathbb{R} \rightarrow \mathbb{R}$ is some rel-valued function of real numbers, and that the indices of interest are

$$I := \{x \in \mathbb{R} : x > 5 \text{ and } 0 \leq g(x) < 5\}.$$

Then $\bigcap_{x \in I} B_x$ can also be expressed as follows:

$$\bigcap_{x \in I} B_x = \bigcap [B_x : x > 5 \text{ and } 0 \leq g(x) < 5] = \bigcap_{x > 5 \text{ and } 0 \leq g(x) < 5} B_x = \bigcap_{\substack{x > 5 \\ 0 \leq g(x) < 5}} B_x. \quad \square$$

Be sure to understand the following example (draw a picture!)

Example 2.10. ★ For $a, b \in \mathbb{R}$, let $Q_{(a,b)} := \{(x, y) \in \mathbb{R}^2 : |x - a| \leq 3/2, |y - b| \leq 3/2\}$. Thus, $Q_{(a,b)}$ is the square in the plane with center (a, b) and side length 3. Compute $\bigcap_{(a,b) \in K} Q_{(a,b)}$

and $\bigcup_{(a,b) \in K} Q_{(a,b)}$.

For $K = \{(a, b) \in \mathbb{R}^2 : -1 \leq a, b \leq 1\}$, compute $\bigcap_{(a,b) \in K} Q_{(a,b)}$ and $\bigcup_{(a,b) \in K} Q_{(a,b)}$.

Solution:

Let $U := \bigcap_{(a,b) \in K} Q_{(a,b)}$ and $V := \bigcup_{(a,b) \in K} Q_{(a,b)}$.

Fix $b_0 \in [-1, 1]$ and consider the squares $Q_{(a,b_0)}$ moving from the left ($a = -1$) all the way to the right ($a = +1$). Even $Q_{(-1,b_0)}$ as the leftmost square has x values as big as $1/2$, and $Q_{(1,b_0)}$ as the rightmost square has x values as small as $-(1/2)$. Thus,

$$(x, y) \in \bigcap_{-1 \leq a \leq 1} Q_{(a,b_0)} \Leftrightarrow \left[-\frac{1}{2} \leq x \leq \frac{1}{2} \text{ and } b_0 - \frac{3}{2} \leq y \leq b_0 + \frac{3}{2} \right].$$

Likewise, if we now also move the squares vertically from $b = -1$ to $b = 1$, then the y values of points in the intersection are exactly those that satisfy $-(1/2) \leq y \leq 1/2$. Thus,

$$U = \{(x, y) : |x| \leq 1/2 \text{ and } |y| \leq 1/2\}.$$

One sees in likewise fashion that the points in the union V are exactly those with x values and y values between $-1 - (3/2) = -5/2$ and $1 + (3/2) = 5/2$. Thus,

$$V = \{(x, y) : |x| \leq 5/2 \text{ and } |y| \leq 5/2\}. \quad \square$$

We finish this section with two very useful propositions. The first one (De Morgan) you already have encountered for two sets (see Proposition 2.3 on p.2.3).¹⁵

Proposition 2.5 (De Morgan's Law for sequences of sets). *Let $(A_n)_n$ be a finite or infinite sequence of subsets of a set Ω . Then the complement of the union is the intersection of the complements, and the complement of the intersection is the union of the complements:*

$$(2.22) \quad (a) \quad \left(\bigcup_k A_k \right)^c = \bigcap_k A_k^c; \quad (b) \quad \left(\bigcap_k A_k \right)^c = \bigcup_k A_k^c;$$

PROOF:

Not very complicated, but we skip it ■

Note that the order of the sequencing does not matter for De Morgan and the next proposition.

Proposition 2.6 (Distributivity of unions and intersections). *Let $(A_n)_n$ be a finite or infinite sequence of sets and let B be a set. Then*

$$(2.23) \quad \bigcup_j (B \cap A_j) = B \cap \bigcup_j A_j,$$

$$(2.24) \quad \bigcap_{j \in I} (B \cup A_j) = B \cup \bigcap_j A_j.$$

PROOF: ■

2.5 Cartesian Products

We next define cartesian products of sets. Those mathematical objects generalize rectangles

$$[a_1, b_1] \times [a_2, b_2] = \{(x, y) : x, y \in \mathbb{R}, a_1 \leq x \leq b_1 \text{ and } a_2 \leq y \leq b_2\}$$

¹⁵Matter of fact, both propositions extend to arbitrary families.

and quads

$$[a_1, b_1] \times [a_2, b_2] \times [a_3, b_3] = \{(x, y, z) : x, y, z \in \mathbb{R}, a_1 \leq x \leq b_1, a_2 \leq y \leq b_2 \text{ and } a_3 \leq z \leq b_3\}.$$

which you certainly have encountered in multivariable calculus.

Definition 2.24 (Cartesian Product). Let X and Y be two sets The set

$$(2.25) \quad X \times Y := \{(x, y) : x \in X, y \in Y\}$$

is called the **cartesian product** of X and Y . We write X^2 as an abbreviation for $X \times X$.

Note that the order is important: (x, y) and (y, x) are different unless $x = y$.

This definition generalizes to more than two sets as follows:

Let X_1, X_2, \dots, X_n be sets. The set

$$(2.26) \quad X_1 \times X_2 \cdots \times X_n := \{(x_1, x_2, \dots, x_n) : x_j \in X_j \text{ for each } j = 1, 2, \dots, n\}$$

is called the cartesian product of X_1, X_2, \dots, X_n .

We write X^n as an abbreviation for $X \times X \times \cdots \times X$.

Example 2.11. In your multivariable calculus course you have learned about twodimensional vectors and threedimensional vectors. Convenient notations would often be

$$(2.27) \quad (x, y) \in \mathbb{R}^2, \quad (a, b) \in \mathbb{R}^2, \quad (x, y, z) \in \mathbb{R}^3, \quad (a, b, c) \in \mathbb{R}^3.$$

Note that those vectors are elements of the cartesian products $\mathbb{R}^2 = \mathbb{R} \times \mathbb{R}$ $\mathbb{R}^3 = \mathbb{R} \times \mathbb{R} \times \mathbb{R}$.

In general, any finite list of real numbers $(\beta_1, \beta_2, \dots, \beta_m)$ is an element of \mathbb{R}^m which we call an m -dimensional **vector** of real numbers.

$$(8, -3, 0, 4, -7)$$

is a 5-dimensional vector of Integers. Since integers are special cases of rational numbers which themselves are also real numbers, this vector is an element of each one of $\mathbb{Z}^5, \mathbb{Q}^5, \mathbb{R}^5$.

The notation used in (2.27) does not scale for higher dimensional vectors, in particular, if the dimension is arbitrary. On the other hand, $(\beta_1, \beta_2, \dots, \beta_m)$ is very suitable. But this is very lengthy notation, so we use the symbol for the subscripted components (that's β) and write an arrow on top to indicate that we are dealing with a vector.¹⁶

We will use as much as possible this arrow notation for vectors. Here are some examples.

$$\vec{x} = (x_1, x_2, \dots, x_n), \quad \vec{b} = (b_1, b_2, b_3, b_4), \quad \vec{Z} = (Z_1, Z_2, \dots, Z_d).$$

Assuming that each subscripted item belongs to \mathbb{R} we have $\vec{x} \in \mathbb{R}^n, \vec{b} \in \mathbb{R}^4, \vec{Z} \in \mathbb{R}^d$.

Notational conveniences for vectors: Unless something else is stated, we will always assume the following. If X is a nonempty set (usually, X is a set of numbers),

¹⁶We borrow that notation from physics.

$\vec{x} \in X^n$ is shorthand for $\vec{x} = (x_1, x_2, \dots, x_n) \in X^n$ (i.e., $x_j \in X$ for $j = 1, 2, \dots, n$.)

We also extend this convention to the case $X_1 \times \cdots \times X_n$ with potentially different sets X_j . This is best explained by example. Having pairs of numbers $a_i < b_i$ for $i = 1, 2, \dots, d$,

$\vec{y} \in]a_1, b_1] \times \cdots \times]a_d, b_d]$ is shorthand for
 $\vec{y} = (y_1, y_2, \dots, y_d)$, where $a_i < y_i \leq b_i$ for $i = 1, \dots, d$. \square

Example 2.12. Cartesian products occur in a natural manner in probability theory when one models the outcomes of repeated experiments.

- (a) If the experiment is three rolls of a die, then the set

$$\Omega = ([1, 6]_{\mathbb{Z}})^3 = \{1, 2, 3, 4, 5, 6\}^3$$

is a natural container for the outcomes of this experiment. For example, $(4, 2, 6) \in \Omega$ is the outcome of having rolled a 4 followed by a 2 followed by a 6.

- (b) n tosses of a coin ($n \in \mathbb{N}$) are modeled as follows. Let H stand for Heads and T for Tails. Then let

$$\Omega = \{H, T\}^n$$

For example, if $n = 5$, then $(H, H, T, H, T) \in \Omega$ models the outcome of having tossed Heads followed by Heads followed by Tails followed by Heads followed by Tails. This example demonstrates that cartesian products are also defined for sets that do not necessarily consist of numbers \square

Here is an abstract example.

Example 2.13. The graph Γ_f of a function with domain X and codomain Y (see def.2.20) is a subset of the cartesian product $X \times Y$. \square

Proposition 2.7. Let X_1, X_2, X_n be finite, nonempty sets. Then,

The size of the cartesian product is the product of the sizes of its factors, i.e.,

$$(2.28) \quad |X_1 \times X_2 \times \cdots \times X_n| = |X_1| \cdot |X_2| \cdot |X_3| \cdots |X_n|.$$

PROOF:

Case $n = 2$: This is trivial for two sets, since the proposition simply states that a matrix (a rectangular grid) of m rows and n columns possesses mn entries.

Case $n = 3$: For three sets X_1, X_2, X_3 , we arrange the $|X_1| \cdot |X_2|$ entries of $X_1 \times X_2$ into a single row. In other words, we consider the members $(x_i^{(1)}, x_j^{(2)}, x_k^{(3)})$ of $X_1 \times X_2 \times X_3$ as members $((x_i^{(1)}, x_j^{(2)}), x_k^{(3)})$

of $(X_1 \times X_2) \times X_3$. We apply the result for two sets to the cartesian product of $X_1 \times X_2$ and X_3 and obtain

$$|X_1 \times X_2 \times X_3| = |(X_1 \times X_2) \times X_3| = |X_1 \times X_2| \cdot |X_3| = |X_1| \cdot |X_2| \cdot |X_3|.$$

We repeat this procedure for $n = 3, 4, 5, \dots$ sets.

Case n : We arrange the elements of $X_1 \times X_2 \times \dots \times X_{n-1}$ into a single row and

interpret each $(x_1, \dots, x_n) \in X_1 \times X_n$ as $((x_1, \dots, x_{n-1}), x_n) \in (X_1 \times \dots \times X_{n-1}) \times X_n$.

Thus, the sets $X_1 \times \dots \times X_n$ and $(X_1 \times \dots \times X_{n-1}) \times X_n$ have the same size. We know from the prior step, case $n - 1$, that $|X_1 \times \dots \times X_{n-1}| = |X_1| \cdot \dots \cdot |X_{n-1}|$. Hence,

$$\begin{aligned} |X_1 \times \dots \times X_n| &= |(X_1 \times \dots \times X_{n-1}) \times X_n| = (|X_1 \times \dots \times X_{n-1}|) \cdot |X_n| \\ &= (|X_1| \cdot \dots \cdot |X_{n-1}|) |X_n| = |X_1| \cdot |X_2| \cdot |X_3| \cdot \dots \cdot |X_n|. \blacksquare \end{aligned}$$

2.6 A Few Things You May Not Have Learned in Calculus

Definition 2.25 (Absolute Convergence). ★

We say that an infinite series $\sum a_j (a_j \in \mathbb{R})$ is **absolutely convergent**, if

$$\sum_{j=1}^{\infty} |a_j| = |a_1| + |a_2| + |a_3| + \dots < \infty, \quad \square$$

Theorem 2.1.

If the series $\sum a_j (a_j \in \mathbb{R})$ is absolutely convergent, then the following holds true:

- (a) The series $\sum a_j$ itself converges, i.e., there is $a \in \mathbb{R}$ such that $\sum_{j=1}^{\infty} a_j = a$,
- (b) ANY rearrangement $\sum_{j=1}^{\infty} a_{n_j} = a_{n_1} + a_{n_2} + \dots$ converges to the same limit as $\sum a_j$.

Here, a **rearrangement**, of a sequence $(a_j)_{j \in \mathbb{N}}$ or series $\sum a_j$ if its members are rearranged into a sequence $(b_j)_{j \in \mathbb{N}}$ or series $\sum b_j$:

There are indices $n_j \in \mathbb{N}$ such that

$$b_1 = a_{n_1}, \quad b_2 = a_{n_2}, \quad b_3 = a_{n_3}, \quad \dots$$

and those indices satisfy the following:

- (1) They are distinct: $i \neq j \Rightarrow n_i \neq n_j$.
- (2) They leave no gaps in the set \mathbb{N} of all indices: For each $k \in \mathbb{N}$ there is $j \in \mathbb{N}$ such that $k = n_j$.

¹⁷ ★ We could have expressed (1) and (2) by stating that the assignment $j \mapsto n_j$ is a bijection $\mathbb{N} \rightarrow \mathbb{N}$. (See Definition 2.16 (Surjective, injective and bijective functions) on p.28.)

PROOF: See your calculus book. ■

Remark 2.16. ★ This remark might seem very strange to you. Assume that the series $\sum a_j$ is convergent, but not absolutely convergent: There is some $a \in \mathbb{R}$ such that $\sum_{j=1}^{\infty} a_j = a$, but $\sum_{j=1}^{\infty} |a_j| = \infty$. The following is known as Riemann's rearrangement theorem: Pick any $-\infty \leq b \leq \infty$. The terms a_j can be rearranged in such a way that the rearranged sequence, call it $\sum_{j=1}^{\infty} a_{n_j}$, converges to b . In other words, you can jumble the terms such that the limit is π . Some other rearrangement yields 0 as the limit, for yet another, the series converges to $-\sqrt{e^{30}}$, ...
□

We make the following blanket assumption.

2.7 Exercises for Ch.2

2.7.1 Exercises for Sets

Exercise 2.1. Prove (2.12) of prop.2.2 on p.20.

Exercise 2.2. Prove the set identities of prop.2.1.

Exercise 2.3. Prove that for any three sets A, B, C it is true that $(A \setminus B) \setminus C = A \setminus (B \cup C)$.

Hint: use De Morgan's formula (2.13.a). ■

Exercise 2.4. Let $X = \{x, y, \{x\}, \{x, y\}\}$. True or false?

- a. $\{x\} \in X$ c. $\{\{x\}\} \in X$ e. $y \in X$ g. $\{y\} \in X$
b. $\{x\} \subseteq X$ d. $\{\{x\}\} \subseteq X$ f. $y \subseteq X$ h. $\{y\} \subseteq X$ □

For the subsequent exercises refer to Definition 2.10 on p.22 of the size $|A|$ of a set A and to Definition 2.24 on p.35 of Cartesian products.

Exercise 2.5. Find the size of each of the following sets:

- a. $A = \{x, y, \{x\}, \{x, y\}\}$ c. $C = \{u, v, v, v, u\}$ e. $E = \{\sin(k\pi/2) : k \in \mathbb{Z}\}$
b. $B = \{1, \{0\}, \{1\}\}$ d. $D = \{3z - 10 : z \in \mathbb{Z}\}$ f. $F = \{\pi x : x \in \mathbb{R}\}$ □

Exercise 2.6. Let $X = \{x, y, \{x\}, \{x, y\}\}$ and $Y = \{x, \{y\}\}$. True or false?

- a. $x \in X \cap Y$ c. $x \in X \cup Y$ e. $x \in X \setminus Y$ g. $x \in X \Delta Y$
b. $\{y\} \in X \cap Y$ d. $\{y\} \in X \cup Y$ f. $\{y\} \in X \setminus Y$ h. $\{y\} \in X \Delta Y$ □

Exercise 2.7. Let $X = \{1, 2, 3, 4\}$ and let $Y = \{x, y\}$.

- a. What is $X \times Y$? c. What is $|X \times Y|$? e. Is $(x, 3) \in X \times Y$? g. Is $3 \cdot x \in X \times Y$?
b. What is $Y \times X$? d. What is $|X \times Y|$? f. Is $(x, 3) \in Y \times X$? h. Is $2 \cdot y \in Y \times X$? □

Exercise 2.8. Let $X = \{8\}$. What is $2^{(2^X)}$?

Exercise 2.9. Let $A = \{1, \{1, 2\}, 2, 3, 4\}$ and $B = \{\{2, 3\}, 3, \{4\}, 5\}$. Compute the following.

- a. $A \cap B$ b. $A \cup B$ c. $A \setminus B$ d. $B \setminus A$ e. $A \Delta B$ □

Exercise 2.10. Let A, X be sets such that $A \subseteq X$ and let $x \in X$. Prove the following:

- a. If $a \in A$ then $A = (A \setminus \{a\}) \uplus \{a\}$.
- b. If $a \notin A$ then $A = (A \uplus \{a\}) \setminus \{a\}$.

□

2.8 Addenda to Ch.2 – EMPTY

EMPTY – EMPTY – EMPTY – EMPTY – EMPTY

EMPTY – EMPTY – EMPTY – EMPTY – EMPTY

EMPTY – EMPTY – EMPTY – EMPTY – EMPTY

2.9 Blank Page after Ch.2

This page is intentionally left blank!

3 The Probability Model

3.1 Probability Spaces

In Section 1.2 (A First Look at Probability) we looked at several examples which motivated us to give a preliminary definition of probability as a function (we called it a probability measure),

$$P : 2^\Omega \longrightarrow [0, 1]$$

which assigns to each element A of the power set 2^Ω of a given set Ω ¹⁸ a number $P(A)$, also written as $P[A]$, between zero and one, such that

- (a) $P(\emptyset) = 0$ and $P(\Omega) = 1$. Here \emptyset denotes the empty set which contains no elements.
- (b) If the subsets A, B of Ω are disjoint, then probability is **additive**:

$$P\left(A \uplus B\right) = P(A) + P(B).$$

Note that additivity holds for three disjoint sets $A, B, C \in 2^\Omega$ since,

$$(*) \quad P\left(A \uplus B \uplus C\right) = P\left[\left(A \uplus B\right) \uplus C\right] = P\left(A \uplus B\right) + P(C) = P(A) + P(B) + P(C).$$

From (*) you get additivity for four disjoint $A, B, C, D \in 2^\Omega$ since,

$$\begin{aligned} P\left(A \uplus B \uplus C \uplus D\right) &= P\left[\left(A \uplus B \uplus C\right) \uplus D\right] \\ &= P\left(A \uplus B \uplus C\right) + P(D) = P[A] + P[B] + P[C] + P[D]. \end{aligned}$$

Now that you have additivity for four disjoint sets, you get it by the same method for five, and then for six, ... and thus, for any finite number of disjoint subsets A_1, \dots, A_n of Ω .

But we are not satisfied since it has proven extremely fruitful to replace (b) with the stronger condition

- (b') If $(A_n)_{n \in \mathbb{N}}$ is a (infinite!) sequence of disjoint subsets of Ω , then probability is " **σ -additive**":¹⁹

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

Unfortunately, this might come with a trade-off. Consider the following example.

Example 3.1. A point located somewhere at $]-\infty, 0[$ starts moving to the right at a constant velocity and is stopped at random somewhere in the unit interval $[0, 1]$ in the following sense: It is stopped just as likely in the left half, $[0, \frac{1}{2}]$, as in the right half, $[\frac{1}{2}, 1]$. More generally, for any $n \in \mathbb{N}$, it is stopped equally likely in each one of the intervals $[\frac{k-1}{n}, \frac{k}{n}]$ ($k = 1, 2, \dots, n$).

- It should be obvious that the only reasonable probability measure on $\Omega := [0, 1]$ is

$$(3.1) \quad P : [0, 1] \rightarrow [0, 1]; \quad [\alpha, \beta] \mapsto \beta - \alpha, \quad \text{where } 0 \leq \alpha \leq \beta \leq 1,$$

- since it is the only one that assigns probabilities proportionate to interval length (including $P([\alpha, \alpha]) = 0$ for intervals of length zero) and also satisfies $P(\Omega) = 1$.

¹⁸ $2^\Omega = \{ \text{all subsets of } \Omega \}$. See Definition 2.8 (Power set) on p.21.

¹⁹ σ ("sigma") is a greek letter. See the appendices for a complete list.

- Unfortunately, it has been proven ²⁰ that no σ -additive function that satisfies those properties exists on the entire power set of $[0, 1]$.
- The only way out of this dilemma without sacrificing σ -additivity is to relax the condition that $P(A)$ must exist for ALL $A \subseteq \Omega$. \square

Since we want to keep σ -additivity, we must define probability as a function

$$P : \mathfrak{F} \longrightarrow [0, 1], \quad \text{where } \mathfrak{F} \text{ is a suitable subset of } 2^\Omega,$$

which satisfies $P(\emptyset) = 0$ and $P(\Omega) = 1$ and

$$P\left(\biguplus_{k=1}^{\infty} A_k\right) = \sum_{k=1}^{\infty} P(A_k) \quad \text{for disjoint } A_1, A_2, \dots \in \mathfrak{F}.$$

Before we further explore this issue, we briefly remark on probability measures which are like the one described by 3.1 of the last example.

To answer the question what conditions a useful domain \mathfrak{F} for a probability measure P should satisfy, it helps to remember De Morgan's Law for finite or infinite sequences of sets. See Proposition 2.5 on p.34. Also, the following proposition which shows how to rewrite any countable union (finite or infinite) as a DISJOINT union will be relevant.

Proposition 3.1 (Rewrite unions as disjoint unions). *Let $(A_j)_{j \in \mathbb{N}}$ be a sequence of sets which all are contained within the universal set Ω . Let*

$$B_n := \bigcup_{j=1}^n A_j = A_1 \cup A_2 \cup \dots \cup A_n \quad (n \in \mathbb{N}),$$

$$C_1 := A_1 = B_1, \quad C_{n+1} := A_{n+1} \setminus B_n \quad (n \in \mathbb{N}).$$

Then

- (a) The sequence $(B_j)_j$ is increasing: $m < n \Rightarrow B_m \subseteq B_n$.
- (b) For each $n \in \mathbb{N}$, $\bigcup_{j=1}^n A_j = \bigcup_{j=1}^n B_j$.
- (c) The sets C_j are mutually disjoint and $\bigcup_{j=1}^n A_j = \biguplus_{j=1}^n C_j$.
- (d) The sets C_j ($j \in \mathbb{N}$) form a partition of the set $\bigcup_{j=1}^{\infty} A_j$.

PROOF: ★ (a) and (b) are trivial. For the proof of (c) and (d), convince yourself that

$$C_n = A_n \setminus (A_1 \cup A_2 \cup \dots \cup A_{n-1}).$$

Thus, C_n precisely contains those elements of A_n that have not previously been encountered! \blacksquare

We return to the question what the domain \mathfrak{F} of a probability should satisfy.

²⁰such a proof is outside the scope of these notes.

If A has a probability $P(A)$, then A^c should have probability $1 - P(A)$. Since probabilities can only be assigned to elements of \mathfrak{F} , we want

$$(A) \quad A \in \mathfrak{F} \Rightarrow A^c \in \mathfrak{F}.$$

If $A_n \in \mathfrak{F}$ are pairwise disjoint, then $\biguplus_{j=1}^{\infty} A_j$ should have probability $\sum_{j=1}^{\infty} P(A_j)$. Since probabilities can only be assigned to elements of \mathfrak{F} , we want

$$A_n \in \mathfrak{F} \text{ disjoint} \Rightarrow \biguplus_{j=1}^{\infty} A_j \in \mathfrak{F}.$$

Since we have seen that any union of a sequence of events can be written as a disjoint union, we need more than the above. We really want

$$(B) \quad A_n \in \mathfrak{F} \text{ arbitrary} \Rightarrow \bigcup_{j=1}^{\infty} A_j \in \mathfrak{F}.$$

Also, it is very reasonable to demand that $P(\emptyset) = 0$ for the impossible event which contains no potential outcomes, i.e., the empty set. It is just as reasonable to ask that $P(\Omega) = 1$ for the sure event, Ω , since it contains all potential outcomes. Thus, we ask that

$$(C) \quad \emptyset \in \mathfrak{F} \quad \text{and} \quad \Omega \in \mathfrak{F}.$$

All this leads to the definition of σ -algebras as suitable domains for probability measures.

Definition 3.1 (σ -algebra). Let Ω be a nonempty set and $\mathfrak{F} \subseteq 2^\Omega$ a collection of subsets of Ω , such that

- (a) $A \in \mathfrak{F} \Rightarrow A^c \in \mathfrak{F}$.
 (b) $A_n \in \mathfrak{F}$ arbitrary $\Rightarrow \bigcup_{j=1}^{\infty} A_j \in \mathfrak{F}$.
 (c) $\emptyset \in \mathfrak{F}$.

Then we call \mathfrak{F} a σ -algebra.

\mathfrak{F} is also called a σ -field, but this is considered old-fashioned terminology. \square

Proposition 3.2. σ -algebras \mathfrak{F} satisfy the following.

- (a) $\Omega \in \mathfrak{F}$.
 (b) $A_1, A_2, \dots, A_n \in \mathfrak{F} \Rightarrow A_1 \cup A_2 \cup \dots \cup A_n \in \mathfrak{F}$.
 (c) Let $n \in \mathbb{N}$ and $A_1, A_2, \dots \in \mathfrak{F}$. Let $A = \bigcap_{k=1}^n A_k$ and $B = \bigcap_{k=1}^{\infty} A_k$. Then $A \in \mathfrak{F}$ and $B \in \mathfrak{F}$. \square

PROOF: ★

PROOF of (a): True, since $\Omega = \emptyset^c$ and complements of elements of \mathfrak{F} belong to \mathfrak{F} and $\emptyset \in \mathfrak{F}$.

PROOF of **(b)**: Since any finite list A_1, \dots, A_n can be written as an infinite sequence

$$B_1 = A_1, B_2 = A_2, \dots, B_n = A_n, B_{n+1} = B_{n+2} = \dots = \emptyset$$

and since $B_j \in \mathfrak{F}$ for each $j \in \mathbb{N}$, it follows from Def.3.1**(b)** that $\bigcup_{j=1}^{\infty} B_j \in \mathfrak{F}$. Since

$$\bigcup_{j=1}^n A_j = \bigcup_{j=1}^n A_j \cup \emptyset \cup \emptyset \cup \dots \cup \emptyset = \bigcup_{j=1}^{\infty} B_j$$

it follows that $\bigcup_{j=1}^n A_j \in \mathfrak{F}$. This proves **(b)**.

PROOF of **(c)**: According to De Morgan's laws, any countable intersection can be written as the union of its complements. Thus we automatically get from **(A)** and **(B)** that countable intersections of a sequence in \mathfrak{F} will again belong to \mathfrak{F} .

Here is a detailed argument. For each j let $C_j := A_j^c$. Further, let $C := \bigcup_{j=1}^n C_j$ and $D := \bigcup_{j=1}^{\infty} C_j$.

Since each C_j is the complement of a member of \mathfrak{F} , we have $C_j \in \mathfrak{F}$. Thus, $D \in \mathfrak{F}$ by the definition of \mathfrak{F} , and we have seen in part **(b)** of this proposition that $C \in \mathfrak{F}$.

It follows from De Morgan's laws that $C^c = A$ and $D^c = B$.

Thus, both A, B belong to \mathfrak{F} as complements of elements of \mathfrak{F} . We have shown **(c)**. ■

Definition 3.2 (Probability measures and probability spaces).

Given are a nonempty set Ω with a σ -algebra $\mathfrak{F} \subseteq 2^\Omega$ and a function

$$P : \mathfrak{F} \longrightarrow [0, 1]; \quad A \mapsto P(A) \quad \text{as follows.}$$

$$(3.2) \quad P(\emptyset) = 0,$$

$$(3.3) \quad P(\Omega) = 1,$$

$$(3.4) \quad (A_n)_{n \in \mathbb{N}} \in \mathfrak{F} \text{ disjoint} \Rightarrow P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n=1}^{\infty} P(A_n) = \sum_{n \in \mathbb{N}} P(A_n). \quad (\sigma\text{-additivity})$$

- We call P a **probability measure** or simply a **probability**.
- The triplet $(\Omega, \mathfrak{F}, P)$ is called a **probability space**.
- We often call disjoint events **mutually exclusive events**. □

Notation 3.1 (Sample spaces and sample points).

- We also call a probability space a **sample space** and an outcome a **sample point**.
- we will also call Ω by itself (as opposed to the triplet $(\Omega, \mathfrak{F}, P)$) a probability space or sample space, but sometimes we refer to Ω as the **carrier set** or **carrier** of $(\Omega, \mathfrak{F}, P)$.
- We like to write Ω for the carrier set, \mathfrak{F} for the σ -algebra and P for the probability measure of a probability space, but different notation may be used. For example, there may be a probability space (S, \mathcal{S}, Q) and outcomes s or x or \vec{y} (vector notation).

Remark 3.1. As I noted in Section 1.2 (A First Look at Probability), “sample space” is the statistician’s terminology for a probability space. I will mostly use the term “probability space”, since we usually think of a sample as a list of items that has been picked in some random fashion from an underlying “population”. We will consider probability spaces in this lecture where it would require a huge stretch of the imagination to consider their elements as such samples. Note though that there are occasions where the term “sample space” seems to be superior terminology.

You, my students, may choose whatever notation you prefer.

And more good news: We have introduced σ -algebras to properly deal with the issue that was raised in Example 3.1 on p.41. It won’t be long and we will on only few occasions deal with σ -algebras and usually refer to probability spaces (Ω, P) and (S, P) . \square

Remark 3.2. How do we interpret $P\left(\bigsqcup_{n \in \mathbb{N}} A_n\right) = \sum_{n=1}^{\infty} P(A_n) = \sum_{n \in \mathbb{N}} P(A_n)$ (formula (3.4) for σ -additivity in the definition of a probability measure)? There are two issues.

- (a) What is the meaning of $\bigsqcup_{n \in \mathbb{N}} A_n$ as opposed to $\bigsqcup_{n=1}^{\infty} A_n$?
- (b) What is the meaning of $\sum_{n \in \mathbb{N}} P(A_n)$, as opposed to $\sum_{n=1}^{\infty} P(A_n)$? Does it really not matter in which order we add the terms of an infinite series?

The answer to (a) is easy. Unions are defined without any reference to an order “first A_1 , then A_2 , then A_3, \dots ”, since the definition of $a \in \bigsqcup_{n \in \mathbb{N}} A_n$ is the existence of at least one index i_0 such that $a \in A_{i_0}$. No reference to an ordering is made. The only justification for the notation $\bigsqcup_{n=1}^{\infty} A_n$ is that it looks more familiar. By the way, what was said here about disjoint unions also applies to arbitrary unions and to intersections.

To answer (b), let us assume until this question is settled, that (3.4) has been replaced by

$$P\left(\bigsqcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n).$$

That formula has none of the issues we are trying to resolve.

Note that the series $\sum P(A_n)$ is absolutely convergent.²¹ To see this, let $A := \bigsqcup_{n=1}^{\infty} A_n$. Clearly, $P(A_n) \geq 0$ for all n . Moreover, by (σ -)additivity applied to $A \bigsqcup A^c = \Omega$,

$$P\left(\bigsqcup_{n=1}^{\infty} A_n\right) \leq P(A) + P\left(\bigsqcup_{n=1}^{\infty} A_n\right) = P(\Omega) = 1 < \infty.$$

Since $\sum P(A_n)$ is absolutely convergent, it does indeed not matter how the terms A_n are arranged. See Theorem 2.1 on p.37. \square

In Section 1.2 (A First Look at Probability) we used throws of a die to illustrate the concepts of random actions and their potential outcomes and let this motivate us to give a preliminary definition of

²¹See Definition 2.25 (Absolute Convergence) on p.37.

a probability measure as a function. Now that we have the final definition of a probability measure, we should study some more examples.

Example 3.2. We model k rolls of a fair die ($k \in \mathbb{N}$) as follows. Let

$$\Omega := \{1, 2, 3, 4, 5, 6\}^k = \{(a_1, a_2, \dots, a_k) : a_j = 1, 2, \dots, 6 \text{ for each } j = 1, 2, \dots, k\}.$$

For example, let $k = 5$. then $\omega_1 = (2, 6, 2, 1, 4) \in \Omega$. On the other hand, $\omega_2 = (2, 6, 2, 9, 4) \notin \Omega$, since $a_j = 1, 2, \dots, 6$ is not true for $j = 4$ (because $a_4 = 9$).

Ω is a finite set, and you will learn later that its size is 6^k . Thus, $\Omega = \{\omega_1, \omega_2, \dots, \omega_{6^k}\}$ where, e.g.,

$$\omega_1 = (1, 1, \dots, 1, 1), \omega_2 = (1, 1, \dots, 1, 2), \dots, \omega_{6^k-1} = (6, 6, \dots, 6, 5), \omega_{6^k} = (6, 6, \dots, 6, 6).$$

Since the die is fair, each one of those 6^k elements of Ω should have the same probability $p := P(\{\omega\})$ for all $\omega \in \Omega$. Since $P(\Omega) = 1$ and

$$\Omega = \bigsqcup [\{\omega\} : \omega \in \Omega] = \bigsqcup_{j=1}^{\infty} \{\omega_j\}.$$

is a union of a sequence of disjoint set, we obtain from the σ -additivity of $P(\cdot)$ the following:

$$1 = P(\Omega) = \sum_{j=1}^{6^k} P\{\omega_j\} = 6^k p \Rightarrow p = \frac{1}{6^k}.$$

- So then, how does one define a probability measure $P : \mathfrak{F} \rightarrow [0, 1]$?
- And what is that σ -algebra \mathfrak{F} going to be?

To answer those questions, we define the function $P : 2^\Omega \rightarrow \mathbb{R}$ as follows.

$$(3.5) \quad P(A) := \frac{|A|}{|\Omega|} = \frac{|A|}{6^k}.$$

Observe the following.

- (1) $A \subseteq \Omega \Rightarrow 0 \leq |A| \leq |\Omega| = 6^k \Rightarrow 0 \leq P(A) \leq 1$.
- (2) The empty set has size $|\emptyset| = 0$ and $|\Omega| = 6^k$. Thus, $P(\emptyset) = 0$ and $P(\Omega) = 1$.
- (3) Assume that A_1, A_2, \dots are disjoint subsets of Ω . Since Ω is finite, only finitely many A_j are not empty (THINK!),
- (4) We rearrange the sequence such that the nonempty members will be A_1, A_2, \dots, A_m for some suitable m .
- (5) Then, $A = A_1 \sqcup A_2 \sqcup \dots \sqcup A_m$ is a finite union and disjointness of the $A_j \Rightarrow |A| = |A_1| + |A_2| + \dots + |A_m|$
- (6) Thus, σ -additivity: $P(A) = |A|/6^k = \sum_{j=1}^m (|A_j|/6^k) = \sum_{j=1}^m P(A_j) = \sum_{\text{all } j} P(A_j)$

Last equation: The omitted sets A_{m+1}, A_{m+2}, \dots were empty, thus $P(A_j) = 0/6^k = 0$ for those j .

We obtain from (1) – (6) that $P(A) = |A|/6^k$ is a probability measure on 2^Ω . \square

Example 3.3. One easily sees the generalization of the last example to arbitrary finite sets:

Let Ω be a finite set of size $N := |\Omega| < \infty$. Let the function $P : 2^\Omega \rightarrow \mathbb{R}$ be given as

$$(3.6) \quad P(A) := \frac{|A|}{|\Omega|} = \frac{|A|}{N}.$$

Then everything stated in **(1) – (6)** of **(a)** remains valid if we replace 6^k with N , and this shows that P is a probability measure on 2^Ω . \square

Definition 3.3 (Equiprobability).

Let (Ω, P) be a finite probability space, i.e., $|\Omega| < \infty$. Let $n := |\Omega|$. We say that P has **equiprobable** outcomes and also, that P **satisfies equiprobability**, if

$$(3.7) \quad P(\{\omega\}) = \frac{1}{|\Omega|} \quad (\text{since then } P\{\omega\} \text{ is constant for all } \omega \in \Omega). \quad \square$$

Remark 3.3. The finiteness of Ω was crucial in the last two examples for the following reason.

If Ω is infinite and countable, then $\Omega = \{\omega_1, \omega_2, \dots\}$ can be written as an infinite sequence of distinct(!) members. It is not possible to define a “uniform” probability measure on Ω as we did in parts **(a)** and **(b)**, i.e., a number p such that $P(\omega_j) = p$ for all $j \in \mathbb{N}$. How so?

- (1) p would have to be strictly positive: Otherwise, $P(\Omega) = \sum_j P(\omega_j) = p + p + \dots \leq 0$, but we require $P(\Omega) = 1$.
- (2) Thus, $p > 0$. Thus, $P(\Omega) = \sum_j P(\omega_j) = p + p + \dots = \infty$, but we require $P(\Omega) = 1$. \square

Remark 3.4. We will see that the most important probability measures on the uncountable set \mathbb{R} ²² satisfy $P(x) = 0$ for all $x \in \mathbb{R}$. That is no contradiction to σ -additivity and $P(\mathbb{R}) = 1$, since one cannot write the real numbers as a countable union $\mathbb{R} = \{x_1\} \uplus \{x_1\} \uplus \{x_2\} \uplus \dots$. Obviously, it is no more possible in those cases to determine a probability measure on \mathbb{R} by only listing the probabilities

$P(x)$ of the atomic events $x \in \mathbb{R}$. Rather, P often is characterized by integrals $P([a, b]) = \int_a^b \varphi(t) dt$.

(And if this is the case, we will indeed obtain $P(x) = \int_x^x \varphi(t) dt = 0$ for all x .) \square

Recall for the next theorem that we denote by $A_n \uparrow$ a nondecreasing sequence of events: $i < J \Rightarrow A_i \subseteq A_j$ and by $B_n \downarrow$ a nondecreasing sequence of events: $i < J \Rightarrow B_i \supseteq B_j$. (See Definition 2.20 on p.30.)

Theorem 3.1 (Continuity property of probability measures). ★

²²the so-called distributions of continuous random variables

Let $(\Omega, \mathfrak{F}, P)$ be a probability space. If $A_n, B_n \in \mathfrak{F}$, then the following is true:

$$(3.8) \quad A_n \uparrow \Rightarrow P(A_n) \uparrow P\left(\bigcup_{n \in \mathbb{N}} A_n\right),$$

$$(3.9) \quad B_n \downarrow \Rightarrow P(B_n) \downarrow P\left(\bigcap_{n \in \mathbb{N}} B_n\right).$$

PROOF: We prove (3.8) as follows: Let $A := \bigcup_{j=1}^{\infty} A_j$ and

$$C_1 := A_1, \quad C_{n+1} := A_{n+1} \setminus A_n \quad (n \in \mathbb{N}).$$

Note that $A_n \uparrow \Rightarrow A_n = \bigcup_{j=1}^n A_j$ and thus, $C_{n+1} := A_{n+1} \setminus \left(\bigcup_{j=1}^n A_j\right)$.

According to Proposition 3.1 (Rewrite unions as disjoint unions) on p.42, the sets C_j form a partition of A and we have

$$A_n = \bigsqcup_{j=1}^n C_j, \quad A = \bigsqcup_{j=1}^{\infty} C_j,$$

It follows from the σ -additivity of P that

$$P(A) = P\left(\bigsqcup_{j=1}^{\infty} C_j\right) = \sum_{j=1}^{\infty} P(C_j) = \lim_{n \rightarrow \infty} \sum_{j=1}^n P(C_j) = \lim_{n \rightarrow \infty} P\left(\bigsqcup_{j=1}^n C_j\right) = \lim_{n \rightarrow \infty} P(A_n).$$

This proves (3.8). We use this result to prove (3.9) as follows.

Let $B := \bigcap_{j=1}^{\infty} B_j$. For $n \in \mathbb{N}$, let $A_n := B_n^c$. Further, let $A := \bigcup_{j=1}^{\infty} A_j$. Then $A_n \uparrow$ and it follows from De Morgan that

$$A^c = \left(\bigcup_{j=1}^{\infty} A_j\right)^c = \bigcap_{j=1}^{\infty} A_j^c = \bigcap_{j=1}^{\infty} B_j = B.$$

We apply (3.8) and obtain

$$1 - P(B_n) = P(A_n) \uparrow P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = 1 - P\left[\left(\bigcup_{n \in \mathbb{N}} A_n\right)^c\right] = 1 - P(B).$$

Thus, $P(B_n) \downarrow P(B)$ and this proves (3.9). ■

Definition 3.4 (Discrete probability space).

Assume that the probability space $(\Omega, \mathfrak{F}, P)$ satisfies the following:

- (a) $P(\{\omega\})$ is defined for all $\omega \in \Omega$. In other words, we ask that $\{\omega\} \in \mathfrak{F}$ for all $\omega \in \Omega$.
- (b) There exists a countable subset A^* of Ω such that $\sum_{\omega \in A^*} P\{\omega\} = 1$

Then we call $(\Omega, \mathfrak{F}, P)$ a **discrete probability space**. \square

We will later on talk about discrete and continuous random variables, but note that there is no such thing as a “continuous probability space”.

Remark 3.5. For the interpretation of the summation $\sum_{\omega \in A^*} P\{\omega\}$ we note the following.

- (a) Either A^* is finite and can be written $A^* = \{\omega_1, \omega_2, \dots, \omega_n\}$ for some suitable n .
Then $\sum_{\omega \in A^*} P\{\omega\} = \sum_{j=1}^n P\{\omega_j\}$.
- (b) Or, A^* is infinite and can be written $A^* = \{\omega_j : j \in \mathbb{N}\}$. We reason as in Remark 3.2 on p.45 with $\{\omega_j\}$ in place of A_j and see that the series $\sum P\{\omega_j\}$ is absolutely convergent. Thus, the value of $\sum_{j=1}^n P\{\omega_j\}$ does not depend on how the elements of A^* were sequenced and we can write $\sum_{\omega \in A^*} P\{\omega\}$ for that common value. \square

Theorem 3.2.

Let $(\Omega, \mathfrak{F}, P)$ be a discrete probability space and $A^* \in \mathfrak{F}$ a countable set such that $\sum_{\omega \in A^*} P\{\omega\} = 1$.

Then

- (a) $A^* \in \mathfrak{F}$.
- (b) $P(A^*) = 1$ and thus, $P((A^*)^c) = 0$.
- (c) $P(A) = P(A \cap A^*)$ for all $A \in \mathfrak{F}$.
- (d) $P(A) = \sum_{\omega \in A \cap A^*} P\{\omega\}$ for all $A \in \mathfrak{F}$.
- (e) \star The formula $\tilde{P}(B) := P(B \cap A^*)$ “extends” P to a probability measure \tilde{P} on the entire power set 2^Ω .

PROOF: \star

PROOF of (a): This is true, because $\{\omega\} \in \mathfrak{F}$ for all ω and $A^* = \biguplus_{\omega \in A^*} \{\omega\}$ is a countable union of elements of \mathfrak{F} .

PROOF of (b): By definition, $\sum_{\omega \in A^*} P\{\omega\} = 1$. Since $A^* = \biguplus_{\omega \in A^*} \{\omega\}$, we obtain $P(A^*) = 1$.

Further, $\Omega = A \biguplus (A^*)^c \Rightarrow 1 = P(A^*) + P((A^*)^c) = 1 + P((A^*)^c)$. Thus, $P((A^*)^c) = 0$.

PROOF of (c): From $0 \leq P(A \cap (A^*)^c) \leq P((A^*)^c) = 0$, we obtain $P(A \cap (A^*)^c) = 0$.

From $A = [A \cap A^*] \biguplus [A \cap (A^*)^c]$, we obtain $P(A) = P(A \cap A^*) + P(A \cap (A^*)^c) = P(A \cap A^*)$.

PROOF of **(d)**: $A \cap A^*$ is a subset of A^* , hence, countable. Thus, $P(A \cap A^*) = \sum_{\omega \in A \cap A^*} P\{\omega\}$. We obtain from **(c)** that $P(A) = \sum_{\omega \in A \cap A^*} P\{\omega\}$.

PROOF of **(e)**: Tedious but easy, if one uses **(c)** and distributivity $A^* \cap \bigcup_j A_j = \bigcup_j (A^* \cap A_j)$. ■

Corollary 3.1.

- (a)** If $(\Omega, \mathfrak{F}, P)$ be a discrete probability space, then P is characterized by the probabilities $P\{\omega\}$ of the outcomes ω .
- (b)** Let Ω be some arbitrary, nonempty set. Assume that $(p_j)_j$ is a finite or infinite sequence of real numbers that satisfies
- $p_j \geq 0$ for all j and $\sum_j p_j = 1$
- Further, assume that $(\omega_j)_j$ is a corresponding sequence of distinct elements of Ω , then $(p_j)_j$ defines a discrete probability space $(\Omega, 2^\Omega, P)$ as follows.
- $P(\emptyset) := 0$, $P(A) := \sum_{j: \omega_j \in A} p_j$, for $A \neq \emptyset$. □

PROOF: ★ This follows from Theorem 3.2. The details are left to the reader. ■

Remark 3.6. The probability spaces $(\Omega, \mathfrak{F}, P)$ we will be faced with are in one of the following categories:

- (a)** $(\Omega, \mathfrak{F}, P)$ is a discrete probability space. According to Theorem 3.2(e) on p.49, we may choose $\mathfrak{F} = 2^\Omega$.
- (b)** $\Omega = \mathbb{R}$ and $P(A)$ is known (at a minimum) for intervals such as $[a, b]$ or $]a, b]$ or $[a, b[$ or $]a, b[$.
- (c)** $\Omega = \mathbb{R}^n$ and $P(A)$ is known (at a minimum) for n -dimensional rectangles such as $[a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_n, b_n]$ (cartesian products of onedimensional intervals!)

It is important that we can assign probabilities to Intervals in **(c)** and n -dimensional rectangles in **(d)**, for the following reason.

- (c')** the most important probabilities P defined for sets in \mathbb{R} come with a so called **probability density function** $f : \mathbb{R} \rightarrow [0, \infty[$ which assigns to an interval $]a, b]$ the probability

$$P(]a, b]) = \int_a^b f(u) du.$$

This makes it plausible that the σ -algebra \mathfrak{B} for such P should contain all intervals $]a, b]$.

- (d') Likewise, the most important probabilities P defined for sets in \mathbb{R}^n come with a probability density function $f : \mathbb{R}^n \rightarrow [0, \infty[$ which assigns to an n -dimensional rectangle $]a_1, b_1] \times]a_2, b_2] \times \cdots \times]a_n, b_n]$ the probability

$$\begin{aligned} P(]a_1, b_1] \times]a_2, b_2] \times \cdots \times]a_n, b_n]) &= \int_{a_n}^{b_n} \int_{a_{n-1}}^{b_{n-1}} \cdots \int_{a_1}^{b_1} f(\vec{u}) \, d\vec{u} \\ &= \int_{a_n}^{b_n} \int_{a_{n-1}}^{b_{n-1}} \cdots \int_{a_1}^{b_1} f(u_1, \dots, u_n) \, du_1 \, du_2 \cdots du_{n-1} \, du_n. \end{aligned}$$

Thus, the σ -algebra \mathfrak{B}^n for such P should contain all rectangles $]a_1, b_1] \times]a_2, b_2] \times \cdots \times]a_n, b_n]$

You may have noticed that we could have worked with either of $]a_j, b_j]$, $[a_j, b_j[$, $[a_j, b_j]$ instead of $]a_j, b_j]$, since $\int_a^a \dots da$ is always zero. Nevertheless, it is more convenient to work with intervals that are open on the left and closed on the right. We will see that when we deal with the so-called cumulative distribution functions on \mathbb{R} and \mathbb{R}^n . \square

Theorem 3.3. ★

Let Ω be some arbitrary set and $(\mathfrak{F}_i)_{i \in I}$ a family of σ -algebras on Ω , i.e., $\mathfrak{F}_i \subseteq 2^\Omega$ for each $i \in I$. No assumption is made about the index set other than $I \neq \emptyset$. Thus, this family may consist of finitely many σ -algebras or of entire sequence or even uncountably many σ -algebras.

- Let $\mathfrak{F} := \bigcap_{i \in I} \mathfrak{F}_i$, i.e., $\mathfrak{F} = \{A \subseteq \Omega : A \in \mathfrak{F}_i \text{ for each index } i\}$. Then \mathfrak{F} is a σ -algebra.

This can also be stated as follows.

Any intersection of σ -algebras results in a σ -algebra.

PROOF: Left to the interested reader. \blacksquare

Theorem 3.4. ★

Let Ω be some arbitrary set and $\mathcal{A} \subseteq 2^\Omega$. In other words, each element of \mathcal{A} is a subset of Ω .

- Then there exists a minimal (i.e., smallest) σ -algebra that contains \mathcal{A} .
- Further, this σ -algebra is uniquely determined by \mathcal{A} . This allows us to name it $\sigma\{\mathcal{A}\}$.

PROOF: We obtain $\sigma\{\mathcal{A}\}$ as the intersection of all σ -algebras that contain \mathcal{A} . According to Theorem 3.3, this intersection is a σ -algebra. \blacksquare

Since the minimal σ -algebra that contains \mathcal{A} is uniquely determined by \mathcal{A} , we can make the following definition.

Definition 3.5. ★

Let Ω be some arbitrary set and $\mathcal{A} \subseteq 2^\Omega$. We call $\sigma\{\mathcal{A}\}$ the σ -algebra generated by \mathcal{A} . If \mathcal{A} is of the form $\mathcal{A} = \{\dots\}$, we also write $\sigma\{\dots\}$ for $\sigma\{\{\dots\}\}$. \square

The next definition is marked optional, but note that Borel sets will be mentioned frequently during lecture.

Definition 3.6. ★ We apply the above to the sets of onedimensional and n -dimensional intervals.

- the smallest σ -algebra of subsets of \mathbb{R} which contains all intervals of real numbers. It is denoted \mathfrak{B} .
- the smallest σ -algebra of subsets of \mathbb{R}^n which contains all n -dimensional rectangles. It is denoted \mathfrak{B}^n .

We call \mathfrak{B} and \mathfrak{B}^n the Borel σ -algebras of \mathbb{R} and of \mathbb{R}^n **Borel σ -algebra** and we call their members **Borel sets**.

It is sufficient for this course that you just remember that

- The Borel sets are the sufficiently well behaved sets of \mathbb{R} and \mathbb{R}^n
- The intervals and n -dimensional rectangles are among those sets.
- Only completely weird and useless sets are not Borel. \square

Remark 3.7. (A) Consider the following sets of intervals of real numbers.

$$\mathfrak{I}_1 := \{[a, b] : a < b\}, \quad \mathfrak{I}_2 := \{[a, b[: a < b\},$$

$$\mathfrak{I}_3 := \{]a, b] : a < b\}, \quad \mathfrak{I}_4 := \{]a, b[: a < b\}.$$

One can show that each one of those sets of intervals is big enough to generate the Borel sets of \mathbb{R} : $\mathfrak{B} = \sigma(\mathfrak{I}_1) = \sigma(\mathfrak{I}_2) = \sigma(\mathfrak{I}_3) = \sigma(\mathfrak{I}_4)$.

(B) The above generalizes to n -dimensional space: Let

$$\mathfrak{I}_5 := \{[a_1, b_1] \times]a_2, b_2[\times \cdots \times]a_n, b_n[: a_1 < b_1, a_2 < b_2, \dots, a_n < b_n\},$$

$$\mathfrak{I}_6 := \{[a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_n, b_n] : a_1 < b_1, a_2 < b_2, \dots, a_n < b_n\},$$

$$\mathfrak{I}_7 := \{]a_1, b_1[\times [a_2, b_2] \times \cdots \times [a_n, b_n] : a_1 < b_1, a_2 < b_2, \dots, a_n < b_n\},$$

$$\mathfrak{I}_8 := \{]a_1, b_1[\times]a_2, b_2[\times \cdots \times]a_n, b_n[: a_1 < b_1, a_2 < b_2, \dots, a_n < b_n\},$$

one can show that $\mathfrak{B}^n = \sigma(\mathfrak{I}_5) = \sigma(\mathfrak{I}_6) = \sigma(\mathfrak{I}_7) = \sigma(\mathfrak{I}_8)$. \square

Fact 3.1. ★ For the following, note that the sets $\mathfrak{I}_1, \dots, \mathfrak{I}_8$ were defined in Example 3.7 on p.52.

(a) Let \mathfrak{I} denote any one of the collections of intervals $\mathfrak{I}_1, \dots, \mathfrak{I}_4$. Let $\mathcal{E} := \mathfrak{I} \uplus \mathbb{R}$. Then any function $P_0 : \mathcal{E} \rightarrow [0, 1]$ which satisfies $P_0(\emptyset) = 0$, $P_0(\mathbb{R}) = 1$ and σ -additivity on \mathcal{E} : $E_n \in \mathcal{E}$ disjoint such that $E := \biguplus_{n \in \mathbb{N}} E_n \in \mathcal{E} \Rightarrow P(E) = \sum_{n \in \mathbb{N}} P(E_n)$, can be uniquely extended to a probability measure on \mathfrak{B} , the Borel sets of \mathbb{R} .

(b) Let \mathfrak{I} denote any one of the collections of intervals $\mathfrak{I}_5, \dots, \mathfrak{I}_8$. Let $\mathcal{E} := \mathfrak{I} \uplus \mathbb{R}^n$. Then any function $P_0 : \mathcal{E} \rightarrow [0, 1]$ which satisfies $P_0(\emptyset) = 0$, $P_0(\mathbb{R}^n) = 1$ and σ -additivity on \mathcal{E} : $E_n \in \mathcal{E}$ disjoint such that $E := \biguplus_{n \in \mathbb{N}} E_n \in \mathcal{E} \Rightarrow P(E) = \sum_{n \in \mathbb{N}} P(E_n)$, can be uniquely extended to a probability measure on \mathfrak{B}^n , the Borel sets of \mathbb{R}^n . \square

Remark 3.8. Consider this a continuation of Remark 3.6. We can summarize it as follows. There are essentially only two kinds of probability spaces $(\Omega, \mathfrak{F}, P)$ we are interested in.

- There is a countable subset A^* of Ω such that $\sum_{\omega \in A^*} P(\{\omega\}) = 1$ (discrete probability spaces). Then $\mathfrak{F} = 2^\Omega$, since the above allows us to define $P(A)$ for arbitrary $A \subseteq \Omega$ as

$$P(A) = \sum_{\omega \in A^* \cap A} P(\{\omega\}).$$

- $\Omega = \mathbb{R}$ or $\Omega = \mathbb{R}^n$. Then $\mathfrak{F} = \text{the Borel sets}$.

Now that we understand the structure of the domain \mathfrak{F} of the probability measures P we will be dealing with, there is no more need to keep carrying this baggage with us.

Henceforth, we will, with very few exceptions, do the following.

We will ignore that probability measures cannot always be given on the entire power set 2^Ω (true only we deal with $(\mathbb{R}, \mathfrak{B}, P)$ or $(\mathbb{R}^n, \mathfrak{B}^n, P)$) and that this necessitated us to introduce a σ -algebra \mathfrak{F} as the domain of that probability measure. Accordingly, we will ignore the σ -algebra and talk about

probability spaces (sample spaces) (Ω, P) , rather than $(\Omega, \mathfrak{F}, P)$. \square

Notational conveniences for probabilities:

If we have a set that is written as $\{\dots\}$, i.e., with curly braces as delimiters, then we may write its probability as $P\{\dots\}$ instead of $P(\{\dots\})$. Specifically for singletons $\{\omega\}$, it is OK to write $P\{\omega\}$.

The next theorem lists two important rules to determine probabilities.

Theorem 3.5 (WMS Ch.02.8, Theorem 2.6). *If A and B are two events in a probability space (Ω, P) , then*

$$(3.10) \quad \textit{Additive Law of Probability:} \quad P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

$$(3.11) \quad \textit{Rule of the Complement:} \quad P[A^c] = 1 - P[A].$$

PROOF of (3.10): We apply the σ -additivity of P as follows:

- (1) $A = (A \setminus B) \uplus (A \cap B)$ and $B = (B \setminus A) \uplus (A \cap B)$
 $\Rightarrow P(A) + P(B) = P(A \setminus B) + P(A \cap B) + P(B \setminus A) + P(A \cap B)$
- (2) $A \cup B = (A \setminus B) \uplus (A \cap B) \uplus (B \setminus A)$
 $\Rightarrow P(A \cup B) = P(A \setminus B) + P(A \cap B) + P(B \setminus A)$

Thus, from (1) and (2), $P(A) + P(B) = P(A \cup B) + P(A \cap B)$.

It follows that $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

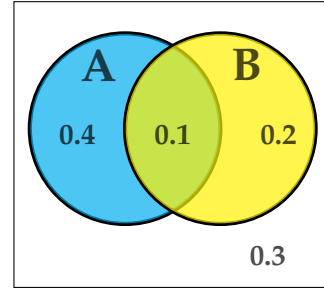
PROOF of (3.11): Immediate from the σ -additivity of P and $\Omega = A \uplus A^c$. \blacksquare

Remark 3.9. If the events A and B are mutually exclusive, i.e., $A \cap B = \emptyset$, then $P[A \cap B] = 0$ and the additive law of probability simply is σ -additivity

$$(3.12) \quad P(A \cup B) = P(A) + P(B). \quad \square$$

Remark 3.10. The additive law of probability is very easy to apply, since all you need is $P(A)$, $P(B)$ and $P(A \cap B)$.

Nevertheless it might be fastest to draw a Venn diagram. Assume you know that $P(A) = 0.5$, $P(B) = 0.3$, $P(A \cap B) = 0.1$. Clearly, $P(A \setminus B) = P(A) - P(A \cap B) = 0.4$ and $P(B \setminus A) = P(B) - P(A \cap B) = 0.2$. It is now immediate that $P(A \cup B) = 0.7$ and we get for free that $P(A \cup B^c) = 0.3$.



The additive law of probability has generalizations for the probability of the union of three or more events.

Theorem 3.6 (Exclusion–Inclusion formula for 3 events). ★

If A_1, A_2, A_3 are events in a probability space (Ω, P) , then

$$(3.13) \quad P(A_1 \cup A_2 \cup A_3) = [P(A_1) + P(A_2) + P(A_3)] - [P(A_1 \cap A_2) + P(A_1 \cap A_3) + P(A_2 \cap A_3)] + P(A_1 \cap A_2 \cap A_3).$$

PROOF: We apply the additive law of probability to the sets A_1 and $A_2 \cup A_3$ and obtain

$$(A) \quad P[A_1 \cup A_2 \cup A_3] = P[A_1] + P[A_2 \cup A_3] - P[A_1 \cap (A_2 \cup A_3)].$$

Next, we apply the additive law of probability to A_2 and A_3 :

$$P[A_2 \cup A_3] = P[A_2] + P[A_3] - P[A_2 \cap A_3].$$

We substitute that in (A) which then reads

$$(B) \quad P[A_1 \cup A_2 \cup A_3] = P[A_1] + P[A_2] + P[A_3] - P[A_2 \cap A_3] - P[A_1 \cap (A_2 \cup A_3)].$$

Since $A_1 \cap (A_2 \cup A_3) = (A_1 \cap A_2) \cup (A_1 \cap A_3)$, (see (2.23) on p.34: distributivity of unions and intersections), it follows from (B) that

$$(C) \quad P[A_1 \cup A_2 \cup A_3] = P[A_1] + P[A_2] + P[A_3] - P[A_2 \cap A_3] - P[(A_1 \cap A_2) \cup (A_1 \cap A_3)].$$

Finally, we apply the additive law of probability to the sets $A_1 \cap A_2$ and $A_1 \cap A_3$:

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

Here is the general formula for any number of events.

Theorem 3.7 (Exclusion–Inclusion formula). ★

If A_1, A_2, \dots, A_n are events in a probability space (Ω, P) , then

$$(3.14) \quad \begin{aligned} P(A_1 \cup A_2 \cup \dots \cup A_n) &= \sum_i P(A_i) - \sum_{i < j} P(A_i \cap A_j) \\ &+ \sum_{i < j < k} P(A_i \cap A_j \cap A_k) - \dots + (-1)^{n+1} \cdot P(A_1 \cap A_2 \cup \dots \cap A_n). \end{aligned}$$

PROOF: Will not be given here. ■

Remark 3.11. This remark is preliminary.

(A) Randomness specifically:

- (1) Random number generator of a statistics package: Generate a random a number $0 \leq x < 1$ with a precision of k decimals (can have big k like $k = 25$. For such a high precision we can model the potential outcomes Ω as the continuum $[0, 1[$).
- (2) Roll a die: $|\Omega| = 6$
- (3) Roll a die 3 times: $|\Omega| = 6^3$
- (4) 20 coin tosses: $|\Omega| = 2^{20} \approx 10^6$ since $2^{10} = 1,024 \approx 10^3$.
- (5) 10^9 coin tosses: $|\Omega| = 2^{10^9} = 2^{10 \cdot 10^8} = (2^{10})^{10^8} \approx (10^3)^{10^8} = 10^{3 \cdot 10^8}$
- (6) A selection of n items from a population is a sample of size n .

(B) A supreme being decides to pick “this” ω . This pick seems random to us since we do not know what choice this being will make. □

3.2 Conditional Probability and Independent Events

Definition 3.7 (Conditional probability).

Given are a probability space (Ω, \mathcal{F}, P) and two events $A, B \in \mathcal{F}$. We call

$$(3.15) \quad P(A | B) := \begin{cases} \frac{P(A \cap B)}{P(B)}, & \text{if } P(B) > 0, \\ \text{undefined}, & \text{if } P(B) = 0, \end{cases}$$

(read: “probability of A given B ” or “probability of A conditioned on B ”) the **conditional probability** of the event A , given that the event B has occurred. □

Theorem 3.8.

Given are a probability space $(\Omega, \mathfrak{F}, P)$ and an event $B \in \mathfrak{F}$ such that $P(B) > 0$. Then

$$(3.16) \quad P(\cdot | B) : \mathfrak{F} \longrightarrow [0, 1]; \quad A \mapsto P(A | B)$$

is another probability measure on (Ω, \mathfrak{F}) .

In other words, $P(\cdot | B)$ satisfies (3.2) – (3.4) of Definition 3.2 (Probability measures and probability spaces) on p.44.

PROOF: First, it follows from $\emptyset \subseteq A \cap B \subseteq B$ that $P(A \cap B)/P(B) \geq 0$ and $P(A \cap B)/P(B) \leq 1$.

This shows that $P(\cdot | B)$ indeed takes values between 0 and 1.

PROOF of (3.2): Since $P(\emptyset \cap B) = 0$, $P(\emptyset | B) = 0/P(B) = 0$.

PROOF of (3.3): Since $\Omega \cap B = B$, $P(\Omega | B) = P(\Omega \cap B)/P(B) = P(B)/P(B) = 1$.

PROOF of (3.4): Assume that $(A_n)_{n \in \mathbb{N}} \in \mathfrak{F}$ is a sequence of disjoint events. Then, for $i \neq j$,

$$(A_i \cap B) \cap (A_j \cap B) \subseteq A_i \cap A_j = \emptyset.$$

Thus, the sequence $(A_n \cap B)_{n \in \mathbb{N}}$ also is mutually disjoint. Further, by (2.23) on p.34,

$$\bigsqcup_{n \in \mathbb{N}} (B \cap A_n) = B \cap \bigsqcup_{n \in \mathbb{N}} A_n.$$

It follows from this and the σ -additivity of P that

$$\begin{aligned} P\left(\bigsqcup_{n \in \mathbb{N}} A_n | B\right) &= \frac{P(B \cap \bigsqcup_{n \in \mathbb{N}} A_n)}{P(B)} = \frac{P(\bigsqcup_{n \in \mathbb{N}} (B \cap A_n))}{P(B)} \\ &= \frac{\sum_{n \in \mathbb{N}} P(B \cap A_n)}{P(B)} = \sum_{n \in \mathbb{N}} \frac{P(B \cap A_n)}{P(B)} = \sum_{n \in \mathbb{N}} P(A_n | B). \end{aligned}$$

We have shown that $P(\cdot | B)$ is σ -additive and this proves (3.4). ■

It is immediate from the definition of $P(A | B)$ that

$$P(A \cap B) = P(A | B) \cdot P(B).$$

This formula is referred to by WMS as the **multiplicative law of probability**. It can be extended to three events as follows.

Proposition 3.3. *If $(\Omega, \mathfrak{F}, P)$ is a probability space and $A, B, C \in \mathfrak{F}$, then*

$$(3.17) \quad P(A \cap B \cap C) = P(A | B \cap C) \cdot P(B | C) \cdot P(C).$$

PROOF:

$$P(A \cap B \cap C) = P(A | B \cap C) \cdot P(B \cap C) = P(A | B \cap C) \cdot P(B | C) \cdot P(C). \blacksquare$$

The multiplicative law of probability generalizes to arbitrarily many sets as follows.

Proposition 3.4 (Multiplicative Law of Probability for n events).

If (Ω, \mathcal{F}, P) is a probability space, $n \in \mathbb{N}$ and $A_1, \dots, A_n \in \mathcal{F}$, then

$$(3.18) \quad P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1 | A_2 \cap \dots \cap A_n) \cdot P(A_2 | A_3 \cap \dots \cap A_n) \cdots \\ \cdots P(A_{n-2} | A_{n-1} \cap A_n) P(A_{n-1} | A_n) P(A_n).$$

PROOF:

It is easier to work with the reverse sequence $A_n \cap A_{n-1} \cap \dots \cap A_1$ instead of $A_1 \cap A_2 \cap \dots \cap A_n$. Repeated use of $P(U \cap V) = P(U | V)P(V)$ with $U = A_j$ and $V = A_{j-1} \cap \dots \cap A_1$ yields

$$\begin{aligned} &P(A_n \cap A_{n-1} \cap \dots \cap A_1) \\ &= P(A_n | A_{n-1} \cap \dots \cap A_1) P(A_{n-1} \cap \dots \cap A_1) \\ &= P(A_n | A_{n-1} \cap \dots \cap A_1) P(A_{n-1} | A_{n-2} \cap \dots \cap A_1) P(A_{n-2} \cap \dots \cap A_1) \\ &= \dots \\ &= P(A_n | A_{n-1} \cap \dots \cap A_1) P(A_{n-1} | A_{n-2} \cap \dots \cap A_1) \cdots P(A_3 | A_2 \cap A_1) P(A_2 | A_1) P(A_1). \blacksquare \end{aligned}$$

Definition 3.8 (Two independent events).

Given are a probability space (Ω, \mathcal{F}, P) and two events $A, B \in \mathcal{F}$. We say that A and B are **independent** if

$$(3.19) \quad P(A \cap B) = P(A) \cdot P(B). \quad \square$$

Independence of three events is not defined as you may have guessed from that last definition.

Definition 3.9 (Three independent events). Given are a probability space (Ω, \mathcal{F}, P) and three events $A, B, C \in \mathcal{F}$. We say that A, B and C are **independent** if

$$(3.20) \quad \begin{aligned} P(A \cap B \cap C) &= P(A) \cdot P(B) \cdot P(C), \\ P(A \cap B) &= P(A) \cdot P(B), \\ P(A \cap C) &= P(A) \cdot P(C), \\ P(B \cap C) &= P(B) \cdot P(C). \quad \square \end{aligned}$$

We can state (3.20) as follows. It must be true for any subsequence of events that the probability of the intersection equals the product of the probabilities of the individual events.

Remark 3.12. It is possible to construct a probability measure P and events A, B, C such that $P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$ and $P(A \cap B) \neq P(A) \cdot P(B)$ \square

Definition 3.9 shows us how to generalize independence to any number of events.

Definition 3.10 (Finitely many independent events).

Given are a probability space (Ω, \mathcal{F}, P) , $n \in \mathbb{N}$ and events $A_1, A_2, \dots, A_n \in \mathcal{F}$. We say that A_1, A_2, \dots, A_n are **independent** if, for ANY subselection of indices

$$1 \leq j_1 < j_2 < \dots < j_k \leq n,$$

it is true that

$$(3.21) \quad P(A_{j_1} \cap A_{j_2} \cap \dots \cap A_{j_k}) = P(A_{j_1}) \cdot P(A_{j_2}) \cdot \dots \cdot P(A_{j_k}). \quad \square$$

Finally, we define independence for infinitely many events.

Definition 3.11 (Independent events – the general case).

Given are a probability space (Ω, \mathcal{F}, P) and a sequence of events $A_1, A_2, \dots \in \mathcal{F}$. We say that this sequence is **independent** if, for ANY FINITE subselection of distinct indices $j_1, j_2, \dots, j_k \in \mathbb{N}$, it is true that

$$(3.22) \quad P(A_{j_1} \cap A_{j_2} \cap \dots \cap A_{j_k}) = P(A_{j_1}) \cdot P(A_{j_2}) \cdot \dots \cdot P(A_{j_k}). \quad \square$$

Remark 3.13. Note that the number k in Definition 3.10 and Definition 3.11 is not fixed. \square

We did not really define independence for any collection of infinitely many events, only for a sequence, i.e., a countable collection of events. The truly general case deals with families (see Definition 2.22 on p.32) of events

Definition 3.12 (Independence of uncountably many events). ★

Given are a probability space (Ω, \mathcal{F}, P) and a family $(A_i)_{i \in I}$ of events $A_i \in \mathcal{F}$. Here I denotes an arbitrary set of indices. We say that this family is **independent** if, for ANY FINITE subselection of distinct indices $i_1, i_2, \dots, i_k \in I$, it is true that

$$(3.23) \quad P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \cdot P(A_{i_2}) \cdot \dots \cdot P(A_{i_k}). \quad \square$$

The next theorem is marked optional, but it is just as easy to remember as the corollary that follows it.

Theorem 3.9. ★

Given are a probability space (Ω, \mathcal{F}, P) and a family $(A_i)_{i \in I}$ of independent events $A_i \in \mathcal{F}$. Here I denotes an arbitrary set of indices. Then we have the following:

If some or all of the A_i are replaced by their complement A_i^c , then the resulting family of events also is independent.

In other words, for each $i \in I$, let B_i be either A_i or A_i^c . Then independence of $(A_i)_{i \in I}$ implies that of $(B_i)_{i \in I}$.

PROOF: Utilizes advanced probabilistic methods that are outside the scope of this course ■

Note that the following corollary is NOT marked as optional!

Corollary 3.2.

Given are a $(\Omega, \mathfrak{F}, P)$ is a probability space, $n \in \mathbb{N}$ and independent events $A_1, \dots, A_n \in \mathfrak{F}$. If some or all of the A_i are replaced by their complement A_i^c , then the resulting list of events also is independent.

In other words, for each $i = 1, 2, \dots, n$, let B_i be either A_i or A_i^c . Then independence of A_1, \dots, A_n implies that of B_1, \dots, B_n .

PROOF: ★

(A): The case $n = 2$ shows the essence of the proof: For convenience, let $B := A_2^c$. First, we show that A_1 and B are independent.

$$\begin{aligned} A_1 &= (A_1 \cap A_2) \dot{\cup} (A_1 \cap B) \Rightarrow P(A_1) = P(A_1 \cap A_2) + P(A_1 \cap B) \\ &= P(A_1) \cdot P(A_2) + P(A_1 \cap B) \\ \Rightarrow P(A_1 \cap B) &= P(A_1) \cdot (1 - P(A_2)) = P(A_1) \cdot P(B). \end{aligned}$$

Thus, A_1 and A_2^c are independent. Since intersection is commutative ($E \cap E' = E' \cap E$), it follows that A_1^c and A_2 also are independent.

Knowing that A_1^c and A_2 are independent, we can apply the proof above to those two independent events and obtain that A_1^c and A_2^c are independent. This finishes the proof for $n = 2$

(B): For general n , let A_1, \dots, A_n be independent. For convenience, let $B := A_1 \cap \dots \cap A_{n-1}$.

Since $P(B \cap A_n) = P(A_1 \cap \dots \cap A_n) = P(A_1) \dots P(A_n) = P(B) \cdot P(A_n)$, B and A_n are independent. We have shown in **(A)** that B and A_n^c are independent, too.

We argue as in **(A)** and conclude from the commutativity of “ \cap ” that replacing any A_j with its complement, i.e., fixing an index j_1 and defining $B_j := A_j$ for $j \neq j_1$ and $B_{j_1} := A_{j_1}^c$, that B_1, \dots, B_n are independent. In other words, replacing just one event with its complement maintains independence.

We apply this to the events $C_j := B_j$ for $j \neq j_2$ and $C_{j_2} := B_{j_2}^c$, where we assume that $j_2 \neq j_1$. The result is that C_1, \dots, C_n also are independent.

At this point we know that replacing $k = 1$ or $k = 2$ events with their complements maintains independence. We apply this to the events $D_j := C_j$ for $j \neq j_3$ and $D_{j_3} := C_{j_3}^c$, where we assume that $j_3 \notin \{j_1, j_2\}$. The result is that D_1, \dots, D_n also are independent.

At this point we know that replacing $k \leq 3$ events with their complements maintains independence. We repeat the above with $k = 4$, then with $k = 5, \dots$, then with $k = n$. This completes the proof. ■

Next, we examine connections between conditional probabilities and independence.

Theorem 3.10.

Given are a probability space (Ω, \mathcal{F}, P) and two events $A, B \in \mathcal{F}$ such that $P(B) > 0$. Then

$$(3.24) \quad A \text{ and } B \text{ are independent} \Leftrightarrow P(A | B) = P(A).$$

PROOF of “ \Rightarrow ”:

Since A and B are independent and $P(B) > 0$,

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A) \cdot P(B)}{P(B)} = P(A).$$

PROOF of “ \Leftarrow ”:

Since $P(A | B) = P(A)$ and $P(B) > 0$,

$$P(A) \cdot P(B) = P(A | B) \cdot P(B) = \frac{P(A \cap B)}{P(B)} \cdot P(B) = P(A \cap B). \blacksquare$$

Corollary 3.3.

If (Ω, \mathcal{F}, P) is a probability space and $A, B \in \mathcal{F}$ such that $P(A) > 0$ and $P(B) > 0$. Then

$$(3.25) \quad A \text{ and } B \text{ are independent} \Leftrightarrow P(A | B) = P(A) \Leftrightarrow P(B | A) = P(B).$$

PROOF: Obvious \blacksquare

3.3 Preimages and Indicator Functions

Introduction 3.1. The major part of this course will be about functions $\omega \mapsto f(\omega)$ which assign the outcomes (= elements) ω of a probability space to items $f(\omega)$ which are usually numbers or vectors of numbers. In other words, the codomain will usually be (a subset of) \mathbb{R} or \mathbb{R}^n . We illustrate this with the following example.

Let the probability space (Ω, P) ²³ represent the outcomes of two rolls of a fair die:

- $\Omega = \{1, 2, \dots, 6\}^2$. Interpret $\omega = (\omega_1, \omega_2)$ as die₁ yields ω_1 , die₂ yields ω_2 .²⁴
 \square Thus, $\omega = (5, 2)$ represents the outcome of die₁ giving a 5 and die₂ giving a 2.
- Probability measure P is determined by $P(\omega_1, \omega_2) = 1/|\Omega| = 1/36$. See Example 3.1 on p.50.

Consider the function which associates with each outcome (ω_1, ω_2) the sum of the throws, i.e.,

- $Y : \Omega \rightarrow \{2, 3, 4, \dots, 11, 12\}$; $(\omega_1, \omega_2) \mapsto Y((\omega_1, \omega_2)) := \omega_1 + \omega_2$.

Get used to the notation! WMS loves to use the letters (X, Y, Z) for function names.

We will create a probability measure P' on $\Omega' := \{2, 3, 4, \dots, 11, 12\}$, the codomain of the function Y .

²³As promised, no more σ -algebra unless absolutely necessary!

²⁴We often prefer to write ω rather than $\vec{\omega}$ if the the symbol Ω is involved, even if the elements are vectors.

- Since Ω' is countable, it suffices to specify $P'(\{2\}), P'(\{3\}), \dots, P'(\{12\})$. (Again, Example 3.1.)
- Define $P'(\{10\}) := P(\{(\omega_1, \omega_2) \in \Omega : Y(\omega_1, \omega_2) = 10\}) = P(\{(4, 6), (5, 5), (6, 4)\}) = 1/12$. This is the probability that the sum of the throws is 10!
- In general, for $\omega' \in \Omega'$, define $P'(\{\omega'\}) := P(\{(\omega_1, \omega_2) \in \Omega : Y(\omega_1, \omega_2) = \omega'\})$. This is the probability that the sum of the throws is ω' !
- One can show quite easily²⁵ that, if $B \subseteq \Omega'$, then

$$(3.26) \quad B \subseteq \Omega' \Rightarrow P'(B) = P(\{\omega \in \Omega : Y(\omega) \in B\}). \quad (\text{We wrote } \omega \text{ for } (\omega_1, \omega_2).)$$

This is the probability that the sum of the throws is in B !

- We have created a probability measure $P'(B)$ on the codomain of Y by assigning P , the original probability on the domain Ω , to the set

$$\{\omega \in \Omega : Y(\omega) \in B\}$$

of all those arguments $\omega \in \Omega$ which are mapped by Y into B .

That makes those sets so important that they warrant their own definition. \square

Since the following definition is of interest not only for probabilistic topics, we will switch from the function notation $Y : \Omega \rightarrow \Omega'$ to the more familiar $f : X \rightarrow Y$.

Definition 3.13.

Let X, Y be two nonempty sets. Let $f : X \rightarrow Y$ and $B \subseteq Y$. Then

$$(3.27) \quad f^{-1}(B) := \{x \in X : f(x) \in B\}$$

is a subset of X which we call the **preimage** of B under f . \square

Remark 3.14. ★

- If we vary $B \subseteq Y$, i.e., $B \in 2^Y$, we can think of the preimage as a function $2^Y \rightarrow 2^X$ (since $f^{-1}(B) \in 2^X$).
- The symbol f^{-1} is the same as that for the ordinary inverse function $f^{-1}(y) = x$, **if this inverse function exists!**
- $f^{-1}(B)$ exists for any choice of $X, Y, f : X \rightarrow Y$, and $B \subseteq Y$, even if the inverse function does not exist!

As an example, let

$$f : \mathbb{R} \rightarrow [-1, \infty[; \quad f(x) = x^2.$$

If there was an inverse function, then it would have to assign to EACH $y \in [-1, \infty[$ a UNIQUE $x \in \mathbb{R}$ (that x would be $f^{-1}(y)$) such that $f(x) = y$. But such is not the case:

²⁵with the help of Proposition 3.6 (f^{-1} is compatible with all basic set ops) further down, on p.63

- If $y = -0.5$, then there is no $x \in \mathbb{R}$ such that $x^2 = y$
- If $y = 10$, then there are too many $x \in \mathbb{R}$ such that $x^2 = y$.
Both $x = \sqrt{10}$ and $x = -\sqrt{10}$ satisfy $x^2 = 10$.
- Note that, for the preimages we obtain $f^{-1}(\{-0.5\}) = \emptyset$
and $f^{-1}(\{10\}) = \{-\sqrt{10}, \sqrt{10}\}$. Coincidence?

For a more extreme example, consider

$$g : [0, \infty[\rightarrow \mathbb{R}; \quad g(x) = \sin(x).$$

If $B_1 = [5, 10]$, $B_2 = \{0\}$, what are $g^{-1}(B_1)$ and $g^{-1}(B_2)$? So, does each $y \in \mathbb{R}$ have a unique $x \in [0, \infty[$ such that $g(x) = y$?

For an even more extreme example, consider

$$h : \mathbb{R} \rightarrow \mathbb{R}; \quad h(x) = 2\pi.$$

If $B_1 = [5, 10]$, $B_2 = \{2\pi\}$, $B_3 = [-500, 5]$, what are $h^{-1}(B_j)$ ($j = 1, 2, 3$)? Again, does each $y \in \mathbb{R}$ have a unique $x \in [0, \infty[$ such that $h(x) = y$? \square

Notational conveniences I:

If we have a set that is written as $\{\dots\}$ then we may write $f^{-1}\{\dots\}$ instead of $f^{-1}(\{\dots\})$. Specifically for singletons $\{y\}$ such that $y \in Y$, it is OK to write $f^{-1}\{y\}$. You also are allowed to write $f^{-1}(y)$ instead of $f^{-1}\{y\}$, even though this author thinks that it is not a good idea to confound elements y and subsets $\{y\}$ of Y .

VERY IMPORTANT: Work the following examples closed book and then check that your solutions are correct!

Example 3.4 (Preimages). Let $f : \mathbb{R} \rightarrow \mathbb{R}; \quad f(x) = x^2$. Determine

- a. $f^{-1}(]-4, -2])$, b. $f^{-1}([1, 2])$, c. $f^{-1}([5, 6])$, d. $\{-4 < f < -2 \text{ or } 1 \leq f \leq 2 \text{ or } 5 \leq f < 6\}$.

Solution:

- a. $f^{-1}(]-4, -2]) = \{x \in \mathbb{R} : x^2 \in]-4, -2]\} = \{-4 < f < -2\} = \emptyset$.
b. $f^{-1}([1, 2]) = \{x \in \mathbb{R} : x^2 \in [1, 2]\} = \{1 \leq f \leq 2\} = [-\sqrt{2}, -1] \cup [1, \sqrt{2}]$.
c. $f^{-1}([5, 6]) = \{x \in \mathbb{R} : x^2 \in [5, 6]\} = \{5 \leq f \leq 6\} = [-\sqrt{6}, -\sqrt{5}] \cup [\sqrt{5}, \sqrt{6}]$.
d. $\{-4 < f < -2 \text{ or } 1 \leq f \leq 2 \text{ or } 5 \leq f < 6\} = f^{-1}(]-4, -2[\cup [1, 2] \cup [5, 6])$
 $= \{x \in \mathbb{R} : x^2 \in]-4, -2[\text{ or } x^2 \in [1, 2] \text{ or } x^2 \in [5, 6]\}$
 $= [-\sqrt{2}, -1] \cup [1, \sqrt{2}] \cup [-\sqrt{6}, -\sqrt{5}] \cup [\sqrt{5}, \sqrt{6}]$. \square

Example 3.5 (Preimages). Let $f : \mathbb{R} \rightarrow \mathbb{R}; \quad f(x) = x^2$. Determine

- a. $f^{-1}(]-4, 2])$, b. $f^{-1}([1, 3])$, c. $\{-4 < f < -2 \text{ and } 1 \leq f \leq 3\}$.

Solution:

- a. $f^{-1}(]-4, 2]) = \{x \in \mathbb{R} : x^2 \in]-4, 2]\} = \{x \in \mathbb{R} : -4 < x^2 < 2\} =]-2, 2[$.
b. $f^{-1}([1, 3]) = \{x \in \mathbb{R} : x^2 \in [1, 3]\} = \{x \in \mathbb{R} : 1 \leq x^2 \leq 3\} = [-\sqrt{3}, -1] \cup [1, \sqrt{3}]$.
c. $\{-4 < f < -2 \text{ and } 1 \leq f \leq 3\} = f^{-1}(]-4, -2[\cap [1, 3])$
 $= \{x \in \mathbb{R} : x^2 \in]-4, -2[\text{ and } x^2 \in [1, 3]\}$
 $= \{x \in \mathbb{R} : 1 \leq x^2 < 2\} =]-\sqrt{2}, -1] \cup [1, \sqrt{2}[$. \square

Proposition 3.5. *Some simple properties:*

$$(3.28) \quad f^{-1}(\emptyset) = \emptyset$$

$$(3.29) \quad B_1 \subseteq B_2 \subseteq Y \Rightarrow f^{-1}(B_1) \subseteq f^{-1}(B_2) \quad (\text{monotonicity of } f^{-1}\{\dots\})$$

$$(3.30) \quad f^{-1}(Y) = X \quad \text{always!}$$

PROOF of 3.29:

We show that $x \in f^{-1}(B_1) \Rightarrow f^{-1}(B_1) \subseteq f^{-1}(B_2)$ as follows.

$$x \in f^{-1}(B_1) \stackrel{(a)}{\Rightarrow} f(x) \in B_1 \stackrel{(b)}{\Rightarrow} f(x) \in B_2 \stackrel{(c)}{\Rightarrow} x \in f^{-1}(B_2)$$

In the above, (a) and (c) state the definition of a preimage and (b) follows from $B_1 \subseteq B_2$

The proof of 3.28 and 3.29 is left as an exercise. ■

Remark 3.15 (Notational conveniences II:).

In probability theory the following notation is also very common:

$$\{f \in B\} := f^{-1}(B), \quad \{f = y\} := f^{-1}\{y\}.$$

Let \mathcal{R} be either of $\mathbb{Z}, \mathbb{Q}, \mathbb{R}$. Assume that the codomain of f is considered a subset of \mathcal{R} . Let $a, b \in \mathcal{R}$ such that $a < b$. We write $\{a \leq f \leq b\} := f^{-1}([a, b]_{\mathcal{R}})$, $\{a < f < b\} := f^{-1}(]a, b[_{\mathcal{R}})$, $\{a \leq f < b\} := f^{-1}([a, b[_{\mathcal{R}})$, $\{a < f \leq b\} := f^{-1}(]a, b]_{\mathcal{R}})$, $\{f \leq b\} := f^{-1}(]-\infty, b]_{\mathcal{R}})$, etc. □

Example 3.6. In the introduction we were examining

- $P'(\{10\}) = P(\{(\omega_1, \omega_2) \in \Omega : Y(\omega_1, \omega_2) = 10\})$.
This can be written as $P'(\{10\}) = P(Y^{-1}\{10\}) = P\{Y = 10\}$
- $P'(\{\omega'\}) = P(\{(\omega_1, \omega_2) \in \Omega : Y(\omega_1, \omega_2) = \omega'\})$.
This can be written as $P'(\{\omega'\}) = P(Y^{-1}\{\omega'\}) = P\{Y = \omega'\}$.
- $P'(B) = P(\{\omega \in \Omega : Y(\omega) \in B\})$.
This can be written as $P'(B) = P(Y^{-1}(B)) = P\{Y \in B\}$.

It is very important that you **remember the first three** of the five formulas of the next proposition.

Proposition 3.6 (f^{-1} is compatible with all basic set ops). *Assume that X, Y be nonempty, $f : X \rightarrow Y$, J is an arbitrary index set.²⁶ Further assume that $B \subseteq Y$ and that $B_j \subseteq Y$ for all j . Then*

²⁶If you have problems with the concept of a family, think of J as a set of integers which are bounded below, i.e., that J is the index set of a finite or infinite sequence or subsequence of sets

$$(3.31) \quad f^{-1}\left(\bigcap_{j \in J} B_j\right) = \bigcap_{j \in J} f^{-1}(B_j)$$

$$(3.32) \quad f^{-1}\left(\bigcup_{j \in J} B_j\right) = \bigcup_{j \in J} f^{-1}(B_j)$$

$$(3.33) \quad f^{-1}(B^c) = (f^{-1}(B))^c$$

$$(3.34) \quad B_1 \cap B_2 = \emptyset \Rightarrow f^{-1}(B_1) \cap f^{-1}(B_2) = \emptyset.$$

$$(3.35) \quad f^{-1}(B_1 \setminus B_2) = f^{-1}(B_1) \setminus f^{-1}(B_2)$$

$$(3.36) \quad f^{-1}(B_1 \Delta B_2) = f^{-1}(B_1) \Delta f^{-1}(B_2)$$

Note that (3.34) implies that the preimages of a disjoint family form a disjoint family.

PROOF:  MF330 notes, ch.8 ■

Proposition 3.7 (Preimages of function composition). Let X, Y, Z be arbitrary, nonempty sets. Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ and $h : X \rightarrow Z$ the composition

$$h(x) = g \circ f(x) = g(f(x)).$$

Let $U \subseteq X$ and $W \subseteq Z$. Then

$$(3.37) \quad (g \circ f)^{-1} = f^{-1} \circ g^{-1}, \text{ i.e., } (g \circ f)^{-1}(W) = f^{-1}(g^{-1}(W)) \text{ for all } W \subseteq Z.$$

PROOF:  MF330 notes, ch.8 ■

Try to understand the above with a simple example, such as $X = Y = \mathbb{R}$,

$f(x) = 3x - 1, g(y) = y^2$, and $W = [0, 1], W = \{-10\}, W = \{10\}$ (three different choices for W).

Indicator functions often are a great notational convenience, for example, when dealing with functions that are defined differently in two or more parts of the domain.

Definition 3.14 (indicator function for a set). Let Ω be a nonempty set and $A \subseteq \Omega$. Let $1_A : \Omega \rightarrow \{0, 1\}$ be the function defined as

$$(3.38) \quad 1_A(\omega) := \begin{cases} 1 & \text{if } \omega \in A, \\ 0 & \text{if } \omega \notin A. \end{cases}$$

1_A is called the **indicator function** of the set A .²⁷ □

²⁷In abstract algebra this is often called the **characteristic function** of A . Some authors write χ_A or $\mathbb{1}_A$ instead of 1_A .

Example 3.7. The so-called density function for the exponential distribution with parameter $\beta > 0$ is

$$f(y) = \begin{cases} \frac{1}{\beta} e^{-y/\beta}, & 0 \leq y < \infty, \\ 0, & \text{elsewhere.} \end{cases}$$

This can also be written as $f(y) = \frac{1}{\beta} e^{-y/\beta} 1_{[0, \infty[}(y)$.

Proposition 3.8. Let A, B, C be subsets of Ω . Then

$$(3.39) \quad A \subseteq B \Rightarrow 1_A \leq 1_B,$$

$$(3.40) \quad 1_{A \cup B} = \max(1_A, 1_B),$$

$$(3.41) \quad 1_{A \cap B} = \min(1_A, 1_B),$$

$$(3.42) \quad 1_{A^c} = 1 - 1_A,$$

$$(3.43) \quad 1_{A \uplus B} = 1_A + 1_B \quad (A, B \text{ disjoint})$$

PROOF: The proof is an easy exercise.

3.4 Random Variables and their Probability Distributions

Introduction 3.2. We continue with an observation we made in the introduction 3.1 to Section 3.3 (Preimages, p.60). There,

- $\Omega = \{1, 2, \dots, 6\}^2$ and $\vec{\omega} = (\omega_1, \omega_2)$ represents a potential (two-number) outcome of two rolls of a fair die, i.e., $P(\{\vec{\omega}\}) = 1/|\Omega| = 1/36$.
- We defined the function $Y : \Omega \rightarrow \Omega' := \{2, 3, 4, \dots, 11, 12\}$; $\vec{\omega} \mapsto Y(\vec{\omega}) := \omega_1 + \omega_2$, which associates with $\vec{\omega} = (\omega_1, \omega_2)$ the sum of the two rolls.
- This function lead to a probability measure P' on Ω' by means of formula (3.26):

$$B \subseteq \Omega' \Rightarrow P'(B) = P\{\vec{\omega} \in \Omega : Y(\vec{\omega}) \in B\}.$$

Observe that the set Ω' has been transformed into a probability space, (Ω', P') .

- With preimage notation and the notational shortcuts of Remark 3.15 on p.63, this can also be written as

$$P'(B) = P(Y^{-1}(B)) = P\{Y \in B\}.$$

These formulas can be written for an arbitrary probability space (Ω, P) , an arbitrary nonempty set Ω' , and an arbitrary function $Y : \Omega \rightarrow \Omega'$. Actually, that is not entirely true, but it will be true for the situations we will deal with in this class. If you are curious, read this optional footnote. ²⁸ \square

The next theorem and the subsequent definitions are very important.

²⁸ \star We have to recall that there really is a σ -algebra \mathfrak{F} on Ω and that $P(A)$ only exists if $A \in \mathfrak{F}$. What if $B \subseteq \Omega'$ does not have a nice preimage, i.e., $\{Y \in B\} \notin \mathfrak{F}$? The only way out is not to allow arbitrary $B \in 2^{\Omega'}$, but **(a)** to also require a σ -algebra \mathfrak{F}' on the codomain Ω' , which **(b)** is so “small” that $B \in \mathfrak{F}' \Rightarrow Y^{-1}(B) \in \mathfrak{F}$; or, if you prefer, \mathfrak{F} must be so “big” that $B \in \mathfrak{F}' \Rightarrow Y^{-1}(B) \in \mathfrak{F}$. There is a name for triplets $[Y, \mathfrak{F}, \mathfrak{F}']$ which satisfy this relationship. The function Y is called **measurable** with respect to \mathfrak{F} and \mathfrak{F}' or $(\mathfrak{F}, \mathfrak{F}')$ -**measurable**. None of this will be an issue in this course!

Theorem 3.11.

Let (Ω, P) be a probability space, Ω' a nonempty set, and $Y : \Omega \rightarrow \Omega'$ a function. Then the formula

$$(3.44) \quad P_Y(B) := P\{Y \in B\} \quad (B \subseteq \Omega')$$

defines a probability measure on Ω' .

PROOF: ★ It follows from $\{Y \in \emptyset\} = \emptyset$ and $\{Y \in \Omega'\} = \Omega$, that

$$P_Y(\emptyset) = P(\emptyset) = 0 \quad \text{and} \quad P_Y(\Omega') = P(\Omega) = 1.$$

Let $B \subseteq \Omega'$. From (3.33) on p.64, we obtain

$$P_Y(B^c) = P\{Y \in B^c\} = P(Y^{-1}(B^c)) = P([Y^{-1}(B)]^c) = 1 - P(Y^{-1}(B)) = 1 - P_Y(B).$$

To prove σ -additivity of P_Y , we apply (3.32) to the index set \mathbb{N} of a sequence of disjoint subsets B_1, B_2, \dots of Ω' . Let $B := B_1 \uplus B_2 \uplus B_3 \uplus \dots$. Then

$$P_Y(B) = P(Y^{-1} \left(\biguplus_{j \in \mathbb{N}} B_j \right)) = P \left(\bigcup_{j \in \mathbb{N}} Y^{-1}(B_j) \right)$$

By (3.34), the sets $Y^{-1}(B_j)$ are disjoint. Thus,

$$P_Y(B) = P \left(\biguplus_{j \in \mathbb{N}} Y^{-1}(B_j) \right) = \sum_{j \in \mathbb{N}} P(Y^{-1}(B_j)) = \sum_{j \in \mathbb{N}} P_Y(B_j).$$

This proves σ -additivity. ■

Definition 3.15 (Probability Distribution).

Let (Ω, P) be a probability space, Ω' a nonempty set, and $Y : \Omega \rightarrow \Omega'$ a function. Then the probability measure P_Y on Ω' which is given by

$$(3.45) \quad P_Y(B) := P\{Y \in B\} \quad (B \subseteq \Omega')$$

is called the **probability distribution** or just the **distribution** of Y with respect to P . Very often the probability space (Ω, P) is fixed for a long stretch. We then simply talk about the probability distribution of Y , without referring to P . □

Definition 3.16 (Random Variables and Random Vectors). Let (Ω, P) be a probability space and let $n \in \mathbb{N}$.

Let $B \subseteq \mathbb{R}$. A function

$$Y : \Omega \longrightarrow B; \quad \omega \mapsto Y(\omega)$$

is called a **random variable** (in short, **r.v.** or **rv**), on $(\Omega, \mathfrak{F}, P)$. Let $B' \subseteq \mathbb{R}^n$. A function

$$\vec{X} = (X_1, X_2, \dots, X_n) : \Omega \longrightarrow B'; \quad \omega \mapsto \vec{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$$

is called a **random vector** on $(\Omega, \mathfrak{F}, P)$.
 If there is a countable subset $B^* = \{y_1, y_2, \dots\}$ of B such that $\sum_j P_Y\{y_j\} = 1$ (i.e., $P\{Y \notin B^*\} = 0$), we call Y a **discrete random variable**. Likewise, if there is a countable subset B'^* of B' such that $P\{\vec{X} \notin B'^*\} = 0$, we call \vec{X} a **discrete random vector**. \square

Note that random variables and vectors which have a countable range are discrete. Also, if you found the footnote at the end of the introduction interesting, have a look at this (optional) one,²⁹

Remark 3.16. In many instances the exact nature of the codomain B of a random variable Y is unimportant. Of course it must be a set of numbers, i.e., $B \subseteq \mathbb{R}$, and it must be big enough to accommodate all function values $Y(\omega)$, i.e., $Y(\omega) \subseteq B$.³⁰ Thus, here is some **good news**.

We often will just say something like “Let Y be a random variable on Ω ” or “Let Y be a discrete random vector on Ω ” and not even mention the codomain of Y . \square

Not all interesting functions on a probability space take values in \mathbb{R} or \mathbb{R}^n . Here is an example.

Example 3.8. The following describes a (unnecessarily complicated) way to simulate n tosses of a fair coin. Let $\Omega := [0, 1[$, where we represent the real number $\omega \in \Omega$ as a decimal $0.d_1d_2d_3$ with infinitely many decimal digits. If necessary, we append infinitely many zeroes to the right. For example, we write $0.25000\dots$ for the number $1/4$. We write H for Heads and T for Tails and define the following function on (Ω, P) .

$$\vec{X} : \Omega \rightarrow \{H, T\}^n$$

- $X_1(\omega) = H$ if d_1 is even, T else.
- $X_2(\omega) = H$ if d_2 is even, T else.
-
- $X_n(\omega) = H$ if d_n is even, T else.

Since $P_{\vec{X}}(\vec{x}) = 1/2^n$ for each $\vec{x} \in \{H, T\}^n$, each combination of a total of n Heads and Tails has the same chance to occur. That is our understanding of a fair coin. \square

Considering that last example, it seems awkward [not] to call a function $\Omega \rightarrow \Omega'$ from a probability space (Ω, P) to a set Ω' a random variable only because its function values are not numbers. We give a name to such functions of randomness.

²⁹ ★ Technically speaking, Y must be $(\mathfrak{F}, \mathfrak{B})$ -measurable and \vec{X} must be $(\mathfrak{F}, \mathfrak{B}^n)$ -measurable. In other words, one must be able to assign probabilities to all preimages of Borel sets. Again, none of this will be an issue in this course!

³⁰It only matters when we need the inverse function $\omega = Y^{-1}(y)$ of $y = Y(\omega)$. (Do not confuse inverse function and preimage, just because they use the same symbol Y^{-1} !) Then $Y^{-1}(y)$ must make sense for all $y \in B$ and that requires that B is minimal: $B = Y(\Omega)$. The same thought also applies to random vectors.

Definition 3.17 (Random element). Let $(\Omega, \mathfrak{F}, P)$ be a probability space and Ω' a nonempty set. We call a function $X : \Omega \rightarrow \Omega'$ a **random element**, also: a **random item**, on Ω . \square

Remark 3.17. We can phrase Theorem 3.11 and the subsequent Definition 3.15 as follows.

All random elements X on a probability space $(\Omega, \mathfrak{F}, P)$ have a distribution

$$P_X(B) = P\{X \in B\} = P(X^{-1}(B)) \quad (B \subseteq \Omega'). \quad \square$$

Since an element x of the domain of a function f (an argument) is assigned to only one function value $y = f(x)$, one should expect that a function of a discrete random element should again be discrete. This is the assertion of the next proposition and the corollary that follows it.

Proposition 3.9. ★ Let $X : (\Omega, P) \rightarrow \Omega'$ be a random element and $g : \Omega' \rightarrow \mathbb{R}$. Let Z be the random variable $\omega \mapsto Z(\omega) := g(X(\omega))$. Let $B^* \in \Omega'$ such that $P_X(B^*) = 1$ and let $C^* := \{g(x) : x \in B^*\}$. Then $P_Z(C^*) = 1$.

PROOF: Let

$$(A) \quad A_1 := \{\omega \in \Omega : Z(\omega) \notin C^*\} = \{\omega \in \Omega : g(X(\omega)) \notin C^*\}.$$

$$(B) \quad \text{Then } \tilde{\omega} \in X^{-1}(B^*) \Leftrightarrow X(\tilde{\omega}) \in B^* \Rightarrow Z(\tilde{\omega}) = g(X(\tilde{\omega})) \in g(B^*) = C^*$$

Here, “ \Leftrightarrow ” follows from the definition of X^{-1} . From (A) + (B) we see that $A_1 \cap X^{-1}(B^*) = \emptyset$.

$$(C) \quad \text{Thus, } A_1 \subseteq [X^{-1}(B^*)]^c.$$

$$(D) \quad \text{Since } P[X^{-1}(B^*)] = P_X(B^*) = 1 \quad (\text{by definition of } B^*),$$

we obtain from (C) that $P(A_1) = 0$ and then, from (A), that

$$(E) \quad P_Z(C^*) = P\{\omega \in \Omega : Z(\omega) \in C^*\} = P(A_1^c) = 1. \quad \blacksquare$$

Corollary 3.4. Let $X : (\Omega, P) \rightarrow \Omega'$ be a random element and $g : \Omega' \rightarrow \mathbb{R}$. Further, let Z be the random variable $g \circ X : \omega \mapsto Z(\omega) = g(X(\omega))$. In other words, Z is the composition of g with X . Then

- (a) If $\omega \mapsto X(\omega)$ only assumes finitely many (distinct) values x_1, \dots, x_n , then $\omega \mapsto Z(\omega)$ only assumes finitely many values z_1, \dots, z_m (and $m \leq n$).
- (b) If $\omega \mapsto X(\omega)$ only assumes an infinite sequence of (distinct) values (x_j) , then $\omega \mapsto Z(\omega)$ assumes a countable set of function values. (This set forms a finite or infinite sequence. (See Definition 2.21 (Countable and uncountable sets) on p.31).
- (c) If X is a discrete random element, then $Z = g(X)$ is a discrete random variable.

PROOF of (a): ★ The potential function values of Z are

$$z'_1 := g(x_1), z'_2 := g(x_2), \dots, z'_n := g(x_n)$$

If g is not injective, there may be duplicate z'_j which must be removed. Thus, Z assumes at m distinct values for some suitable $m \leq n$. We rename them z_1, \dots, z_m .

PROOF of (b): ★ The potential function values of Z the members of the sequence $z'_j = g(x_j)$, where $j \in \mathbb{N}$. Removing the duplicates leaves us with a finite or infinite subsequence of distinct items z_j and those form the countable set of all function values of Z .

PROOF of (c): ★ Since X is discrete, there is a countable set $B^* \subseteq \Omega'$ such that $P_X(B^*) = 1$.

We have seen in the proof of (b) that a function g transports countably many arguments b^* into countably many function values $c^* = g(b^*)$. Thus, the set $C^* := \{g(b^*) : b^* \in B^*\}$ is countable.

It follows from Proposition 3.9 on p.68 that $P_Z(C^*) = 1$. Since C^* is countable, Z is discrete. ■

Remark 3.18. Consider the following of a philosophical rather than mathematical nature. Not all mathematicians agree with it.

I like to think of a probability space (Ω, P) as a seat of randomness in the following sense. Some all-powerful supreme being or supreme force of nature, let's call it \mathcal{SB} , decides to pick "this" particular $\omega_0 \in \Omega$. As a result, all random elements X, Y, Z, \dots are invoked with ω_0 as argument, resulting in the outcomes $X(\omega_0), Y(\omega_0), Z(\omega_0), \dots$. With this interpretation it makes a lot of sense to talk about functions on (Ω, P) as **random** elements since, when we interpret $\omega \in \Omega$ as "randomness",

$$x = X(\omega) \text{ simply means that } x \text{ is a function of randomness.}$$

Only \mathcal{SB} knows what ω_0 will be picked. But if we know, say, the distribution P_X of a random variable X , then we can at least quantify the likelihood that \mathcal{SB} chose an ω such that $17.8 \leq X(\omega) \leq 21.3$. It will be $P_X([17.8, 21.3]) = P\{17.8 \leq X \leq 21.3\}$. □

Often it only is the distribution of a random element with values in a set Ω' that matters and there may be many different choices of probability space plus random element which result in that same probability measure on Ω' . We illustrate that with two more settings for the modeling of the distribution of n tosses of a fair coin on the space $\{H, T\}^n$. See Example 3.8. We fix $n = 3$ since this example illustrates all essential points.

Example 3.9. (a) Let $\Omega_1 := \{0, 1\}^3$ with the probability measure $P\{(a, b, c)\} = 1/|\Omega_1| = 1/8$.

Let $Y_1 : \Omega_1 \rightarrow \{H, T\}^3$ the random element that changes each H into a 1 and each T into a 0. For example, $Y_1(1, 0, 1) = (H, T, H)$ and $Y_1(0, 0, 1) = (T, T, H)$.

Then P_{Y_1} is the same probability measure as $P_{\bar{X}}$ of Example 3.8, since both assign the number $1/8$ to each element of $\{H, T\}^3$.

(b) Let $\Omega_2 := \{H, T\}^3$ with the probability measure $P\{(a, b, c)\} = 1/|\Omega_2| = 1/8$. (Same as in (a), except that now a, b, c represent either of H or T rather than 0 or 1.)

Let $Y_2 : \Omega_2 \rightarrow \{H, T\}^3$ be the **identity** (also, **identity function**) on Ω_2 . That is the "do nothing" function which assigns each element of a set to itself, i.e., $Y_2(\omega) = \omega$ for all $\omega \in \Omega_2$.

Clearly, P_{Y_2} also assigns probability $P_{Y_2}(\{\omega\}) = 1/8$ to each element of $\{H, T\}^3$.

(c) Let $\Omega_3 := \{H, T\}^3 \times \{1, 2, 3, 4\}$ with the probability measure $P\{(a, b, c, d)\} = 1/|\Omega_3| = 1/32$. (Same as in (a), except that now a, b, c represent either of H or T rather than 0 or 1.)

Let $Y_3 : \Omega_3 \rightarrow \{H, T\}^3$ be the function defined as $Y_3(a, b, c, d) := (a, b, c)$. We compute the distribu-

tion P_{Y_3} for the outcomes (a, b, c) of the probability space $(\{H, T\}^3, P_{Y_3})$.

$$\begin{aligned}(a, b, c) \in Y_3 &\Rightarrow P_{Y_4}\{(a, b, c, d)\} = P\{Y_4 = (a, b, c, d)\} \\ &= P\{(a, b, c, 1), (a, b, c, 2), (a, b, c, 3), (a, b, c, 4)\} = 4(1/32) = 1/8.\end{aligned}$$

We have obtained in this example and Example 3.9 the probability P' which models three tosses of a fair coin, i.e., $P'\{(a, b, c)\} = 1/8$ for each $(a, b, c) \in \{H, T\}^3$, as the distribution of four different random elements \vec{X}, Y_1, Y_2, Y_3 which were defined on four different probability spaces. Thus, you have multiple choices of probability spaces and random items to model a distribution. you will hopefully agree that Y_1 and Y_2 are much better choices than \vec{X} and Y_3 . \square

4 Combinatorial Analysis

In many important cases we find ourselves in the situation of Example 3.2 on p.46, where we have a finite probability space (Ω, P) , in which each outcome $\omega \in \Omega$ has equal probability

$$P\{\omega\} = \frac{1}{|\Omega|}$$

and thus, for each event $A \subset \Omega$,

$$P(A) = \frac{|A|}{|\Omega|}.$$

Hence, all we need to determine $P(A)$, is the knowledge of how to count the elements of Ω and of A . Combinatorial analysis, also called **combinatorics**, is a branch of mathematics that provides us with tools to accomplish that task.

4.1 The Multiplication Rule

The first result is known under names such as the basic principle of counting ([8] Ross, Sheldon M.: A First Course in Probability, 3rd edition) and the mn rule (WMS text).

Theorem 4.1 (Multiplication rule).

(A) Assume that two actions A and B are performed such that

- the first one has m outcomes, $\{a_1, a_2, \dots, a_m\}$,
- the second one has n outcomes $\{b_1, b_2, \dots, b_n\}$ for each outcome of the first one.
- Then the number of combined outcomes (a_i, b_j) is mn .

(B) Generalization. Assume that k actions A_1, \dots, A_k are performed such that

- action A_1 has n_1 outcomes, $\{a_1^{(1)}, a_2^{(1)}, \dots, a_{n_1}^{(1)}\}$,
- action A_2 has n_2 outcomes, $\{a_1^{(2)}, a_2^{(2)}, \dots, a_{n_2}^{(2)}\}$ for each outcome of A_1 ,
- action A_3 has n_3 outcomes, $\{a_1^{(3)}, a_2^{(3)}, \dots, a_{n_3}^{(3)}\}$ for each combined outcome (x_1, x_2) , where x_1 is one of the A_1 -outcomes and x_2 is one of the A_2 -outcomes,

- action A_k has n_k outcomes, $\{a_1^{(k)}, a_2^{(k)}, \dots, a_{n_k}^{(k)}\}$ for each combined outcome $(x_1, x_2, \dots, x_{k-1})$, where each x_j is one of the A_j -outcomes, i.e., x_j is one of $a_1^{(j)}, \dots, a_{n_j}^{(j)}$.
- Then there are $n_1 \cdot n_2 \cdot \dots \cdot n_k$ combined outcomes (x_1, x_2, \dots, x_k) .
Here, each x_j is one of the n_j outcomes $a_1^{(j)}, \dots, a_{n_j}^{(j)}$ of A_j .

PROOF: We identify the actions with their outcomes, i.e., we define

$$A_j = \{a_1^{(j)}, \dots, a_{n_j}^{(j)}\}, \quad \text{for } j = 1, 2, \dots, k.$$

Now, the multiplication rule merely states that $|A_1 \times A_2 \times \dots \times A_n| = |A_1| \cdot |A_2| \cdot \dots \cdot |A_n|$, and this is true according to (2.28) on p.36. ■

Example 4.1 (Ross-prob-thy-3ed Example 2c). How many 7–digit license plates can be created if the first three are letters (CAPS) and the last four are digits?

Answer: $26^3 \cdot 10^4 = 175,760,000$ \square

Example 4.2 (Ross-prob-thy-3ed Example 2e). How many different 7–digit license plates can be created if the first three are letters (CAPS) and the last four are digits and none of those symbols can be repeated?

Answer: $26 \cdot 25 \cdot 24 \cdot 10 \cdot 9 \cdot 8 \cdot 7 = 78,624,000$ \square

Example 4.3. How many 7–digit license plates can be created if the first three are letters (CAPS) and the last four are digits and none of the letters can be repeated?

Answer: $26 \cdot 25 \cdot 24 \cdot 10^4 = 26 \cdot 600 \cdot 10^4 = 15,600 \cdot 10^4 = 15,600,000$. \square

Example 4.4 (Ross-prob-thy-3ed Example 2d). If $|\Omega| = n$, how many different functions $\psi : \Omega \rightarrow \{0, 1\}$, i.e., how many functions on Ω that can only take the values 0 and 1, do exist?

Answer: If $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$, then

- we have 2 choices for the $\psi(\omega_1)$ selection.
- For each of those there are 2 choices for the $\psi(\omega_2)$ selection.
- For each of those $\psi(\omega_1), \psi(\omega_2)$ selections there are 2 choices for the $\psi(\omega_3)$ selection.
- -----
- For each of those $\psi(\omega_1), \dots, \psi(\omega_{n-1})$ selections there are 2 choices for the $\psi(\omega_n)$ selection.

So we have $2 \cdot 2 \cdots 2 = 2^n$ selections. \square

Example 4.5. If $|\Omega| = n$, how many subsets of Ω , including \emptyset and Ω , do exist?

Answer: If $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$, any subset $A \subseteq \Omega$ can be uniquely represented by an element $\vec{d} = \vec{d}(A) = (d_1, d_2, \dots, d_n)$ of $\{0, 1\}^n$ as follows:

- $d_j = 1 \Leftrightarrow \omega_j \in A$ and $d_j = 0 \Leftrightarrow \omega_j \notin A$.

The assignment $F : A \mapsto \vec{d}(A)$ between the subsets of Ω and $\{0, 1\}^n$ is injective:

- If $A' \subseteq \Omega$ such that $\vec{d}(A) = \vec{d}(A')$, then $\omega \in A \Leftrightarrow \omega \in A'$, i.e., $A = A'$.

F also is surjective: if $\vec{d}(d_1, d_2, \dots, d_n) \in \{0, 1\}^n$, then

- $B := \{\omega_j : d_j = 1\}$ (a subset of Ω) which satisfies $F(B) = \vec{d}$.

Thus, F is a bijection. We illustrate this with the following example. Let $\Omega := \{\omega_1, \omega_2, \omega_3, \omega_4\}$.

- $A_1 = \{\omega_2, \omega_3\} \Rightarrow F(A_1) = (0, 1, 1, 0)$. Also, $F^{-1}(0, 1, 1, 0) = \{\omega_j : d_j = 1\} = \{\omega_2, \omega_3\} = A_1$.
- $A_2 = \{\omega_4\} \Rightarrow F(A_2) = (0, 0, 0, 1)$. Also, $F^{-1}(0, 0, 0, 1) = \{\omega_j : d_j = 1\} = \{\omega_4\} = A_2$

Since F is a bijection, there are as many subsets of Ω as there are vectors

$\vec{d}(A) = (d_1, d_2, \dots, d_n)$ of zeros and ones of length n . And how many are those?

- we have 2 choices for d_1 : either $d_1 = 0$ or $d_1 = 1$.
- For each of those choices: either $d_2 = 0$ or $d_2 = 1$.
- -----
- For each of those 2^{n-1} choices $[d_j = 0 \text{ or } d_j = 1 \text{ (} j = 1, 2, \dots, n-1 \text{)}]$: either $d_n = 0$ or $d_n = 1$.

Thus, we have $2 \cdot 2 \cdots 2 = 2^n$ choices. \square

4.2 Permutations

Definition 4.1 (WMS Ch.02.6, Definition 2.7 - Permutation).

An ordered arrangement of r distinct objects is called a **permutation** of size r . The number of ways of ordering n distinct objects taken r at a time will be designated by the symbol P_r^n .

\square

Theorem 4.2 (WMS Ch.02.6, Theorem 2.2).

$$(4.1) \quad P_r^n = n(n-1)(n-2) \cdots (n-r+1) = \frac{n!}{(n-r)!}.$$

Here, $n!$ (“ n factorial”) is defined as follows.

$$(4.2) \quad n! = \begin{cases} n(n-1) \cdots 2 \cdot 1, & \text{if } n \in \mathbb{N}, \\ 1, & \text{if } n = 0. \end{cases}$$

PROOF: We can consider each permutation as the result of the following actions A_1, \dots, A_r .

- A_1 is the selection of the first item. Since all n items are available for selection, A_1 has n outcomes.
- A_2 is the selection of the second item. Since one item was already selected and duplicates are not allowed, only $n-1$ items are available for selection. Thus, A_2 has $n-1$ outcomes.
- -----
- A_r is the selection of item r . Since $r-1$ items have been previously selected and duplicates are not allowed, only $n-(r-1) = n-r+1$ items are available for selection. Thus, A_r has $n-r+1$ outcomes.

It follows from the multiplication rule that there are $n(n-1) \cdots (n-r+1)$ different ways to select r items without repeating a selection, i.e., of obtaining a permutation of size r of those n items. \blacksquare

Example 4.6 (WMS Ch.02.8, Example 2.8). The names of 3 employees are to be randomly drawn, without replacement, from a bowl containing the names of 30 employees of a small company. The person whose name is drawn first receives \$100, and the individuals whose names are drawn second and third receive \$50 and \$25, respectively. How many sample points are associated with this experiment?

Solution: Because the prizes awarded are different, the number of sample points is the number of ordered arrangements of $r = 3$ out of the possible $n = 30$ names. Thus, the number of sample points in S is

$$P_3^{30} = \frac{30!}{27!} = (30)(29)(28) = 24,360. \quad \square$$

Example 4.7. Jenny has collected 20 post cards, all of them different: 4 from France, 2 from Peru, 8 from Japan, 6 from Kenia. She wants to place them into 4 numbered boxes according to their country of origin.

(A) Jenny consider two arrangements different if, say, Esteban’s card takes a different spot in the Peru box, but she does not care whether the Peru cards end up in box #1 or #2 or #3 or #4. How many different arrangements are possible?

Answer:

- 4 choices for France card #1,
- 3 choices for France card #2 (into the same box),
- 2 choices for France card #3 (into the same box),
- 1 choice for France card #4 (into the same box).
- Thus, there are $4!$ choices for the France cards.
- For each one of those $4!$ choices we obtain in a similar manner that there are $2!$ choices for Peru.
- For each one of those $4! \cdot 2!$ choices we obtain in a similar manner that there are $8!$ choices for Japan.
- For each one of those $4! \cdot 2! \cdot 8!$ choices we obtain in a similar manner that there are $6!$ choices for Kenia.

Thus, $4! \cdot 2! \cdot 8! \cdot 6!$ different arrangements are possible.

(B) As before, Jenny considers two arrangements different if, say, Esteban’s card takes a different spot in the Peru box. But this time it also matters in which box a country’s cards are placed.. How many different arrangements are possible now?

Answer: There are $4!$ permutations of the 4 boxes. This amounts to $4!$ rearrangements of each choice made in (A). Thus, $4! \cdot 2! \cdot 8! \cdot 6! \cdot 4!$ arrangements are possible. \square

4.3 Combinations, Binomial and Multinomial Coefficients

In Example 4.5 on p.72, a simple application of the multiplication rule showed the following:

If Ω is a set of finite size, then its powerset 2^Ω (i.e., the set of all subsets of Ω), has size $|2^\Omega| = 2^{|\Omega|}$.

A related question would be the following:

- How many subsets of Ω have size k ?

Examining how many permutations of size k can be obtained from the elements $\omega_1, \omega_2, \dots, \omega_n$ might not be a bad idea, since permutations of distinct items remain free of duplicates, just as we require for (sub-)sets. But rearrangements of the order in which the elements $\omega_{n_1}, \omega_{n_2}, \dots, \omega_{n_k}$ of such a subset lead to different permutations although the subset remains the same, since the order of the elements of a set is disregarded.

Thus, we must divide P_k^n , the number of permutations of size k of the elements of Ω , by the number of rearrangements that one can obtain from a given set of its members. Since that number is P_k^k , we have obtained the following result.

Theorem 4.3.

Let $0 \leq k \leq n$. A set of size n has

$$\frac{n!}{k!(n-k)!} \cdot$$

subsets of size k .

PROOF: There are $P_k^n = n(n-1) \cdots (n-k+1)$ permutations of size k that can be obtained from the n (distinct!) elements $\omega_1, \omega_2, \dots, \omega_n$ of Ω . Let $A := \{\omega_{n_1}, \omega_{n_2}, \dots, \omega_{n_k}\}$ be such a permutation.

There are $P_k^k = k!$ rearrangements of $\omega_{n_1}, \omega_{n_2}, \dots, \omega_{n_k}$. Since order does not matter in sets (and their subsets), each one of those $k!$ permutations forms one and the same set A .

To say this differently, the number P_k^n was obtained by counting each size k subset $k!$ times.

Thus, we must divide P_k^n by P_k^k to obtain the number of subsets of size k . We obtain

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

This proves the theorem. ■

Selections of size k from a collection of n distinct objects disregarding the order in which those k items were selected (as is the case when selecting a subset of size k from a set of size $n \geq k$), are so important when counting is involved that they deserve a name of their own. For the following see also WMS Ch.02.6, Definition 2.8.

Definition 4.2 (Number of combinations).

We call the number of selections of size k from a collection of n distinct items when the order in which those k items were selected is ignored, the **number of combinations of n objects taken k at a time**. We write $\binom{n}{k}$ for this number. □

Remark 4.1.

- (a) Some texts also use the symbol C_k^n instead of $\binom{n}{k}$. This is considered outdated terminology.
- (b) We emphasize that both are true: $\binom{n}{k}$
 = number of selections of size k from n distinct items when disregarding order
 = number of subsets of size k of a set of size n . □

Theorem 4.4.

Given are n items of which n_1 are alike, n_2 are alike, \dots , n_r are alike ($n_1 + \dots + n_r = n$). Then the number of distinguishable arrangements of those n items is

$$\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! n_2! \cdots n_r!}.$$

PROOF:

- We tag the group 1 items as $x_1^{(1)}, x_2^{(1)}, \dots, x_{n_1}^{(1)}$,
- the group 2 items as $x_1^{(2)}, x_2^{(2)}, \dots, x_{n_2}^{(2)}$,
- -----
- the group r items as $x_1^{(r)}, x_2^{(r)}, \dots, x_{n_r}^{(r)}$,

to make all n items artificially distinguishable. We have learned that there are $n!$ permutations.

When we only keep the superscripts that indicate the group but we remove the subscripts, since in truth items belonging the same group cannot be distinguished, there will be a lot less arrangements that are distinct.

To fix the ideas, assume that group 2 has 4 members and we have an arrangement

$$\text{Arr \#1: } * * * x_3^{(2)} * * * * * x_2^{(2)} x_4^{(2)} * * * * x_1^{(2)} * *$$

and that we have another arrangement

$$\text{Arr \#2: } * * * x_1^{(2)} * * * * * x_4^{(2)} x_2^{(2)} * * * * x_3^{(2)} * *$$

where all items that do not belong to group 2 (the ones marked “*”) occupy the same column in both arrangements. To put it differently, we obtained Arr #2 from Arr #1 by permuting the items in group 2 and leaving all other items in place.

In total there are $n_2! = 4! = 24$ such permutations. Let us consider one of them as special. For example, this one,

$$\text{Arr \#5: } * * * x_1^{(2)} * * * * * x_2^{(2)} x_3^{(2)} * * * * x_4^{(2)} * *$$

where the group 2 items are arranged, left to right, in increasing order of their subscripts.

We go through all $n!$ permutations and discard all those where the group 2 items are ordered differently from $x_1^{(2)}, x_2^{(2)}, x_3^{(2)}, x_4^{(2)}$.

$$\text{Then only } \frac{n!}{n_2!} \text{ arrangements remain,}$$

but for those the artificial distinction which was introduced by the subscripts is gone in group 2.

We repeat the above procedure to those survivors, but for group 1. We discard all those where the group 1 items are not ordered $x_1^{(1)}, x_2^{(1)}, \dots, x_{n_1}^{(1)}$.

$$\text{Then only } \frac{n!}{n_2! n_1!} \text{ arrangements remain,}$$

but for those the artificial distinction which was introduced by the subscripts is gone in groups 1 and 2.

We keep going with the remaining groups.

$$\text{Then only } \frac{n!}{n_1! n_2! \cdots n_r!} \text{ arrangements remain,}$$

but for those the artificial distinction which was introduced by the subscripts is gone in all r groups.

It follows that there are $n! / (n_1! n_2! \cdots n_r!)$ different arrangements if we cannot distinguish the items belonging to the same group. ■

Example 4.8. How many distinct permutations are there of the word SHANANANANA

Answer: We designate Groups 1–4 according to the letters S, H, A, N.

Then $n_1 = n_2 = 1, n_3 = 5, n_4 = 4$. Further, $n = 1 + 1 + 5 + 4 = 11$. Thus, there are

$$\frac{11!}{5! \cdot 4! \cdot 1! \cdot 1!} = \frac{11 \cdot 10 \cdot 9 \cdot 8 \cdot 7 \cdot 6}{4 \cdot 3 \cdot 2} = \frac{11 \cdot 10 \cdot 9 \cdot 8 \cdot 7 \cdot 6}{3(4 \cdot 2)} = 11 \cdot 10 \cdot 9 \cdot 7 \cdot 2 = 13,860$$

distinguishable arrangements of the word SHANANANANA. □

Definition 4.3 (Multinomial coefficients).

The numbers

$$(4.3) \quad \binom{n}{n_1 n_2 \cdots n_r} = \frac{n!}{n_1! n_2! \cdots n_r!}.$$

that appear in Theorem 4.4 are called **multinomial coefficients**. If $r = 2$, then there is some integer $0 \leq k \leq n$ such that $n_1 = k$ and $n_2 = n - k$. We write

$$(4.4) \quad \binom{n}{k} := \frac{n!}{k!(n-k)!} \quad \text{for} \quad \binom{n}{k, n-k}$$

and speak of **binomial coefficients**. Convention: We define $\binom{n}{k} := 0$ for $k > n$. \square

The next theorem explains the appropriateness of the previous definition.

Theorem 4.5.

Let $r, n \in \mathbb{N}$ such $r \leq n$ and $x_1, x_2, \dots, x_r \in \mathbb{R}$. Then

$$(4.5) \quad (x_1 + x_2 + \cdots + x_r)^n = \sum_{\substack{n_1, \dots, n_r \geq 0 \\ n_1 + \cdots + n_r = n}} \binom{n}{n_1, n_2, \dots, n_r} x_1^{n_1} x_2^{n_2} \cdots x_r^{n_r}.$$

In particular, if $r = 2$, we obtain the **binomial theorem**:

$$(x_1 + x_2)^n = \sum_{j=0}^n \binom{n}{j} x_1^j x_2^{n-j}.$$

PROOF:

First, we show that the case $n = 2$ follows from 4.5.

Since $n_1, n_2 \geq 0$ and $n_1 + n_2 = n \Rightarrow 0 \leq n_1 \leq n$ and $n_2 = n - n_1$, writing j for n_1 yields the binomial theorem formula.

To prove the first formula, We start by "multiplying out" the product

$$(x_1 + x_2 + \cdots + x_r)^n = (x_1 + x_2 + \cdots + x_r)(x_1 + x_2 + \cdots + x_r) \cdots (x_1 + x_2 + \cdots + x_r)$$

and obtain in the resulting expansion terms of the form

$$a_1 \cdot a_2 \cdots a_n \quad \text{such that each factor } a_j \text{ is either } x_1 \text{ or } x_2 \dots \text{ or } x_r.$$

In the following we consider the sizes n_1, n_2, \dots, n_r as fixed

Note that it is not possible to obtain two selections

$$\vec{a} = (a_1, a_2, \dots, a_n) \quad \text{and} \quad \vec{b} = (b_1, b_2, \dots, b_n) \quad \text{such that} \quad a_j = b_j \quad \text{for all } j.$$

The reason: We multiply out the n factors $(x_1 + \cdots + x_r)$ in such a way that for no two of the resulting products we picked the same variable x_i in each one of those n factors $(x_1 + \cdots + x_r)$

But then the following is true if we consider such a selection as a word $a_1 a_2 \dots a_n$ where each letter is one of x_1 or $x_2 \dots$ or x_r . Any two of those words are distinguishable even though some or all of the letters x_i can occur multiple times.

For example, if $n = 7, n_1 = 2, n_2 = 3, n_3 = 2$ and we write X for x_1, Y for x_2, Z for x_3 , we have this situation.

The word $YXZZYYX$ is formed only once. But of course, we obtain other words with the same sizes n_j , e.g. the rearrangement $ZYXZYXY$ which is distinguishable from the first word.

Thus, in the general case, there are as many terms in the expansion of $(x_1 + x_2 + \dots + x_r)^n$ containing each symbol x_j exactly n_j times as there are distinguishable “words” that contain each x_j exactly n_j times. According to Theorem 4.4, there are

$$\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! n_2! \dots n_r!}.$$

such terms. Since this is the number of times the product $x_1^{n_1} x_2^{n_2} \dots x_r^{n_r}$ occurs in the expansion of $(x_1 + x_2 + \dots + x_r)^n$, it follows that

$$(x_1 + x_2 + \dots + x_r)^n = \sum_{\substack{n_1, \dots, n_r \geq 0 \\ n_1 + \dots + n_r = n}} \binom{n}{n_1, n_2, \dots, n_r} x_1^{n_1} x_2^{n_2} \dots x_r^{n_r}. \blacksquare$$

Theorem 4.6.

Given are n distinct items and r distinct bins of fixed sizes n_1, n_2, \dots, n_r such that $n_1 + \dots + n_r = n$. Then the number of distinguishable placements of the n items into those r bins, when disregarding the order in which the items were placed into any one of those bins, is

$$\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! n_2! \dots n_r!}.$$

The proof is given after the following example which will help clarify how to interpret Theorem 4.6.

Example 4.9. Given are a list of $n = 7$ items and $r = 3$ bins as follows.

- The 7 items are a, b, c, d, e, f, g .
- Bin 1 has size 2, bin 2 has size 3, bin 3 has size 2 (thus $n = 2 + 3 + 2 = 7$).
- Arr #1: bin 1 has b, c , bin 2 has e, a, g , bin 3 has f, d
- Arr #2: bin 1 has c, b , bin 2 has a, g, e , bin 3 has d, f
- Arr #3: bin 1 has b, d , bin 2 has a, g, e , bin 3 has c, f
- Then Arr #1 and Arr #2 are considered the same since each bin contains the same items. Only their order is different.
- On the other hand, both Arr #1 and Arr #2 both are considered different from Arr #3 since, e.g., bin 1 contains item d for #3, but bin 1 does not contain item d for the other two arrangements. \square

PROOF of Theorem 4.6:

The proof is very similar to that of Theorem 4.4, so we keep the discussion brief.

- For each one of the $n!$ permutations of all n items, there are $n_1! - 1$ others which possess the same n_1 elements in bin 1, only differently ordered, but have exactly the same item at each other of the remaining $n - n_1$ spots. Removing those duplicates leaves us with $n!/n_1!$ arrangements.
- Of those $n!/n_1!$ arrangements, there are $n_2! - 1$ others which possess the same n_2 elements in bin 2, only differently ordered, but have exactly the same item at each other of the remaining $n - n_1 - n_2$ spots. Removing those duplicates leaves us with $n!/(n_1!n_2!)$ arrangements.
- -----
- Having removed the duplicates from bins 1 through $k - 1$, we are left with $\frac{n!}{n_1! \cdots n_{k-1}!}$ arrangements. For each one of those there are $n_k! - 1$ others which possess the same n_k elements in bin k , only differently ordered. Removing those duplicates leaves us with $\frac{n!}{n_1! \cdots n_k!}$ arrangements.
- For any two surviving arrangements the following is true: There is at least one bin that does not contain the same elements (possibly rearranged) for both those arrangements.

to make all n items artificially distinguishable. We have learned that there are $n!$ permutations.

When we only keep the superscripts that indicate the group but we remove the subscripts, since in truth items belonging the same group cannot be distinguished, there will be a lot less arrangements that are distinct.

To fix the ideas, assume that group 2 has 4 members and we have an arrangement

$$\text{Arr \#1: } * * * x_3^{(2)} * * * * * x_2^{(2)} x_4^{(2)} * * * * x_1^{(2)} * *$$

and that we have another arrangement

$$\text{Arr \#2: } * * * x_1^{(2)} * * * * * x_4^{(2)} x_2^{(2)} * * * * x_3^{(2)} * *$$

where all items that do not belong to group 2 (the ones marked “*”) occupy the same column in both arrangements. To put it differently, we obtained Arr #2 from Arr #1 by permuting the items in group 2 and leaving all other items in place.

In total there are $n_2! = 4! = 24$ such permutations. Let us consider one of them as special. For example, this one,

$$\text{Arr \#5: } * * * x_1^{(2)} * * * * * x_2^{(2)} x_3^{(2)} * * * * x_4^{(2)} * *$$

where the group 2 items are arranged, left to right, in increasing order of their subscripts.

We go through all $n!$ permutations and discard all those where the group 2 items are ordered differently from $x_1^{(2)}, x_2^{(2)}, x_3^{(2)}, x_4^{(2)}$.

$$\text{Then only } \frac{n!}{n_2!} \text{ arrangements remain,}$$

but for those the artificial distinction which was introduced by the subscripts is gone in group 2.

We repeat the above procedure to those survivors, but for group 1. We discard all those where the group 1 items are not ordered $x_1^{(1)}, x_2^{(1)}, \dots, x_{n_1}^{(1)}$.

$$\text{Then only } \frac{n!}{n_2! n_1!} \text{ arrangements remain,}$$

but for those the artificial distinction coming from the subscripts is gone in groups 1 and 2.

We keep going with the remaining groups....

In the end only $\frac{n!}{n_1! n_2! \cdots n_r!}$ arrangements remain,

but for those the artificial distinction which was introduced by the subscripts is gone in all r groups. It follows that there are $n! / (n_1! n_2! \cdots n_r!)$ different arrangements if we cannot distinguish the items belonging to the same group. ■

Proposition 4.1.

(A) There are $\binom{n-1}{r-1}$ distinct integer-valued vectors $\vec{x} = (x_1, x_2, \dots, x_r)$ such that

$$x_1 + x_2 + \cdots + x_r = n \quad \text{and} \quad x_i > 0, i = 1, \dots, r.$$

(B) There are $\binom{n+r-1}{r-1}$ distinct integer-valued vectors $\vec{y} = (y_1, y_2, \dots, y_r)$ such that

$$y_1 + y_2 + \cdots + y_r = n \quad \text{and} \quad y_i \geq 0, i = 1, \dots, r.$$

PROOF of (A):

Each such equation corresponds to an arrangement of n symbols \otimes which denote the numbers $1, 2, \dots, n$ in sequence, and $r - 1$ bars $|$ which are places in-between those symbols, in such a way, that no two bars are adjacent. For example, the arrangement

$$\bullet \bullet | \bullet \bullet \bullet | \bullet \bullet \bullet$$

expresses the equation $2 + 4 + 3 = 7$. In the general case, one or zero bars can be placed in the $n - 1$ gaps between the n bullets:

(A) $\bullet \otimes \bullet \otimes \bullet \otimes \bullet \otimes \bullet \otimes \bullet \otimes \bullet \otimes \dots \otimes \bullet \otimes \bullet$

Thus, there are as many different integer equations as there are ways to select $r - 1$ of those $n - 1$ gaps for the $r - 1$ bars. This number is $\binom{n-1}{r-1}$.

FIRST PROOF of (B):

An equation $\sum_{j=1}^r y_j = n; y_j \geq 0$ of part (B) becomes an equation $\sum_{j=1}^r x_j = n + r; x_j > 0$ of part (A), by setting $x_j := y_j + 1$.

In reverse, equation $\sum_{j=1}^r x_j = n + r; x_j > 0$ of part (A) becomes an equation $\sum_{j=1}^r y_j = n; y_j \geq 0$ of part (B), by setting $y_j := x_j - 1$.

We have shown in (A) that there are $\binom{n+r-1}{r-1}$ different equations of the form $\sum_{j=1}^r x_j = n + r; x_j > 0$.

Thus, there also that many of the form $\sum_{j=1}^r y_j = n; y_j \geq 0$. This proves (B).

ALTERNATE PROOF of **(B)**: We add two more placeholders \otimes for the separating bars. One to the left of the leftmost bullet and another to the right of the rightmost bullet. The condition $y_j \geq 0$ instead of $x_j > 0$ implies that each one of those placeholders can be occupied by as few as zero bars and as many as all $r - 1$ bars. To put it differently, any combination of bullets and bars is admissible. We create a tagged list of $n + r - 1$ distinct placeholders for both bullets and bars and select $r - 1$ of them for the bars. Obviously, the order of the bars does not matter. Thus there are $\binom{n+r-1}{r-1}$ such selections. ■

Consider the issue of distributing n distinct items into r distinct bins where bin _{j} contains $0 \leq n_j \leq n$ items and the n_j are allowed to vary for different selections. (But of course, $n_1 + \dots + n_r = n$.)

Then each such selection corresponds to an integer vector $\vec{n} = (n_1, \dots, n_r)$ which is a solution of the equation $\sum_{j=1}^r n_j = n; n_j \geq 0$.

If we demand in addition that each bin contains at least one item, then each such selection corresponds to an integer vector $\vec{n} = (n_1, \dots, n_r)$ which is a solution of the equation $\sum_{j=1}^r n_j = n; n_j > 0$.

We obtain from Proposition 4.1 the following.

Proposition 4.2.

(A) There are $\binom{n-1}{r-1}$ ways to select n indistinguishable items into r distinct bins such that each bin contains at least one item.

(B) There are $\binom{n+r-1}{r-1}$ ways to select n indistinguishable items into r distinct bins.

PROOF: This follows from from Proposition 4.1. ■

Example 4.10. Mother Jones’ cookies and the stars & bars examples:

- How many ways are there to give 10 cookies to 4 kids if each one gets at least one cookie? **A:** There are $\binom{10-1}{4-1} = (9 \cdot 8 \cdot 7)/(3 \cdot 2 \cdot 1) = 84$ ways.
- How many ways are there to separate 6 stars by two bars into three parts, if one or more of those parts may contain zero stars? **A:** There are $\binom{6+3-1}{3-1} = (8 \cdot 7)/(2 \cdot 1) = 28$ ways. □

Here is another example that employs binomial coefficients.

Example 4.11 (Ross-prob-thy-3ed Example 4c). Given are n antennas of which d are defective. They will be arranged in a linear order and will relay signals. This chain will not function if two or more defective items are placed next to each other.

How many ways are there to arrange the antennas so that we obtain a functioning arrangement?

Answer: We denote the $n - d$ working antennas by the \otimes symbol, separate them by bullets \bullet and add one \bullet each to the left of the leftmost and to the right of the rightmost.



Then the functioning relays are precisely those where one or zero defective antennas are placed at each one of those \bullet spots. Each such placement corresponds to a selection of size d of those $n - d + 1$

bullets: The selected spots will get a defective antenna and nothing will happen to the others.

Thus, there are $\binom{n-d+1}{d}$ functioning arrangements. \square

Example 4.12. A lottery is held among N participants. There are K drawings in which a prize is given away. ($K < N$). In each drawing, each participant has an equal chance of obtaining the prize. (Thus, it is possible, though unlikely, that one single person walks away with all K prizes.) Amanda is one of the participants. What is the probability that she will walk away with exactly k prizes? Of course, ($k \leq K$).

Solution:

- (a) There are N different selections for drawing #1.
- (b) Each one of those has N selections for drawing #2. Thus, there are N^2 different ways to distribute the first two prizes
- (c) Each one of those N^2 has N selections for drawing #3. Thus, there are N^3 different ways to distribute the first 3 prizes
- (d) Thus, there are N^K different ways to distribute all K prizes

It follows that the sample space Ω has size N^K . Since all drawings are done at random, all outcomes $\omega \in \Omega$ are equally likely. Thus, $P\{\omega\} = 1/(N^K)$ for all ω . Note that an outcome $\omega \in \Omega$ is of the form

(\star) $\omega = (i_1, i_2, \dots, i_K) : \text{ prize 1 goes to person } i_1, \dots \text{ prize } K \text{ goes to person } i_K$

- Let $A := \{ \text{Jane gets exactly } k \text{ prizes} \}$.

Assume that the outcomes ω and ω' are as follows:

- ω : participant i_1 gets prize j_1 and i_2 gets prize j_2
- ω' : participant i_1 gets prize j_2 and i_2 gets prize j_1
- There is no difference how other $K - 2$ prizes were awarded.

Even though order matters, we only are able to distinguish the outcomes ω and ω' if j_1 and j_2 are given to different persons. Otherwise all K slots of both ω and ω' are identical, i.e., $\omega = \omega'$.

Thus, there are (only) as many different ways to give k of the K prizes to Jane as there are ways to select k of K items DISREGARDING ORDER. That number is $\binom{K}{k}$.

Next, consider that each one of those $\binom{K}{k}$ ways of designing k of the K slots of an outcome ω to Jane must be complemented by filling each one of the remaining $K - k$ slots with one of the other $N - 1$ participants. This time we CANNOT DISREGARD ORDER. See the discussion above concerning the outcomes ω and ω' .

- We repeat the reasoning of (a) – (d) to $N - 1$ instead of N choices for those $K - k$ instead of k drawings and see that there are $(N - 1)^{K-k}$ possible selections.
- The event A consists all outcomes obtained by matching any one of those $(N - 1)^{K-k}$ selections with any one of the $\binom{K}{k}$ ways of allocating k prizes to Jane.
- By the multiplication rule, $|A| = \binom{K}{k} (N - 1)^{K-k}$.
- Since all outcomes are equally likely, $P(A) = \frac{|A|}{|\Omega|} = \frac{\binom{K}{k} (N - 1)^{K-k}}{N^K}$. \square

We summarize the results of Theorem 4.4, Theorem 4.6, Proposition 4.1, and Proposition 4.2.

Remark 4.2. The multinomial coefficients

$$\binom{n}{n_1 n_2 \cdots n_k} = \frac{n!}{n_1! n_2! \cdots n_k!}.$$

of Definition 4.3 appear in the following settings:

- Distinct selections of n items of which n_1 are alike, n_2 are alike, ..., n_k are alike. Example: different rearrangements of the word "BANANA".
- They are coefficients in the expansion of $(x_1 + x_2 + \cdots + x_k)^n$.
- Distinct selections of n items into k distinct bins of fixed sizes n_1, \dots, n_k . That is the WMS definition in their Theorem 2.3 of Ch.02.6.
- Subdividing n indistinguishable items into k partitions, where the sizes n_1, \dots, n_k of those partitions are allowed to vary for different subdivisions. Example: number of integer valued vectors (n_1, \dots, n_k) such that $n_1, \dots, n_k \geq 0$ and $\sum_j n_j = n$. \square

5 More on Probability

This chapter corresponds to material found in WMS ch.2

5.1 Total Probability and Bayes Formula

Theorem 5.1 (Total Probability and Bayes Formula ³¹).

Assume that $\{B_1, B_2, \dots\}$ is a partition of Ω and that $A \subseteq \Omega$. such that $P(B_j) > 0$ for all j . Then

$$(5.1) \quad P(A) = \sum_{j=1}^{\infty} P(A | B_j) P(B_j).$$

$$(5.2) \quad P(B_j | A) = \frac{P(A | B_j)P(B_j)}{\sum_{i=1}^{\infty} P(A | B_i) P(B_i)}. \quad (\text{Bayes formula})$$

Note that the above also covers finite partitions $\{B_1, B_2, \dots, B_k\}$ of Ω : apply the formulas with

$$B_{k+1} := B_{k+2} := \dots := 0.$$

PROOF: Since $(B_j)_j$ partitions Ω $(A \cap B_j)_j$ partitions A . Thus, $A = \bigsqcup_j (A \cap B_j)$. Thus,

$$P(A) = \sum_{j=1}^{\infty} P(A \cap B_j). = \sum_{j=1}^{\infty} P(A | B_j) P(B_j).$$

This proves (5.1). To prove (5.2), we apply to its right-hand side the already proven (5.1). We obtain

$$\frac{P(A | B_j)P(B_j)}{\sum_{i=1}^{\infty} P(A | B_i) P(B_i)} = \frac{P(A | B_j)P(B_j)}{P(A)} = \frac{P(A \cap B_j)}{P(A)} = P(B_j | A). \quad \blacksquare$$

When working with conditional probabilities, in particular when one wants to apply the Bayes formula, it often is convenient to work with tree diagrams. This is demonstrated in the next example.

Example 5.1. It has been established that 40% of all jobs for college graduates are in the technology sector. Of those college graduates who work in technology, one quarter enjoys listening to classical music. Of those college graduates who hold other kinds of jobs, one out of three enjoys listening to classical music.

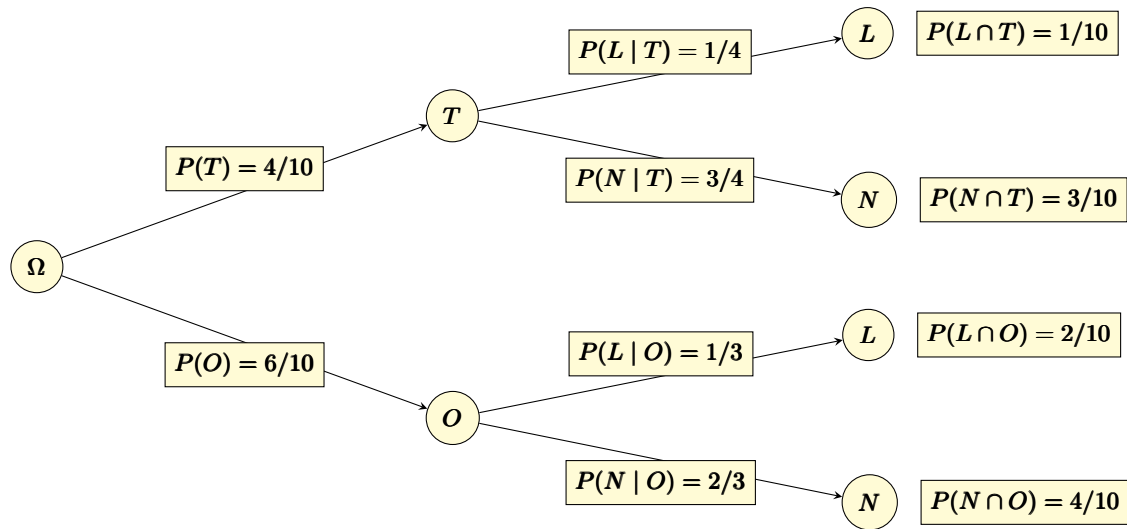
- What is the probability that Pedro neither works in technology, nor listens to classical music?
- Harry does not listen to classical music. How likely is it that he works in technology?
- Jane says that she likes classical music. What is the probability that she works in technology?

Solution: We use the following abbreviations:

T: Works in technology O: "Other": does not work in technology
L: Listens to classical music N: Does not listen to classical music

³¹Thomas Bayes (1702 - 1761) was an English clergyman and mathematician.

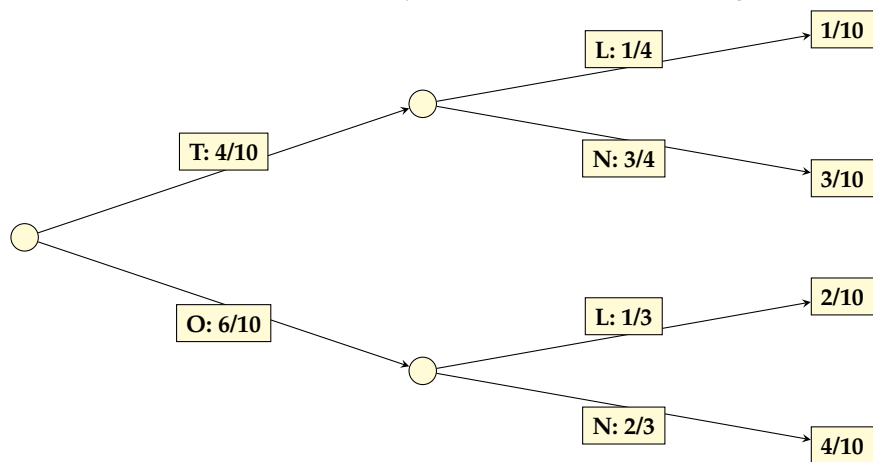
The information available to us is sufficient to draw the following tree diagram:



A line segment that connects two nodes indicates conditioning of the right side on the left side. For example, the node that connects T and N signifies that the event N is conditioned on the event T . $P(L | T)$, the corresponding conditional probability, is attached to the line segment. Note that this is also true for the two line segments that emanate from Ω , since $P(T) = P(T | \Omega)$ and $P(O) = P(O | \Omega)$. Note that T and O partition Ω and the same is true for L and N .

Tree diagrams can be very convenient because the probability of an intersection is obtained by multiplying the two probabilities to the left. For example, $P(T \cap N) = (4/10)(3/4) = 3/10$.

Not all the notation is necessary to work with such a diagram. Here is a pared down version:



Let us now discuss the answers to the three problems posed above

- (a) What is the probability that Pedro neither works in technology, nor listens to classical music?
 - This is the ordinary (no conditioning) probability $P(O \cap N) = 4/10$.

- (b) Harry works in technology. How likely is it that he does not listen to classical music?
 - We are conditioning on the event T and want to compute $P(N | T)$. The diagram shows that $P(N | T) = 3/4$.
- (c) Jane says that she likes classical music. What is the probability that she works in technology?
 - We are asking for the conditional probability $P(T | L)$.

This is a reverse conditioning (Bayes formula problem). The tree diagram makes it easy to find all the probabilities involved:

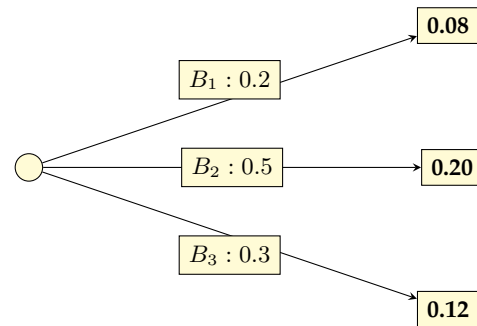
- $P(T | L) = P(T \cap L)/P(L)$.
- $P(T \cap L) = 1/10$ and $P(L) = P(O \cap L) + P(T \cap L) = (2 + 1)/10 = 3/10$.
- Thus, $P(T | L) = (1/10)/(3/10) = 1/3$.

We continue with some general remarks concerning tree diagrams.

It should be clear how to generalize such diagrams. One can condition at each stage on more than just two events. For example, Let us assume the following. In stage 1, we “condition” Ω on $\Omega = A_1 \uplus A_2 \uplus A_3$, In stage 2, we condition A_2 on $\Omega = B_1 \uplus B_2 \uplus B_3 \uplus B_4$. If

$$P(A_2) = 0.4,$$

then the resulting tree fragment is to the right.



Because $\Omega = \uplus_j B_j$, it is always true that

$$\sum_j P(B_j | A_k) = \frac{\sum_j P(B_j \cap A_k)}{P(A_k)} = \frac{P(A_k)}{P(A_k)} = 1$$

Thus, the sum of the conditional probabilities over all line segment that emanate from a given node is 1. In the tree excerpt above: that node is $A_k = A_2$ and the sum of the conditional probabilities is

$$P(B_1 | A_2) + P(B_2 | A_2) + P(B_3 | A_2) = 0.2 + 0.5 + 0.3 = 1. \quad \square$$

5.2 Random Sampling and Urn Models With and Without Replacement

The following definition is **PRELIMINARY** and will be amended in Definition 5.2 (Sampling as a Random element) below (see p.88).

Definition 5.1.

- (a) We call the action of picking n items x_1, x_2, \dots, x_n from a collection of N items a **sampling action of size n** . Alternatively, we also use the phrases **sampling process** and **sampling procedure**. Here, $n \in \mathbb{N}$ and $N \in \mathbb{N}$ or $N = \infty$.
- (b) We call the specific outcome of such a sampling action (the list x_1, x_2, \dots, x_n) a **realization** of that sampling action. \square
- (c) In yet another instance of notational abuse, both the sampling action and an outcome of this action (a realization) will be referred to as a **sample** of size n if this does not lead to any confusion. \square

Example 5.2. Each of the following can be considered sampling actions.

- (a) Drawing blindfolded a ball from an urn that contains N balls $n = 5$ times in a row recording each time the outcome and then replacing the ball (putting it back).
- (b) Drawing blindfolded $n = 5$ balls from an urn that contains N balls in one fell swoop, i.e., not replacing any of the balls
- (c) Rolling a die twice in a row and recording the outcome.
- (d) Selecting in a random fashion $n = 2,000$ persons from all persons eligible to vote without replacement, i.e., we want a sample of n distinct voters. Note that N is huge when compared to n .
- (e) Same as (d), but we only record their voting preference, their annual income and their age and discard all other data.
- (f) Same as (e), but we only record their annual income.
- (g) The random numbers generator of a computer creates a sample of n numbers such that they are uniformly distributed on the interval $[0, 1]$.³² (Computers can do that!) Since there are infinitely many such numbers and the computer can generate any one of them,³³ $N = \infty$.
- (h) A factory mass-produces an item, e.g., screws, at a huge rate per hour. Quality control randomly picks $n = 50$ every hour and checks for defective items. Since the number N of screws from which the sample is obtained is so huge, one can, for all practical purposes, act as if $N = \infty$. (This will considerably simplify the mathematics involved in computing, e.g., the probability that such a sample contains 5 or more defective items) if the rate of defectives is supposed to be 3.5%.
- (i) We write down the numbers 1, 2, ;10. This **deterministic sampling action**. is very boring for a course called “Probability Theory”, because no randomness is involved. Nevertheless, Definition 5.1 encompasses deterministic sampling. \square

Remark 5.1.

- (a) We only are interested in sampling actions that involve randomness. In other words, if there is a set U such that $x_j \in U$ for all j , our sampling action can be modeled, for fixed n , as a random element $\vec{X} : (\Omega, P) \rightarrow U^n$. Since deterministic actions also are (constant) random elements, deterministic sampling actions are also covered by Definition 5.1.
- (b) Since the “population” from which each item $x_j = X_j(\omega)$ is sampled is the set U from (a), it is possible to implement $\Omega := U^N$ as the carrier set of the probability space (Ω, P) . In other words, we could narrow things down to $\vec{X} : (U^N, P) \rightarrow U^n$. Matter of fact, you will be as specific as you can when trying to find the formula or just the particular number that solves a given problem.
- (c) But there are advantages to refer to an unspecified probability space (Ω, P) when dealing with the general theory. A good example are the theorems and definitions about expectation and variance in MF Chapter 6 (Discrete Random Variables and Random Elements) where going into specific settings would hinder rather than help the understanding. \square

Here is the promised amended version of Definition 5.1.

³²“uniformly distributed” means that the proportion of numbers x_j that fall within the interval $0 \leq a < b \leq 1$ is (approximately) $b - a$.

³³in theory, since there is no such thing as “infinitely many”) in our physical reality

Definition 5.2 (Sampling as a Random element). Let (Ω, P) be a probability space. Let $U \neq \emptyset$ be a collection of N items ($N \in \mathbb{N}$ or $N = \infty$), which we can think of as the “population of interest”. Let $n \in \mathbb{N}$ (so $n < \infty$), such that $n \leq N$.

- (a) Let $\vec{X} : (\Omega, P) \rightarrow U^n$ be a random element with codomain U^n . If we interpret \vec{X} as the action of picking n items

$$\vec{x} = x_1, x_2, \dots, x_n = X(\vec{\omega}) = X_1(\omega), X_2(\omega), \dots, X_n(\omega)$$

from U , then we call \vec{X} a **sampling action of size n** . Alternatively, we also use the phrases **sampling process** and **sampling procedure**.

- (b) We call a specific outcome (the list $\vec{x} = (x_1, x_2, \dots, x_n)$) a **realization** of that sampling action. See Example 1.5 on p.11.
- (c) Both the sampling action and an outcome of this action (a realization) are called a **sample** of size n if the context makes it clear what is being discussed.
- (d) If there is a specific $\vec{x}^* \in U^n$ such that $P\{\vec{X} = \vec{x}^*\} = 1$, (this certainly is the case if $\vec{X}(\omega) = \vec{x}^*$ for all $\omega \in \Omega$), then we call both the sampling action \vec{X} and the realization \vec{x}^* a **deterministic sample**. \square

Remark 5.2.

- (a) You may wonder about the difference between a U^n -valued random element and a sample of n items which are picked from a population U . The answer: Mathematically speaking, there is no difference whatsoever. It is the interpretation that matters!
- (b) Going back to using the terms probability space and sample space interchangeably, this author likes to think not of (Ω, P) , but only of $(U^n, P_{\vec{X}})$ as a sample space. The reason is that the latter hosts the potential outcomes of the sampling action \vec{X} . (And yes, the probability measure $P_{\vec{X}}$ on that sample space is the distribution of \vec{X}).
- (c) Do those individual sample picks X_j happen with or without replacement? In other words, can the same $x \in U$ be picked more than once or are for a fixed ω all outcomes distinct? The answer: The definition does not say. This must always be explicitly stated or known from the context.
- (d) Consider items (d) and (h) of Example 5.2. If $N \gg n$, then the computational differences between selecting the sample with or without replacement are so small that we can assume sampling with replacement even if the sampled items are not returned to the population after each pick. This often simplifies the computational effort involved. \square

Remark 5.3. We switch focus to the role of proper randomization when picking a sample.

- (a) Picking a small size sample that allows us to make inferences to the population from which it was drawn, can require a lot of thought. The budget available for collecting that sample is often limited and will limit the methods available. Of course, a smaller sample will cost less than a bigger one if the procedure to collect the data is the same in both cases.

We fix $n \in \mathbb{N}$. What will make the sample representative of the population, i.e.

- what guarantees that the composition of the sample mirrors that of the population?

It certainly will not help if the sample has, e.g., 90% students if the population of interest only has 20%. So, we can fix that by establishing quota and restrict the proportion of students to 20%. Of course, there is also the ethnic composition of the population that we want mirrored in the sample. And there is income distribution, gender and 5,000 or more attributes for which we want to maintain close to identical proportions reasonably well.

- (b) Clearly, a practical limit to the number of ways a (hopefully small) can be partitioned into “strata” is reached quickly, so we must look for an alternative way to obtain a sample that is not biased in favor of value a , say “is male” of attribute B (here: gender), when compared to the proportion in the population. And we need this for all important a and B .
The solution is to make the sample selection as random as possible. If we pick the first item at random, i.e., with the same chance $\frac{1}{N}$, then pick #2 at random from the remaining $N - 1$, then pick #3 at random from the remaining $N - 2$, and finally pick # n at random from the remaining $N - n + 1$ items, then this degree of randomness should prevent any kind of gross distortion (bias) in the sample.
- (c) So then, that means that every item has equal chance of being selected, doesn't it? **The answer is NO.** Rather, any collection $\vec{x} = x_1, \dots, x_n$ should have the same chance of being selected as any other collection $\vec{x}' = x'_1, \dots, x'_n$. By the way, we know that probability:
- If we do not worry about the order in which the n distinct items were selected, then there are $\binom{N}{n}$ different selections and that probability must be $1/\binom{N}{n}$.
 - If order does matter and we deal with permutations, then the answer is $1/P_n^N$.
- (d) Would the above requirement be the same as simply asking that each item in the population has the same probability, $1/N$, of being selected? Next comes a counterexample. \square

Example 5.3. We have a population of $N = 600$ students. 100 of them are freshmen, 100 of them are sophomores, 100 of them are juniors, 100 of them are seniors, 100 of them are first year graduate students, the others are second year graduate students.

A sample of $n = 100$ will be selected as follows. A fair die is rolled. If the outcome is 1, all freshmen will be selected, On a 2, all sophomores will be selected, On a 6, all second year graduate students will be selected.

- In the resulting sample each student has the same probability $1/6$ of being selected.
- But only 6 of the possible $\binom{600}{100}$ possible outcomes have a non-zero chance (of $1/6$ each) of being selected: Those where each student belongs to the same group as all the others! \square

There is a special name for the ideal kind of samples (with respect to randomness of the selection). Note that the following definition is tied to sampling without replacement!

Definition 5.3 (Simple Random Sample).

- (a) We call a sampling action of size n ($n \in \mathbb{N}$) from a population of size $N < \infty$ a **simple random sampling action**, in brief, an **SRS action**, if there are no duplicates allowed (i.e., we sample without replacement) and each of the potential outcomes has equal chance of being selected.
- (b) As in Definition 5.2 (Sampling as a Random element), we call both an SRS action and a realization of this action a **simple random sample of size n** . (Briefly, an **SRS**.) \square

Definition 5.4 (Urn models). SRS requires that a single item is selected with equal probability $|U| = 1/N$. When abstracting from the specifics, this boils down to being blindfolded and selecting, without replacement, n well shuffled balls from an urn containing N numbered balls. Some authors also use the scenario of tickets in a box rather than balls in an urn.

- (a) An **urn model without replacement** describes a mechanism by which a blindfolded person selects a fixed number of balls from an urn in which the balls have been well mixed. Note that the resulting realizations will contain no duplicates.
- (b) An **urn model with replacement** describes a mechanism by which a blindfolded person selects a fixed number of balls from an urn as follows.
 - (1) The balls are well mixed.
 - (2) A ball is picked and the outcome is recorded.
 - (3) The ball is put back into the urn.
 - (4) Steps (1) through (3) are repeated until all n balls have been selected. \square

More material may be added to this section at a later time.

6 Discrete Random Variables and Random Elements

This chapter corresponds to material found in WMS ch.3

6.1 Probability Mass Function and Expectation

We start with a trivial observation.

Proposition 6.1. *A real-valued function of a random element is a random variable.*

PROOF: Let $X : (\Omega, P) \rightarrow \Omega'$ be a random element on a probability space (Ω, P) and $g : \Omega' \rightarrow \mathbb{R}$ be a real-valued function. Then $\omega \mapsto g(X(\omega))$ is a real-valued function of ω , hence it is a random variable. ■

Definition 6.1 (Probability mass function).

For a discrete random element X on (Ω, P) , define

$$(6.1) \quad p(x) := p_X(x) := P_X\{x\} = P\{X = x\}.$$

We call p_X the **probability mass function** (WMS: **probability function**) for X . We also write **PMF** for probability mass function. □

Theorem 6.1.

If p_X is the probability mass function of a discrete random element X , then

$$(6.2) \quad 0 \leq p_X(x) \leq 1; \quad \text{for all } x$$

$$(6.3) \quad \sum_{x \text{ s.t. } p_X(x) > 0} p_X(x) = 1$$

Proof: See WMS ch.3. ■

Next, we elaborate on the meaning of $\sum_{x \text{ s.t. } p_X(x) > 0} \dots$

We make the following blanket assumption.

Assumption 6.1 (All series are absolutely convergent).

Unless explicitly stated otherwise, all sequences are either known to be absolutely convergent or assumed to be absolutely convergent. In particular, if $p_X(x)$ is the probability mass function of a discrete random element X which takes values in a set Ω' , $g : \Omega' \rightarrow \mathbb{R}$ is a real-valued function and x_n is a sequence in Ω' , then the series $\sum g(x_j)p_X(x_j)$ is absolutely convergent. □

Remark 6.1. Assume that $p_X(x)$ is the probability mass function of a discrete random element X with values in a set Ω' . Then there exists a countable set $\Omega^* \subseteq \Omega'$ such that $P_X(\Omega^*) = 1$. Thus, the probability mass function $p_X(\cdot)$ of X satisfies

$$p_X(x) = 0 \quad \text{for all } x \in (\Omega^*)^c.$$

Let $g : \Omega' \rightarrow \mathbb{R}$ be a real-valued function. Clearly,

$$g(x) \cdot p_X(x) = 0 \quad \text{for all } x \in (\Omega^*)^c.$$

Ω^* being countable means that $\Omega^* = \{x_1, x_2, \dots\}$ for some finite or infinite sequence x_j . All that follows is trivial in the finite case, so let us confine ourselves to the infinite case $\Omega^* = \{x_j : j \in \mathbb{N}\}$.

For $j \in \mathbb{N}$, let $a_j := g(x_j)p_X(x_j)$. By Assumption 6.1 on p.91, the series $\sum a_j$ is absolutely convergent. Hence, its value does not depend on the ordering of the elements of Ω^* . Thus, we are justified to write

$$\sum_{x \in \Omega^*} g(x)p_X(x) \quad \text{rather than} \quad \sum_{j=1}^{\infty} g(x_j)p_X(x_j).$$

We go a step further. Since $g(x)p_X(x) = 0$ for $x \notin \Omega^*$, we can omit “ $x \in \Omega^*$ ” and write either of the following:

$$\begin{aligned} \sum_x g(x)p_X(x) &= \sum_{x \in \Omega'} g(x)p_X(x) = \sum_{x \in \Omega^*} g(x)p_X(x) \\ (6.4) \qquad &= \sum_{p_X(x) > 0} g(x)p_X(x) = \sum_{j=1}^{\infty} g(x_j)p_X(x_j). \end{aligned}$$

Choosing $g(x) = 1$, we can express probabilities involving X as follows. If $B \subseteq \Omega'$, then

$$(6.5) \quad P\{X \in B\} = P_X(B) = \sum_{x \in B} p_X(x) = \sum_{x \in \Omega^* \cap B} p_X(x) = \sum_{x \in B, p_X(x) > 0} p_X(x). \quad \square$$

Example 6.1. Johnny may choose 2 cookies from a plate with 4 chocolate cookies and 3 oatmeal cookies. We write CC for chocolate cookies and OC for oatmeal cookies. Johnny has no preference and picks two cookies at random.

Let $Y :=$ number of CC chosen by Johnny. Find the PMF $p_Y(y)$ for Y .

Solution:

Note that you were not given the domain (sample space) (S, P) of the random variable Y . There is no need to specify it completely. It suffices to know that, since Johnny can choose 2 of the 7 cookies in $\binom{7}{2}$ ways,

$$(1) \quad |S| = \binom{7}{2} = \frac{7 \cdot 6}{2!} = 21. \quad \text{Since selection was at random, } P\{s\} = \frac{1}{21} \quad \text{for all } s \in S.$$

The codomain can be any set of numbers that contains 0, 1, 2, because $p_Y(y) = P\{Y = y\} = 0$ for all other numbers y . Thus, our task is to compute $p_Y(0)$, $p_Y(1)$, $p_Y(2)$.

(2) Each selection of y CCs comes with a selection of $2 - y$ OCs

Thus, there are $\binom{4}{y} \cdot \binom{3}{2-y}$ ways to select y CCs and $2 - y$ OCs. ($y = 0, 1, 2$)

$$\begin{aligned}
 \text{(3)} \quad p_Y(0) &= \frac{\binom{4}{0} \cdot \binom{3}{2}}{21} = \frac{3}{3 \cdot 7} = \frac{1}{7}, \\
 p_Y(1) &= \frac{\binom{4}{1} \cdot \binom{3}{1}}{21} = \frac{4 \cdot 3}{3 \cdot 7} = \frac{4}{7}, \\
 p_Y(2) &= \frac{\binom{4}{2} \cdot \binom{3}{0}}{21} = \frac{(4 \cdot 3)/2}{3 \cdot 7} = \frac{2}{7}. \quad \square
 \end{aligned}$$

Whereas a PMF is defined for any discrete random element Y , the next definition needs that the values of Y are numbers.

Definition 6.2 (WMS Ch.03.2, Definition 3.4).

Let Y be a discrete random variable with probability mass function $p_Y(y)$. Then

$$E[Y] := \sum_y y p_Y(y) = \sum_y y P\{Y = y\},$$

is called the **expected value**, also **expectation** or **mean** of Y . \square

Remark 6.2.

A strict definition of $E[Y]$ would explicitly require that the sum $\sum_y y \cdot p_Y(y)$ is absolutely convergent, i.e.,

$$\sum_y |y| p_Y(y) < \infty.$$

The reason: Only absolute convergence of a series guarantees that its value does not depend on the order in which the terms are added. As in WMS and according to Assumption 6.1 on p.91, we will quietly assume that absolute convergence is satisfied for all random variables for which the expected value is used. \square

Proposition 6.2. ★ Let A_1, A_2, \dots, A_n a list of mutually disjoint events in a probability space (Ω, P) . Let $y_1, y_2, \dots, y_n \in \mathbb{R}$. Then

$$(6.6) \quad E \left[\sum_{j=1}^n y_j 1_{A_j} \right] = \sum_{j=1}^n y_j P(A_j).$$

PROOF: Let $Y := \sum_{j=1}^n y_j 1_{A_j}$; let $A := \bigsqcup_{j=1}^n A_j$. We may assume that $A = \Omega$, since we can add the zero term $0 \cdot 1_{A^c}$ to Y if $A^c \neq \emptyset$.

We further may assume that all numbers y_1, \dots, y_n are distinct for the following reason. Assume for example, that $y_{n_1} = y_{n_2} = y_{n_k} = y'$ and that this is the complete list of indices n_j such that $y_{n_j} = y'$.

We define $A' := A_{n_1} \uplus A_{n_2} \uplus \dots \uplus A_{n_k}$. Since

$$\sum_{j=1}^k y_{n_j} 1_{A_{n_j}} = \sum_{j=1}^k y' \cdot 1_{A_{n_j}} = y' \sum_{j=1}^k 1_{A_{n_j}} = y' \cdot 1_{A_{n_1} \uplus \dots \uplus A_{n_k}} = y' \cdot 1_{A'},$$

we can replace those terms with duplicate y' -values with the single term $y' \cdot 1_{A'}$.

We repeat this procedure with all y -values, even if they occur even once. This way we can write

$$(6.7) \quad Y = \sum_{j=1}^m y'_j 1_{A'_j}, \quad \text{where } \Omega = \biguplus_{i=1}^m A'_i \text{ and all } y'_i \text{ are distinct.}$$

In such a representation of Y , the distinctness of the y'_i implies that

$$Y(\omega) = y'_i \Leftrightarrow \omega \in A'_i \Leftrightarrow \{Y = y'_i\} = A'_i.$$

In particular, $P\{Y = y'_i\} = P(A'_i)$. Thus,

$$(6.8) \quad E[Y] = E \left[\sum_{i=1}^m y'_i 1_{A'_i} \right] = \sum_y y' P\{Y = y'\} = \sum_{i=1}^m y'_i P\{Y = y'_i\} = \sum_{i=1}^m y'_i P(A'_i).$$

In the last step of the proof we bring back the duplicate y -values. As above, we assume that

$y_{n_1} = y_{n_2} = \dots = y_{n_k} = y'_i$ and $A'_i := A_{n_1} \uplus A_{n_2} \uplus \dots \uplus A_{n_k}$. Then

$$y'_i P(A'_i) = y'_i P \left(\biguplus_{j=1}^k A_{n_j} \right) = y'_i \sum_{j=1}^k P(A_{n_j}) = \sum_{j=1}^k y_{n_j} P(A_{n_j}).$$

We substitute this result in (6.8) and obtain $E[Y] = \sum_{i=1}^m \sum_{j=1}^k y_{n_j} P(A_{n_j})$.

Since $\sum_{i=1}^m$ is the summation over all complete groups of equal y -values and each $\sum_{j=1}^k$ sums over all items in that group, that double sum equals $\sum_{j=1}^n y_j P(A_{n_j})$. Thus, $E[Y] = \sum_{j=1}^n y_j P(A_{n_j})$.

This proves the proposition. ■

Theorem 6.2.

Let Y be a discrete random variable and $g : \mathbb{R} \rightarrow \mathbb{R}$; $y \mapsto g(y)$ be a real-valued function. Then the random variable $g \circ Y : \omega \mapsto g(Y(\omega))$ has the following expected value:

$$(6.9) \quad E[g(Y)] = \sum_{\text{all } y} g(y) p_Y(y) = \sum_{\text{all } y} g(y) P\{Y = y\}.$$

PROOF: We give the proof assuming that Y takes only finitely many distinct values y_1, y_2, \dots, y_n .³⁴

³⁴As an aside, note that $y \mapsto g(y)$ need not be defined for all $y \in \mathbb{R}$. It suffices that the domain of g contains $Y(\Omega) = \{Y(\omega) : \omega \in \Omega\}$. (The range of the function Y ; see Definition 2.14 on p.27.)

Let $\{z_1, z_2, \dots, z_m\}$ denote the set of all distinct function values $g(y_i), i = 1, \dots, n$. In general, $m \leq n$ rather than $m = n$, because it is possible for one or more of the arguments y to have the same function value $g(y)$.

For $j = 1, \dots, m$, let

$$I_j := \{i \in [1, n] : g(y_i) = z_j\}$$

denote the set of all those indices i such that g assigns y_i to the same function value z_j . Note that

- (1) each I_j contains at least one index.
- (2) The index sets I_j form a partition of the indices i for the arguments y_i of g :

$$(A) \quad [1, n] = I_1 \uplus I_2 \uplus \dots \uplus I_m.$$

For $i = 1, \dots, n$ and $j = 1, \dots, m$, let

$$(B) \quad B_i := \{Y = y_i\} = \{\omega \in \Omega : Y(\omega) = y_i\}; \quad C_j := \{Z = z_j\} = \{\omega \in \Omega : Z(\omega) = z_j\}.$$

Since $\omega \in C_j \Leftrightarrow Z(\omega) = z_j \stackrel{(B)}{\Leftrightarrow} Y(\omega) = y_i$ for some $i \in I_j \Leftrightarrow \omega \in \biguplus_{i \in I_j} B_i$, it follows that

$$(C) \quad C_j = \biguplus_{i \in I_j} B_i.$$

We have for Y and Z the representations

$$(D) \quad Z(\omega) = \sum_{j=1}^m z_j 1_{\{Z=z_j\}}(\omega) = \sum_{j=1}^m z_j 1_{C_j}(\omega) \stackrel{(C)}{=} \sum_{j=1}^m z_j 1_{\biguplus_{i \in I_j} B_i}(\omega) = \sum_{j=1}^m z_j \sum_{i \in I_j} 1_{B_i}(\omega).$$

Here the last equation holds because the indicator function of a disjoint union is the sum of the indicator functions. That is a triviality which has been noted in (3.43) on p.65.

Since $g(y_i) = \text{const} = z_j$ for all $i \in I_j$, we can rewrite that last sum as

$$(E) \quad Z(\omega) = \sum_{j=1}^m \sum_{i \in I_j} g(y_i) 1_{B_i}(\omega) \stackrel{(A)}{=} \sum_{i=1}^n g(y_i) 1_{B_i}(\omega).$$

We conclude from (D) and (E) that $E[Z] = E\left[\sum_{i=1}^n g(y_i) 1_{B_i}\right]$.

Finally, we apply Proposition 6.2 on p.93 and obtain, since $B_i = \{Y = y_i\}$,

$$E[Z] = \sum_{i=1}^n g(y_i) P(B_i) = \sum_{i=1}^n g(y_i) P\{Y = y_i\}. \blacksquare$$

The following corresponds to WMS Theorems 3.4 and 3.5.

Theorem 6.3.

Let $c \in \mathbb{R}$, Y be a discrete random variable and $g_1, g_2, g_n : \mathbb{R} \rightarrow \mathbb{R}$ be a list of n real-valued functions. Then

$$(6.10) \quad E[c] = c \quad \text{and} \quad E[cY] = cE[Y],$$

$$(6.11) \quad E[cg_j(Y)] = cE[g_j(Y)].$$

Further, the random variable

$$\sum_{j=1}^n g_j \circ Y : \Omega \rightarrow \mathbb{R}; \quad \omega \mapsto \sum_{j=1}^n g_j(Y(\omega))$$

has the following expected value:

$$(6.12) \quad E \left[\sum_{j=1}^n g_j \circ Y \right] = \sum_{j=1}^n E[g_j \circ Y].$$

PROOF: Let Z denote the random variable $Z = c : \omega \mapsto c$, then

$$P\{Z = z\} = \begin{cases} 1, & \text{if } z = c, \\ 0, & \text{if } z \neq c. \end{cases}$$

Thus, $E[Z] = \sum_{z: P_Z\{z\} > 0} z \cdot P_Z\{z\} = c \cdot 1 = c$. This proves the first half of (6.10).

For the proof of the second half, note that $c = 0$ implies $cY = 0$. Thus, $E[cY] = cE[Y]$ becomes $E[0] = 0$, and we covered that case already. So we may assume that $c \neq 0$.

Let $Y' := cY$ and $y' := cy$. Then $Y(\omega) = y \Leftrightarrow Y'(\omega) = y'$. Thus, $P\{Y' = y'\} = P\{Y = \frac{y'}{c}\}$. Thus,

$$\begin{aligned} E[cY] &= E[Y'] = \sum_{y'} y' \cdot P\{Y' = y'\} = \sum_{y'} y' \cdot P\{Y = \frac{y'}{c}\} \\ &= \sum_y c \cdot \frac{y'}{c} \cdot P\{Y = \frac{y'}{c}\} = c \sum_y y \cdot P\{Y = y\} = c \cdot E[Y]. \end{aligned}$$

This proves the second half of (6.10). We apply this formula with $g_j(Y)$ in place of Y and (6.11) follows.

Finally, we apply Theorem 6.4 with Y_j in place of $g_j \circ Y$.³⁵ This results in (6.12). ■

The following cannot be found in the WMS text.

³⁵The proof of that theorem does not make use of this current one.

Theorem 6.4.

Let $Y_1, Y_2, \dots, Y_n : \Omega \rightarrow \mathbb{R}$ be discrete random variables which all are defined on the same probability space (Ω, P) ($n \in \mathbb{N}$). Then the random variable

$$\sum_{j=1}^n Y_j : \Omega \longrightarrow \mathbb{R}; \quad \omega \mapsto \sum_{j=1}^n Y_j(\omega)$$

has the following expected value:

$$(6.13) \quad E \left[\sum_{j=1}^n Y_j \right] = \sum_{j=1}^n E[Y_j].$$

In other words, the expectation of the sum is the sum of the expectations.

PROOF: ★ There are finite or infinite sequences $x_i, y_j \in \mathbb{R}$ as follows. Let $A_i := \{X = x_i\}$ and $B_j := \{Y = y_j\}$. Then the A_i are disjoint, the B_j are disjoint, and $A_* := (\biguplus_i A_i)^c, B_* := (\biguplus_j B_j)^c$ have probability zero. We may assume that $X(\omega) = 0$ for $\omega \in A_*$ and $Y(\omega) = 0$ for $\omega \in B_*$, since that does not change any assertions that are based on probabilities, such as the taking of expected values: Being able to discard the expressions $\biguplus_i A_i$ and $\biguplus_j B_j$ considerably simplifies the proof. For example, this assumption allows us to write, without having to exclude any $\omega \in \Omega$,

$$(A) \quad X(\omega) = \sum_i x_i 1_{\{X=x_i\}}(\omega), \quad Y(\omega) = \sum_j y_j 1_{\{Y=y_j\}}(\omega).$$

If $P\{X = 0\} > 0$, then we include 0 as one of the x_i and if $P\{Y = 0\} > 0$, then we include 0 as one of the y_j . We do so even though $0 \cdot 1_{\{X=0\}} = 0 \cdot 1_{\{Y=0\}} = 0$ contributes nothing to those sums, since then

$$A_i := \{X = x_i\}, \quad B_j := \{Y = y_j\}_j, \quad C_{i,j} := A_i \cap B_j$$

form partitions $\biguplus_i A_i = \biguplus_j B_j = \biguplus_{i,j} C_{i,j} = \Omega$ of Ω . Moreover, for each i, j ,

$$A_i = \biguplus_k C_{i,k} \quad \text{and} \quad B_j = \biguplus_k C_{k,j},$$

$$(B) \quad \text{which implies} \quad 1_{A_i} = \sum_k 1_{C_{i,k}} \quad \text{and} \quad 1_{B_j} = \sum_k 1_{C_{k,j}}.$$

$$\text{Since} \quad X \stackrel{(A)}{=} \sum_i x_i 1_{A_i} \stackrel{(B)}{=} \sum_{i,j} x_i 1_{C_{i,j}} \quad Y \stackrel{(A)}{=} \sum_j y_j 1_{B_j} \stackrel{(B)}{=} \sum_{i,j} y_j 1_{C_{i,j}}$$

$$\text{and thus,} \quad X + Y = \sum_{i,j} x_i 1_{C_{i,j}} + \sum_{i,j} y_j 1_{C_{i,j}} = \sum_{i,j} (x_i + y_j) 1_{C_{i,j}},$$

it follows from Prop.6.2 on p.93, that

$$(C) \quad E[X] = \sum_{i,j} x_i P(C_{i,j}), \quad E[Y] = \sum_{i,j} y_j P(C_{i,j}), \quad E[X + Y] = \sum_{i,j} (x_i + y_j) P(C_{i,j}).$$

We conclude the proof as follows:

$$E[X + Y] \stackrel{\text{(C)}}{=} \sum_{i,j} (x_i + y_j)P(C_{i,j}) = \sum_{i,j} x_i P(C_{i,j}) + \sum_{i,j} y_j P(C_{i,j}) \stackrel{\text{(C)}}{=} E[X] + E[Y]. \blacksquare$$

Remark 6.3.

- (1) The last theorem encompasses all variants of Theorem 6.3. For example, (6.12) follows with $Y_j = g_j \circ Y$.
- (2) The reason that many texts on an undergraduate probability theory do not list this theorem is that the proof, though elementary, is very tedious and requires working with the PMF of the random element $\vec{Y} = (Y_1, \dots, Y_n)$, given by

$$p_{\vec{Y}}(\vec{y}) = P\{Y_1 = y_1, \dots, Y_n = y_n\} \quad \square$$

Variance and standard deviation of a random variable indicate how strongly its distribution is concentrated around its expected value.

Definition 6.3 (Variance and standard deviation of a random variable).

Y be a random variable. The **variance** of Y is defined as the expected value of $(Y - E[Y])^2$. In other words,

$$(6.14) \quad \text{Var}[Y] := \sigma_Y^2 := E[(Y - E[Y])^2].$$

We call $SD(Y) := \sigma_Y := \sqrt{\text{Var}[Y]}$ The **standard deviation** of Y . \square

Theorem 6.5.

If Y is a discrete random variable, then

$$\text{Var}[Y] = E[Y^2] - (E[Y])^2.$$

PROOF:

$$\begin{aligned} \text{Var}[Y] &= E[(Y - E[Y])^2] = E(Y^2 - (2E[Y])Y + (E[Y])^2) \\ &= E(Y^2) - 2E[Y]E[Y] + (E[Y])^2 = E(Y^2) - (E[Y])^2. \blacksquare \end{aligned}$$

Theorem 6.6.

Let Y be a discrete random variable and $a, b \in \mathbb{R}$. Then

$$(6.15) \quad \text{Var}[aY + b] = a^2 \text{Var}[Y].$$

In other words, shifting a random variable by b , leaves its variance unchanged and multiplying it by a constant multiplies its variance by the square of that constant.

PROOF: We prove this by first showing that, for random variables Y and Y' ,

$$\text{Var}[aY] = a^2\text{Var}[Y] \quad \text{and} \quad \text{Var}[Y' + b] = \text{Var}[Y']$$

The assertion then follows from replacing Y' with aY .

We obtain from (6.10) that

$$\text{Var}[aY] = E[a^2Y^2] - (E[aY])^2 = a^2E[Y^2] - (aE[Y])^2 = a^2(E[Y^2] - (E[Y])^2) = a^2\text{Var}[Y].$$

To prove that $\text{Var}[Y' + b] = \text{Var}[Y']$, we observe that for any random variable Z and constant c , $E[Z + a] = E[Z] + E[a] = E[Z]$. Thus,

$$\begin{aligned} \text{Var}[Y' + b] &= E\left[\left((Y' + b) - E[Y' + b]\right)^2\right] \\ &= E\left[\left((Y' + b) - (E[Y'] + b)\right)^2\right] = E\left[\left(Y' - E[Y']\right)^2\right] = \text{Var}[Y']. \quad \blacksquare \end{aligned}$$

Remark 6.4. Since $\sqrt{a^2} = -a$ for negative numbers a ,

$$(6.16) \quad \sigma(aY) = |a|\sigma(Y). \quad \square$$

The following cannot be found in the WMS text.

Theorem 6.7 (Bienaymé formula).

Let $Y_1, Y_2, \dots, Y_n : \Omega \rightarrow \mathbb{R}$ be independent discrete random variables which all are defined on the same probability space (Ω, P) ($n \in \mathbb{N}$). Here we take the naive definition of independence: The outcomes of any Y_k are not influenced by the outcomes of the other Y_j . We will give a formulation of independence in terms of probabilities in a later chapter. Then

$$(6.17) \quad \text{Var}\left[\sum_{j=1}^n Y_j\right] = \sum_{j=1}^n \text{Var}[Y_j].$$

In other words, for independent random variables, the variance of the sum is the sum of the variances.

PROOF: Not given here. \blacksquare

Remark 6.5. The independence is necessary, otherwise there are counterexamples:

If $Y_1 = Y_2 = Y$ for some random variable Y , then

$$\text{Var}[Y + Y] = \text{Var}[2Y] = 4\text{Var}[Y] \neq \text{Var}[Y] + \text{Var}[Y]. \quad \square$$

6.2 Bernoulli Variables and the Binomial Distribution

Definition 6.4 (iid sequences).

Let $X_1, X_2, \dots (\Omega, P) \rightarrow \Omega'$ be a sequence of random elements. We speak of an **independent and identically distributed sequence**, in short, an **iid sequence** of random elements, if

- (1) the X_j are independent. Here we take the naive definition of independence: The outcomes of any X_k are not influenced by the outcomes of the other X_j . We will give a formulation of independence in terms of probabilities in a later chapter.
- (2) All random elements have the same distribution:
 - $P_{X_1}(B) = P_{X_2}(B) = P_{X_3}(B) = \dots$ for all j and all $B \subseteq \Omega'$.
 - Note that this can also be written
 $P\{X_1 \in B\} = P\{X_2 \in B\} = P\{X_3 \in B\} = \dots$ for all j and all $B \subseteq \Omega'$.
 - If the X_j are discrete random elements, identical distribution translates to identical PMFs $p_{X_1}(x) = p_{X_2}(x) = p_{X_3}(x) = \dots$ for all j and all $x \in \Omega'$. \square

Definition 6.5 (Bernoulli items and variables).

Let X be a binary random element on a probability space (Ω, P) , i.e., a random element which only assumes two outcomes, such as

- S (success) or F (failure)
- T (true) or F (false)
- Y (Yes) or N (No)
- 1 or 0

We call X a **Bernoulli random element**, or a **Bernoulli trial**.

- We call $p := P\{X = \text{success}\}$ the **success probability** and $q := 1 - p$, i.e., $q = P\{X = \text{failure}\}$, the **failure probability** of the Bernoulli trial.
- If a Bernoulli trial X has outcomes 1 and 0, then we call X a **Bernoulli variable** or a **0–1 encoded Bernoulli trial**.
- We call an iid sequence of Bernoulli trials a **Bernoulli sequence**. \square

Remark 6.6.

(a) The entire distribution of a Bernoulli trial is determined by the value of its success probability.

(b) Note that the definition of a Bernoulli sequence $(X_j)_j$ implies that

- (1) the X_j are independent
- (2) each X_j has the same success and failure probabilities. We write p and q for those numbers.

(c) Unless stated otherwise, we interpret the value 0 of a 0–1 encoded Bernoulli trial as failure and the value 1 as success. \square

Theorem 6.8 (Expected value and variance of a 0–1 encoded Bernoulli trial).

Let X be a 0–1 encoded Bernoulli trial with $p := P\{X = 1\}$. Then

$$(6.18) \quad E[X] = p \quad \text{and} \quad \text{Var}[X] = pq.$$

PROOF:

$$E[X] = 0q + 1 \cdot p = p.$$

For the variance, $Var[X] = E[X^2] - (E[X])^2 = E[X^2] - p^2$. Further,

$$E[X^2] = 0^2 \cdot q + 1^2 \cdot p = p.$$

Hence, $Var[X] = p - p^2 = p(1 - p) = pq$. ■

Definition 6.6 (Binomial Distribution).

Let $n \in \mathbb{N}$ and $0 \leq p \leq 1$. Let Y be a random variable with probability mass function

$$(6.19) \quad p_Y(y) = \binom{n}{y} p^y q^{n-y}.$$

Then we say that Y has a **binomial distribution**, with parameters n and p or, in short, a **binom(n, p) distribution**. We also say that Y is binom(n, p). □

Remark 6.7. How does one see that p_Y of (6.19) satisfies $p_Y(y) \geq 0$ for all y and $\sum_y p_Y(y) = 1$, i.e., it really is a probability mass function?

- $p_Y(y) \geq 0$ is true, since $p, q, \binom{n}{y} \geq 0$.
- We apply the binomial theorem (see Theorem 4.5) to $(p + q)^n$ and obtain

$$1 = 1^n = (p + q)^n = \sum_{j=0}^n \binom{n}{j} p^j q^{n-j}. \quad \square$$

Theorem 6.9.

Let X_1, X_2, X_n be a Bernoulli sequence of size n with success probability p . Let Y be the number of successes in that sequence, i.e., $Y(\omega) =$ number of indices j such that $X_j(\omega) = S$.

- Then Y is binom(n, p).

PROOF: Clearly,

$$Y(\omega) = y \Leftrightarrow \begin{cases} X_j(\omega) = S & \text{for } y \text{ indices } j, \\ X_j(\omega) = F & \text{for } n - y \text{ indices } j. \end{cases}$$

Let $\vec{x} := (x_1, \dots, x_n)$ a vector that consists of y components S and $n - y$ components F . For such an arrangement \vec{x} of y successes and $n - y$ failures, let n_1, n_2, n_y denote the indices for which $X_{n_j} = S$ and m_1, m_2, m_{n-y} those indices for which $X_{m_j} = F$. Further, let $A(\vec{x})$ denote the event

$$A(\vec{x}) := \{X_1 = x_1, X_2 = x_2, \dots, X_n = x_n\}.$$

Then independence of the Bernoulli trials X_j and thus, of the events $\{X_j = x_j\}$, yields

$$\begin{aligned} P(A(\vec{x})) &= P(\{X_1 = x_1\} \cap \dots \cap \{X_n = x_n\}) = P\{X_1 = x_1\} \cdot P\{X_2 = x_2\} \cdots P\{X_n = x_n\} \\ \text{(A)} \quad &= P\{X_{n_1} = S\} \cdots P\{X_{n_y} = S\} \cdot P\{X_{m_1} = F\} \cdots P\{X_{m_{n-y}} = F\} = p_y \cdot q^{n-y}. \end{aligned}$$

There are as many different vectors \vec{x} with y successes and $n - y$ failures as there are ways to form different lists of size n consisting of y items S and $n - y$ items F . That number is $\binom{n}{y}$.

We observe that the events $A(\vec{x})$ and $A(\vec{x}')$ are disjoint for different \vec{x} and \vec{x}' , since this means that there is at least one index j such that either $x_j = S$ and $x'_j = F$ or the other way around.

Let us assume that $x_j = S$ and $x'_j = F$. If $\omega \in A(\vec{x})$, then $X_j(\omega) = S$. But then $\omega \notin A(\vec{x}')$, since then $X_j(\omega)$ would have to be F . Thus, $A(\vec{x}) \cap A(\vec{x}') = \emptyset$. The case that $x_j = F$ and $x'_j = S$ is handled in the same fashion. Since

$$\{Y = y\} = \bigcup_{\vec{x}} A(\vec{x})$$

where \vec{x} assumes all $\binom{n}{y}$ arrangements of y successes and $n - y$ failures, it follows that

$$P\{Y = y\} = \sum_{\vec{x}} A(\vec{x})P(A(\vec{x})) \stackrel{(A)}{=} \binom{n}{y} p^y q^{n-y}.$$

This last expression equals the PMF of a binom(n, p) distribution and this concludes the proof. ■

Theorem 6.10 (Expected value and variance of a binom(n, p) variable).

Let Y be a binom(n, p) variable. Then

$$(6.20) \quad E[Y] = np \quad \text{and} \quad \text{Var}[Y] = npq.$$

PROOF: Let X_1, \dots, X_n be an iid list of 0–1 encoded Bernoulli trials with $p := P\{X = 1\}$. Let $Y' := \sum_{j=1}^n X_j$. according to Theorem 6.8, Theorem 6.4 on p.96, and, since the X_j are independent, Theorem 6.7 (Bienaymé formula) on p.99,

$$E[Y'] = \sum_{j=1}^n E[X_j] = np \quad \text{and} \quad \text{Var}[Y'] = \sum_{j=1}^n \text{Var}[X_j] = npq.$$

Further, $Y' = y \Leftrightarrow$ exactly y of the X_j have outcome y . Thus, Y' denotes the number of successes of those Bernoulli trials. According to Theorem 6.9 on p.101, Y' has a binom(n, p) distribution.

Since expected value and variance of a discrete random variable are determined by its PMF,

$$E[Y] = E[Y'] = np \quad \text{and} \quad \text{Var}[Y] = \text{Var}[Y'] = npq. \quad \blacksquare$$

6.3 Geometric + Negative Binomial + Hypergeometric Distributions

Definition 6.7 (Geometric distribution).

A random variable Y is said to have a **geometric distribution** with parameter $0 \leq p \leq 1$ or, in short, a **geom(p) distribution**, if its probability mass functions is as follows:

$$(6.21) \quad p_Y(y) = q^{y-1} p, \quad \text{for } y = 1, 2, 3, \dots \quad \square$$

Theorem 6.11. Let $X_1, X_2, \dots : (\Omega, P) \rightarrow \{S, F\}$ be an infinite Bernoulli sequence with success probability $0 \leq p \leq 1$.

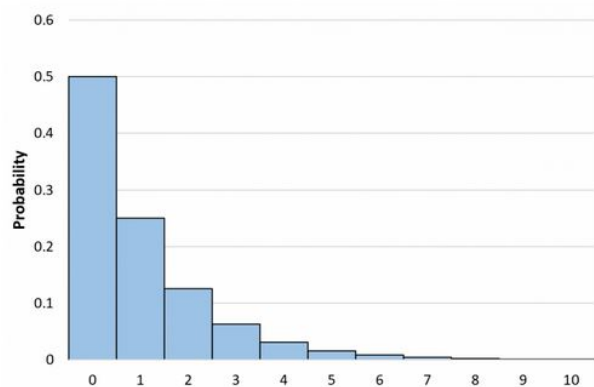
Let $T(\Omega, P) \rightarrow \mathbb{N}$ be the random variable

$$T(\omega) := \begin{cases} \text{smallest integer } k > 0 \text{ such that } X_k(\omega) = S \text{ if such a } k \text{ exists,} \\ \infty, & \text{else.} \end{cases}$$

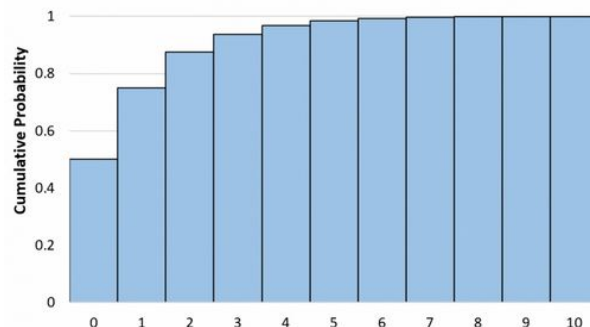
- Then T is $\text{geom}(p)$.

PROOF: Since $T(\omega) = n \Leftrightarrow X_1(\omega) = X_2(\omega) = \dots = X_{n-1}(\omega) = F$ and $X_n(\omega) = S$ and the independence of the X_i implies that the events $\{X_1 = F\}, \{X_2 = F\}, \dots, \{X_{n-1} = F\}, \{X_n = S\}$, are independent, we obtain

$$\begin{aligned} P\{X_1 = F, X_2 = F, \dots, X_{n-1} = F, X_n = S\} &= P\{X_1 = F\} \cap \dots \cap \{X_{n-1} = F\} \cap \{X_n = S\} \\ &= P\{X_1 = F\} \cdot P\{X_2 = F\} \cdot \dots \cdot P\{X_{n-1} = F\} \cdot P\{X_n = S\} = q^{n-1} p. \blacksquare \end{aligned}$$



6.1 (Figure). PMF for $\text{geom}(0.5)$.



6.2 (Figure). CDF for $\text{geom}(0.5)$.

Remark 6.8. In Theorem ?? we wrote $T(\omega)$ rather than the usual $Y(\omega)$ for the following reason. If we interpret the index j of the Bernoulli trial X_j as the point in time when the j th trial takes place, then $\omega \mapsto T(\omega)$ represents a **random time**, the time at which the first success happens. \square

Theorem 6.12 (WMS Ch.03.5, Theorem 3.8).

If Y is a $\text{geom}(p)$ random variable, then

$$E[Y] = \frac{1}{p}, \quad \text{and} \quad \text{Var}[Y] = \frac{q}{p^2}.$$

PROOF:

A: Expectation $E[Y]$:

One can obtain the derivative of the series $\sum_{y=1}^{\infty} q^y$ by differentiating it term-by-term. Since

$$\frac{d}{dq} q^y = yq^{y-1},$$

it follows that

$$(A) \quad \frac{d}{dq} \left(\sum_{y=1}^{\infty} q^y \right) = \sum_{y=1}^{\infty} yq^{y-1}.$$

We use (A) as follows.

$$\begin{aligned} E(Y) &= \sum_{y=1}^{\infty} yp_Y() = \sum_{y=1}^{\infty} yq^{y-1}p = p \sum_{y=1}^{\infty} yq^{y-1} \stackrel{(A)}{=} p \frac{d}{dq} \left(\sum_{y=1}^{\infty} q^y \right) \\ &= p \frac{d}{dq} \left(\frac{q}{1-q} \right) = p \frac{1 \cdot (1-q) - q(-1)}{(1-q)^2} = p \frac{1}{p^2} = \frac{1}{p}. \end{aligned}$$

B: Variance $Var[Y]$:³⁶

[6] Kargin, Vladislav: BU Lecture Notes for the Introduction to Probability Course

We compute the variance by again interchanging differentiation and summation. It follows from

$$\frac{d^2}{dq^2} q^y = y(y-1)q^{y-2},$$

that

$$(B) \quad \frac{d^2}{dq^2} \left(\sum_{y=1}^{\infty} q^y \right) = \sum_{y=2}^{\infty} y(y-1)q^{y-2} = \frac{1}{pq} \sum_{y=2}^{\infty} y(y-1)q^{y-1} \cdot p.$$

We use (B) as follows.

$$\begin{aligned} E[Y(Y-1)] &= \sum_{y=1}^{\infty} y(y-1)p_Y() = \sum_{y=2}^{\infty} y(y-1)q^{y-1}p = pq \sum_{y=2}^{\infty} y(y-1)q^{y-2} \\ &\stackrel{(B)}{=} pq \frac{d^2}{dq^2} \left(\sum_{y=2}^{\infty} q^y \right) = pq \cdot \frac{d^2}{dq^2} \left(\sum_{y=0}^{\infty} q^y \right) = pq \cdot \frac{d^2}{dq^2} \left(\frac{1}{1-q} \right) \\ &= pq \cdot \frac{d}{dq} \left(\frac{-1}{(1-q)^2} \right) = pq \cdot \frac{2}{p^3} = \frac{2q}{p^2}. \blacksquare \end{aligned}$$

Since $Var[Y] = E[Y^2] - (E[Y])^2 = E[Y(Y-1)] + E[Y] - \left(\frac{1}{p}\right)^2$, we conclude that

$$\begin{aligned} Var[Y] &= (E[Y^2] - E[Y]) - \left(\frac{1}{p}\right)^2 + E[Y] = E[Y(Y-1)] - \left(\frac{1}{p}\right)^2 + \frac{1}{p} \\ &= \frac{2q}{p^2} - \frac{1}{p^2} + \frac{p}{p^2} = \frac{2q - (1-p)}{p^2} = \frac{q}{p^2}. \blacksquare \end{aligned}$$

³⁶Source: [6] Kargin, Vladislav: BU Lecture Notes for the Introduction to Probability Course

Definition 6.8 (Negative binomial distribution). ★

A random variable Y has a **negative binomial distribution** with parameters p and r if

$$(6.22) \quad p_Y(y) = \binom{y-1}{r-1} p^r q^{y-r}, \quad \text{where } r \in \mathbb{N}, \quad y = r, r+1, r+2, \dots, \quad 0 \leq p \leq 1. \quad \square$$

This last definition has been marked as ★, so you are not expected to recall p_Y from memory. In contrast, the next theorem is NOT optional.

Theorem 6.13. Let $X_1, X_2, \dots : (\Omega, P) \rightarrow \{S, F\}$ be an infinite Bernoulli sequence with success probability $0 \leq p \leq 1$.

Let $t_1 < t_2 < \dots$ be the subsequence of those indices at which a success happens. In other words,

$$X_n(\omega) = \begin{cases} S = \text{success} & \text{if } n \text{ is one of } t_1, t_2, \dots, \\ F = \text{failure}, & \text{else.} \end{cases}$$

Two points to note:

- There will be different subsequences t_1, t_2, \dots for different arguments $\omega \in \Omega$. In other words, we are dealing with a sequence of random variables(!)

$$t_1 = T_1(\Omega), \quad t_2 = T_2(\Omega), \quad t_3 = T_3(\Omega), \quad \dots$$

- It is possible that we are dealing with an ω for which there are only 18 successes in the entire (infinite) sequence $X_1(\omega), X_2(\omega), \dots$. In this case, we define $T_{19}(\omega) = T_{20}(\omega) = \dots = \infty$. More generally, if $r \in \mathbb{N}$ and the sequence $X_1(\omega), X_2(\omega), \dots$ has less than r successes, we define

$$T_r(\omega) := \infty.$$

Now that we have defined $T_r = T_r(\omega)$, we are ready to state the theorem.

- The random variable T_r has a negative binomial distribution with parameters p and r .

PROOF: See the introductory remarks of WMS Chapter 3.6 before Definition 3.9. ■

Remark 6.9. If we think of the indices n of the sequence X_n as points in time, we can interpret the random variables T_1, T_2, \dots as follows.

- T_r is the time of the r th success in the underlying Bernoulli sequence X_n . □

Theorem 6.14. ★

If the random variable Y is negative binomial with parameters p and r ,

$$E[Y] = \frac{r}{p} \quad \text{and} \quad \text{Var}[Y] = \frac{r(1-p)}{p^2}.$$

PROOF: Not given here. ■

Definition 6.9 (Hypergeometric distribution).

A random variable Y has a **hypergeometric distribution** with parameters N , R and n if its PMF is

$$(6.23) \quad p_Y(y) = \frac{\binom{R}{y} \binom{N-R}{n-y}}{\binom{N}{n}},$$

where the nonnegative integers N , R , n and y are subject to the following conditions:

- $y \leq n$
- $y \leq R$
- $n - y \leq N - R$ □

Remark 6.10. For the following you should review Section 5.2 (Random Sampling and Urn Models With and Without Replacement).

The hypergeometric distribution provides the mathematical model for drawing SRS samples of size n from a population of size N where each item in that population is classified as either S (success) or F (failure).

In contrast to the scenarios involving the binomial, geometric and negative binomial distributions, those n picks X_1, X_2, \dots, X_n do NOT constitute a Bernoulli sequence since SRS sampling is sampling without replacement and the X_j will neither be independent nor have the same success probability across all j .

Rather, we must model this kind of sampling with an urn model without replacement. See Definition 5.4 (Urn models) on p.90. It simplifies matters greatly that we are only interested in success or failure of each sample pick, since this means that we can model our population as N well-mixed balls in an urn, of which R are labeled S and the remaining $N - R$ are labeled F . Picking the SRS sample of size n from the population then is modeled by picking a sample of size n without replacement from that urn. □

Theorem 6.15.

- Given is an urn which contains N well-mixed balls of two colors, Red and Black. We assume that R are Red and thus, the remaining $N - R$ are Black.
- A sample of size n is drawn without replacement from that urn, according to Definition 5.4(a).

Let the random variable Y denote the number of Red balls in that sample. Then Y is hypergeometric with parameters N , R and n . In other words, its PMF is

$$p_Y(y) = \frac{\binom{R}{y} \binom{N-R}{n-y}}{\binom{N}{n}}.$$

PROOF: We give here a very skeletal proof. For more detail consult WMS Chapter 3.7.

We are not interested in the order in which those Red balls were picked, so our probability space Ω will be that of all combinations of size n that can be selected from N balls. Thus,

$$|\Omega| = \binom{N}{n}.$$

$p_Y(y)$ is the probability of selecting exactly y Red balls in the sample of size n . Such a selection is obtained by partitioning the N balls into the heap of all R red balls, the heap of all $N - R$ Black balls and then proceeding as follows.

Conceptually we pick one of the $\binom{R}{y}$ possible selections of y items from the R red balls and then complementing it with one of the $\binom{N-R}{n-y}$ possible selections of the remaining $n - y$ items from the $N - R$ black balls. By Theorem 4.1 (multiplication rule of combinatorial analysis) on p.71, there are $\binom{R}{y} \cdot \binom{N-R}{n-y}$ such selections. It follows that

$$p_Y(y) = P\{Y = y\} = \frac{\binom{R}{y} \cdot \binom{N-R}{n-y}}{\binom{N}{n}}.$$

It follows that Y is hypergeometric with parameters N , R and n . ■

Theorem 6.16 (WMS Ch.03.7, Theorem 3.10). ★

Let Y be a hypergeometric random variable with parameters N , R and n . Then

$$(6.24) \quad E[Y] = \frac{nR}{N} \quad \text{and} \quad \text{Var}[Y] = n \left(\frac{R}{N} \right) \left(\frac{N-R}{N} \right) \left(\frac{N-n}{N-1} \right).$$

PROOF: We reproduce here the plausibility argument given by WMS in their “proof” of WMS Theorem 3.10.

Since we consider picking an R -item as a success, the above formulas read with $p := \frac{R}{N}$ and $q = 1 - p = \frac{N-R}{N}$ as follows:

$$E[Y] = n \cdot p \quad \text{and} \quad \text{Var}[Y] = n \cdot p \cdot q \left(\frac{N-n}{N-1} \right).$$

Except for the factor $(N - n)/(N - 1)$

those are expectation and variance of the binom($n, R/n$) distribution. Note for the

$$\text{correction factor } \frac{N - n}{N - 1}, \quad \text{that} \quad \lim_{N \rightarrow \infty} \frac{N - n}{N - 1} = 1.$$

This reflects the fact that, if N is huge in comparison to n , drawing from an urn with or without replacement yields, up to a rounding error, the same probabilities. ■

6.4 The Poisson Distribution

We start out with the simple observation that $e^x = \sum_{j=0}^{\infty} \frac{x^j}{j!}$ for any $x \in \mathbb{R}$.

Proposition 6.3. *Let $\lambda > 0$. Then the function $p(y) := e^{-\lambda} \frac{\lambda^y}{y!}$ defines a probability mass function on $[0, \infty[\mathbb{Z} = \{0, 1, 2, \dots\}$.*

PROOF: Obviously, $p(y) \geq 0$ for all y .

To show that $\sum_y p(y) = 1$, we apply the formula $e^x = \sum_{j=0}^{\infty} \frac{x^j}{j!}$, which is true for any $x \in \mathbb{R}$, with $x = \lambda$ and $j = y$. ■

This simple proposition enables us to make the following definition.

Definition 6.10 (Poisson variable).

Let Y be a random variable and $\lambda > 0$. We say that Y has a **Poisson probability distribution** with parameter λ , in short, Y is **poisson**(λ), if its probability mass function is

$$p_Y(y) = \frac{\lambda^y}{y!} e^{-\lambda}, \quad \text{for } y = 0, 1, 2, \dots, \quad \square$$

We follow WMS Chapter 3.8 to show what phenomena can be modeled by a Poisson variables

Proposition 6.4. *Given is some event of interest, E .*

- (1) *We define a random variable Y which counts how often E happen in a “unit”. We leave it open whether this unit is a time interval (maybe a minute or a year) or a subset of d -dimensional space ($d = 1, 2, 3$). Let us write A for that unit.*
 - *Example: Y is the number of car accidents that happen in Binghamton during a day (unit of time),*
 - *Example: Y is the number of typos on a randomly picked page of these lecture notes (“page” is a twodimensional unit – square inches).*
- (2) *Given $n \in \mathbb{N}$, we subdivide the unit (A) into n parts of equal size. Let*

$$\vec{X}^{(n)} := X_1^{(n)}, X_2^{(n)}, X_n^{(n)},$$

where $X_j^{(n)}$ = the number of times that E happens in subunit j .

- Assume that for all big enough, FIXED n ,
 - the $X_j^{(n)}$ are independent
 - for each j , $P\{X_j^{(n)} = 0 \text{ or } 1\} = 1$: E (i.e., the event of interest) happens at most once in such a small subunit
 - $p_n := P\{X_j^{(n)} = 1\}$ is constant in j ($j = 1, 2, \dots, n$)
 - $\lambda := n \cdot p_n$ is constant in n : For large enough k , $k p_k = (k+1)p_{k+1} = (k+2)p_{k+2} = \dots = \lambda$.

Given these assumptions, the following is true:

- The random variable $Y^{(n)} := X_1^{(n)} + X_2^{(n)} + \dots + X_n^{(n)}$ is $\text{binom}(n, p_n)$ for large n .
- The $\text{binom}(n, p_n)$ probability mass functions $p_{Y^{(n)}}$ converge to that of a $\text{poisson}(\lambda)$ variable:

$$(6.25) \quad \lim_{n \rightarrow \infty} p_{Y^{(n)}}(y) = \lim_{n \rightarrow \infty} \binom{n}{p} p^y (1-p)^{n-y} = e^{-\lambda} \cdot \frac{\lambda^y}{y!}, \quad \text{for } y = 0, 1, 2, \dots,$$

PROOF: We follow WMS:

Recall that $\lambda = np$. Thus,

$$\begin{aligned} \binom{n}{p} p^y (1-p)^{n-y} &= \frac{n(n-1) \cdots (n-y+1)}{y!} \left(\frac{\lambda}{n}\right)^y \left(1 - \frac{\lambda}{n}\right)^{n-y} \\ (\star) \quad &= \frac{\lambda^y}{y!} \left(1 - \frac{\lambda}{n}\right)^n \frac{n(n-1) \cdots (n-y+1)}{n^y} \left(1 - \frac{\lambda}{n}\right)^{-y} \\ &= \left(\frac{\lambda^y}{y!}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-y} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{y-1}{n}\right). \end{aligned}$$

From calculus we obtain $\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda}$. Further,

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-y} = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right) = \lim_{n \rightarrow \infty} \left(1 - \frac{2}{n}\right) = \dots = \lim_{n \rightarrow \infty} \left(1 - \frac{y-1}{n}\right) = 1.$$

We take limits in (\star) and obtain

$$\lim_{n \rightarrow \infty} \binom{n}{p} p^y (1-p)^{n-y} = \left(\frac{\lambda^y}{y!}\right) e^{-\lambda}. \quad \blacksquare$$

Theorem 6.17 (WMS Ch.03.8, Theorem 3.11).

A $\text{poisson}(\lambda)$ random variable has expectation and variance λ . In other words,

$$(6.26) \quad E[Y] = \text{Var}[Y] = \lambda.$$

A. PROOF of $E[Y] = \lambda$:

$$E(Y) = \sum_y y p_Y(y) = \sum_{y=0}^{\infty} y \frac{\lambda^y e^{-\lambda}}{y!} = \sum_{y=1}^{\infty} y \frac{\lambda^y e^{-\lambda}}{y!} = \lambda \sum_{y=1}^{\infty} \frac{\lambda^{y-1} e^{-\lambda}}{(y-1)!}.$$

In the last equation we used $y!/y = (y-1)!$. We write $k = y - 1$ for the index variable and obtain

$$E(Y) = \lambda \sum_{k=0}^{\infty} \frac{\lambda^k e^{-\lambda}}{k!} = \lambda \sum_{k=0}^{\infty} p(k),$$

where $p(k) = \frac{\lambda^k e^{-\lambda}}{k!}$ is the PMF of a poisson(λ) random variable. Thus, $\sum_{k=0}^{\infty} p(k) = 1$ and it follows that $E[Y] = \lambda$.

B. PROOF of $Var[Y] = \lambda$:

$$\text{Observe that } y^2 e^{-\lambda} \frac{\lambda^y}{y!} = e^{-\lambda} \cdot \frac{y^2 \lambda \lambda^{y-1}}{y!} = (\lambda e^{-\lambda}) \frac{y \lambda^{y-1}}{(y-1)!} = (\lambda e^{-\lambda}) \frac{1}{(y-1)!} \frac{d}{d\lambda} (\lambda^y)$$

Since we can interchange summation and differentiation, this yields

$$\begin{aligned} E[Y^2] &= \sum_{y=0}^{\infty} y^2 e^{-\lambda} \frac{\lambda^y}{y!} = \sum_{y=1}^{\infty} y^2 e^{-\lambda} \frac{\lambda^y}{y!} = (\lambda e^{-\lambda}) \sum_{y=1}^{\infty} \frac{d}{d\lambda} \left(\frac{\lambda \cdot \lambda^{y-1}}{(y-1)!} \right) \\ &= (\lambda e^{-\lambda}) \frac{d}{d\lambda} \left(\lambda \sum_{y=1}^{\infty} \frac{\lambda^{y-1}}{(y-1)!} \right) = (\lambda e^{-\lambda}) \frac{d}{d\lambda} \left(\lambda \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} \right). \end{aligned}$$

Since $\sum_{y=0}^{\infty} \frac{\lambda^y}{y!} = e^\lambda$, this implies $E[Y^2] = (\lambda e^{-\lambda}) \frac{d}{d\lambda} (\lambda e^\lambda) = \lambda e^{-\lambda} (e^\lambda + \lambda e^\lambda) = \lambda + \lambda^2$.

By part **A**, $E[Y] = \lambda$. Thus, $E[Y^2] = \lambda + (E[Y])^2$. We finally obtain

$$Var[Y] = E[Y^2] - (E[Y])^2 = (\lambda + (E[Y])^2) - (E[Y])^2 = \lambda. \blacksquare$$

We refer to the WMS text for examples of random variables with a Poisson distribution.

6.5 Moments, Central Moments and Moment Generating Functions

Unless something different is stated, Y is a random variable $Y : (\Omega, P) \rightarrow \mathbb{R}$ on some probability space (Ω, P) .

$$\mu = E[Y], \quad \sigma^2 = Var[Y], \quad \sigma = \sqrt{Var[Y]},$$

denote expectation, variance and standard deviation of Y .

Definition 6.11 (k th Moment).

If Y is a random variable and $k \in \mathbb{N}$,

$$(6.27) \quad \mu'_k := E[Y^k]$$

is called the k th **moment** of Y . μ'_k also is referred to as the k th **moment of Y about the origin**. \square

Note in particular that the first moment of Y is the expectation of Y and that

$$\mu'_2 = \text{Var}[Y] + E[Y]^2.$$

Another useful moment of a random variable is one taken about its mean.

Definition 6.12 (k th Central Moment).

If Y is a random variable and $k \in \mathbb{N}$,

$$(6.28) \quad \mu_k := E[(Y - E[Y])^k] = E[(Y - \mu)^k]$$

is called the k th **central moment** of Y . μ_k also is referred to as the k th **moment of Y about its mean**. \square

Proposition 6.5 (The moments determine the distribution). ★ Under fairly slight assumptions the following is true for two random variables Y_1 and Y_2 .

$$\text{If } E[Y_1^k] = E[Y_2^k] \text{ for } k = 1, 2, 3, \dots, \text{ then } P_{Y_1} = P_{Y_2}.$$

In other words, the distribution of a random variable is uniquely determined by its moments.

PROOF: Beyond the scope of these lecture notes. \blacksquare

Next we associate with a random variable Y which is a function $\omega \mapsto Y(\omega)$ a function $t \mapsto m_Y(t)$ of a real variable t . It allows us to generate all moments μ'_k of Y by computing its k th derivative at $t = 0$. Since $m_Y(t)$ determines in this way all moments of Y and since those in turn determine P_Y ,³⁷ $m_Y(t)$ uniquely determines the entire distribution of Y .

Definition 6.13 (Moment-generating function).

Let Y be a random variable for which one can find $\delta > 0$ (no matter how small), such that

$$(6.29) \quad m(t) := m_Y(t) := E[e^{tY}] \quad \text{is finite for } |t| < \delta.$$

Then we say that Y has **moment-generating function**, in short, **MGF**, $m_Y(t)$. \square

Theorem 6.18. The following is WMS Ch.03.9, Theorem 3.12.

Let Y be a random variable with MGF $m_Y(t)$ and $k \in \mathbb{N}$. Then its k th moment is obtained as the k th derivative of $m_Y(\cdot)$, evaluated at $t = 0$:

$$(6.30) \quad \mu'_k = m^{(k)}(0) = \left. \frac{d^k m(t)}{dt^k} \right|_{t=0}.$$

³⁷See Proposition 6.5

PROOF: We write $m(t)$ for $m_Y(t)$. From the series expansion $e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$, we obtain

$$\begin{aligned} m(t) &= E[e^{tY}] = E\left[\sum_{k=0}^{\infty} \frac{t^k Y^k}{k!}\right] = \sum_{k=0}^{\infty} \frac{t^k}{k!} E[Y^k] = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mu'_k \\ &= 1 + t\mu'_1 + \frac{t^2}{2!}\mu'_2 + \frac{t^3}{3!}\mu'_3 + \dots \end{aligned}$$

Taking derivatives repeatedly,

$$\begin{aligned} m^{(1)}(t) &= \mu'_1 + \frac{2t}{2!}\mu'_2 + \frac{3t^2}{3!}\mu'_3 + \dots & \Rightarrow m^{(1)}(0) &= \mu'_1 + 0 + 0 + \dots, \\ m^{(2)}(t) &= \mu'_2 + \frac{2t}{2!}\mu'_3 + \frac{3t^2}{3!}\mu'_4 + \dots & \Rightarrow m^{(2)}(0) &= \mu'_2 + 0 + 0 + \dots, \\ \dots & \dots & & \\ m^{(k)}(t) &= \mu'_k + \frac{2t}{2!}\mu'_{k+1} + \frac{3t^2}{3!}\mu'_{k+2} + \dots & \Rightarrow m^{(k)}(0) &= \mu'_k + 0 + 0 + \dots \end{aligned}$$

In summary,

$$m^{(1)}(0) = \mu'_1, \quad m^{(2)}(0) = \mu'_2, \quad \dots, \quad m^{(k)}(0) = \mu'_k. \quad \blacksquare$$

Technical note: The existence of the MGF of Y allowed us to compute the derivative of a series as the sum of the derivatives.

You find the next proposition as Example 3.23 in WMS Ch.3.9.

Proposition 6.6. ★ *If Y is a poisson(λ) random variable ($\lambda > 0$), its MGF is*

$$(6.31) \quad m_Y(t) = e^{\lambda(e^t-1)}.$$

PROOF: For this proof, we abbreviate **(A)** $\tilde{\lambda} := \lambda e^t$.

Note that the Taylor expansion $e^x = \sum_{j=0}^{\infty} \frac{x^j}{j!}$ yields, with x and j replaced by $\tilde{\lambda}$ and y ,

$$\mathbf{(B)} \quad e^{\tilde{\lambda}} = \sum_{y=0}^{\infty} \frac{\tilde{\lambda}^y}{y!}.$$

$$\begin{aligned} \text{Then, } m_Y(t) &= E(e^{tY}) = \sum_{y=0}^{\infty} e^{ty} p(y) = \sum_{y=0}^{\infty} e^{ty} \frac{\lambda^y e^{(-\lambda)}}{y!} \\ &= \sum_{y=0}^{\infty} (e^t)^y \lambda^y \frac{e^{-\lambda}}{y!} = \sum_{y=0}^{\infty} \frac{(\lambda e^t)^y e^{-\lambda}}{y!} \stackrel{\mathbf{(A)}}{=} e^{-\lambda} \sum_{y=0}^{\infty} \frac{\tilde{\lambda}^y}{y!} \\ &\stackrel{\mathbf{(B)}}{=} e^{-\lambda} e^{\tilde{\lambda}} \stackrel{\mathbf{(A)}}{=} e^{(-1)\lambda} e^{\lambda e^t} = e^{\lambda(-1+e^t)} = e^{\lambda(e^t-1)}. \quad \blacksquare \end{aligned}$$

- The subsection titled “The Tchebysheff Inequality” which was at this location has been integrated into subsection 7.8 (Inequalities for Probabilities)

6.6 Exercises for Ch.6

Exercise 6.1. If the random variable Y has expectation $E[Y] = -2$ and standard deviation $\sigma_Y = 2$, what is $E[(Y + 3)^2]$?

Answer: Since $E[Y^2] = \text{Var}[Y] + (E[Y])^2 = (\sigma_Y)^2 + (-2)^2 = 8$,

$$E[(Y + 3)^2] = E[Y^2] + 6E[Y] + 9 = 8 - 12 + 9 = \boxed{5} \blacksquare$$

Exercise 6.2. If the random variable Y has the PMF

$$p_Y(-2) = 0.13, p_Y(0) = 0.24, p_Y(1) = 0.18, p_Y(2) = 0.45,$$

- (a) compute $E[Y]$
- (b) compute $\text{Var}[Y]$
- (c) compute σ_Y

Answer (the numeric computations might have errors):

- (a) $E[Y] = \sum_y y \cdot p_Y(y) = (-2)(0.13) + 0(0.24) + 1(0.18) + 2(0.45) = 0.82$
- (b) $\text{Var}[Y] = \sum_y (y - E[Y])^2 \cdot p_Y(y)$
 $= (-2 - 0.82)^2(0.13) + (0 - 0.82)^2(0.24) + (1 - 0.82)^2(0.18) + (2 - 0.82)^2(0.45) = 1.8276$
- (c) $\sigma_Y = \sqrt{\text{Var}[Y]} = \sqrt{1.8276} \approx 1.3513888$

Exercise 6.3. Let Y be a 0–1 encoded Bernoulli variable with $P\{Y = 1\} = p$.

- (a) Compute its MGF
- (b) Use the MGF method to compute the n th moment about the origin, $E[X^n]$

Answer:

$$(a) \quad M_Y(t) = E[e^{tY}] = e^{0t} \cdot q + e^{1t} \cdot p = \boxed{q + pe^t}$$

(b) The derivatives of $M_Y(t)$ are

$$M'_Y(t) = (q + pe^t)' = pe^t, M''_Y(t) = (pe^t)' = pe^t, \dots, M_Y^{(n)}(t) = pe^t, \dots,$$

Thus, $E[Y^n] = \mu'_n = M_Y^{(n)}(0) = pe^0 = \boxed{p}$ for all n .

(c) We use the results of (b) to compute the variance:

$$\text{Var}[Y] = E[Y^2] - (E[Y])^2 = \mu'_2 - (\mu'_1)^2 = p - p^2 = (1 - p)p = \boxed{pq} \blacksquare$$

Exercise 6.4. Let Y be a binom(n, p) variable. Use the MGF method to verify that $E[Y] = np$ and $\text{Var}[Y] = npq$.

Answer: Since the PMF of Y is $p_Y(y) = \binom{n}{y} p^y q^{n-y}$,

$$\begin{aligned} M_Y(t) &= E[e^{tY}] = \sum_{y=0}^n e^{ty} \binom{n}{y} p^y q^{n-y} = \sum_{y=0}^n \binom{n}{y} (e^t)^y p^y q^{n-y} \\ &= \sum_{y=0}^n \binom{n}{y} (pe^t)^y q^{n-y} = (pe^t + q)^n \end{aligned}$$

Here we obtained the last equation by applying the binomial theorem,

$$(a + b)^n = \sum_{j=0}^n \binom{n}{j} a^j b^{n-j},$$

with $a = pe^t$ and $b = q$.

$$\begin{aligned} M_Y(t)' &= npe^t(pe^t + q)^{n-1}, \\ M_Y(t)'' &= npe^t(pe^t + q)^{n-1} + n(n-1)(pe^t)^2(pe^t + q)^{n-2}. \end{aligned}$$

Thus,

$$\begin{aligned} E[Y] &= M_Y(0)' = np, \\ E[Y^2] &= M_Y(0)'' = np + n(n-1)p^2. \end{aligned}$$

It follows that

$$\begin{aligned} \text{Var}[Y] &= E[Y^2] - (E[Y])^2 = \\ &= E[Y^2] - M_Y(0)'' = np + n(n-1)p^2 - n^2p^2 = npq. \blacksquare \end{aligned}$$

7 Continuous Random Variables

7.1 Cumulative Distribution Function of a Random Variable

The material found in this section does not make any references to continuous random variables.

Definition 7.1 (Cumulative Distribution Function).

Let Y denote any random variable (it need not be discrete). The **distribution function** of Y , also called its **cumulative distribution function** or **CDF (cumulative distribution function)**, is defined as follows.

$$(7.1) \quad F(y) := F_Y(y) := P\{Y \leq y\} \quad \text{for } y \in \mathbb{R}. \quad \square$$

Example 7.1. Let Y be a binom(2, 1/4) random variable, i.e., $n = 2$ and $p = 1/4$. Compute $F_Y(y)$.

Solution: The probability mass function for Y is

$$p_Y(y) = \binom{2}{y} \left(\frac{1}{4}\right)^y \left(\frac{3}{4}\right)^{2-y}.$$

Thus,

$$p_Y(0) = \frac{1}{16}, \quad p_Y(1) = 2 \left(\frac{1}{4}\right) \left(\frac{3}{4}\right) = \frac{6}{16}, \quad p_Y(2) = \frac{9}{16}.$$

It follows that

- $y < 0 \Rightarrow F_Y(y) = P_Y(\emptyset) = 0.$
- $0 \leq y < 1 \Rightarrow F_Y(y) = p_Y(0) = 1/16.$
- $1 \leq y < 2 \Rightarrow F_Y(y) = p_Y(0) + p_Y(1) = 7/16.$
- $y \geq 2 \Rightarrow F_Y(y) = p_Y(0) + p_Y(1) + p_Y(2) = 1.$

Note that F_Y is constant on intervals A of \mathbb{R} if $p_Y(a) = 0$ for all $a \in A$. \square

Theorem 7.1 (Properties of a Cumulative Distribution Function).

If $F_Y(y)$ is the cumulative distribution function of a random variable Y , then

- (1) $F_Y(-\infty) = \lim_{y \rightarrow -\infty} P(Y \leq y) = 0.$
- (2) $F_Y(\infty) = \lim_{y \rightarrow \infty} P(Y \leq y) = 1.$
- (3) $F_Y(y)$ is a nondecreasing function of y . In other words, if $y_1 < y_2$, then $F_Y(y_1) \leq F_Y(y_2)$. See Definition 2.20 on p.30.

PROOF:

The proof of (1) and (2) follows from

It follows from $-\infty < Y(\omega) < \infty$ that

$$\bigcap_{y \in \mathbb{R}} \{Y \leq y\} = \bigcap_{n \in \mathbb{N}} \{Y \leq -n\} = \emptyset$$

$$\bigcup_{y \in \mathbb{R}} \{Y \leq y\} = \bigcup_{n \in \mathbb{N}} \{Y \leq n\} = \Omega$$

We apply Theorem 3.1 (Continuity property of probability measures) on p.47 and obtain

$$F_Y(-\infty) = \lim_{n \rightarrow \infty} P\left(\bigcap_{y \in \mathbb{R}} \{Y \leq y\}\right) = P(\emptyset) = 0,$$

$$F_Y(\infty) = \lim_{n \rightarrow \infty} P\left(\bigcup_{y \in \mathbb{R}} \{Y \leq y\}\right) = P(\Omega) = 1.$$

Obvious from $P \geq 0$ and $y_1 < y_2 \Rightarrow \{Y \leq y_2\} = \{Y \leq y_1\} \uplus \{y_1 < Y \leq y_2\}$, since this implies

$$F(y_2) = P\{Y \leq y_2\} = P\{Y \leq y_1\} + P\{y_1 < Y \leq y_2\} \geq P\{Y \leq y_1\} = F(y_1). \blacksquare$$

Remark 7.1. ★ There is a fourth property that is satisfied by all CDFs:

$y \mapsto F_Y(y)$ is **right continuous** at all arguments y .

This means the following. if y is approached from the right by a sequence y_n such as $y_n = y + \frac{1}{n}$ or $y_n = y(1 + e^{-n})$, then

$$\lim_{n \rightarrow \infty} F(y_n) = F(y). \square$$

7.2 Continuous Random Variables and Probability Density Functions

Definition 7.2 (Continuous random variable).

We call a random variable Y with distribution function $F_Y(y)$ **continuous**, if $F_Y(y)$ is continuous, for all arguments y . \square

Proposition 7.1. Let Y be a continuous random variable with CDF $F_Y(y)$. Then its distribution gives zero probability to all singletons $\{a\}$ ($a \in \mathbb{R}$). Also, it gives the same probability to an interval with endpoints $-\infty < a < b < \infty$, regardless whether a and/or b do or do not belong to that interval. In other words,

$$(7.2) \quad a \in \mathbb{R} \Rightarrow P\{Y = a\} = P_Y\{a\} = 0,$$

$$(7.3) \quad -\infty < a < b < \infty \Rightarrow P\{a < Y < b\} = P\{a \leq Y < b\}$$

$$= P\{a < Y \leq b\} = P\{a \leq Y \leq b\}.$$

PROOF: Since $\{a\} \subseteq]a - \frac{1}{n}, a]$ and $]a - \frac{1}{n}, a] =]-\infty, a] \setminus]-\infty, a - \frac{1}{n}]$ (set difference),

$$P\{Y = a\} \leq P\{a - \frac{1}{n} < Y \leq a\} = P\{Y \leq a\} - P\{Y \leq a - \frac{1}{n}\} = F_Y(a) - F_Y\left(a - \frac{1}{n}\right).$$

F_Y is continuous at a , in particular, F_Y is continuous from the left at a . Thus,

$$\lim_{n \rightarrow \infty} F_Y\left(a - \frac{1}{n}\right) = F_Y(a).$$

It follows that $P\{Y = a\} = F_Y(a) - F_Y(a) = 0$. This proves (7.2).

This result, plus additivity of probability measures, plus

$$[a, b] =]a, b[\cup \{a\} \cup \{b\}, \quad [a, b] = [a, b[\cup \{b\}, \quad [a, b] =]a, b] \cup \{a\},$$

show that (7.3) holds. ■

A lot more can be done with a CDF that is not only continuous but has a continuous derivative. We make the following blanket assumption.

Assumption 7.1 (All continuous random variables have a differentiable CDF). Unless explicitly stated otherwise, all continuous random variables are assumed to satisfy the following:

The first derivative $\frac{dF_Y}{dy}$ of F_Y exists and is continuous except for, at most, a finite number of points in any finite interval.

All cumulative distribution functions for continuous random variables that we deal with in this course satisfy this assumption. □

This last assumption allows us to make the following definition.

Definition 7.3 (Probability density function).

Let Y be a continuous random variable with CDF $F_Y(y)$. For all arguments y where the derivative $F'_Y(y) = \frac{dF_Y(y)}{dy}$ exists, we define

$$f(y) := f_Y(y) := \frac{dF_Y(y)}{dy}.$$

We call f_Y the **probability density function** or, in short, the **PDF** of the continuous random variable Y . □

Theorem 7.2.

Let Y be a continuous random variable with CDF $F_Y(y)$ and PDF $f_Y(y)$.

(1) If $a, b \in \mathbb{R}$ and $a < b$, then

$$(7.4) \quad P\{a < Y \leq b\} = F_Y(b) - F_Y(a) = \int_a^b f(y) dy.$$

(2) $f_Y(y) \geq 0$ for $-\infty < y < \infty$.

(3) $\int_{-\infty}^{\infty} f_Y(y) dy = 1$.

PROOF: (1) is the fundamental theorem of calculus. Of course, we interpret $\int_a^b f(y)dy$ as follows. Assume that some of the points y at which $f'_Y(y)$ does not exist fall within the interval $[a, b]$. Our assumption guarantee that there are only finitely many such y , say,

$$a \leq y_1 < y_2 < \cdots < y_k \leq b.$$

Then, by the definition of integrals,

$$\int_a^b f(y)dy = \int_a^{y_1} f(y)dy + \int_{y_1}^{y_2} f(y)dy + \cdots + \int_{y_k}^b f(y)dy.$$

(2) and (3) are obvious. ■

The following is the reverse of Theorem 7.2.

Theorem 7.3. Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ satisfy the following:

- (1) ψ is integrable: $\int_a^b \psi(x)dx$ exists for $a < b$.
- (2) $\psi(x) \geq 0$ for $-\infty < x < \infty$.
- (3) $\int_{-\infty}^{\infty} \psi(x)dx = 1$.
- Then, $Q\{a < Y \leq b\} := \int_a^b \psi(x)dx$ defines a probability measure Q on Ω .

PROOF:

Out of scope for this course: Some advanced tools are needed to show is the σ -additivity of Q . ■

Remark 7.2. We combine (7.3) and (7.4) and obtain the following for a continuous random variable Y with PDF $f_Y(y)$: If $a, b \in \mathbb{R}$ and $a < b$, then

$$(7.5) \quad \begin{aligned} P\{a < Y < b\} &= P\{a \leq Y \leq b\} = P\{a \leq Y < b\} \\ &= P\{a < Y \leq b\} = \int_a^b f(y)dy. \quad \square \end{aligned}$$

The next definition applies to any random variable, be it continuous or discrete or neither. It is based on the following elementary observation.

Remark 7.3. ★ Assume that Y is a random variable with CDF $F_Y(y)$. For $0 < p < 1$, let

$$A_p := \{\alpha \in \mathbb{R} : F_Y(\alpha) \geq p\}.$$

Note that the function $y \mapsto F_Y(y)$ is nondecreasing.

- It is obvious that $[\alpha < \alpha' \text{ and } F_Y(\alpha) \geq p] \Rightarrow F_Y(\alpha') \geq p$.
- In other words, $[\alpha < \alpha' \text{ and } \alpha \in A_p] \Rightarrow \alpha' \in A_p$.
- In other words, A_p is an interval that stretches all the way to $+\infty$: There must be some real number β such that $A_p =]\beta, \infty[$ or $A_p = [\beta, \infty[$.

We see that $\beta \in A_p$ and thus, $A_p = [\beta, \infty[$, as follows. Let $\beta_n := \beta + \frac{1}{n}$.

- Since $\beta_n \in A_p$, $F_Y(\beta_n) \geq p$. Since F_Y is right continuous,³⁸ $F_Y(\beta) = \lim_{n \rightarrow \infty} F_Y(\beta_n)$.
- Thus, $F_Y(\beta) \geq p$. Thus, $\beta \in A_p$. Thus, $A_p = [\beta, \infty[$.
- Since $A_p = \{\alpha \in \mathbb{R} : F_Y(\alpha) \geq p\}$ and $A_p = [\beta, \infty[$, β is the smallest element of A_p , i.e.,

$$\beta = \min\{\alpha \in \mathbb{R} : F_Y(\alpha) \geq p\}.$$

The number β is uniquely determined by p . This allows us to denote it by the symbol ϕ_p . \square

Definition 7.4 (p th quantile).

Let Y denote any random variable and $0 < p < 1$. Let ϕ_p be the number derived in the previous remark, i.e.,

$$(7.6) \quad \phi_p = \min\{\alpha \in \mathbb{R} : F_Y(\alpha) \geq p\}$$

We call ϕ_p the p th **quantile** and also the $100p$ th **percentile** of Y .

Moreover, we call $\phi_{0.25}$ the **first quartile**, $\phi_{0.5}$ the **median**, and $\phi_{0.75}$ the **third quartile**, of the random variable Y . \square

Remark 7.4. Remark: How does the definition of the $100p$ th percentile given above correspond to the one experienced in everyday life: the number y_p that divides a list of numeric observations into $100p\%$ of the data being $\leq y_p$ and the remaining data being above y_p ? The connection is as follows.

- Assume that $\vec{y} = (y_1, y_2, \dots, y_K)$ is the list of observations. It may contain duplicates.
- We remove the duplicates and $N \leq K$ distinct values $\omega_1, \omega_2, \dots, \omega_N$ remain.
- We define $\Omega := \{\omega_1, \omega_2, \dots, \omega_N\}$ and $P\{\omega_j\} := \frac{n_j}{K}$ (we divide by K , **not** by $N!$), where $n_j =$ number of times that ω_j occurs in the original list, \vec{y} .
- σ -additivity extends P from the simple events $\{\omega_j\}$ to all events of Ω .
- Since ϕ_p is defined in terms of the CDF F_Y of a random variable Y , we define the following “dummy” random variable on (Ω, P) : $\omega \mapsto Y(\omega) := \omega$ ³⁹

For example, if the list of observations is $\vec{y} = (0, 2, 2, 2, 3, 4, 4, 6, 6, 6, 6, 7, 8, 8, 8)$, then

- $K = 15$, $\Omega = \{0, 2, 3, 4, 6, 7, 8\}$, $N = 7$,
- $P\{0\} = \frac{1}{15}$, $P\{2\} = \frac{3}{15}$, $P\{3\} = \frac{1}{15}$, $P\{4\} = \frac{2}{15}$, $P\{6\} = \frac{4}{15}$, $P\{7\} = \frac{1}{15}$, $P\{8\} = \frac{3}{15}$.
- Then $F_Y(6) = (1 + 3 + 1 + 2)/15 = 7/15$. Thus, $\phi_{7/15} = 6$.
- Also, the percentage of observations with a score of 6 or less is $700/15 \approx 46.667\%$. Hence, a score of 6 corresponds to the 46.667th percentile of \vec{y} . \square

³⁸See Remark 7.1 on p.116.

³⁹This method is more frequently employed in reverse: Given is a function $y \mapsto F(y)$ on the real numbers which satisfies the assumptions of Theorem 7.1 (Properties of a Cumulative Distribution Function) on p.115 and the subsequent Remark 7.1: F is nondecreasing, right-continuous, $F(-\infty) = 0$, $F(\infty) = 1$. We then define $\Omega := \mathbb{R}$ and, for $]a, b] \subseteq \Omega$, $P(]a, b]) := F(b) - F(a)$. σ -additivity extends this to a probability measure on all Borel sets of Ω (i.e., of \mathbb{R}). Now we define the random variable Y on (Ω, P) via $Y(y) := y$. Its CDF F_Y matches F , since,

$$F_Y(y) = P\{Y \leq y\} = P(]-\infty, y]) = F(y) - F(-\infty) = F(y).$$

In other words, Any function F that conforms to Theorem 7.1 and Remark 7.1 can be represented as the CDF F_Y of an appropriate random variable Y .

Example 7.2. Given the toss of a fair coin, let $Y(\omega) = 1$ if Heads and $Y(\omega) = 0$ if Tails come up.

Then Y has PMF $p_Y(0) = p_Y(1) = 1/2$

and CDF $F_Y(y) = 0$ for $y < 0$, $F_Y(y) = 0.5$ for $0 \leq y < 1$, $F_Y(y) = 1$ for $y \geq 1$.

We now easily compute ϕ_p for any $0 < p < 1$ by separately considering the cases

$$\begin{aligned} 0 < p < \frac{1}{2}: & F_Y(\alpha) \geq p \Leftrightarrow \alpha \geq 0. \text{ Thus, } \phi_p = 0. \\ p = \frac{1}{2}: & F_Y(\alpha) \geq \frac{1}{2} \Leftrightarrow \alpha \geq 0. \text{ Thus, } \phi_{1/2} = 0. \\ \frac{1}{2} < p < 1: & F_Y(\alpha) \geq p \Leftrightarrow \alpha \geq 1. \text{ Thus, } \phi_p = 1. \end{aligned}$$

Note that there are only two different ϕ_p values across all $0 < p < 1$: Either $\phi_p = 0$ or $\phi_p = 1$

This example also demonstrates that

$$\min\{\alpha \in \mathbb{R} : F_Y(\alpha) \geq p\}$$

cannot be replaced with the simpler expression

$$\min\{\alpha \in \mathbb{R} : F_Y(\alpha) = p\} :$$

The set $\{\alpha \in \mathbb{R} : F_Y(\alpha) = p\}$ is empty for $0 < p < 1$ unless $p = 0.5$, meaning that the minimum does not even exist! \square

The issues encountered in that last example do not occur if $F_Y(y)$ is a continuous function of y .

Proposition 7.2.

Let Y be a continuous random variable with CDF $F_Y(y)$. Then

$$(7.7) \quad \phi_p = \min\{\alpha \in \mathbb{R} : F_Y(\alpha) = p\}.$$

PROOF: The continuity of F_Y ensures that the sets

$$B_p := \{\alpha \in \mathbb{R} : F_Y(\alpha) = p\}$$

are not empty. The result follows from the fact that the function F_Y is nondecreasing. Further details are omitted. \blacksquare

Remark 7.5. For a continuous random variable Y with PMF $p_Y(y)$, quantiles have the following geometric meaning:

- The p th quantile is that value on the horizontal(!) axis which splits the area under the PMF into $100 \cdot p\%$ to the left and $100(1 - p)\%$ to the right. In particular,
- the median splits the area under the PMF into two halves.
- the first quartile splits the area under the PMF into 25% to the left and 75% to the right.
- the third quartile splits the area under the PMF into 75% to the left and 25% to the right. \square

We also use functional notation $\phi(p)$ for ϕ_p , since this makes what follows easier to understand.

Proposition 7.3.

Let Y be a random variable with an injective CDF $F_Y(y)$. (Note that it is not assumed that F_Y is continuous.) Then

$$(7.8) \quad \phi(F_Y(y)) = y \quad \text{for all } y \in \mathbb{R}$$

PROOF:

Let $p := F_Y(y)$. Since F_Y is nondecreasing, its injectivity means that

$$(7.9) \quad y_1 < y < y_2 \Rightarrow F_Y(y_1) < F_Y(y) < F_Y(y_2)$$

We infer that $\alpha < y$ does not satisfy $F_Y(\alpha) \geq F_Y(y) = p$. Since (see 7.6 on p.119)

$$(7.10) \quad \phi(F_Y(y)) = \min\{\alpha \in \mathbb{R} : F_Y(\alpha) \geq \phi(F_Y(y))\},$$

it follows from (7.10) that $\phi(F_Y(y)) < y$ is not possible. Thus, $\phi(F_Y(y)) \geq y$.

On the other hand, $\alpha = y$ does satisfy $F_Y(\alpha) \geq F_Y(y) = p$ and we just have seen that y is the smallest possible of those α . We apply (7.10) once more and conclude that $\phi(F_Y(y)) = y$. ■

Proposition 7.4.

Let Y be a random variable with a bijective CDF $F_Y : \mathbb{R} \xrightarrow{\sim}]0, 1[$. Then $F_Y(y)$ and $\phi(p)$ are inverse to each other, i.e.,

$$(7.11) \quad \phi(F_Y(y)) = y \text{ for all } y \in \mathbb{R} \quad \text{and} \quad F_Y(\phi(p)) = p \text{ for all } 0 < p < 1.$$

PROOF:

The equation $\phi(F_Y(y)) = y$ was shown in Proposition 7.3. Thus, it only remains to be shown that

$$(7.12) \quad F_Y(\phi(p)) = p \text{ for all } 0 < p < 1.$$

We observe that the bijective and nondecreasing function F_Y is strictly increasing and continuous. It is easy to see that F_Y is strictly increasing: Note that $y_1 < y_2 \Rightarrow F_Y(y_1) \leq F_Y(y_2)$ because F_Y is nondecreasing. Injectivity prohibits $F_Y(y_1) = F_Y(y_2)$. Thus, F_Y is strictly increasing.

It is harder to see that F_Y is continuous:

- If there was a point of discontinuity $y_0 \in \mathbb{R}$ for F_Y , then F_Y being nondecreasing and right-continuous would mean that $F_Y(y_0-) = \lim_{y < y_0, y \rightarrow y_0} F_Y(y) < F_Y(y_0)$.
- Also, F_Y nondecreasing $\Rightarrow F_Y(y) \leq F_Y(y_0-)$ for $y < y_0$ and $F_Y(y) \geq F_Y(y_0)$ for $y \geq y_0$.
- Thus, no $y \in \mathbb{R}$ and $p \in]F_Y(y_0-), F_Y(y_0)[$ satisfies $F_Y(y) = p$, contradicting surjectivity of F_Y .

Since F_Y is continuous, we obtain from Proposition 7.2 on p.120 that

$$(7.13) \quad \phi(p) = \min\{\alpha \in \mathbb{R} : F_Y(\alpha) = p\}.$$

In particular, $\phi(p)$ is an element of the set $\{\alpha \in \mathbb{R} : F_Y(\alpha) = p\}$. Thus, $\phi(p)$ satisfies $F_Y(\phi(p)) = p$. We have shown (7.12). We noted previously that the proposition follows. ■

7.3 Expected Value, Variance and MGF of a Continuous Random Variable

Assumption 7.2 (All continuous random variables have Expectations). **A.** Unless explicitly stated otherwise, all continuous random variables are assumed to possess a probability density function $f_Y(y)$ that satisfies

$$\int_{-\infty}^{\infty} |y|f(y) dy < \infty.$$

This technical condition guarantees the existence of $\int_{-\infty}^{\infty} yf(y)dy$ which is needed to define the expected value of Y .

B. We further assume that, unless specifically stated otherwise, there is a common probability space (Ω, P) for all random variables Y , be they discrete, continuous or neither, are of the form $Y : (\Omega, P) \rightarrow \mathbb{R}$. \square

Definition 7.5 (Expected value of a continuous random variable).

Let Y be a continuous random variable with PDF $f_Y(y)$. We call

$$(7.14) \quad E(Y) := \int_{-\infty}^{\infty} yf_Y(y) dy$$

the **expected value**, also **expectation** or **mean** of Y . \square

We will use the next theorem in the proof of Theorem 7.5 on p.123. The presentation given here follows [4] Ghahramani, Saeed.

Theorem 7.4. ★

Let Y be a continuous random variable with CDF F_Y and PDF f_Y .

Then

$$(7.15) \quad E[Y] = \int_0^{\infty} (1 - F_Y(y)) dy - \int_0^{\infty} F_Y(-y) dy$$

$$(7.16) \quad = \int_0^{\infty} P\{Y > y\} dy - \int_0^{\infty} P\{Y \leq -y\} dy.$$

PROOF: We only need to prove (7.15), since (7.16) follows from the definition of a CDF.

$$\text{Let } A_1 := \{(u', y') : y' < 0, 0 < u' < -y'\}, \quad B_1 := \{(u', y') : u' > 0, y' < -u'\}.$$

Then $u' < -y' \Leftrightarrow y' < -u'$ implies $A_1 = B_1 = \{(u', y') : u' > 0, y' < 0, u' < -y'\}$. Thus,

$$\begin{aligned} \int_{-\infty}^0 \left(\int_0^{-y} du \right) f(y) dy &= \iint_{A_1} f_Y(y) d(u, y) \\ \text{(a)} \quad &= \iint_{B_1} f_Y(y) d(u, y) = \int_0^{\infty} \left(\int_{-\infty}^{-u} f_Y(y) dy \right) du. \end{aligned}$$

$$\text{Let } A_2 := \{(u', y') : y' > 0, 0 < u' < y'\}, \quad B_2 := \{(u', y') : u' > 0, y' > u'\}.$$

Then $A_2 = B_2$, because both denote the set $\{(u', y') : u' > 0, y' > 0, u' < y'\}$. It follows that

$$\begin{aligned} \int_0^\infty \left(\int_0^y du \right) f(y) dy &= \iint_{A_2} f_Y(y) d(u, y) \\ \text{(b)} \qquad \qquad \qquad &= \iint_{B_2} f_Y(y) d(u, y) = \int_0^\infty \left(\int_u^\infty f_Y(y) dy \right) du. \end{aligned}$$

We use **(a)** and **(b)** in the following chain of equations:

$$\begin{aligned} E[Y] &= \int_{-\infty}^\infty y f_Y(y) dy = \int_{-\infty}^0 y f_Y(y) dy + \int_0^\infty y f_Y(y) dy \\ &= - \int_{-\infty}^0 \left(\int_0^{-y} du \right) f_Y(y) dy + \int_0^\infty \left(\int_0^y du \right) f_Y(y) dy \\ &\stackrel{\text{(a),(b)}}{=} - \int_0^\infty \left(\int_{-\infty}^{-u} f_Y(y) dy \right) du + \int_0^\infty \left(\int_u^\infty f_Y(y) dy \right) du. \\ &= - \int_0^\infty F_Y(-u) du + \int_0^\infty (1 - F_Y(u)) du. \end{aligned}$$

The last equation follows from $\int_\alpha^\beta f_Y(y) dy = F_Y(\beta) - F_Y(\alpha)$. ■

Corollary 7.1. ★

Let Y be a nonnegative, continuous random variable with CDF F_Y and PDF f_Y . Then

$$(7.17) \qquad E[Y] = \int_0^\infty (1 - F_Y(y)) dy = \int_0^\infty P\{Y > y\} dy.$$

PROOF: $Y \geq 0$ implies $P\{Y \leq -y\} = 0$ for $0 \leq y < \infty$. Thus, (7.17) follows from (7.15) and (7.16). ■

Quite a few theorems about discrete random variables have continuous counterparts when one replaces probability mass function $p(y)$ with probability density function $f(y)$ and summation over the countably many y for which $p(y) > 0$ with integration over all y . The following theorem corresponds to Theorem 6.2 on p.94. Note that the continuous random variable $\omega \mapsto g(Y(\omega))$ of that theorem is covered by Assumption 7.2 on p.121, i.e., $E[g \circ Y]$ exists.

Theorem 7.5.

Let Y be a discrete or continuous random variable with PDF f_Y and $g : \mathbb{R} \rightarrow \mathbb{R}$; $y \mapsto g(y)$ be a real-valued function. Then the random variable $g \circ Y : \omega \mapsto g(Y(\omega))$ has expectation

$$(7.18) \qquad E[g(Y)] = \int_{-\infty}^\infty g(y) f_Y(y) dy.$$

PROOF: ★ The proof of Theorem 6.2 on p.94 handles the discrete case. So we may assume that Y is a continuous random variable.

According to Proposition 3.7 (Preimages of function composition) on p.64,

$$\begin{aligned}\{g \circ Y > u\} &= (g \circ Y)^{-1}(]u, \infty[) = Y^{-1}(g^{-1}(]u, \infty[)) = \{Y \in g^{-1}(]u, \infty[)\}. \\ \{g \circ Y \leq -u\} &= (g \circ Y)^{-1}(]-\infty, -u]) = Y^{-1}(g^{-1}(]-\infty, -u])) = \{Y \in g^{-1}(]-\infty, -u])\}.\end{aligned}$$

Thus,

$$\begin{aligned}\text{(a)} \quad P\{g \circ Y > u\} &= P\{Y \in g^{-1}(]u, \infty[)\} = P_Y\{g^{-1}(]u, \infty[)\} = P_Y\{y : g(y) > u\} \\ P\{g \circ Y \leq -u\} &= P\{Y \in g^{-1}(]-\infty, -u])\} = P_Y\{g^{-1}(]-\infty, -u])\} = P_Y\{y : g(y) \leq -u\}.\end{aligned}$$

Next, we show that $A_1 = B_1$. Here, we define A_1 and B_1 as follows:

$$\text{(b1)} \quad A_1 := \{(u', y') : 0 < u' < \infty, g(y') > u'\}, \quad B_1 := \{(u', y') : g(y') > 0, 0 < u' < g(y')\},$$

To show $A_1 \subseteq B_1$, let $(u, y) \in A_1$, i.e., $(u, y) \in \{(u', y') : 0 < u' < \infty, g(y') > u'\}$.

- $0 < u$ and $u < g(y)$ yields $g(y) > 0$ and $0 < u < g(y)$. Thus, $(u, y) \in B_1$.

To see that $B_1 \subseteq A_1$, let $(u, y) \in B_1$, i.e., $(u, y) \in \{(u', y') : g(y') > 0, 0 < u' < g(y')\}$.

- Since $0 < u < g(y)$, it follows that $0 < u < \infty$ and $u < g(y)$. Thus, $(u, y) \in A_1$.

$$\text{(c1)} \quad \text{We proved that } A_1 = B_1. \text{ It follows that } \iint_{A_1} f_Y(y) d(t, y) = \iint_{B_1} f_Y(y) d(t, y).$$

On a parallel track, we show that $A_2 = B_2$, where we define A_2 and B_2 as follows:

$$\text{(b2)} \quad A_2 := \{(u', y') : 0 < u' < \infty, g(y') \leq -u'\} \quad B_2 := \{(u', y') : g(y') < 0, 0 < u' \leq -g(y')\}.$$

To show $A_2 \subseteq B_2$, let $(u, y) \in A_2$, i.e., $(u, y) \in \{(u', y') : 0 < u' < \infty, g(y') \leq -u'\}$.

- Since $g(y) \leq -u \Leftrightarrow u \leq -g(y)$ and we also have $0 < u < \infty$, $(u, y) \in A_2$ implies $0 < u \leq -g(y)$.
- To show that also $g(y) < 0$ we observe that $g(y) \leq -u < -0 = 0$.

Finally, to show $B_2 \subseteq A_2$, let $(u, y) \in B_2 = \{(u', y') : g(y') < 0, 0 < u' \leq -g(y')\}$.

- $0 < u < \infty$ is immediate from $0 < u \leq -g(y)$. We still must show that $g(y) \leq -u$.
- To show that also $g(y) < 0$ we observe that $g(y) \leq -u < -0 = 0$. But this is immediate from $0 < u \leq -g(y)$.

$$\text{(c2)} \quad \text{We proved that } A_2 = B_2. \text{ It follows that } \iint_{A_2} f_Y(y) d(t, y) = \iint_{B_2} f_Y(y) d(t, y).$$

We apply (c1) and (c2) to the integrals $\int_0^\infty P\{g \circ Y > u\} du$ and $\int_0^\infty P\{g \circ Y \leq -u\} du$ as follows.

$$\begin{aligned}\int_0^\infty P\{g \circ Y > u\} du &\stackrel{\text{(a)}}{=} \int_0^\infty P\{Y \in g^{-1}(]u, \infty[)\} du = \int_0^\infty P_Y\{g^{-1}(]u, \infty[)\} du \\ &= \int_0^\infty P_Y\{y : u < g(y) < \infty\} du = \int_0^\infty \left(\int_{\{y: u < g(y) < \infty\}} f_Y(y) dy \right) du \\ &\stackrel{\text{(b1)}}{=} \iint_{A_1} f_Y(y) d(t, y) \stackrel{\text{(c1)}}{=} \iint_{B_1} f_Y(y) d(t, y) \stackrel{\text{(b1)}}{=} \int_{\{y: g(y) > 0\}} \left(\int_0^{g(y)} du \right) f_Y(y) dy\end{aligned}$$

Hence, since $\int_0^{g(y)} du = g(y)$,

$$(d1) \quad \int_0^\infty P\{g \circ Y > u\} du = \int_{\{y:g(y)>0\}} g(y) f_Y(y) dy$$

$$\begin{aligned} \int_0^\infty P\{g \circ Y \leq -u\} du &\stackrel{(a)}{=} \int_0^\infty P\{Y \in g^{-1}([-\infty, -u])\} du = \int_0^\infty P_Y\{g^{-1}([-\infty, -u])\} du \\ &= \int_0^\infty P_Y\{y : -\infty < g(y) < -u\} du = \int_0^\infty \left(\int_{\{y:-\infty < g(y) < -u\}} f_Y(y) dy \right) du \\ &\stackrel{(b2)}{=} \iint_{A_2} f_Y(y) d(t, y) \stackrel{(c2)}{=} \iint_{B_2} f_Y(y) d(t, y) \stackrel{(b2)}{=} \int_{\{y:g(y)<0\}} \left(\int_0^{-g(y)} du \right) f_Y(y) dy \end{aligned}$$

Hence, since $\int_0^{-g(y)} du = -g(y)$,

$$(d2) \quad \int_0^\infty P\{g \circ Y \leq -u\} du = - \int_{\{y:g(y)<0\}} g(y) f_Y(y) dy$$

It follows from (d1) and (d2) and Theorem 7.4 on p.122 and

$$\int_{\{y:g(y)=0\}} g(y) f_Y(y) dy = \int_{\{y:g(y)=0\}} 0 f_Y(y) dy = 0,$$

that

$$\begin{aligned} E[g \circ Y] &= \int_0^\infty P\{g \circ Y > u\} du - \int_0^\infty P\{g \circ Y \leq -u\} du \\ &= \int_{\{y:g(y)>0\}} g(y) f_Y(y) dy + \int_{\{y:g(y)<0\}} g(y) f_Y(y) dy \\ &= \int_{\{y:g(y)>0\}} g(y) f_Y(y) dy + \int_{\{y:g(y)<0\}} g(y) f_Y(y) dy + \int_{\{y:g(y)=0\}} g(y) f_Y(y) dy \\ &= \int_{\mathbb{R}} g(y) f_Y(y) dy = \int_{-\infty}^\infty g(y) f_Y(y) dy \quad \blacksquare \end{aligned}$$

The following corresponds to WMS Theorem 4.5.

Theorem 7.6.

Let $c \in \mathbb{R}$, Y be a discrete or continuous random variable and $g_1, g_2, g_n : \mathbb{R} \rightarrow \mathbb{R}$; $y \mapsto g(y)$ be a list of n real-valued functions. Then

$$(7.19) \quad E[c] = c,$$

$$(7.20) \quad E[cg_j(Y)] = cE[g_j(Y)].$$

Further, the random variable

$$\sum_{j=1}^n g_j \circ Y : \Omega \longrightarrow \mathbb{R}; \quad \omega \mapsto \sum_{j=1}^n g_j(Y(\omega))$$

has the following expected value:

$$(7.21) \quad E \left[\sum_{j=1}^n g_j \circ Y \right] = \sum_{j=1}^n E[g_j \circ Y].$$

PROOF: ■

We will not deal in this course with the sums of continuous and discrete random variables, so the next definition is only included for completeness' sake and to allow the formulation of theorems 7.7 and 7.8 below.

Definition 7.6. ★

If Y_1, Y_2, \dots, Y_m is a list of discrete random variables and Y'_1, Y'_2, \dots, Y'_n is a list of continuous random variables, all of which are defined on the same probability space (Ω, P) , then we define

$$(7.22) \quad E \left[\sum_{i=1}^m Y_i + \sum_{j=1}^n Y'_j \right] := \sum_{i=1}^m E[Y_i] + \sum_{j=1}^n E[Y'_j] p. \quad \square$$

The following is the continuous random variables version of Theorem 6.4 on p.96.

Theorem 7.7.

Let $Y_1, Y_2, \dots, Y_n : \Omega \rightarrow \mathbb{R}$ be random variables. (which all are defined on the same probability space (Ω, P) ($n \in \mathbb{N}$ by Assumption 7.2.B). Some may be continuous, others may be discrete. Then the random variable

$$\sum_{j=1}^n Y_j : \Omega \longrightarrow \mathbb{R}; \quad \omega \mapsto \sum_{j=1}^n Y_j(\omega)$$

has the following expected value:

$$(7.23) \quad E \left[\sum_{j=1}^n Y_j \right] = \sum_{j=1}^n E[Y_j].$$

In other words, the expectation of the sum is the sum of the expectations.

PROOF: Not given here. ■

We extend Definition 6.3 on p.98 of the variance and standard deviation of a discrete random vari-

able to the continuous case without modification, i.e.,

$$(7.24) \quad \text{Var}[Y] := \sigma_Y^2 := E[(Y - E[Y])^2],$$

$$(7.25) \quad \sigma_Y := \sqrt{\text{Var}[Y]}.$$

Theorems 6.5, 6.6 6.7 about the variances of discrete random variables have the following counterpart.

Theorem 7.8. *Let Y be a discrete or continuous random variable. Let $Y_1, Y_2, \dots, Y_n : \Omega \rightarrow \mathbb{R}$ be independent random variables (which all are defined on the same probability space (Ω, P) ($n \in \mathbb{N}$ by Assumption 7.2.B)). Some may be continuous, others may be discrete. Further, let $a, b \in \mathbb{R}$. Then*

$$(7.26) \quad \text{Var}[Y] = E[Y^2] - (E[Y])^2,$$

$$(7.27) \quad \text{Var}[aY + b] = a^2 \text{Var}[Y],$$

$$(7.28) \quad \text{Var}\left[\sum_{j=1}^n Y_j\right] = \sum_{j=1}^n \text{Var}[Y_j].$$

PROOF: The proof of (7.26) is the same as for Theorem 6.5 on p.98. The proof of the other formulas is not given here. ■

Remark 7.6. Note that independence of Y_1, \dots, Y_n is required for the validity of (7.28)! □

The moments about the origin μ'_k , the moments about the mean μ_k and the MGF $m_Y(t)$ of a discrete random variable Y , all were defined as expected values. This allows us to use those same definitions for continuous random variables.

Unless something different is stated, Y is a random variable $Y : (\Omega, P) \rightarrow \mathbb{R}$ on some probability space (Ω, P) . Further, $\mu = E[Y]$, $\sigma^2 = \text{Var}[Y]$ and $\sigma = \sqrt{\text{Var}[Y]}$ denote expectation, variance and standard deviation of Y .

Definition 7.7. For $k \in \mathbb{N}$, we define

$$(7.29) \quad \mu'_k := E[Y^k] \quad (\textit{kth moment of } Y \textit{ about the origin})$$

$$(7.30) \quad \mu_k := E[(Y - E[Y])^k] = E[(Y - \mu)^k] \quad (\textit{kth central moment of } Y)$$

$$(7.31) \quad m(t) := m_Y(t) := E[e^{tY}] \quad (\textit{moment-generating function of } Y)$$

As in the discrete case we assume that the expectations defining μ'_k and μ_k exist and that there is some $\delta > 0$ such that $m_Y(t)$ is defined (i.e., finite) for $|t| < \delta$. □

Theorem 6.18 on p.111 remains valid for continuous random variables:

Theorem 7.9.

Let Y be a random variable with MGF $m_Y(t)$ and $k \in \mathbb{N}$. Then its k th moment is obtained as the k th derivative of $m_Y(\cdot)$, evaluated at $t = 0$:

$$(7.32) \quad \mu'_k = m^{(k)}(0) = \left. \frac{d^k m(t)}{dt^k} \right|_{t=0}.$$

PROOF: The proof of Theorem 6.18 can be used without any alterations. ■

Proposition 7.5.

Let Y be a random variable with MGF $m_Y(t)$. Let $a, b \in \mathbb{R}$, $Y' := Y + a$, $Y'' := bY$. Then

$$(7.33) \quad m_{Y'}(t) = e^{ta} m_Y(t),$$

$$(7.34) \quad m_{Y''}(t) = m_Y(bt).$$

PROOF: To prove (7.33), we note that e^{ta} is constant in ω . Thus, $E[e^{ta}W] = e^{ta}E[W]$ for any random variable W . Thus,

$$m_{Y'}(t) = E[e^{t(Y+a)}] = E[e^{tY} e^{ta}] = e^{ta} E[e^{tY}] = e^{ta} m_Y(t).$$

Formula (7.34) follows from

$$m_{Y''}(t) = E[e^{t(bY)}] = E[e^{(tb)Y}] = m_Y(tb). \quad \blacksquare$$

7.4 The Uniform Probability Distribution

Given two real numbers $\theta_1 < \theta_2$, we consider a random variable $Y(\omega)$ that “lives” in the interval $[\theta_1, \theta_2]$, i.e., $P\{\theta_1 \leq Y \leq \theta_2\} = 1$ and has the same likelihood of occurring in any subinterval of same length:

Definition 7.8 (Continuous, uniform random variable).

Let Y be a random variable and $-\infty < \theta_1 < \theta_2 < \infty$. We say that Y has a **continuous uniform probability distribution** with parameters θ_1 and θ_2 — also, that Y is **uniform on** $[\theta_1, \theta_2]$ or $Y \sim \mathbf{uniform}(\theta_1, \theta_2)$ — if Y has probability density function

$$(7.35) \quad f_Y(y) = \begin{cases} \frac{1}{\theta_2 - \theta_1}, & \text{if } \theta_1 \leq y \leq \theta_2, \\ 0, & \text{else. } \square \end{cases}$$

Remark 7.7 (uniform and equiprobable probability measures). Uniform distributions are the equivalent of the distribution of discrete random variables Y that satisfy equiprobability, i.e., their PMF $p_Y(y) = P\{Y = y\}$ is strictly positive only for finitely many numbers y_1, y_2, \dots, y_n and $p_Y(y_j) = 1/n$ for all $j \in [1, n]_{\mathbb{Z}}$. See Definition 3.3 on p.47. □

Theorem 7.10 (WMS Ch.04.4, Theorem 4.6).

If $\theta_1 < \theta_2$ and Y is a uniform random variable with parameters θ_1, θ_2 , then

$$E[Y] = \frac{\theta_1 + \theta_2}{2} \quad \text{and} \quad \text{Var}[Y] = \frac{(\theta_2 - \theta_1)^2}{12}.$$

PROOF: A simple exercise in integrating $\int_{\theta_1}^{\theta_2} y \, dy$ and $\int_{\theta_1}^{\theta_2} y^2 \, dy$. ■

Theorem 7.11.

Assume that Y is a continuous random variable with CDF $F_Y(y)$. Let $U := F_Y(Y)$. Then $U \sim \text{uniform}(0, 1)$.

SIMPLIFIED PROOF under the assumption that the CDF F_Y is a bijection $F_Y : \mathbb{R} \xrightarrow{\sim}]0, 1[$.

The inverse F_Y^{-1} of F_Y satisfies $F_Y^{-1}(F_Y(y)) = y$ for all $y \in \mathbb{R}$. Thus, for $0 < u < 1$,

$$\begin{aligned} F_U(u) &= P\{U \leq u\} = P\{F_Y \circ Y \leq u\} = P\{F_Y^{-1} \circ F_Y \circ Y \leq F_Y^{-1}(u)\} \\ &= P\{Y \leq F_Y^{-1}(u)\} = F_Y(F_Y^{-1}(u)) = u. \end{aligned}$$

We still must handle the cases $u \leq 0$ and $u \geq 1$. We assumed that the codomain of F_Y is $]0, 1[$.

- Thus, $y \in \mathbb{R} \Rightarrow 0 < F_Y(y) < 1$.
- Thus, $\omega \in \Omega \Rightarrow 0 < U(\omega) = F_Y(Y(\omega)) < 1$
 $\Rightarrow [P\{U \leq 0\} = 0 \text{ and } P\{U \leq 1\} = 1] \Rightarrow [F_U(0) = 0 \text{ and } F_U(1) = 1]$.
- Thus, $[u \leq 0 \Rightarrow F_U(u) \leq F_U(0) = 0]$ and $[u \geq 1 \Rightarrow F_U(u) \geq F_U(1) = 1]$.

It follows that F_U is the CDF of a $\text{uniform}(0, 1)$ random variable. Thus, $U \sim \text{uniform}(0, 1)$. ■

GENERAL PROOF ★ (We drop the assumption that F_Y is a bijection $\mathbb{R} \xrightarrow{\sim}]0, 1[$):

This proof follows the one of Theorem 2.1.10 in Casella, Berger [3], but it gives additional detail.

Let $0 < p < 1$ and let

$$(A) \quad G(p) := \min\{y \in \mathbb{R} : F_Y(y) \geq p\}.$$

In other words, $G(p)$ is the p th quantile ϕ_p for the random variable Y . Since G is nondecreasing,

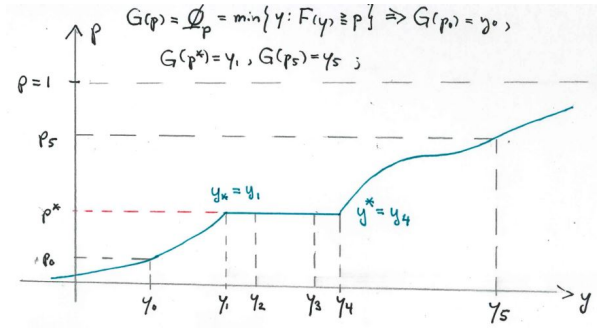
$$(B) \quad F_U(p) = P\{U \leq p\} = P\{F_Y(Y) \leq p\} = P\{G(F_Y(Y)) \leq G(p)\}.$$

The most difficult part of the proof is to show that

$$(C) \quad P\{G(F_Y(Y)) \leq G(p)\} = P\{Y \leq G(p)\}.$$

We consider two different cases.

- **Case 1:** There is a unique y such that $G(p) = y$. In the picture, that would be y_0 for p_0 and y_5 for p_5
- (a) Observe that $G(p) = y \Leftrightarrow p = F_Y(y)$.
- (b) $G(p') < G(p) < G(p'') \Leftrightarrow p' < p < p''$.
- **Case 2:** There are $y_* < y^*$, determined by $G(p) = y \Leftrightarrow y_* < y < y^*$. In the picture, that would be $y_* = y_1$ and $y^* = y_4$ for $F(y) = p$.



7.1 (Figure). non-injective, continuous CDF.

We now show that (C) is true for **Case 1**.

We deduce from (a) and (b) that

$$\begin{aligned} \omega \in \{G(F_Y(Y)) \leq G(p)\} &\Leftrightarrow F_Y(Y(\omega)) \leq G(p) (= F_Y(y)) \\ &\Leftrightarrow Y(\omega) \leq y (= G(p)) \Leftrightarrow \omega \in \{Y \leq G(p)\}. \end{aligned}$$

Taking probabilities shows that (C) is valid, since we obtain

$$P\{G(F_Y(Y)) \leq G(p)\} = P\{Y \leq G(p)\}.$$

Next, we show that (C) is true for **Case 2**.

The picture shows that, if $F_Y(y') = p'$ and $F_Y(y) = p \Leftrightarrow y_* \leq y \leq y^*$, then

- (c) $G(p') < G(p) \Leftrightarrow y' < y_*$; $G(p') = G(p) \Leftrightarrow y_* \leq y' \leq y^*$;
- (d) Thus, $G(p') \leq G(p) \Leftrightarrow y' \leq y^* \Leftrightarrow [y' \leq y_* \text{ or } y_* < y' \leq y^*]$.

Clearly,

$$\omega \in \{G(F_Y(Y)) \leq G(p)\} \Leftrightarrow G(F_Y(Y(\omega))) \leq G(p) (= y_*).$$

We apply (d) with $y' = Y(\omega)$ and $p' = F_Y(Y(\omega))$ and obtain

$$G(F_Y(Y(\omega))) \leq G(p) \Leftrightarrow [Y(\omega) \leq y_* \text{ or } y_* < Y(\omega) \leq y^*].$$

Thus, $\{G(F_Y(Y)) \leq G(p)\} = \{Y \leq y_*\} \cup \{y_* < Y \leq y^*\}$. Taking probabilities,

$$\begin{aligned} P\{G(F_Y(Y)) \leq G(p)\} &= P\{Y \leq y_*\} + P\{y_* < Y \leq y^*\} \\ &= F_Y(y_*) + (F_Y(y^*) - F_Y(y_*)) = F_Y(G(p)) = P\{Y \leq G(p)\}. \end{aligned}$$

Here, the equation next to the last follows from $G(p) = y_*$ and $F_Y(y_*) = G(p) = F_Y(y^*)$.

We have shown that (C) also is true for **Case 2**.

We combine (B) and (C) and obtain

$$(D) \quad F_U(p) = P\{F_Y(Y) \leq p\} = P\{Y \leq G(p)\} = F_Y(G(p)).$$

Our next goal is to show that $F_Y(G(p)) = p$. We break this down into the following steps.

- (1) By **(A)**, $F_Y(G(p)) \geq p$. We now show that also $F_Y(G(p)) \leq p$.
- (2) Let $y_n := G(p) - 1/n$. Then $G(p) = \lim_{n \rightarrow \infty} y_n$.
- (3) $G(p)$ being the smallest y such that $F_Y(y) \geq p$ implies that $F_Y(y_n) < p$.
- (4) Since Y is continuous, $F(y)$ is continuous. Thus, $F_Y(G(p)) = \lim_{n \rightarrow \infty} F_Y(y_n)$.
- (5) Since $F_Y(y_n) < p$ by **(3)**, $\lim_{n \rightarrow \infty} F_Y(y_n) \leq p$, i.e., $F_Y(G(p)) \leq p$. (See **(4)**.)
- (6) We have shown **(1)** and it follows that $F_Y(G(p)) = p$.

It now follows from **(D)** that $P\{U \leq p\} = p$ for any $0 < p < 1$.

The boundary cases $p = 0$ and $p = 1$ are taken into account by extending the definition of $G(p)$ given in **(A)**, which is $G(p) = \min\{y \in \mathbb{R} : F_Y(y) \geq p\}$, as follows.

- Since $F_Y(y) \geq 0$ for all y , it is natural to define $G(0) := -\infty$.
- If there is some y_* such that $F_Y(y_*) = 1$, then **(A)** remains in force for $G(1)$.
- Otherwise, (if $F_Y(y) < 1$ for all y), we define $G(1) := \infty$. ■

Theorem 7.12.

Given are a uniform(0, 1) random variable U and a continuous function $F : \mathbb{R} \rightarrow [0, 1]$ that satisfies the conditions of Theorem 7.1 (Properties of a Cumulative Distribution Function) on p.115:

- F is nondecreasing
- $F(-\infty) := \lim_{y \rightarrow -\infty} F(y) = 0$
- $F(\infty) := \lim_{y \rightarrow \infty} F(y) = 1$

$$(7.36) \quad \text{Let } G : [0, 1] \rightarrow \mathbb{R}; \quad p \mapsto G(p) := \min\{y \in \mathbb{R} : F(y) \geq p\}.$$

Let $Z := G(U)$ be the random variable $\omega \mapsto Z(\omega) := G(U(\omega))$.

Then its CDF matches F . In other words, $F_Z(y) = F(y)$ for all $y \in \mathbb{R}$.

SIMPLIFIED PROOF under the assumption that the F is a bijection $F : \mathbb{R} \xrightarrow{\sim}]0, 1[$.

We first show that G is the inverse of F .

- Since F is both nondecreasing and injective, F is strictly increasing.
- Let $0 < p_0 < 1$ and $y_0 := F^{-1}(p_0)$ or, equivalently, $p_0 = F(y_0)$.
- Let $A := \{y \in \mathbb{R} : F(y) \geq p_0\}$. Since $F(y_0) = p_0 \geq p_0$, it follows that $y_0 \in A$.
- Since F is strictly increasing, $y < y_0 \Rightarrow F(y) < F(y_0) = p_0 \Rightarrow y \notin A$.
- Since $y_0 \in A$ and $y < y_0 \Rightarrow y \notin A$, we conclude that $y_0 = \min(A)$.
- By (7.36), $G(p_0) = \min(A)$. We have shown $G(p_0) = y_0 = F^{-1}(p_0)$ for each $0 < p_0 < 1$.

Let $y \in \mathbb{R}$. Since $G = F^{-1}$, we obtain

$$F_Z(y) = P\{Z \leq y\} = P\{G \circ U \leq y\} = P\{F^{-1} \circ U \leq y\} = P\{U \leq F(y)\} = F(y).$$

The last equation follows from $0 < F(y) < 1$ and $U \sim \text{uniform}(0, 1)$. It follows that $F_Z(y) = F(y)$ for all y , i.e., $F_Z = F$ ■

GENERAL PROOF ★ (We drop the assumption that F_Y is a bijection $\mathbb{R} \xrightarrow{\sim}]0, 1[$):

Let $I := F_Y(\mathbb{R}) = \{F_Y(y) : y \in \mathbb{R}\}$ be the range of F_Y .

- Note that $G(p)$ equals the p th quantile ϕ_p of a random variable with CDF $F(y)$. (See Definition 7.4 on p.119.)
- Further, the continuity of F guarantees that for each $0 < p < 1$ one can find $y \in \mathbb{R}$ such that $F(y) = p$ (and thus, $p \mapsto G(p)$ is injective).

- Thus, I is one of the following intervals:
 - If $0 < F(y) < 1$ for all y , then $I =]0, 1[$
 - If $0 \leq F(y) < 1$ for all y , then $I = [0, 1[$
 - If $0 < F(y) \leq 1$ for all y , then $I =]0, 1]$
 - If $0 \leq F(y) \leq 1$ for all y , then $I = [0, 1]$
- We will refer in this proof to Figure 7.1 on p.130 (non-injective, continuous CDF) in the proof of Theorem 7.11.

We fix $y \in \mathbb{R}$. Let $p := F(y)$. Then

- (a) Since F is continuous and nondecreasing, there are numbers $y_* \leq y_*$ such that $F(\tilde{y}) = p \Leftrightarrow y_* \leq \tilde{y} \leq y_*$.
- (b) Either F is strictly increasing at y and then $y_* = y = y_*$, or F is “flat around y ” and $y_* < y_*$.
- (c) For $p' \in I$, choose y' such that $F(y') = p'$. Then, since $F(y_*) = p$, $p' < p \Leftrightarrow F(y') < p \Leftrightarrow y' < y_*$ and $p' \leq p \Leftrightarrow F(y') \leq p \Leftrightarrow y' \leq y_* \Leftrightarrow G(p') \leq y_*$.
- (d) Further, since F is nondecreasing, G also is nondecreasing. Thus, $p' \leq p \Leftrightarrow G(p') \leq G(p)$. It follows from (c) that $p' \leq p \Leftrightarrow G(p') \leq G(p) \Leftrightarrow y' \leq y_* \Leftrightarrow G(p') \leq y_*$.

Let $\omega \in \Omega$ and $p' := U(\omega)$. Recall that $p = F(y)$. Then

$$G(U(\omega)) \leq y \Leftrightarrow [G(p') \leq G(p)] \stackrel{(d)}{\Leftrightarrow} [p' \leq p] \Leftrightarrow [U(\omega) \leq F(y)].$$

We take probabilities and obtain, since $U \sim \text{uniform}(0, 1)$ implies $P\{U \leq \tilde{p}\} = \tilde{p}$ for $0 \leq \tilde{p} \leq 1$,

$$F_Z(y) = P\{G(U) \leq y\} = P\{U \leq F(y)\} = F(y).$$

To summarize, we have shown that $F_Z(y) = F(y)$ for all $y \in \mathbb{R}$. ■

Remark 7.8. A special case of Theorem 7.12 can be found in WMS Ch.06.3, Example 6.5, which shows how to solve the following problem: Let U be a uniform random variable on the interval $(0, 1)$. Find a transformation $G(U)$ such that $G(U)$ possesses an exponential distribution with mean β . □

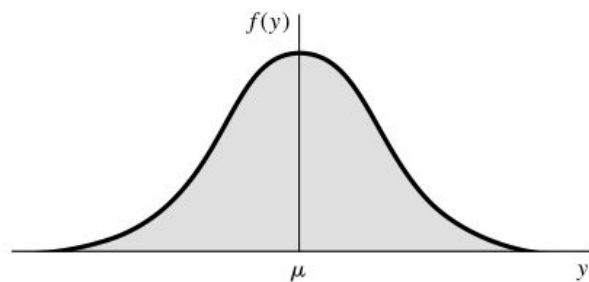
7.5 The Normal Probability Distribution

Many numerical random phenomena yield histograms which are approximately unimodal (a single highest value) and symmetric around the mean μ , like the picture to the right, and they adhere to the **empirical rule**: Approximately

- 68% of the data fall between $\mu \pm 1 \cdot \sigma$
- 95% of the data fall between $\mu \pm 2 \cdot \sigma$
- 99.7% of the data fall between $\mu \pm 3 \cdot \sigma$

Such data are adequately modeled by the normal distribution.

The empirical rule is also known as the **68%–95%–99.7% rule**.



Source: WMS Ch.4.5

Definition 7.9 (Normal random variable).

Let $\sigma > 0$ and $-\infty < \mu < \infty$. We say that a random variable Y has a **normal probability distribution** with mean μ and variance σ^2 if its probability density function is

$$(7.37) \quad f_Y(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(y-\mu)^2/(2\sigma^2)}, \quad (y \in \mathbb{R}). \quad \square$$

We also express that by saying that Y is $\mathcal{N}(\mu, \sigma^2)$. Moreover, we call Y **standard normal** if Y is $\mathcal{N}(0, 1)$.

We will see that $E[Y] = \mu$ and $\text{Var}[Y] = \sigma^2$. This justifies calling the parameters μ and σ^2 the mean and variance of the distribution.

Lemma 7.1.

$$(7.38) \quad (y - \mu)^2 - 2yt\sigma^2 = [y - (\mu + t\sigma^2)]^2 - 2\mu t\sigma^2 - t^2\sigma^4.$$

PROOF: We multiply out the right-hand expression and obtain

$$\begin{aligned} \text{R.S.} &= [y - (\mu + t\sigma^2)]^2 - 2\mu t\sigma^2 - t^2\sigma^4 \\ &= y^2 - 2y(\mu + t\sigma^2) + (\mu^2 + 2\mu t\sigma^2 + t^2\sigma^4) - 2\mu t\sigma^2 - t^2\sigma^4 \\ &= y^2 - 2\mu y - 2yt\sigma^2 + \mu^2 \\ &= (y - \mu)^2 - 2yt\sigma^2 = \text{L.S.} \quad \blacksquare \end{aligned}$$

Proposition 7.6.

Let the random variable Y be $\mathcal{N}(\mu, \sigma^2)$. Then

$$(7.39) \quad m_Y(t) = e^{\mu t + (\sigma^2 t^2)/2}.$$

PROOF:

$$\begin{aligned} m_Y(t) &= \int_{-\infty}^{\infty} e^{yt} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} dy \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\frac{(yt)(2\sigma^2)}{2\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} dy \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2\sigma^2} [(y-\mu)^2 - 2yt\sigma^2]} dy. \end{aligned}$$

We apply Lemma 7.1 and obtain for the exponent the following.

$$\begin{aligned} -\frac{1}{2\sigma^2} [(y - \mu)^2 - 2yt\sigma^2] &= -\frac{1}{2\sigma^2} \{ [y - (\mu + t\sigma^2)]^2 - 2\mu t\sigma^2 - t^2\sigma^4 \} \\ &= -\frac{[y - (\mu + t\sigma^2)]^2}{2\sigma^2} + \frac{1}{2\sigma^2} [2\mu t\sigma^2 + t^2\sigma^4] \\ &= \mu t + \frac{t^2\sigma^2}{2} - \frac{1}{2} \left[\frac{y - (\mu + t\sigma^2)}{\sigma} \right]^2 \end{aligned}$$

It follows that

$$\begin{aligned} m_Y(t) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\mu t + \frac{t^2\sigma^2}{2}} e^{-\frac{1}{2} \left[\frac{y - (\mu + t\sigma^2)}{\sigma} \right]^2} dy \\ &= e^{\mu t + \frac{t^2\sigma^2}{2}} \left[\frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2} \left(\frac{y - (\mu + t\sigma^2)}{\sigma} \right)^2} dy \right]. \end{aligned}$$

The expression in square brackets is the integral $\int_{-\infty}^{\infty} \varphi(y) dy$, where $\varphi(y)$ is the PDF of a normal variable with mean $\mu + t\sigma^2$ and variance σ^2 . Thus, this integral evaluates to 1 and it follows that

$$m_Y(t) = e^{\mu t + \frac{t^2\sigma^2}{2}}. \blacksquare$$

Theorem 7.13 (WMS Ch.04.5, Theorem 4.7).

If Y is a normally distributed random variable with parameters μ and σ , then

$$E[Y] = \mu \quad \text{and} \quad \text{Var}[Y] = \sigma^2.$$

PROOF: We differentiate $m_Y(t) = \exp\{\mu t + \frac{t^2\sigma^2}{2}\}$ twice and obtain

$$\begin{aligned} m'_Y(t) &= (\mu + t\sigma^2) \exp\left\{\mu t + \frac{t^2\sigma^2}{2}\right\}, \\ m''_Y(t) &= (\mu + t\sigma^2)^2 \exp\left\{\mu t + \frac{t^2\sigma^2}{2}\right\} + \sigma^2 \exp\left\{\mu t + \frac{t^2\sigma^2}{2}\right\}. \end{aligned}$$

Thus, the first and second moment about the origin are

$$\begin{aligned} E[Y] &= \mu'_1 = m'_Y(0) = (\mu + 0)e^0 = \mu, \\ E[Y^2] &= \mu'_2 = m''_Y(0) = (\mu + 0)^2 e^0 + \sigma^2 e^0 = \mu^2 + \sigma^2. \end{aligned}$$

Finally,

$$\text{Var}[Y] = E[Y^2] - (E[Y])^2 = \mu^2 + \sigma^2 - \mu^2 = \sigma^2. \blacksquare$$

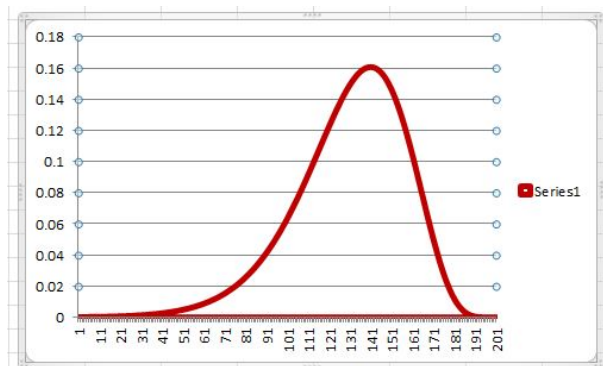
Remark 7.9. The importance of the normal distribution stems from the so called Central Limit Theorem (Theorem 10.13 on p.231), which we will discuss in Chapter 10 (Limit Theorems). It states the following.

- Given is an iid sequence of random variables Y_1, Y_2, \dots with common expectation $\mu := E[Y_j]$ and finite standard deviation $\sigma := \sqrt{\text{Var}[Y_j]} < \infty$ and a standard normal variable Z .
- For $n \in \mathbb{N}$, we define $\bar{Y}_n := \frac{1}{n} \sum_{j=1}^n Y_j = \frac{Y_1 + \dots + Y_n}{n}$ and $Z_n := \frac{\bar{Y}_n - \mu}{\sigma/\sqrt{n}}$.
- An aside: One easily sees from Theorems 7.7 (p.126) and 7.8 that $E[\bar{Y}_n] = \mu$, $\sigma_{\bar{Y}_n} = \sigma/\sqrt{n}$ and thus, $E[Z_n] = 0$, $\text{Var}[Z_n] = 1$.
- The Central Limit Theorem states that for each fixed $z \in \mathbb{R}$, $F_{Z_n}(z)$ converges to $F_Z(z)$.
- In other words, $\lim_{n \rightarrow \infty} P\{Z_n \leq z\} = \lim_{n \rightarrow \infty} F_{Z_n}(z) = F_Z(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$ for all z .

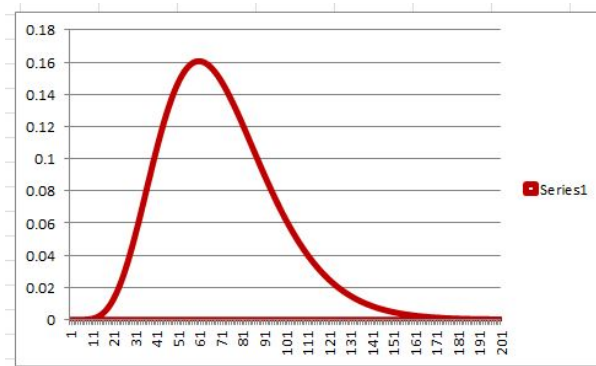
□

7.6 The Gamma Distribution

Whereas the normal distribution is a good fit for histograms which are symmetric, many random phenomena yield **left skewed** (also referred to as **left tailed**) or **right skewed** (also referred to as **right tailed**) histograms which are more appropriately modeled by distributions which themselves also are left or right skewed.



Left skewed distribution



Right skewed distribution

The gamma distribution which we discuss here can be used to generate all kinds of right skewed distributions.

Definition 7.10 (Gamma random variable).

Let $\sigma > 0$ and $-\infty < \mu < \infty$. We say that a random variable Y has a **gamma distribution** with **shape parameter** $\alpha > 0$ and **scale parameter** $\beta > 0$ if its probability density function is

$$(7.40) \quad f_Y(y) = \begin{cases} \frac{y^{\alpha-1} e^{-y/\beta}}{\beta^\alpha \Gamma(\alpha)}, & \text{if } 0 \leq y < \infty, \\ 0, & \text{else,} \end{cases}$$

where $\Gamma(\alpha)$ is the **gamma function**

$$(7.41) \quad \Gamma(\alpha) = \int_0^\infty y^{\alpha-1} e^{-y} dy.$$

We also express that by saying that Y is $\text{gamma}(\alpha, \beta)$. \square

Proposition 7.7. *The gamma function satisfies the following:*

$$(7.42) \quad \Gamma(1) = 1,$$

$$(7.43) \quad \Gamma(\alpha) = (\alpha - 1)\Gamma(\alpha - 1) \quad \text{for all } \alpha > 1,$$

$$(7.44) \quad \Gamma(n) = (n - 1)! \quad \text{for all } n \in \mathbb{N}.$$

PROOF: (7.42) is immediate from $\int_0^\infty e^{-y} dy = -e^{-y} \Big|_0^\infty = 0 - (-1) = 1$.

We obtain (7.43) from integration by parts of $\Gamma(\alpha)$:

$$\begin{aligned}\Gamma(\alpha) &= y^{\alpha-1}(-e^{-y})\Big|_0^\infty + \int_0^\infty (\alpha-1)y^{\alpha-2}e^{-y} dy \\ &= 0 + (\alpha-1) \int_0^\infty y^{(\alpha-1)-1}e^{-y} dy \\ &= (\alpha-1)\Gamma(\alpha-1).\end{aligned}$$

To show (7.44) we observe that repeated application of (7.43) yields

$$\begin{aligned}\Gamma(n) &= (n-1)\Gamma(n-1) \\ &= (n-1)(n-2)\Gamma(n-2) \\ &= (n-1)(n-2)(n-3)\cdots 2\Gamma(2) \\ &= (n-1)(n-2)(n-3)\cdots 2 \cdot 1\Gamma(1).\end{aligned}$$

Since $\Gamma(1) = 1$ by (7.42), it follows that

$$\Gamma(n) = (n-1)(n-2)(n-3)\cdots 2 \cdot 1 = (n-1)!.$$

Proposition 7.8.

If the random variable Y is gamma(α, β) it has MGF

$$(7.45) \quad m_Y(t) = \frac{1}{(1-t\beta)^\alpha} \quad \text{for } t < \frac{1}{\beta}.$$

PROOF: ★ We define

$$(A) \quad \tilde{\beta} := \frac{\beta}{1-t\beta}$$

and observe that $\tilde{\beta} > 0$ for $1-t\beta > 0$, i.e., for $t < 1/\beta$. Further,

$$(B) \quad ty - \frac{y}{\beta} = \frac{(-y+ty\beta)}{\beta} = \frac{-y(1-t\beta)}{\beta} = -y \Big/ \frac{\beta}{(1-t\beta)} = \frac{-y}{\tilde{\beta}}.$$

Thus,

$$\begin{aligned}m_Y(t) &= E(e^{tY}) = \int_0^\infty e^{ty} \left[\frac{y^{\alpha-1}e^{-y/\beta}}{\beta^\alpha\Gamma(\alpha)} \right] dy \\ &= \frac{1}{\beta^\alpha} \int_0^\infty \frac{y^{\alpha-1}}{\Gamma(\alpha)} \exp\left[ty - \frac{y}{\beta}\right] dy \stackrel{(B)}{=} \frac{1}{\beta^\alpha} \int_0^\infty \frac{y^{\alpha-1}e^{-y/\tilde{\beta}}}{\Gamma(\alpha)} dy\end{aligned}$$

Part of (B) is $\frac{-y(1-t\beta)}{\beta} = \frac{-y}{\tilde{\beta}}$. Thus, $(1-t\beta)\tilde{\beta} = \beta$; thus, $\beta^\alpha = (1-t\beta)^\alpha \cdot \tilde{\beta}^\alpha$; thus,

$$m_Y(t) = \frac{1}{(1-t\beta)^\alpha} \cdot \int_0^\infty \frac{y^{\alpha-1}e^{-y/\tilde{\beta}}}{\tilde{\beta}^\alpha\Gamma(\alpha)} dy = \frac{1}{(1-t\beta)^\alpha} \cdot \int_0^\infty \varphi(y) dy.$$

Here, the function $\varphi(y)$ is the PDF of a $\text{gamma}(\alpha, \tilde{\beta})$ random variable. It follows that $\int_0^{\infty} \varphi(y) dy = 1$ and we conclude that $m_Y(t) = 1/(1 - t\beta)^\alpha$. ■

Theorem 7.14 (WMS Ch.04.6, Theorem 4.8).

Let the random variable Y be $\text{gamma}(\alpha, \beta)$ with $\alpha, \beta > 0$. Then

$$E[Y] = \alpha\beta \quad \text{and} \quad \text{Var}[Y] = \alpha\beta^2.$$

PROOF: We obtain those results by differentiating the MGF of Y .

$$\begin{aligned} m_Y(t) &= (1 - \beta t)^{-\alpha} \Rightarrow m'_Y(t) = (-\alpha)(1 - \beta t)^{-\alpha-1}(-\beta) \\ &\Rightarrow m''_Y(t) = (-\alpha)(-\beta)(-\beta)(-\alpha - 1)(1 - \beta t)^{-\alpha-2}. \end{aligned}$$

Thus,

$$\begin{aligned} m'_Y(0) &= (-\alpha)(1 - 0)^{-\alpha-1}(-\beta) = \alpha\beta, \\ m''_Y(0) &= (-\alpha)\beta^2(-\alpha - 1)(1 - 0)^{-\alpha-2} = (-\alpha)^2\beta^2 - (-\alpha)\beta^2 = \alpha^2\beta^2 + \alpha\beta^2. \end{aligned}$$

In other words, $E[Y] = \alpha\beta$ and $E[Y^2] = \alpha\beta^2$. From this,

$$\text{Var}[Y] = E[Y^2] - (E[Y])^2 = (\alpha^2\beta^2 + \alpha\beta^2) - \alpha^2\beta^2 = \alpha\beta^2. \quad \blacksquare$$

Definition 7.11 (Chi-square distribution).

Let $\nu \in \mathbb{N}$. We say that a random variable Y has a **chi-square distribution** with ν **degrees of freedom**, in short, Y is **chi-square with ν df** or Y is **chi-square(ν)**, or Y is $\chi^2(\nu)$, if Y is $\text{gamma}(\nu/2, 2)$. In other words, Y must have a gamma distribution with shape parameter $\nu/2$ and scale parameter 2. □

Theorem 7.15 (WMS Ch.04.6, Theorem 4.9).

A chi-square random variable Y with ν degrees of freedom has expectation and variance

$$E[Y] = \nu \quad \text{and} \quad \text{Var}[Y] = 2\nu.$$

PROOF: This follows from Theorem 7.14 with $\alpha = \nu/2$ and $\beta = 2$. ■

Definition 7.12 (Exponential distribution).

We say that a random variable Y has an **exponential distribution** with parameter $\beta > 0$, in short, Y is **expon(β)**, if it has density function

$$(7.46) \quad f_Y(y) = \begin{cases} \frac{1}{\beta} e^{-y/\beta}, & \text{for } 0 \leq y < \infty, \\ 0, & \text{else. } \square \end{cases}$$

Remark 7.10. In many textbooks exponential random variables are expressed in terms of $\lambda = 1/\beta$. Then its PDF is

$$(7.47) \quad f_Y(y) = \begin{cases} \lambda e^{-\lambda y}, & \text{for } 0 \leq y < \infty, \\ 0, & \text{else. } \square \end{cases}$$

Theorem 7.16.

An exponential random variable Y with parameter β has expectation and variance

$$E[Y] = \beta \quad \text{and} \quad \text{Var}[Y] = \beta^2.$$

PROOF: This follows from Theorem 7.14 with $\alpha = 1$. ■

Proposition 7.9 (Memorylessness of the exponential distribution). *Let Y be an exponential random variable. Let $t > 0$ and $h > 0$. Then*

$$(7.48) \quad P\{Y > t+h \mid Y > t\} = P\{Y > h\}.$$

PROOF: From the definition of conditional probability and

$$\{Y > t+h\} \cap \{Y > t\} = \{Y > t+h\},$$

it follows that

$$P\{Y > t+h \mid Y > t\} = \frac{P\{Y > t+h\}}{P\{Y > t\}}.$$

We obtain

$$P\{Y > t+h\} = \int_{t+h}^{\infty} \frac{1}{\beta} e^{-y/\beta} dy = -\frac{1}{1/\beta} \cdot \frac{1}{\beta} \cdot e^{-y/\beta} \Big|_{t+h}^{\infty} = -e^{-y/\beta} \Big|_{t+h}^{\infty} = e^{-(t+h)/\beta}$$

and

$$P\{Y > t\} = \int_t^{\infty} \frac{1}{\beta} e^{-y/\beta} dy = -e^{-y/\beta} \Big|_t^{\infty} = e^{-t/\beta}.$$

Thus,

$$P\{Y > t+h \mid Y > t\} = \frac{e^{-(t+h)/\beta}}{e^{-t/\beta}} = e^{-h/\beta} = P\{Y > h\}. \quad \blacksquare$$

Remark 7.11. The property (7.48) of an exponential distribution is referred to as the **memoryless property** of the exponential distribution. It also occurs in the geometric distribution. □

7.7 The Beta Distribution

This chapter is merely a summary of the most important material of WMS Chapter 4.7 (The Beta Probability Distribution).

Like the gamma PDF, the beta density function is a two-parameter PDF defined over the closed interval $0 \leq y \leq 1$. y often plays the role of a proportion, such as the proportion of impurities in a chemical product or the proportion of time that a machine is under repair.

Definition 7.13 (Beta distribution). ★

A random variable Y has a **beta probability distribution** with parameters $\alpha > 0$ and $\beta > 0$ if it has density function

$$(7.49) \quad f_Y(y) = \begin{cases} \frac{y^{\alpha-1}(1-y)^{\beta-1}}{B(\alpha, \beta)}, & \text{if } 0 \leq y \leq 1, \\ 0, & \text{else,} \end{cases}$$

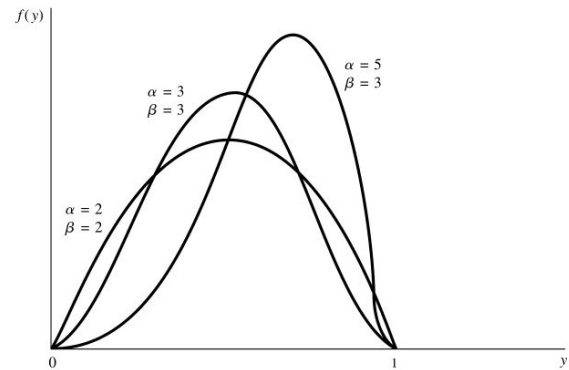
where

$$(7.50) \quad B(\alpha, \beta) = \int_0^1 y^{\alpha-1}(1-y)^{\beta-1} dy = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}.$$

We also express that by saying that Y is beta(α, β). \square

Beta density functions come in a large variety of shapes for different values of α and β . Some of these are shown in the figure to the right.

Note that $0 \leq y \leq 1$ does not restrict the use of the beta distribution. It can be applied to a random variable defined on the interval $c \leq y \leq d$ by means of the transformation $\tilde{y} = (y - c)/(d - c)$ which defines a new variable $0 \leq \tilde{y} \leq 1$ which has the correct domain for the beta density.



Beta density functions. Source: WMS

Theorem 7.17. ★

If Y is a beta-distributed random variable with parameters $\alpha > 0$ and $\beta > 0$, then

$$E[Y] = \frac{\alpha}{\alpha + \beta} \quad \text{and} \quad \text{Var}[Y] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}.$$

PROOF: See the WMS text \blacksquare

7.8 Inequalities for Probabilities

This chapter lists some very useful estimates for probabilities which involve the moments of a random variable. Among them is the Tchebysheff inequality.

Theorem 7.18. ★

Let Y, Z be continuous or discrete random variables and $a > 0$. Assume further that $Y \geq 0$. Then

$$(7.51) \quad P\{Y \geq a\} \leq \frac{E[Y]}{a},$$

$$(7.52) \quad P\{|Z| \geq a\} \leq \frac{E[|Z|^n]}{a^n}.$$

(7.51) is known as the **Markov inequality**

PROOF of (7.51):⁴⁰ We give the proof for continuous random variables. The discrete case is even simpler since it involves summation instead of integration.

Let $f_Y(y)$ be the PDF of Y . We observe the following:

- (a) $Y \geq 0$ implies $y f_Y(y) = 0$ for $-\infty < y < 0$.
- (b) $y f_Y(y) \geq 0$ for $0 \leq y < \infty$.
- (c) $y f_Y(y) \geq a f_Y(y)$ for $a \leq y < \infty$.

Thus,

$$\begin{aligned} E[Y] &= \int_{-\infty}^{\infty} y f_Y(y) dy \stackrel{\text{(a)}}{=} \int_0^{\infty} y f_Y(y) dy = \int_0^a y f_Y(y) dy + \int_a^{\infty} y f_Y(y) dy \\ &\stackrel{\text{(b)}}{\geq} \int_a^{\infty} y f_Y(y) dy \stackrel{\text{(c)}}{\geq} \int_a^{\infty} a f_Y(y) dy = a \int_a^{\infty} f_Y(y) dy = a P\{Y \geq a\}. \end{aligned}$$

We divide by $a > 0$ and obtain (7.51).

PROOF of (7.52): Since $|Z|^n \geq 0$ and $a^n > 0$, we can apply (7.51) with $|Z|^n$ in place of Y and a^n in place of a :

$$(A) \quad P\{|Z|^n \geq a^n\} \leq \frac{E[|Z|^n]}{a^n}.$$

Since the function $x \mapsto x^n$ is (strictly) increasing, $|Z(\omega)|^n \geq a^n \Leftrightarrow |Z(\omega)| \geq a$.

Thus, (A) yields $P\{|Z| \geq a\} \leq E[|Z|^n]/a^n$ and this proves (7.52). ■

The work we have done here allows us to quickly prove the Tchebysheff inequalities in the form listed in WMS Chapter 4.10 (Tchebysheff's Theorem).

Theorem 7.19 (Tchebysheff Inequalities).

Let Y be a random variable with mean $\mu = E[Y]$ and standard deviation $\sigma = \sqrt{\text{Var}[Y]}$. Let $k > 0$. Then

$$(7.53) \quad P\{|Y - \mu| \geq k\sigma\} \leq \frac{1}{k^2},$$

$$(7.54) \quad P\{|Y - \mu| < k\sigma\} \geq 1 - \frac{1}{k^2}.$$

Both (7.53) and (7.54) are known as the **Tchebysheff inequalities**

⁴⁰Source: https://en.wikipedia.org/wiki/Markov%27s_inequality

PROOF: We apply (7.52) with $n = 2$, $Y - \mu$ in place of Z , and $k\sigma$ in place of a . We obtain

$$P\{|Y - \mu| \geq k\sigma\} \leq \frac{E[|Y - \mu|^2]}{(k\sigma)^2} = \frac{E[(Y - \mu)^2]}{(k\sigma)^2} = \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2}.$$

This proves (7.53). Since the event $\{|Y - \mu| < k\sigma\}$ is the complement of the event $\{|Y - \mu| \geq k\sigma\}$, (7.54) follows. ■

Remark 7.12. Some comments about the Tchebysheff inequalities:

- (a) Both inequalities state the same, since the events $\{|Y - \mu| < c\sigma\}$ and $\{|Y - \mu| \geq c\sigma\}$ are complements of each other. We had noted this in the proof of Theorem 7.19.
- (b) The inequalities are not particularly powerful, but consider that they are universally valid, regardless of any particulars concerning Y !
- (c) If we write $a := k\sigma$ and thus, $k = a/\sigma$, we obtain

$$P\{|Y - \mu| < a\} \geq 1 - \frac{\text{Var}[Y]}{a^2} \quad \text{and} \quad P\{|Y - \mu| \geq a\} \leq \frac{\text{Var}[Y]}{a^2}. \quad \square$$

Example 7.3. The screws produced by ACME Co. follow a distribution with a mean of $\mu = 18.40$ mm and a variance of $\sigma^2 = 0.64$ mm². In other words, the length Y of a randomly picked screw (a sample of size 1) has $E[Y] = 18.40$ and $\text{Var}[Y] = 0.64$.

A screw can only be sold if its length is within 17.20 and 19.60 mm. How likely is it that a screw is produced that cannot be sold?

Solution: We observe that $E[Y] = 18.40$ is the midpoint of the interval $[17.20, 19.60]$ and that

- a screw cannot be sold $\Leftrightarrow Y(\omega) \notin [17.20, 19.60] \Leftrightarrow |Y(\omega) - E[Y]| > (17.20, 19.60)/2 = 1.2$.

We solve

$$k\sigma = |Y - E[Y]| = 1.2, \quad \text{i.e.,} \quad \sqrt{0.64}k = 0.8k = 1.2,$$

for k and obtain $k = 1.2/0.8 = 3/2$. Thus, $k^2 = 9/4$.

Tchebysheff's inequality (7.54) then yields the following upper bound for the probability of obtaining a sample with a difference $\bar{Y}(\omega) - \mu$ as large as or even larger than the one we have sampled:

$$P\{|Y - \mu| > k\sigma\} \leq \frac{1}{k^2} = 4/9.$$

This example demonstrates the low quality of the bounds that we obtain from Tchebysheff's inequalities. For example, let us assume we know that Y follows a normal distribution, i.e.,

$$Y \sim \mathcal{N}(\mu = 18.40, \sigma^2 = 0.64),$$

then we can deduce from the empirical rule (the 68%–95%–99.7% rule)⁴¹ that

$$\begin{aligned} 0.32 &= 1 - 0.68 \approx P\{|Y - \mu| > 1 \cdot \sigma\} \\ &\geq P\{|Y - \mu| > 1.5\sigma\} \\ &\geq P\{|Y - \mu| > 2\sigma\} \approx 1 - 0.95 = 0.05. \end{aligned}$$

⁴¹see the introduction to subch.7.5: The Normal Probability Distribution

Thus, higher precision calculations show that the more likely event of $Y(\omega)$ not being within one standard deviation of 18.40 mm only has a probability of 0.32, substantially less than our overly generous estimate of $4/9 = 0.44\bar{4}$ for the less likely event of being within 1.5 standard deviations. By the way, the exact figure (in the case of $Y \sim \mathcal{N}(18.40, 0.64)$) is $P\{|Y - \mu| > 1.5\sigma\} \approx 0.1336$. This is less than one third of the Tchebysheff estimate. \square

Example 7.4. It has been established some time ago that the data in the population of interest follow a distribution with a mean of $\mu = 18.40$. In other words, a random pick Y (a sample of size 1) from that population has $E[Y] = 18.40$. There have been concerns that the composition of the population has changed significantly and μ with it. An SRS (simple random sample) is drawn from that population and mean and variance are estimated from the realization of this sample as

$$\bar{Y}(\omega) = 17.60 \quad \text{and} \quad S^2(\omega) = 6.25.^{42}$$

Is the deviation of $\bar{Y}(\omega)$ from μ big enough to discard $\mu = 18.40$ and go through the process of establishing a new population mean?

Solution: We use $S^2 = 6.25$ for $\sigma^2 = \text{Var}[Y]$. Then $\sigma = \sqrt{6.25} = 2.5$. We solve

$$k\sigma = |\bar{Y} - E[Y]|, \quad \text{i.e.,} \quad 0.25k = |17.60 - 18.40| = 0.8,$$

for k and obtain $k = 3.2$. Thus, $k^2 = 10.24$. Since $E[\bar{Y}] = E[Y]$ it follows from Tchebysheff's inequality (7.54) that the probability of obtaining a sample with a difference $\bar{Y}(\omega) - E[Y]$ as large as or even larger than the one of the sample we have drawn, is

$$P\{|Y - \mu| < k\sigma\} \geq 1 - \frac{1}{k^2} = 1 - \frac{1}{10.24} = 0.902344.$$

This probability is very large and shows that our sample mean $\bar{Y} = 17.60$ does not contradict the assumption that the population mean 18.40 \square

⁴² $\bar{Y} = 17.60$ is the so called sample mean (see Example 8.5: Variance of the sample mean on p.163) and

$S^2 = S^2(\omega) = \frac{1}{n-1} \left(\sum_{j=1}^n (Y_j(\omega) - \bar{Y}(\omega))^2 \right)$ is the so called sample variance which will be introduced in subchapter 10.3 (Sampling Distributions) of Chapter 10(Limit Theorems). See Definition 10.4: Sample variance on p.225.

8 Multivariate Probability Distributions

Like the previous chapter, this one is extremely skeletal in nature. It contains very few examples. You are reminded again that you must work through the corresponding chapters in the WMS text. In this case, that would be WMS Chapter 5 (Multivariate Probability Distributions).

8.1 Multivariate CDFs, PMFs and PDFs

Assumption 8.1 (Comma separation denotes intersection). We will follow the following convention for the notation of events that are generated by random variables or random elements $X, Y, Z \dots$

Separating commas are to be interpreted as “and” and not as “or”. Thus, for example,

$$\begin{aligned} \{X \in B, Y = \alpha, 5 \leq Z < 8\} &= \{X \in B \text{ and } Y = \alpha \text{ and } 5 \leq Z < 8\} \\ &= \{X \in B\} \cap \{Y = \alpha\} \cap \{5 \leq Z < 8\}. \quad \square \end{aligned}$$

Definition 8.1 (Joint cumulative distribution function).

Given are two random variables Y_1 and Y_2 . No assumption is made whether they are discrete or continuous. We call

$$(8.1) \quad F(y_1, y_2) := F_{Y_1, Y_2}(y_1, y_2) := P(Y_1 \leq y_1, Y_2 \leq y_2), \quad \text{where } y_1, y_2 \in \mathbb{R},$$

the **joint cumulative distribution function** or **bivariate cumulative distribution function** or **joint CDF** or **joint distribution function** of Y_1 and Y_2 . \square

Theorem 8.1.

Let Y_1 and Y_2 be random variables with joint CDF $F_{Y_1, Y_2}(y_1, y_2)$. Further, assume that $\vec{a} := (a_1, a_2) \in \mathbb{R}^2$ and $\vec{b} := (b_1, b_2) \in \mathbb{R}^2$ satisfy $a_1 < b_1$ and $a_2 < b_2$. Then,

$$(8.2) \quad F_{Y_1, Y_2}(-\infty, -\infty) = F_{Y_1, Y_2}(-\infty, y_2) = F_{Y_1, Y_2}(y_1, -\infty) = 0.$$

$$(8.3) \quad F_{Y_1, Y_2}(\infty, \infty) = 1,$$

$$(8.4) \quad \begin{aligned} P\{a_1 < Y_1 \leq b_1, a_2 < Y_2 \leq b_2\} &= F_{Y_1, Y_2}(b_1, b_2) - F_{Y_1, Y_2}(a_1, b_2) \\ &\quad - F_{Y_1, Y_2}(b_1, a_2) + F_{Y_1, Y_2}(a_1, a_2), \end{aligned}$$

$$(8.5) \quad F_{Y_1, Y_2}(b_1, b_2) - F_{Y_1, Y_2}(a_1, b_2) - F_{Y_1, Y_2}(b_1, a_2) + F_{Y_1, Y_2}(a_1, a_2) \geq 0,$$

PROOF:

(8.2) follows from

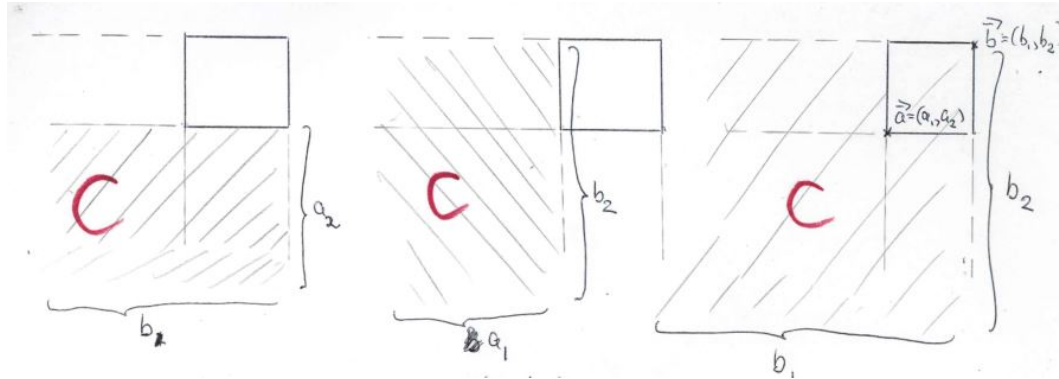
$$P\{Y_1 < -\infty\} = P\{Y_2 < -\infty\} = 0.$$

(8.3) follows from

$$P\{Y_1 < \infty, Y_2 < \infty\} = P(\Omega) = 1.$$

(8.5) is immediate from (8.4).

Finally, for the proof of (8.4), we see from the three pictures below the following:



- $P\{a_1 < Y_1 \leq b_1, a_2 < Y_2 \leq b_2\} \hat{=}$ black rectangle in the upper right corner
- $F_{Y_1, Y_2}(b_1, b_2) \hat{=}$ shaded area in the right drawing
- $F_{Y_1, Y_2}(b_1, a_2) \hat{=}$ shaded area (below black rectangle) in the left drawing
- $F_{Y_1, Y_2}(b_1, a_2) \hat{=}$ shaded area (to left of black rectangle) in the middle drawing
- $F_{Y_1, Y_2}(a_1, a_2) \hat{=}$ area marked with a red C

The expression $F_{Y_1, Y_2}(b_1, b_2) - F_{Y_1, Y_2}(a_1, b_2) - F_{Y_1, Y_2}(b_1, a_2)$ would correspond to the black rectangle, except that we subtracted the red C area twice. We add $F_{Y_1, Y_2}(a_1, a_2)$ to compensate. ■

Definition 8.2 (Joint probability mass function).

Let Y_1 and Y_2 be discrete random variables. We call

$$(8.6) \quad p(y_1, y_2) := p_{Y_1, Y_2}(y_1, y_2) := P\{Y_1 = y_1, Y_2 = y_2\}, \quad \text{where } y_1, y_2 \in \mathbb{R},$$

the **joint probability mass function** or **bivariate probability mass function** or **joint PMF** of Y_1 and Y_2 . □

Just as in the univariate case, $p_{Y_1, Y_2}(y_1, y_2)$ assigns nonzero probabilities to only finitely or countably many pairs of values (y_1, y_2) . As in the univariate case, by definition,

$$\sum_{(y_1, y_2) \in B} p_{Y_1, Y_2}(y_1, y_2) = \sum_{\substack{(y_1, y_2) \in B, \\ p_{Y_1, Y_2}(y_1, y_2) > 0}} p_{Y_1, Y_2}(y_1, y_2).$$

Proposition 8.1 (WMS Ch.05.2, Theorem 5.1).

If Y_1 and Y_2 are discrete random variables with joint PMF $p_{Y_1, Y_2}(y_1, y_2)$, then

- (1) $p_{Y_1, Y_2}(y_1, y_2) \geq 0$ for all $y_1, y_2 \in \mathbb{R}$,
- (2) $\sum_{y_1, y_2} p_{Y_1, Y_2}(y_1, y_2) = 1$.
- (3) $F_{Y_1, Y_2}(y_1, y_2) = \sum_{u_1 \leq y_1, u_2 \leq y_2} p_{Y_1, Y_2}(u_1, u_2) = \sum_{u_1 \leq y_1} \sum_{u_2 \leq y_2} p_{Y_1, Y_2}(u_1, u_2)$.

PROOF: Obvious. ■

Definition 8.3 (Jointly continuous random variables).

Let Y_1 and Y_2 be random variables with joint CDF $F(y_1, y_2)$. We call Y_1 and Y_2 **jointly continuous** if $F(y_1, y_2)$ is a continuous function of both arguments. □

Assumption 8.2 (Jointly continuous random variables have PDFs). We will follow the following convention for the notation of events that are generated by random variables or random elements $X, Y, Z \dots$

We assume for all jointly continuous random variables Y_1 and Y_2 that $\frac{\partial^2 F_{Y_1, Y_2}}{\partial y_1 \partial y_2}$ exists and is continuous except for $(y_1, y_2) \in B^*$, where the set $B^* \subseteq \mathbb{R}^2$ satisfies that $B^* \cap B$ is finite for any bounded subset $B \in \mathbb{R}^2$ (bounded sets are those contained in a circle with sufficiently large radius).

This assumption guarantees for all $y_1, y_2 \in \mathbb{R}$, when we write f_{Y_1, Y_2} for $\frac{\partial^2 F_{Y_1, Y_2}}{\partial y_1 \partial y_2}$, that

$$\begin{aligned}
 (8.7) \quad F_{Y_1, Y_2}(y_1, y_2) &= \int_{-\infty}^{y_1} \int_{-\infty}^{y_2} f_{Y_1, Y_2}(u_1, u_2) \, du_2 \, du_1 \\
 &= \int_{-\infty}^{y_2} \int_{-\infty}^{y_1} f_{Y_1, Y_2}(u_1, u_2) \, du_1 \, du_2 . \\
 &= \iint_{]-\infty, y_1[\times]-\infty, y_2]} f_{Y_1, Y_2}(u_1, u_2) \, du_1 \, du_2 . \quad \square
 \end{aligned}$$

Definition 8.4 (WMS Ch.05.2, Definition 5.3).

Let Y_1 and Y_2 be continuous random variables with joint distribution function $F(y_1, y_2)$ and second derivative $f_{Y_1, Y_2}(y_1, y_2) = \frac{\partial^2 F_{Y_1, Y_2}}{\partial y_1 \partial y_2}(y_1, y_2)$. We call $f_{Y_1, Y_2}(y_1, y_2)$ the **joint probability density function** or **joint PDF** of Y_1 and Y_2 . □

Theorem 8.2.

Let Y_1 and Y_2 be jointly continuous random variables with joint PDF $f_{Y_1, Y_2}(y_1, y_2)$, then

- (1) $f_{Y_1, Y_2}(y_1, y_2) \geq 0$ for all y_1, y_2 .
- (2) $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{Y_1, Y_2}(y_1, y_2) \, dy_1 \, dy_2 = 1$.

PROOF: An easy consequence of Theorem 8.1 on p.143. ■

8.2 Marginal and Conditional Probability Distributions

Definition 8.5 (Marginal distribution of two random variables).

Let $\vec{Y} = (Y_1, Y_2)$ be a vector of two random variables with joint distribution
 $(B_1, B_2) \mapsto P_{Y_1, Y_2}(B_1, B_2) = P\{Y_1 \in B_1, Y_2 \in B_2\}$, where $B_1, B_2 \subseteq \mathbb{R}$.

We call the probability measures

$$(8.8) \quad Q_1 : B_1 \mapsto P_{Y_1, Y_2}(B_1, \mathbb{R}) \quad \text{and} \quad Q_2 : B_2 \mapsto P_{Y_1, Y_2}(\mathbb{R}, B_2)$$

the **marginal distributions** of $\vec{Y} = (Y_1, Y_2)$. \square

Proposition 8.2.

The marginal distributions of $\vec{Y} = (Y_1, Y_2)$ are the distributions P_{Y_1} and P_{Y_2} of the coordinates Y_1 and Y_2 . In other words, $Q_1 = P_{Y_1}$ and $Q_2 = P_{Y_2}$

PROOF: Since, $Y_1(\omega) \in B \Leftrightarrow Y_1(\omega) \in B$ and $Y_2(\omega) \in \mathbb{R}$ holds for all $B \subseteq \mathbb{R}$, we obtain

$$Q_1(B) = P_{Y_1, Y_2}(B, \mathbb{R}) = P\{Y_1 \in B, Y_2 \in \mathbb{R}\} = P\{Y_1 \in B\} = P_{Y_1}(B), \quad \text{whenever } B \subseteq \mathbb{R}.$$

Thus, $Q_1 = P_{Y_1}$. We obtain in a similar fashion from $Y_2(\omega) \in B \Leftrightarrow Y_1(\omega) \in \mathbb{R}$ and $Y_2(\omega) \in B$, that

$$Q_2(B) = P_{Y_2}(B), \quad \text{for all } B \subseteq \mathbb{R}. \quad \blacksquare$$

Henceforth, we will retire the symbols Q_1, Q_2 and denote the marginal distributions of $\vec{Y} = (Y_1, Y_2)$ by P_{Y_1} and P_{Y_2} .

Definition 8.5 translates for discrete random variables, whose distribution is determined by their joint PMF and for continuous random variables, whose distribution is determined by their joint PDF, to the following.

Definition 8.6 (Marginal PMF and PDF).

(a) Let Y_1 and Y_2 be discrete random variables with joint PMF $p_{Y_1, Y_2}(y_1, y_2)$. We call

$$(8.9) \quad p_{Y_1}(y_1) = \sum_{\text{all } y_2} p_{Y_1, Y_2}(y_1, y_2) \quad \text{and} \quad p_{Y_2}(y_2) = \sum_{\text{all } y_1} p_{Y_1, Y_2}(y_1, y_2)$$

the **marginal probability mass functions** or **marginal PMFs** of Y_1 and Y_2 .

(b) Let Y_1 and Y_2 be continuous random variables with joint PDF $f_{Y_1, Y_2}(y_1, y_2)$. We call

$$(8.10) \quad f_{Y_1}(y_1) = \int_{-\infty}^{\infty} f_{Y_1, Y_2}(y_1, y_2) dy_2 \quad \text{and} \quad f_{Y_2}(y_2) = \int_{-\infty}^{\infty} f_{Y_1, Y_2}(y_1, y_2) dy_1.$$

the **marginal density functions** or **marginal PDFs** of Y_1 and Y_2 . \square

Remark 8.1. We recall Definition 3.7 of $P(A | B)$, the probability of the event A conditioned on the event B , which is defined for $P(B) > 0$ as

$$P(A | B) = \frac{P(A \cap B)}{P(B)}.$$

We also recall that, if $P(B) > 0$, the set function $A \mapsto P(A | B)$ is a probability measure on Ω . See Theorem 3.8 on p.55. We replace the general events A and B with events $\{Y_1 = y_1\}$ and $\{Y_2 = y_2\}$ and obtain, if $P\{Y_2 = y_2\} > 0$,

$$(8.11) \quad P\{Y_1 = y_1 | Y_2 = y_2\} = \frac{P\{Y_1 = y_1, Y_2 = y_2\}}{P\{Y_2 = y_2\}}.$$

As we always do for conditional probabilities, we interpret (8.11) as the probability that the random variable Y_1 equals y_1 , given that Y_2 equals y_2 .

Not much can be done with formula (8.11) for continuous random variables Y_1 and Y_2 , because $P\{Y_2 = y_2\} = 0$ for all $y_2 \in \mathbb{R}$; but it shows us how to define conditional PMFs for discrete random variables. \square

Definition 8.7 (Conditional probability mass function).

Let Y_1 and Y_2 be discrete random variables with joint PMF $p_{Y_1, Y_2}(y_1, y_2)$ and marginal PMFs $p_{Y_1}(y_1)$ and $p_{Y_2}(y_2)$. Then we call

$$(8.12) \quad p_{Y_1|Y_2}(y_1 | y_2) := \begin{cases} P\{Y_1 = y_1 | Y_2 = y_2\}, & \text{if } P\{Y_2 = y_2\} > 0, \\ \text{undefined}, & \text{if } P\{Y_2 = y_2\} = 0, \end{cases}$$

the **conditional probability mass function** or the **conditional PMF** of Y_1 given Y_2 .

Likewise, we call

$$(8.13) \quad p_{Y_2|Y_1}(y_2 | y_1) := \begin{cases} P\{Y_2 = y_2 | Y_1 = y_1\}, & \text{if } P\{Y_1 = y_1\} > 0, \\ \text{undefined}, & \text{if } P\{Y_1 = y_1\} = 0, \end{cases}$$

the **conditional PMF** of Y_2 given Y_1 . \square

Remark 8.2. Note that conditional PMFs can be expressed in terms of joint PMF and marginal PMFs:

$$(8.14) \quad p_{Y_1|Y_2}(y_1 | y_2) = \frac{p_{Y_1, Y_2}(y_1, y_2)}{p_{Y_2}(y_2)} \quad \text{if } p_{Y_2}(y_2) > 0,$$

$$(8.15) \quad p_{Y_2|Y_1}(y_2 | y_1) = \frac{p_{Y_1, Y_2}(y_1, y_2)}{p_{Y_1}(y_1)} \quad \text{if } p_{Y_1}(y_1) > 0. \quad \square$$

The author does not think that there is much use for the next definition (WMS Ch.05.3, Definition 5.6) because all jointly continuous random variables come with PDF

$$f_{Y_1, Y_2}(y_1, y_2) = \frac{\partial^2 F_{Y_1, Y_2}}{\partial y_1 \partial y_2}(y_1, y_2).$$

It is included only for the sake of completeness.

Definition 8.8. \star Let Y_1 and Y_2 be two jointly continuous random variables. Then,

$$(8.16) \quad F_{Y_1|Y_2}(y_1 | y_2) := P(Y_1 \leq y_1 | Y_2 = y_2) := \int_{-\infty}^{y_1} \frac{f_{Y_1, Y_2}(u_1, y_2)}{f_{Y_2}(y_2)} du_1$$

defines the **conditional distribution function** or **conditional CDF** of Y_1 given $Y_2 = y_2$. \square

Definition 8.9 (Conditional probability density function).

Let Y_1 and Y_2 be continuous random variables with joint PDF $f_{Y_1|Y_2}(y_1, y_2)$ and marginal densities $f_{Y_1}(y_1)$ and $f_{Y_2}(y_2)$. Then we call

$$(8.17) \quad f_{Y_1|Y_2}(y_1 | y_2) := \begin{cases} \frac{f_{Y_1, Y_2}(y_1, y_2)}{f_{Y_2}(y_2)}, & \text{if } f_{Y_2}(y_2) > 0, \\ \text{undefined}, & \text{if } f_{Y_2}(y_2) = 0, \end{cases}$$

the **conditional probability density function** or the **conditional PDF** of Y_1 given Y_2 .

Likewise, we call

$$(8.18) \quad f_{Y_2|Y_1}(y_2 | y_1) := \begin{cases} \frac{f_{Y_1, Y_2}(y_1, y_2)}{f_{Y_1}(y_1)}, & \text{if } f_{Y_1}(y_1) > 0, \\ \text{undefined}, & \text{if } f_{Y_1}(y_1) = 0, \end{cases}$$

the **conditional PDF** of Y_2 given Y_1 . \square

8.3 Independence of Random Variables and Discrete Random Elements

Introduction 8.1. Let $X_1, X_2 : (\Omega, P) \rightarrow \Omega'$ be two random elements (recall that they are random variables if $\Omega' = \mathbb{R}$). Not all events $A \subseteq \Omega$ are meaningful for X_1 and X_2 . Rather, only **events generated by** X_1 and by X_2 , i.e., events of the form $\{X_1 \in B_1\}$ and $\{X_2 \in B_2\}$ for suitable $B_1, B_2 \subseteq \Omega'$ will matter.

Since independence of two events A_1 and A_2 is defined by $P(A_1 \cap A_2) = P(A_1)P(A_2)$, the proper way to define independence of X_1 and X_2 seems to be

$$(8.19) \quad P\{X_1 \in B_1, X_2 \in B_2, \} = P\{X_1 \in B_1\} \cdot P\{X_2 \in B_2, \} \quad \text{for all relevant } B_1, B_2 \subseteq \Omega'.$$

What are the relevant sets B_j ? We answer that question for discrete random elements (hence, also for discrete random variables) and for continuous random variables.

(a) Assume that $X : (\Omega, P) \rightarrow \Omega'$ is a discrete random element with PMF $p_X(x)$. In other words, there is a countable $\Omega^* \subseteq \Omega'$ such that, for any $B \subseteq \Omega'$,

$$P\{X \in B\} = P_X(B) = \sum_{x \in \Omega^* \cap B} p_X(x) = \sum_{x \in B} p_X(x) = \sum_{x \in B} P\{X = x\}.$$

These equations show that the distribution of X is determined by the events $\{X = x\}$. Thus, the relevant sets for X are of the form $B = \{x\}$.

(b) Assume that Y is a continuous random variable on (Ω, P) with PDF $f_Y(y)$. Then the probabilities for the events that matter, the events $\{a < Y \leq b\}$ where $a < b$, are

$$P\{a < Y \leq b\} = \int_a^b f_Y(y) dy.$$

(See (7.4) in theorem 7.2 on p.117.) This equation shows that the distribution of Y is determined by the probability density function $f_Y(y)$. Thus, the relevant sets for Y are the intervals $B =]a, b]$.⁴³

⁴³Since $P\{X = a\} = 0$ for all $a \in \mathbb{R}$, it does not matter whether we do or do not include the end points. See Proposition 7.1 on p.116.

In summary, we could define independence of discrete random elements X_1 and X_2 as

$$P\{X_1 = x_1, X_2 = x_2, \} = P\{X_1 = x_1\} \cdot P\{X_2 = x_2, \} \quad \text{for all } x_1, x_2 \in \Omega'.$$

Equivalently, this can be expressed as

$$(8.20) \quad p_{X_1, X_2}(x_1, x_2) = p_{X_1}(x_1) \cdot p_{X_2}(x_2) \quad \text{for all } x_1, x_2 \in \Omega'.$$

Moreover, independence of continuous random variables Y_1 and Y_2 could be defined as

$$P\{a < X_1 \leq b, c < X_2 \leq d\} = P\{a < X_1 \leq b\} \cdot P\{c < X_2 \leq d\} \quad \text{for all } a < b \text{ and } c < d.$$

Equivalently, this can be expressed as

$$(8.21) \quad \int_a^b \int_c^d f_{Y_1, Y_2}(y_1, y_2) dy_2 dy_1 = \int_a^b f_{Y_1}(y_1) dy_1 \cdot \int_c^d f_{Y_2}(y_2) dy_2 \quad \text{for all } a < b \text{ and } c < d.$$

The CDF (cumulative distribution function) $F_Y(y)$ gives us for both discrete and continuous random variables (but we must exclude discrete random elements) a unified way to express what was stated in **(a)** and **(b)** as follows.

In the discrete case **(a)** we have

$$P\{Y = y\} = P\{Y \leq y\} - P\{Y < y\} = F_Y(y) - F_Y(y-).$$

Here $F_Y(y-) = \lim_{a < y, a \rightarrow y} F_Y(a)$ is the left-hand limit of the (monotone) function $F_Y(\cdot)$ at y .

In the continuous case **(b)** we have

$$P\{a < Y \leq b\} = P\{Y \leq b\} - P\{Y \leq a\} = F_Y(b) - F_Y(a).$$

In both cases, independence of Y_1 and Y_2 can now be defined as

$$(8.22) \quad F_{Y_1, Y_2}(y_1, y_2) = F_{Y_1}(y_1) \cdot F_{Y_2}(y_2) \quad \text{for all } y_1, y_2 \in \mathbb{R}. \quad \square$$

We make (8.22) the basis for the definition of independence of random variables.

Definition 8.10 (Independent random variables).

Let Y_1 and Y_2 be random variables with CDFs $F_{Y_1}(y_1)$ and $F_{Y_2}(y_2)$ and with joint CDF $F_{Y_1, Y_2}(y_1, y_2)$. We call Y_1 and Y_2 **independent** if

$$(8.23) \quad F_{Y_1, Y_2}(y_1, y_2) = F_{Y_1}(y_1) \cdot F_{Y_2}(y_2) \quad \text{for all } y_1, y_2 \in \mathbb{R}.$$

If Y_1 and Y_2 are not independent, we call them **dependent**.

We must treat discrete random elements separately since there are no CDFs.

Let X_1 and X_2 be discrete random elements with PMFs $p_{X_1}(x_1)$ and $p_{X_2}(x_2)$ and with joint PMF $p_{X_1, X_2}(x_1, x_2)$. We call X_1 and X_2 **independent** if

$$(8.24) \quad p_{X_1, X_2}(x_1, x_2) = p_{X_1}(x_1) \cdot p_{X_2}(x_2) \quad \text{for all } x_1, x_2 \in \mathbb{R}.$$

If X_1 and X_2 are not independent, we call them **dependent**. \square

Theorem 8.3 (Functions of independent random variables are independent).

Let $\vec{Y} = (Y_1, \dots, Y_k) : (\Omega, P) \rightarrow \mathbb{R}$ be a vector of k independent random variables and $h_j : \mathbb{R} \rightarrow \mathbb{R}$.

- Then the random variables $h_1 \circ Y_1, \dots, h_k \circ Y_k$ also are independent.

PROOF: We recall (3.37) of Proposition 3.7 (Preimages of function composition) on p.64: Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ and $W \subseteq Z$. Then

$$(A) \quad (g \circ f)^{-1} = f^{-1} \circ g^{-1}, \text{ i.e., } (g \circ f)^{-1}(W) = f^{-1}(g^{-1}(W)).$$

We use this twice in the following calculations.

$$\begin{aligned} P\{h_j \circ Y_j \in B_j, (j = 1, \dots, n)\} &= P\{(h_j \circ Y_j)^{-1}(B_j), (j = 1, \dots, n)\} \\ &\stackrel{(A)}{=} P\{Y_j^{-1} \circ h_j^{-1}(B_j), (j = 1, \dots, n)\} = P\{Y_j \in h_j^{-1}(B_j), (j = 1, \dots, n)\}. \end{aligned}$$

Since the Y_j are independent, the product rule holds. We obtain

$$\begin{aligned} P\{h_j \circ Y_j \in B_j, (j = 1, \dots, n)\} &= \prod_j P\{Y_j \in h_j^{-1}(B_j)\} = \prod_j P\{Y_j^{-1} \circ h_j^{-1}(B_j)\} \\ &\stackrel{(A)}{=} \prod_j P\{\prod_j P\{(h_j \circ Y_j)^{-1}(B_j)\}\} = \prod_j P\{h_j \circ Y_j \in B_j\}. \quad \blacksquare \end{aligned}$$

Theorem 8.4 (WMS Ch.05.4, Theorem 5.4).

If Y_1 and Y_2 are discrete random variables with joint PMF $p_{Y_1, Y_2}(y_1, y_2)$ and marginal PMFs $p_{Y_1}(y_1)$ and $p_{Y_2}(y_2)$, then

$$(8.25) \quad Y_1, Y_2 \text{ are independent} \Leftrightarrow p_{Y_1, Y_2}(y_1, y_2) = p_{Y_1}(y_1) \cdot p_{Y_2}(y_2) \text{ for all } y_1, y_2 \in \mathbb{R}.$$

If Y_1 and Y_2 are continuous random variables with joint PDF $f_{Y_1, Y_2}(y_1, y_2)$ and marginal PDFs $f_{Y_1}(y_1)$ and $f_{Y_2}(y_2)$, then

$$(8.26) \quad Y_1, Y_2 \text{ are independent} \Leftrightarrow f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) \cdot f_{Y_2}(y_2) \text{ for all } y_1, y_2 \in \mathbb{R}.$$

PROOF: We only prove here the \Rightarrow directions of (8.25) (8.26). The proof of the opposite direction is left as an exercise to the reader.

We apply (??) of Theorem ?? on p.?? and (8.23) of Definition 8.10 (Independent random variables) on p.149 as follows.

$$\begin{aligned} &P\{a_1 < Y_1 \leq y_1, a_2 < Y_2 \leq y_2\} \\ &\stackrel{(9.18)}{=} F_{Y_1, Y_2}(y_1, y_2) - F_{Y_1, Y_2}(a_1, y_2) - F_{Y_1, Y_2}(y_1, a_2) + F_{Y_1, Y_2}(a_1, a_2) \\ &\stackrel{(8.23)}{=} F_{Y_1}(y_1)F_{Y_2}(y_2) - F_{Y_1}(a_1)F_{Y_2}(y_2) - F_{Y_1}(y_1)F_{Y_2}(a_2) + F_{Y_1}(a_1)F_{Y_2}(a_2) \\ (A) \quad &= (F_{Y_1}(y_1) - F_{Y_1}(a_1))(F_{Y_2}(y_2) - F_{Y_2}(a_2)) = P\{a_1 < Y_1 \leq y_1\} \cdot P\{a_2 < Y_2 \leq y_2\} \end{aligned}$$

For discrete Y_1 and Y_2 , we obtain with $a_1 = y_1^-$ and $a_2 = y_2^-$,

$$\begin{aligned} p_{Y_1, Y_2}(y_1, y_2) &= P\{y_1^- < Y_1 \leq y_1, y_2^- < Y_2 \leq y_2\} \\ &\stackrel{(A)}{=} P\{y_1^- < Y_1 \leq y_1\} \cdot P\{y_2^- < Y_2 \leq y_2\} = p_{Y_1}(y_1) \cdot p_{Y_2}(y_2). \end{aligned}$$

For continuous Y_1 and Y_2 , we obtain,

$$\begin{aligned} \int_{a_1}^{y_1} \int_{a_2}^{y_2} f_{Y_1, Y_2}(u_1, u_2) du_1 du_2 &= P\{a_1 < Y_1 \leq y_1, a_2 < Y_2 \leq y_2\} \\ &\stackrel{(A)}{=} P\{a_1 < Y_1 \leq y_1\} \cdot P\{a_2 < Y_2 \leq y_2\} = \int_{a_1}^{y_1} f_{Y_1}(u_1) du_1 \cdot \int_{a_2}^{y_2} f_{Y_2}(u_2) du_2 \end{aligned}$$

We differentiate with respect to y_1 and y_2 and obtain $f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) f_{Y_2}(y_2)$. ■

The next theorem will be generalized in Theorem 8.10 on p.157. There Y_1 and Y_2 will be replaced with functions $g(Y_1)$ and (Y_2) .

Theorem 8.5.

If Y_1 and Y_2 are independent random variables, then

$$(8.27) \quad E[Y_1 \cdot Y_2] = E[Y_1] \cdot E[Y_2].$$

PROOF: We show the proof for continuous Y_1 and Y_2 . Since $f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) \cdot f_{Y_2}(y_2)$,

$$\begin{aligned} E[Y_1 Y_2] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y_1 y_2 f_{Y_1, Y_2}(y_1, y_2) dy_1 dy_2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y_1 y_2 f_{Y_1}(y_1) f_{Y_2}(y_2) dy_1 dy_2 \\ &= \int_{-\infty}^{\infty} y_2 \left[\int_{-\infty}^{\infty} y_1 f_{Y_1}(y_1) dy_1 \right] f_{Y_2}(y_2) dy_2 = \int_{-\infty}^{\infty} y_2 E[Y_1] f_{Y_2}(y_2) dy_2 \\ &= E[Y_1] \int_{-\infty}^{\infty} y_2 f_{Y_2}(y_2) dy_2 = E[Y_1] E[Y_2]. \end{aligned}$$

The proof for discrete random variables is obtained by employing $p_{Y_1, Y_2}(y_1, y_2) = p_{Y_1}(y_1) \cdot p_{Y_2}(y_2)$ and replacing integration with summation. ■

Theorem 8.6 (WMS Ch.05.4, Theorem 5.5).

Let the continuous random variables Y_1 and Y_2 have a joint PDF $f_{Y_1, Y_2}(y_1, y_2)$ that is strictly positive if and only if there are constants $a < b$ and $c < d$ such that

$$f_{Y_1, Y_2}(y_1, y_2) > 0 \quad \Leftrightarrow \quad a \leq y_1 \leq b \quad \text{and} \quad c \leq y_2 \leq d.$$

$$(8.28) \quad \text{Then } Y_1, Y_2 \text{ are independent} \quad \Leftrightarrow \quad f_{Y_1, Y_2}(y_1, y_2) = g_1(y_1) \cdot g_2(y_2)$$

for suitable nonnegative functions $g_1, g_2 : \mathbb{R} \rightarrow \mathbb{R}$ such that the only argument of g_1 is y_1 and the only argument of g_2 is y_2 .

PROOF:

The \Rightarrow direction is trivially true: Choose the marginal densities f_{Y_1} and f_{Y_2} for g_1 and g_2 .

PROOF of \Leftarrow : From $f(y_1, y_2) = g_1(y_1)g_2(y_2)$, we obtain for the marginal densities,

$$\begin{aligned} \text{(A)} \quad f_{Y_1}(y_1) &= \int_{-\infty}^{\infty} f(y_1, y_2) dy_2 = \int_{-\infty}^{\infty} g_1(y_1)g_2(y_2) dy_2 = g_1(y_1) \int_{-\infty}^{\infty} g_2(y_2) dy_2 = \alpha g_1(y_1), \\ f_{Y_2}(y_2) &= \int_{-\infty}^{\infty} f(y_1, y_2) dy_1 = \int_{-\infty}^{\infty} g_1(y_1)g_2(y_2) dy_1 = g_2(y_2) \int_{-\infty}^{\infty} g_1(y_1) dy_1 = \beta g_2(y_2), \end{aligned}$$

Here, the constants $\alpha = \int_{-\infty}^{\infty} g_2(y_2) dy_2$ and $\beta = \int_{-\infty}^{\infty} g_1(y_1) dy_1$ satisfy

$$\begin{aligned} \text{(B)} \quad \alpha\beta &= \int_{-\infty}^{\infty} g_2(y_2) dy_2 \cdot \int_{-\infty}^{\infty} g_1(y_1) dy_1 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g_1(y_1)g_2(y_2) dy_1 dy_2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{Y_1, Y_2}(y_1, y_2) dy_1 dy_2 = 1. \end{aligned}$$

We conclude that

$$f_{Y_1, Y_2}(y_1, y_2) \stackrel{\text{(B)}}{=} \alpha\beta f_{Y_1, Y_2}(y_1, y_2) = \alpha\beta g_1(y_1)g_2(y_2) = (\alpha g_1(y_1))(\beta g_2(y_2)) \stackrel{\text{(A)}}{=} f_{Y_1}(y_1)f_{Y_2}(y_2).$$

We have proved independence. ■

Example 8.1 (Buffon's needle). The plane is segmented by parallel lines into strips of width $d > 0$. A needle of length $\lambda < d$ is dropped at random onto the plane. What is the probability that the line will intersect one of those parallel lines?

Solution: A needle that is dropped on the plane uniquely determines a right-angled triangle as follows:

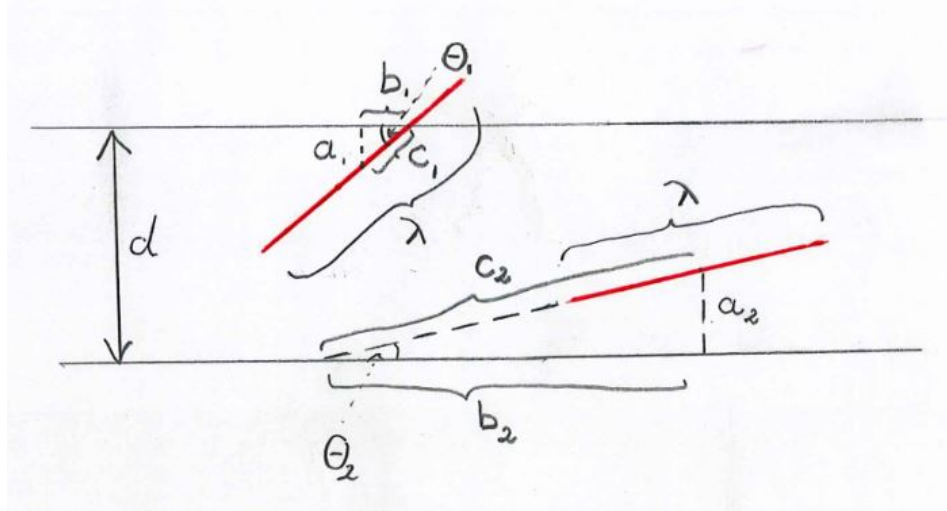
- Leg #1 is perpendicular to the parallels. It extends from the midpoint of the needle to the nearest parallel line. Its length is denoted a .
- Its hypotenuse of length c is on the same line as the needle. Thus, it extends from the midpoint of the needle to the point of intersection with that parallel line.
- Leg #2 is located on that parallel line. Its length is denoted b .

We denote the angle formed by the hypotenuse and leg #2 by θ . Thus,

$$\text{(A)} \quad \sin(\theta) = \frac{a}{c}, \text{ thus, } c = \frac{a}{\sin(\theta)}.$$

$$\text{(B)} \quad \text{The needle intersects the (nearest) parallel} \Leftrightarrow c < \lambda/2 \stackrel{\text{(A)}}{\Leftrightarrow} \frac{a}{\sin(\theta)} < \lambda/2.$$

8.1 (Figure).
Buffon's needle



In Figure 8.1, the triangle on the left satisfies **(B)**:

- $c_1 < \lambda/2$ means that the NE part of the needle extends past the nearest parallel.

On the other hand, the one on the right does not satisfy **(B)**:

- $c_2 > \lambda/2$ means that the SW end of the needle does not reach the nearest parallel.

Note that the triangle created by the random position of the needle is uniquely determined by the two random variables

$$\omega \mapsto A(\omega) := \text{length of leg \#2,}$$

$$\omega \mapsto \Theta(\omega) := \text{angle between leg \#1 and the hypotenuse.}$$

Let $\Gamma \subseteq \Omega$ be the event that the needle intersects with a parallel line. We have seen that

$$\omega \in \Gamma \stackrel{\text{(B)}}{\iff} \frac{A(\omega)}{\sin(\Theta(\omega))} < \frac{\lambda}{2} \iff (A(\omega), \Theta(\omega)) \in B,$$

where

$$B = \left\{ (a, \theta) \in]0, d/2[\times]0, \pi[: \frac{a}{\sin(\theta)} < \frac{\lambda}{2} \right\} = \left\{ (a, \theta) \in]0, d/2[\times]0, \pi[: a < \frac{\lambda}{2} \cdot \sin(\theta) \right\}.$$

Here, the constraint $0 < a < d/2$ results from the fact that the midpoint of the needle has a distance of at most $d/2$ from the nearest parallel. Thus, the length $A(\omega)$ of leg #2 cannot exceed $d/2$.

The randomness of the needle toss ensures that

- $A \sim \text{uniform}(0, \lambda/2)$
- $\Theta \sim \text{uniform}(0, \pi)$
- A and Θ are independent.

It follows that the joint PDF of (A, Θ) is

$$f_{A,\Theta}(a, \theta) = f_A(a) \cdot f_\Theta(\theta) = \begin{cases} \frac{2}{d\pi}, & \text{if } 0 < a < \frac{d}{2}, 0 \leq \theta \leq \pi, \\ 0, & \text{elsewhere.} \end{cases}$$

We obtain the probability that a randomly tossed needle intersects one of the parallel lines as

$$\begin{aligned} P(\Gamma) &= P\{(A, \Theta) \in B\} = \iint_B f_{A,\Theta}(a, \theta) da d\theta \\ &= \int_0^\pi \int_0^{(\lambda/2)\sin(\theta)} \frac{2}{d\pi} da d\theta = \frac{\lambda}{d\pi} \int_0^\pi \sin(\theta) d\theta = \frac{\lambda}{d\pi} (-\cos \theta) \Big|_0^\pi = \frac{2\lambda}{d\pi}. \quad \square \end{aligned}$$

8.4 The Multivariate Uniform Distribution

In this section we extend uniform distribution of Chapter 7.4 (The Uniform Probability Distribution) to regions in two- and three-dimensional space.

Definition 8.11 (Continuous, uniform random variable).

(A) Let $\vec{Y} = (Y_1, Y_2)$ be a twodimensional random vector of continuous random variables with a joint PDF $f_{\vec{Y}}(y_1, y_2)$ that satisfies the following:

- There is a constant $c > 0$ such that either $f_{\vec{Y}}(y_1, y_2) = c$ or $f_{\vec{Y}}(y_1, y_2) = 0$.

Let $C := \{(y_1, y_2) \in \mathbb{R}^2 : f_{\vec{Y}}(y_1, y_2) > 0\}$. Then we say that \vec{Y} has a **continuous uniform probability distribution** on C . \square

(B) Let $\vec{Y} = (Y_1, Y_2, Y_3)$ be a threedimensional random vector of continuous random variables with a joint PDF $f_{\vec{Y}}(y_1, y_2, y_3)$ that satisfies the following:

- There is a constant $d > 0$ such that either $f_{\vec{Y}}(y_1, y_2, y_3) = d$ or $f_{\vec{Y}}(y_1, y_2, y_3) = 0$.

Let $D := \{(y_1, y_2, y_3) \in \mathbb{R}^3 : f_{\vec{Y}}(y_1, y_2, y_3) > 0\}$. Then we say that \vec{Y} has a **continuous uniform probability distribution** on D . \square

Remark 8.3. The constants c and d of the previous definition are uniquely determined as follows:

(A) In the twodimensional case,

$$\iint_{\mathbb{R}^2} f_{\vec{Y}}(y_1, y_2) dy_1 dy_2 = 1 \Rightarrow c = 1 / \iint_C dy_1 dy_2.$$

In other words, c is the reciprocal of the area of C .

(B) In the threedimensional case,

$$\iiint_{\mathbb{R}^3} f_{\vec{Y}}(y_1, y_2, y_3) dy_1 dy_2 dy_3 = 1 \Rightarrow d = 1 / \iiint_D dy_1 dy_2 dy_3.$$

Thus, d is the reciprocal of the volume of D .

(C) It should be obvious how to generalize uniform distribution to n -dimensional random vectors:

Let $\vec{Y} = (Y_1, \dots, Y_n)$ be an n -dimensional random vector of continuous random variables with a joint PDF $f_{\vec{Y}}(\vec{y})$ that satisfies the following:

- There is a constant $e > 0$ such that either $f_{\vec{Y}}(\vec{y}) = e$ or $f_{\vec{Y}}(\vec{y}) = 0$.

Let $E := \{\vec{y} \in \mathbb{R}^n : f_{\vec{Y}}(\vec{y}) > 0\}$. Then we say that \vec{Y} has a **continuous uniform probability distribution** on E .

Similarly to the cases $n = 2$ and $n = 3$, we obtain that e is the reciprocal of the (n -dimensional) volume of E : $e = 1/e'$, where

$$e' := \int \cdots \int_{\vec{y} \in E} d\vec{y} \quad \square$$

Example 8.2. (a) What is the uniform density on $C := C_1 \uplus C_2$, where

$$C_1 := \{\vec{y} \in \mathbb{R}^2 : y_1 < 0, 0 \leq y_2 \leq e^{y_1}\}, \quad C_2 := \{\vec{y} \in \mathbb{R}^2 : 0 \leq y_1 \leq 2, 0 \leq y_2 \leq 1\}?$$

Note that C_1 has area $\int_{-\infty}^0 e^{y_1} dy_1 = 1$ and C_2 , a rectangle of with 2 and height 1, has area 2. Thus, C has area 3 and thus, $c = 1/3$. It follows that

$$f_{\vec{Y}}(\vec{y}) = \begin{cases} \frac{1}{3}, & \text{if } y_1 < 0, 0 \leq y_2 \leq e^{y_1}, \text{ or } 0 \leq y_1 \leq 2, 0 \leq y_2 \leq 1, \\ 0, & \text{else.} \end{cases}$$

(b) Determine the uniform density on

$$D := \{\vec{y} \in \mathbb{R}^3 : y_1 > 0, y_2 > 0, y_3 > 0, y_1^2 + y_2^2 + y_3^2 \leq 1\}.$$

Since $\text{Vol}(D)$, the volume of D , is one eighth of $(4/3)\pi$, the volume of the unit sphere, we obtain

$$d = \frac{1}{\text{Vol}(D)} = \frac{8}{(4/3)\pi} = \frac{6}{\pi}.$$

Thus,

$$f_{\vec{Y}}(\vec{y}) = \begin{cases} \frac{6}{\pi}, & \text{if } y_1 > 0, y_2 > 0, y_3 > 0, y_1^2 + y_2^2 + y_3^2 \leq 1, \\ 0, & \text{else. } \square \end{cases}$$

8.5 The Expected Value of a Function of Several Random Variables

In this section we must work with vectors (x_1, x_2, \dots, x_k) of fixed, but arbitrary dimension k , where each component x_j is a real number and thus, $(x_1, x_2, \dots, x_k) \in \mathbb{R}^k$. Since it is extremely space consuming to repeatedly write such lengthy objects, we remind you of the “arrow notation” that was introduced in Example 2.11 on p.35.

Notation 8.1 (Arrow notation for vectors).

- We write \vec{x} as an abbreviation for a vector (x_1, x_2, \dots, x_n) . The dimension n is either explicitly stated or known from the context.
- If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function of n real numbers and $U = [a_1, b_1] \times \dots \times [a_n, b_n]$ is an n -dimensional rectangle, we write

$$\int_A f(\vec{x}) d\vec{x} = \int_{a_1}^{b_1} \dots \int_{a_2}^{b_2} \int_{a_1}^{b_1} f(x_1, x_2, \dots, x_n) dy_1 dy_2 \dots dy_n$$

Note that all integrands that occur in this course are so well behaved that the order in which those n integrations take place can be switched around, just as you remember it in the cases $n = 2$ and $n = 3$ from multidimensional calculus.

- Let $a_1 < b_1, a_2 < b_2, \dots, a_n < b_n$ for some $n \in \mathbb{N}$. Then $\vec{y} \in [a_1, b_1] \times \dots \times [a_n, b_n]$ denotes the following: $\vec{y} = (y_1, y_2, \dots, y_n)$ and $a_i < y_i \leq b_i$ for $i = 1, \dots, n$.

Here are some examples.

- (a) $\vec{z} \in \mathbb{R}^m$ means: $\vec{z} = (z_1, z_2, \dots, z_m)$ and $z_j \in \mathbb{R}$ for all j .
 (b) If $f : \mathbb{R}^k \rightarrow \mathbb{R}$, then $g(\vec{y})$ means: $f(y_1, \dots, y_k)$.
 (c) If $g : \mathbb{R}^d \rightarrow \mathbb{R}$, then $g(\vec{Y})$ means: $g(Y_1, \dots, Y_d)$; $g(\vec{Y}(\omega))$ means: $g(Y_1(\omega), \dots, Y_d(\omega))$.
 (d) If $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$, then $E[\psi(\vec{Y})]$ means: $E[\psi(Y_1, \dots, Y_n)]$.

Definition 8.12 (Expected value of $g(\vec{Y})$).

(a) Let $k \in \mathbb{N}$ and let $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ be a vector of discrete random variables on a probability space (Ω, P) with PMF $p_{\vec{Y}}(\vec{y})$. Further, let $g : \mathbb{R}^k \rightarrow \mathbb{R}$ be a function of k real numbers y_1, y_2, \dots, y_k . Then

$$(8.29) \quad E[g(\vec{Y})] = E[g(Y_1, Y_2, \dots, Y_k)] := \sum_{y_1, y_2, \dots, y_k} \cdots \sum g(\vec{y}) p_{\vec{Y}}(\vec{y})$$

is called the **expected value** or **mean** of the random variable $g(\vec{Y})$. As usual, the sum on the right is countable summation over those $\vec{y} = (y_1, y_2, \dots, y_k)$ for which $p_{\vec{Y}}(\vec{y}) \neq 0$.

(b) Let $k \in \mathbb{N}$ and let $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ be a vector of continuous random variables on a probability space (Ω, P) with PDF $f_{\vec{Y}}(\vec{y})$. Further, let $h : \mathbb{R}^k \rightarrow \mathbb{R}$ be a function of k real numbers y_1, y_2, \dots, y_k . Then

$$(8.30) \quad E[h(\vec{Y})] = E[h(Y_1, Y_2, \dots, Y_k)] := \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} h(\vec{y}) f_{\vec{Y}}(\vec{y}) d\vec{y}$$

is called the **expected value** or **mean** of the random variable $g(\vec{Y})$.

See Notations 8.1 (Arrow notation for vectors) for an explanation of $\int \cdots d\vec{y}$.

As in the onedimensional case, we only are allowed to say that $E[g(\vec{Y})]$ exists if $\sum \cdots \sum |g(y_1, \dots, y_k)| p(y_1, \dots, y_k)$ is finite and that $E[h(\vec{Y})]$ exists if $\int \cdots \int |g(y_1, \dots, y_k)| f(y_1, \dots, y_k) dy_1 \dots dy_k$ is finite. The functions g and h we deal with in this course will always satisfy that assumption. \square

Example 8.3. As an example of the power of that definition, we give here the proof that

$$E[Y_1 + \cdots + Y_n] = E[Y_1] + \cdots + E[Y_n].$$

Let $h(\vec{y}) := y_1 + \cdots + y_n$. Then, by definition 8.12,

$$E[h(\vec{Y})] = \int_{\mathbb{R}^n} (y_1 + \cdots + y_n) f_{\vec{Y}}(\vec{y}) d\vec{y} = \sum_{j=1}^n \int_{\mathbb{R}^n} y_j f_{\vec{Y}}(\vec{y}) d\vec{y}.$$

Let $\vec{y} := (y_1, \dots, y_{j-1}, y_{j+1}, \dots, y_n)$. Then $\int (\cdots) d\vec{y} = \int (\cdots) d\vec{y} dy_j$ because the order of integration can be switched. Since y_j is constant with respect to \vec{y} ,

$$\int_{\mathbb{R}^n} y_j f_{\vec{Y}}(\vec{y}) d\vec{y} = \int_{\mathbb{R}} \left(\int_{\mathbb{R}^{n-1}} y_j f_{\vec{Y}}(\vec{y}) d\vec{y} \right) dy_j = \int_{-\infty}^{\infty} y_j \left(\int_{\mathbb{R}^{n-1}} f_{\vec{Y}}(\vec{y}) d\vec{y} \right) dy_j.$$

The inner integral “integrates out” all variables except y_j from the PDF of \vec{Y} . Thus, it is the marginal PDF f_{Y_j} of Y_j . It follows from $E[Y_j] = \int_{-\infty}^{\infty} y_j f_{Y_j} dy_j$ that

$$E[h(\vec{Y})] = \sum_{j=1}^n \int_{\mathbb{R}^n} y_j f_{\vec{Y}}(\vec{y}) d\vec{y} = \sum_{j=1}^n \int_{-\infty}^{\infty} y_j f_{Y_j} dy_j = \sum_{j=1}^n E[Y_j]. \quad \square$$

We list here the theorems of WMS Chapter 5.6 (Special Theorems) that detail the rules for evaluating expectations. For the remainder of this section we assume that Y_1, Y_2, \dots are random variables on a common probability space (Ω, P)

Theorem 8.7 (WMS Ch.05.6, Theorem 5.6).

$$(8.31) \quad c \in \mathbb{R} \Rightarrow E[c] = c.$$

PROOF: Trivial. ■

Theorem 8.8 (WMS Ch.05.6, Theorem 5.7).

Let $c \in \mathbb{R}$ and $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. Then the random variable $g(Y_1, Y_2)$ satisfies

$$(8.32) \quad E[cg(Y_1, Y_2)] = cE[g(Y_1, Y_2)].$$

PROOF: Trivial. ■

Theorem 8.9 (WMS Ch.05.6, Theorem 5.8).

Let $g_1, g_2, \dots, g_k : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\vec{Y} := (Y_1, \dots, Y_n)$. Then the random variables $g_j(\vec{Y})$ ($j = 1, \dots, k$) satisfy

$$(8.33) \quad \begin{aligned} E[g_1(\vec{Y}) + g_2(\vec{Y}) + \dots + g_k(\vec{Y})] \\ = E[g_1(\vec{Y})] + E[g_2(\vec{Y})] + \dots + E[g_k(\vec{Y})]. \end{aligned}$$

PROOF: We proved in Example 8.3 on p.156 that $E[\sum_j U_j] = \sum_j E[U_j]$ for discrete or continuous random variables U_1, \dots, U_k . We apply this formula to $U_j := g_j(\vec{Y})$ and the theorem follows. ■

The next theorem generalizes Theorem 8.5 on p.151. That one stated that, for independent random variables, the expectation of the product is the product of the expectations.

Theorem 8.10.

Let $g, h : \mathbb{R} \rightarrow \mathbb{R}$ be functions of a single variable and assume that the random variables Y_1 and Y_2 are independent. Then the random variables $g(Y_1)$ and $h(Y_2)$ also are independent and they satisfy

$$(8.34) \quad E[g(Y_1)h(Y_2)] = E[g(Y_1)]E[h(Y_2)].$$

PROOF: We give the proof for the continuous case only. It is the WMS proof without any alterations. The proof for the discrete case is similar.

Let $f_{Y_1, Y_2}(y_1, y_2)$ denote the joint PDF of Y_1 and Y_2 . Independence of Y_1 and Y_2 yields

$$f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) f_{Y_2}(y_2).$$

The product $g(Y_1)h(Y_2)$ is a function $\varphi(Y_1, Y_2)$ of Y_1 and Y_2 . Hence, by Definition 8.12 (Expected value of $g(\vec{Y})$) on p.156,

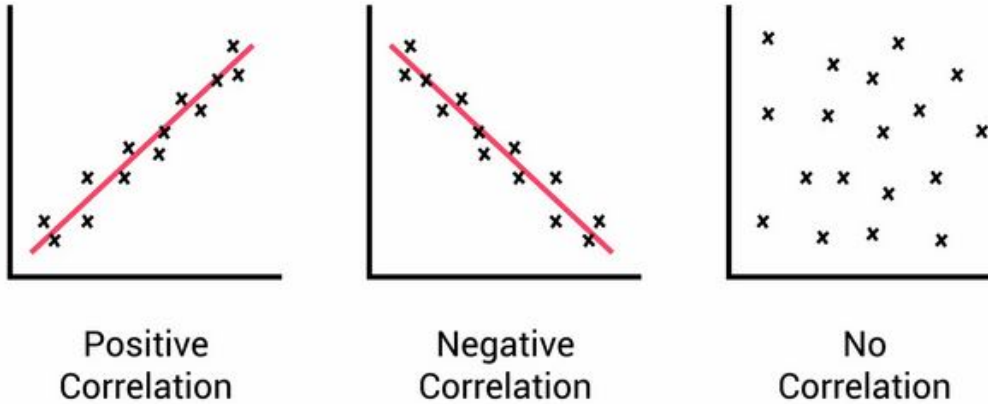
$$\begin{aligned} E[g(Y_1)h(Y_2)] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(y_1)h(y_2)f_{Y_1, Y_2}(y_1, y_2) dy_2 dy_1 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(y_1)h(y_2)f_{Y_1}(y_1) f_{Y_2}(y_2) dy_2 dy_1 \\ &= \int_{-\infty}^{\infty} g(y_1)f_{Y_1}(y_1) \left[\int_{-\infty}^{\infty} h(y_2)f_{Y_2}(y_2) dy_2 \right] dy_1 \\ &= \int_{-\infty}^{\infty} g(y_1)f_{Y_1}(y_1) E[h(Y_2)] dy_1 \\ &= E[h(Y_2)] \int_{-\infty}^{\infty} g(y_1)f_{Y_1}(y_1) dy_1 = E[g(Y_1)] E[h(Y_2)]. \end{aligned}$$

The proof of the independence of $g \circ Y_1$ and $h \circ Y_2$ is based on a characterization of the independence if random elements X_i which involves $\sigma\{X_i\}$, the sigma algebras generated by each X_i . it is omitted here. ■

8.6 The Covariance of Two Random Variables

Introduction 8.2. If we examine how two random variables Y_1 and Y_2 relate to each other, we can consider among other issues the following:

- If the values of Y_1 increase, will the values of Y_2 , on average, also tend to increase? One says in this case that Y_1 and Y_2 have **positive correlation**.
- Or will the values of Y_2 , on average, tend to decrease as the values of Y_1 increase? One says in this case that Y_1 and Y_2 have **negative correlation**.
- Or will the values of Y_2 , on average, have neither increasing nor falling tendency as the values of Y_1 increase? One says in this case that Y_1 and Y_2 have **zero correlation** or that they are **uncorrelated**.
- What if Y_1 and Y_2 are independent? We should expect in that case that Y_1 and Y_2 are uncorrelated.



One can associate with Y_1 and Y_2 a number ρ , their which measures the strength of their correlation. More precisely, it measures the strength of the linear association between Y_1 and Y_2 and whether that association is of an increasing or decreasing nature. ρ is defined in terms of the covariance of Y_1 and Y_2 and this will be the topic of the current section. \square

In this entire section, we consider two random variables Y_1 and Y_2 on a probability space (Ω, P) . As usual, we denote mean and standard deviation

$$\mu_j := E[Y_j], \quad \sigma_j := \sqrt{\text{Var}[Y_j]}, \quad \text{for } j = 1, 2.$$

Definition 8.13 (Covariance).

The **covariance** of Y_1 and Y_2 is

$$(8.35) \quad \text{Cov}[Y_1, Y_2] = E[(Y_1 - E[Y_1])(Y_2 - E[Y_2])] = E[(Y_1 - \mu_1)(Y_2 - \mu_2)]. \quad \square$$

Remark 8.4. $\text{Cov}[Y_1, Y_2]$ has the following properties:

- (a) The larger the absolute value of the covariance of Y_1 and Y_2 , the greater the linear dependence between Y_1 and Y_2 .
- (b) $\text{Cov}[Y_1, Y_2] > 0$ indicates that, on average, Y_1 increases as Y_2 increases.
- (c) $\text{Cov}[Y_1, Y_2] < 0$ indicates that, on average, Y_1 decreases as Y_2 increases.
- (d) $\text{Cov}[Y_1, Y_2] = 0$ indicates that, on average, Y_1 remains constant as Y_2 increases. It is a peculiarity of the statistician's lingo that this kind of linear relationship, even if it is very strong, is defined to be as **NO linear relationship** between Y_1 and Y_2 .
- (e) If we consider $10Y_1$ instead of Y_1 and $10Y_2$ instead of Y_2 the correlation changes by a factor of $10^2 = 100$: $\text{Cov}[10Y_1, 10Y_2] = 100\text{Cov}[Y_1, Y_2]$. This is not convenient in many situations and one defines a standardized correlation by relating Y_1 and Y_2 to their variances. This will be done in the next definition. \square

Definition 8.14 (Correlation coefficient).

The **correlation coefficient**, of Y_1 and Y_2 is

$$(8.36) \quad \rho = \frac{\text{Cov}(Y_1, Y_2)}{\sigma_1 \sigma_2} \quad \square$$

We say that Y_1 and Y_2 have **positive correlation** if $\rho > 0$, (i.e., if $\text{Cov}(Y_1, Y_2) > 0$), they have **negative correlation** if $\rho < 0$, (i.e., if $\text{Cov}(Y_1, Y_2) < 0$), and that they have **zero correlation** or that they are **uncorrelated** if $\rho = 0$, (i.e., if $\text{Cov}(Y_1, Y_2) = 0$).

Proposition 8.3. *The correlation coefficient satisfies the inequality*

$$(8.37) \quad -1 \leq \rho \leq 1 \quad \square$$

PROOF: Omitted ■

The next formula often makes it easier to compute the covariance.

Theorem 8.11.

$$(8.38) \quad \text{Cov}[Y_1, Y_2] = E[(Y_1 - \mu_1)(Y_2 - \mu_2)] = E[Y_1 Y_2] - E[Y_1]E[Y_2].$$

PROOF: Since $E[U + V] = E[U] + E[V]$ and $E[cU] = cE[U]$ and $E[c] = c$ for all random variables U, V and numbers c ,

$$\begin{aligned} \text{Cov}[Y_1, Y_2] &= E[(Y_1 - \mu_1)(Y_2 - \mu_2)] \\ &= E(Y_1 Y_2 - \mu_1 Y_2 - \mu_2 Y_1 + \mu_1 \mu_2) \\ &= E[Y_1 Y_2] - \mu_1 E[Y_2] - \mu_2 E[Y_1] + \mu_1 \mu_2 \\ &= E[Y_1 Y_2] - \mu_1 \mu_2 - \mu_2 \mu_1 + \mu_1 \mu_2 = E[Y_1 Y_2] - \mu_1 \mu_2. \quad \blacksquare \end{aligned}$$

Theorem 8.12.

Independent random variables are uncorrelated.

PROOF: By Theorem 8.5 on p.151, independent random variables Y_1 and Y_2 satisfy $E[Y_1 Y_2] = E[Y_1]E[Y_2]$. Together with (8.38), we obtain

$$\text{Cov}[Y_1, Y_2] = E[Y_1 Y_2] - E[Y_1]E[Y_2] = 0. \quad \blacksquare$$

Example 8.4 (Uncorrelated, but not independent). The following simple example shows two discrete random variables Y_1 and Y_2 which are uncorrelated, but they are not independent.

We obtain from the joint PMF $p(y_1, y_2)$ of Y_1 and Y_2 , shown at the right, that

$$E[Y_1] = (-1)\frac{1}{4} + 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{4} = 0,$$

$$E[Y_2] = (-1)\frac{1}{2} + 1 \cdot \frac{1}{2} = 0,$$

$$\begin{aligned} E[Y_1 Y_2] &= (-1)(-1)0 + 0(-1)\frac{1}{2} + (1)(-1)0 \\ &\quad + (-1)(1)\frac{1}{4} + 0 \cdot 1 \cdot 0 + 1 \cdot 1 \cdot \frac{1}{4} = 0. \end{aligned}$$

| | Y_2 | |
|-------|-------|-----|
| Y_1 | -1 | 1 |
| -1 | 0 | 1/4 |
| 0 | 1/2 | 0 |
| 1 | 0 | 1/4 |

Thus, $E[Y_1Y_2] = E[Y_1]E[Y_2] = 0$ and Y_1 and Y_2 are uncorrelated. On the other hand, $p(-1, -1) = 0$, whereas $p_{Y_1}(-1) \cdot p_{Y_2}(-1) = \frac{1}{4} \cdot \frac{1}{2} \neq 0$. Thus, Y_1 and Y_2 are not independent. \square

Definition 8.15 (Linear function). ★

Let $n \in \mathbb{N}$. We call a function $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$; $\vec{x} = (x_1, \dots, x_n) \mapsto \varphi(\vec{x})$, a **linear function**, of x_1, \dots, x_n , if there are constants $a_1, \dots, a_n \in \mathbb{R}$ such that

$$(8.39) \quad \varphi(\vec{x}) = a_1x_1 + a_2x_2 + \dots + a_nx_n = \sum_{j=1}^n a_jx_j. \quad \square$$

Remark 8.5. Note that if $\vec{Y} = (Y_1, \dots, Y_n)$ is a vector of random variables, then the function φ of (8.39) defines a random variable $V = \varphi(\vec{Y}) = \sum_{j=1}^n a_jY_j$. \square

Theorem 8.13 (WMS Ch.05.8, Theorem 5.12). Let $\vec{X} = X_1, \dots, X_m$ and $\vec{Y} = Y_1, \dots, Y_n$ be random variables on a probability space (Ω, P) . For $i = 1, \dots, m$ and $j = 1, \dots, n$, let $\xi_i := E(X_i)$ and $\eta_j := E(Y_j)$. Further, let

$$U := \sum_{i=1}^m a_iX_i \quad \text{and} \quad V := \sum_{j=1}^n b_jY_j,$$

where $\vec{a} = (a_1, a_2, \dots, a_m)$ and $\vec{b} = (b_1, b_2, \dots, b_n)$ are two constant vectors. Then

$$(8.40) \quad E[U] = \sum_{i=1}^m a_i\xi_i,$$

$$(8.41) \quad \text{Var}[U] = \sum_{i=1}^m a_i^2 \text{Var}[X_i] + 2 \sum_{1 \leq i < j \leq m} a_i a_j \text{Cov}[X_i, X_j].$$

$$(8.42) \quad \text{Cov}[U, V] = \sum_{i=1}^m \sum_{j=1}^n a_i b_j \text{Cov}[X_i, Y_j].$$

In (8.41), $\sum_{1 \leq i < j \leq m} \dots$ refers to summation over all pairs (i, j) satisfying $i < j$.

PROOF: The theorem consists of three parts, of which (8.40) follows directly from Theorems 8.8 and 8.9.

Proof of (8.41): From the definition of variance we obtain

$$\begin{aligned}
 \text{Var}[U] &= E[U - E[U]]^2 = E\left[\sum_{i=1}^n a_i X_i - \sum_{i=1}^n a_i \xi_i\right]^2 = E\left[\sum_{i=1}^n a_i (X_i - \xi_i)\right]^2 \\
 &= E\left[\sum_{i=1}^n a_i^2 (X_i - \xi_i)^2 + \sum_{\substack{i=1 \\ j=1 \\ i \neq j}}^n \sum_{i=1}^n a_i a_j (X_i - \xi_i)(X_j - \xi_j)\right] \\
 &= \sum_{i=1}^n a_i^2 E[(X_i - \xi_i)^2] + \sum_{\substack{i=1 \\ j=1 \\ i \neq j}}^n \sum_{i=1}^n a_i a_j E[(X_i - \xi_i)(X_j - \xi_j)].
 \end{aligned}$$

By the definitions of variance and covariance, we have

$$E[(X_i - \xi_i)^2] = \text{Var}[X_i] \quad \text{and} \quad E[(X_i - \xi_i)(X_j - \xi_j)] = \text{Cov}[X_i, X_j].$$

Thus,

$$\text{Var}[U] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + \sum_{\substack{i=1 \\ j=1 \\ i \neq j}}^n \sum_{i=1}^n a_i a_j \text{Cov}[X_i, X_j].$$

We apply symmetry $\text{Cov}[X_i, X_j] = \text{Cov}[X_j, X_i]$ to the double summation and obtain

$$\text{Var}[U] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + 2 \sum_{1 \leq i < j \leq n} a_i a_j \text{Cov}[X_i, X_j].$$

We have shown (8.41). To prove (8.42), we proceed in a similar fashion: We have

$$\begin{aligned}
 \text{Cov}[U, V] &= E[(U - E[U])(V - E[V])] \\
 &= E\left[\left(\sum_{i=1}^m a_i X_i - \sum_{i=1}^m a_i \xi_i\right) \left(\sum_{j=1}^n b_j X_j - \sum_{j=1}^n b_j \eta_j\right)\right] \\
 &= E\left[\left(\sum_{i=1}^m a_i (X_i - \xi_i)\right) \left(\sum_{j=1}^n b_j (Y_j - \eta_j)\right)\right]
 \end{aligned}$$

$$\begin{aligned}
 \text{Thus, } \text{Cov}[U, V] &= E\left[\sum_{i=1}^m \sum_{j=1}^n a_i b_j (X_i - \xi_i)(Y_j - \eta_j)\right] \\
 &= \sum_{i=1}^m \sum_{j=1}^n a_i b_j E[(X_i - \xi_i)(Y_j - \eta_j)] \\
 &= \sum_{i=1}^m \sum_{j=1}^n a_i b_j \text{Cov}[X_i, Y_j]. \quad \blacksquare
 \end{aligned}$$

Remark 8.6. Note the following about Theorem 8.13:

- (a) Neither CDFs, PMFs or PDFs were needed to prove the theorem. Thus, the proof applies to both discrete and continuous random variables.
- (b) Since $Cov[Y_i, Y_i] = Var[Y_i]$, (8.41) is a particular version of (8.42). \square

We are now in a position to prove (7.28) of Theorem 7.8 on p.127 (and thus, also (6.16) of Theorem 6.4 on p.99) Those formulas state that, for independent random variables, the variance of the sum equals the sum of the variances. Even better, independence can be replaced with the weaker assumption of correlation zero. (See Theorem 8.12.)

Corollary 8.1 (Bienaymé formula for uncorrelated variables). ★

Let $Y_1, Y_2, \dots, Y_n : \Omega \rightarrow \mathbb{R}$ be uncorrelated random variables (which all are defined on the same probability space (Ω, P)) ($n \in \mathbb{N}$). Then

$$(8.43) \quad Var \left[\sum_{j=1}^n Y_j \right] = \sum_{j=1}^n Var[Y_j].$$

PROOF: Since Y_1, \dots, Y_n are uncorrelated, $Cov[Y_i, Y_j] = 0$ for $1 \leq i, j \leq n$ and $i \neq j$. We employ (8.41) on p.161 with $a_1 = a_2 = \dots = a_n = 1$ and obtain

$$Var \left[\sum_{i=1}^n Y_i \right] = \sum_{i=1}^n Var[Y_i] + 2 \sum_{1 \leq i < j \leq n} Cov[Y_i, Y_j] = \sum_{i=1}^n Var[Y_i] + 0. \blacksquare$$

Example 8.5 (Variance of the sample mean ⁴⁴). This example belongs thematically to Section 5.2 (Random Sampling and Urn Models With and Without Replacement). We model SRS sampling from a population to infer statistical knowledge about it as follows.

- The population is represented by a probability space (Ω, P) and the statistical knowledge we are interested in is part of the distribution of a random variable Y on (Ω, P) .
- Picking at random an item from the population is modeled as the outcome $Y(\omega)$ of an invocation of Y .
- Picking an SRS sample of size n from the population is modeled as the n outcomes $\vec{Y}(\omega) = (Y_1(\omega), \dots, Y_n(\omega))$ of n independent random variables Y_1, \dots, Y_n which have the same distribution as Y . In other words, the Y_j are a (finite) iid sequence in the sense of Definition 6.4 on p.100.
- Of course, that last point is an idealization, since independent sample picks correspond to sampling with replacement, whereas SRS models to sampling without replacement. See Definitions 5.3 on p.89 and 5.4 about SRS and urn models. On the other hand, the computational differences between results based on sampling with and without replacement are of practical insignificance if the sample size is small when compared to the population size. ⁴⁵

In this example we specifically consider the mean of the population data.

⁴⁴This is a modified version of WMS, Example 5.27.

⁴⁵See parts (c) and (d) of Remark 5.2 on p.88.

- It seems natural to model this mean by the mean of Y , i.e., the expectation $\mu = E[Y]$ of Y .
- So that's it then. $E[Y]$ is the answer we are looking for. Well, it would be if we only knew the distribution of Y or, at least, $E[Y]$.
- But we don't! We "defined" Y as the action of taking a single random pick from the population, and that is the extent of our knowledge of Y .
- This is why we introduced the vector \vec{Y} of n iid sample picks. The randomness and independence of Y_1, \dots, Y_n should make the specific sample \vec{y} that consists of the outcomes $y_j = Y_j(\omega)$ representative of the population. Thus, its **sample mean** $\bar{y} = \bar{Y}(\omega)$ which is obtained by averaging the sample data, i.e.,

$$\bar{Y}(\omega) = \frac{Y_1(\omega) + Y_2(\omega) + \dots + Y_n(\omega)}{n},$$

should result in a good estimate of the population mean.

All of the above serves as motivation for the following setup. Let Y_1, Y_2, \dots, Y_n be independent random variables with common expectation $E[Y_j] = \mu$ and variance $Var[Y_j] = \sigma^2$ ($j = 1, \dots, n$). Let

$$(8.44) \quad \bar{Y} := \frac{1}{n} \sum_{j=1}^n Y_j.$$

It follows from (8.40) on p.161 and Corollary 8.1 on p.163 that

$$E[\bar{Y}] = E\left[\frac{1}{n} \sum_{j=1}^n Y_j\right] = \frac{1}{n} E\left[\sum_{j=1}^n Y_j\right] = \frac{1}{n} \sum_{j=1}^n E[Y_j] = \frac{1}{n} (n\mu) = \mu,$$

$$Var[\bar{Y}] = Var\left[\frac{1}{n} \sum_{j=1}^n Y_j\right] = \frac{1}{n^2} Var\left[\sum_{j=1}^n Y_j\right] = \frac{1}{n^2} \sum_{j=1}^n Var[Y_j] = \frac{1}{n^2} (n\sigma^2) = \frac{\sigma^2}{n}.$$

We infer from those two formulas the following.

Recall that the purpose of \bar{Y} is to serve as an **estimator** for the following population parameter: The population mean, which is the mean of anyone of the sample picks $\mu = E[Y_j]$.

The significance of the formula $E[\bar{Y}] = \mu$ is as follows

- The expected value of this estimator equals the parameter it is meant to estimate.

An estimator with that property is referred to as an **unbiased estimator**.

Now to the formula $Var[\bar{Y}] = \sigma^2/n$. We use it to compare the standard deviations

$$\sigma_{Y_j} = \sqrt{Var[Y_j]} \quad \text{and} \quad \sigma_{\bar{Y}} = \sqrt{Var[\bar{Y}]}$$

of a single pick Y_j and the average \bar{Y} of n such independent picks. Note that the standard deviation of a random variable U is a measure for its concentration about its expected value. (And the same is true for its variance.) A small σ_U signifies that most outcomes $U(\omega)$ are in close vicinity of $E[U]$. Thus, $\sigma_{\bar{Y}}$ is a measure for the lack of precision with which \bar{Y} estimates $E[\bar{Y}] = \mu$.

- In the extreme case of a sample of size 1, i.e., $n = 1$, that lack of precision is σ .
- For $n = 100$, that lack of precision goes down to $\frac{\sigma}{10}$. Thus, precision has improved by a factor of 10.
- Generally speaking, increasing the sample size by the factor K (and spending all that time and money doing so) does not reward us with a proportionate improvement of the precision of the estimate \bar{Y} . It only increases by the factor \sqrt{K} . \square

8.7 Conditional Expectations and Conditional Variance

8.7.1 The Conditional Expectation With Respect to an Event

We will start with a definition of the conditional expectation $E[Y | B]$ of a random variable Y where conditioning happens with respect to an event $B \subseteq \Omega$. This definition is usually not taught in an undergraduate level course on probability theory for the following reason: It cannot be extended, in the case of continuous random variables Y and \tilde{Y} , to $E[Y | \tilde{Y} = \tilde{y}]$, i.e., conditioning on \tilde{Y} having a fixed outcome \tilde{y} .

All that follows in this subsection is based on Theorem 3.8 on p.55 which states the following: If (Ω, P) is a probability space and $B \subseteq \Omega$ is an event that satisfies $P(B) > 0$, then the function $Q(\cdot)$, defined as $Q(A) := P(A | B)$ for $A \subseteq \Omega$, is a probability measure on Ω .⁴⁶

Assumption 8.3.

In all of this subsection we deal with a fixed probability space (Ω, P) and a fixed event $B \subseteq \Omega$ that satisfies $P(B) > 0$. We further assume that $Q(\cdot)$ is the probability measure

$$(8.45) \quad A \mapsto Q(A) := P(A | B), \quad \text{where } A \subseteq \Omega.$$

The symbols X, X_1, X_2, \dots denote random elements and X, X_1, X_2, \dots denote random variables on Ω . We need not be specific about whether we mean (Ω, P) or (Ω, Q) , because the definition of random element and random variable does not involve the probability measure, only the carrier space Ω . \square

Remark 8.7. The following mathematical triviality allows us to translate much that we have done with random variables in connection with P to their analogues with respect to $Q = P(\cdot | B)$.

- All definitions, propositions and theorems in which an unspecified probability measure P is involved can be reformulated by replacing P with Q .

Here is a list (certainly not complete) of many such concepts.

- cumulative distribution function, • probability mass function
- probability density function • joint CDF • joint PMF • joint PDF
- expectation • variance • moments • moment generating function

⁴⁶To be exact, there also was a σ -algebra \mathcal{F} and we had to assume that $B \in \mathcal{F}$ and that $P(A)$ is defined only for $A \in \mathcal{F}$. This in turn implies that $Q(A) = P(A | B)$ only is defined for arguments $A \in \mathcal{F}$. We do not mention \mathcal{F} since we decided to avoid dealing with σ -algebras whenever possible.

BEWARE: The above does not apply to cases where a specific probability measure is considered. An example for this would be, e.g., Proposition 7.9 on p.138 (memorylessness of the exponential distribution). Here the probability measure is an exponential distribution P_Y .

We will elaborate on some of the items in that bulleted list in the next remark. \square

Remark 8.8. In the following, the phrase “ Q -.....” serves as an abbreviation for the lengthier “..... with respect to Q ”.

- (a) The Q -CDF of a random variable Y is $F_Y^Q(y) = Q\{Y \leq y\} = P\{Y \leq y \mid B\}$.
- (b) The Q -PMF of a discrete random element⁴⁷ X is $p_X^Q(x) = Q\{X = x\} = P\{X = x \mid B\}$.
- (c) Assume that the derivative $f_Y^Q(y) = \frac{dF_Y^Q(y)}{dy}$ of the Q -CDF of a random variable Y exists and is continuous except for at most finitely many y in any finite interval. Then Y is a Q -continuous random variable with Q -PDF $f_Y^Q(y)$.⁴⁸
- (d) We skip joint Q -CDFs and joint Q -PDFs and only elaborate on the joint Q -PMF. of two random elements X_1 and X_2 . It is, as one should expect, defined as $p_{X_1, X_2}^Q(x_1, x_2) = Q\{X_1 = x_1, X_2 = x_2\} = P\{X_1 = x_1, X_2 = x_2 \mid B\}$.
- (e) The Q -expected value of a discrete random variable Y is $E^Q[Y] = \sum_y y \cdot p_Y^Q(y) = \sum_y y \cdot P\{Y = y \mid B\}$. (\sum_y is over all y where $p_Y^Q(y) > 0$.)
- (f) The Q -expectation of a continuous random variable Y is $E^Q[Y] = \int_{-\infty}^{\infty} y \cdot f_Y^Q(y) dy$.
- (g) The Q -variance of a random variable Y is $Var^Q[Y] = E^Q[(Y - E^Q[Y])^2]$.
- (h) The Q -MGF of a random variable Y is $m_Y^Q(t) = E^Q[e^{tY}]$.

For expectations of functions of random variables we skip the case of one or two random variables and proceed directly to the case of a vector $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ of random variables. (See Definition 8.12 on p.156.)

- (i) If the Y_j are discrete and $g : \mathbb{R}^k \rightarrow \mathbb{R}$, then $E^Q[g(\vec{Y})] = \sum_{y_1, y_2, \dots, y_k} g(\vec{y}) p_{\vec{Y}}^Q(\vec{y})$.
- (j) If the Y_j are continuous and $h : \mathbb{R}^k \rightarrow \mathbb{R}$, then $E^Q[h(\vec{Y})] = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} h(\vec{y}) f_{\vec{Y}}^Q(\vec{y}) d\vec{y}$. \square

Here are some of the theorems we get for free because we have shown them for any probability measure. Again, BEWARE: We made the assumption $P(B) > 0$!

Theorem 8.14.

If $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ is a vector of k discrete or Q -continuous random variables, then

$$(8.46) \quad E^Q \left[\sum_{j=1}^n Y_j \right] = \sum_{j=1}^n E^Q[Y_j].$$

⁴⁷Since $P\{X = x\} \cap B \leq P\{X = x\}$, $P\{X = x\} = 0$ implies $Q\{X = x\} = 0$. Thus, any P -discrete random element also is Q -discrete.

⁴⁸There may be some reasonably general and simple conditions that guarantee Y being Q -continuous from being P -continuous, but this author is not aware of them.

PROOF: This follows from Theorem 8.14 on p.166. ■

Theorem 8.15. *If Y is a discrete or Q -continuous random variable and $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ is a vector of k Q -independent discrete or Q -continuous random variables, then*

$$(8.47) \quad \text{Var}^Q[Y] = E^Q[Y^2] - (E^Q[Y])^2,$$

$$(8.48) \quad \text{Var}^Q[aY + b] = a^2 \text{Var}^Q[Y],$$

$$(8.49) \quad \text{Var}^Q \left[\sum_{j=1}^n Y_j \right] = \sum_{j=1}^n \text{Var}^Q[Y_j].$$

PROOF: This follows from Theorem 7.8 on p.127. ■

There is an issue with that last theorem. Not just with the proof, but with the assumptions that were made. How is Q -independence defined for random variables, or even for events A_1, A_2, A_k ? The answer is, of course, that we apply all previously made definitions of independence of two or more events or random variables, replacing the original probability measure P with Q .

The following theorem about the Q -independence of two events is worthwhile mentioning.

Theorem 8.16.

Let $A_1, A_2, B \subseteq \Omega$ be three events such that $P(A_1) > 0, P(A_2) > 0$ and $P(B) > 0$. Then

$$(8.50) \quad \begin{aligned} & \text{(a)} \quad P(A_1 \cap A_2 | B) = P(A_1 | B) \cdot P(A_2 | B) \\ \Leftrightarrow & \text{(b)} \quad P(A_1 | A_2 \cap B) = P(A_1 | B) \\ \Leftrightarrow & \text{(c)} \quad P(A_2 | A_1 \cap B) = P(A_2 | B). \end{aligned}$$

In other words, if A_i and A_j are independent with respect to “just” conditioning on B , then “further” conditioning of A_i on both A_j and B has no effect. Here, either $i = 1, j = 2$ or $i = 2, j = 1$.

PROOF: Since (a) is asymmetrical in A_1 and A_2 and (c) is obtained from (b) by switching the roles of A_1 and A_2 , it suffices to prove (a) \Leftrightarrow (b).

PROOF that (a) \Rightarrow (b):

$$\begin{aligned} P(A_1 | A_2 \cap B) &= \frac{P(A_1 \cap A_2 \cap B)}{P(A_2 \cap B)} = \frac{P(A_1 \cap A_2 \cap B)}{P(B)} \cdot \frac{P(B)}{P(A_2 \cap B)} \\ &= P(A_1 \cap A_2 | B) \cdot \frac{1}{P(A_2 | B)} \stackrel{\text{(a)}}{=} P(A_1 | B) \cdot P(A_2 | B) \cdot \frac{1}{P(A_2 | B)} \\ &= P(A_1 | B). \end{aligned}$$

PROOF that (b) \Rightarrow (a):

$$\begin{aligned} P(A_1 \cap A_2 | B) &= \frac{P(A_1 \cap A_2 \cap B)}{P(B)} = \frac{P(A_1 \cap A_2 \cap B)}{P(A_2 \cap B)} \cdot \frac{P(A_2 \cap B)}{P(B)} \\ &= P(A_1 | A_2 \cap B) \cdot P(A_2 | B) \stackrel{\text{(b)}}{=} P(A_1 | B) \cdot P(A_2 | B). \quad \blacksquare \end{aligned}$$

8.7.2 The Conditional Expectation w.r.t a Random Variable or Random Element

Remark 8.9. ★ We mentioned at the beginning of the previous subsection 8.7.1 (The Conditional Expectation With Respect to an Event), that conditioning with respect to an event B constitutes a dead end street. This is the reason why the material has been marked as ★ (optional). Now let us give the reason.

As far as modeling reality by means of probability theoretical concepts is concerned, the primary interest of conditioning is being able to assume during certain calculations of the probability involving a random element X_1 , that another random element X_2 has as its outcome a fixed value x_2 . Thus, we typically are interested in

- $P\{X_1 \in B_1 \mid X_2 = x_2\}$, where x_2 is some fixed outcome that can be attained by X_2 .

Having stated the issue in the most general terms, we will restrict ourselves for the remainder of this remark to random variables Y_1 and Y_2 rather than working with random elements. This will allow us to contrast discrete and continuous random variables.

The method of subsection 8.7.1 (using the probability measure $Q(A) = P\{A \mid \tilde{Y} = \tilde{y}\}$) will actually work if we condition on specific values of a discrete random variable Y_2 . This is so because we only are interested in those outcomes y_2 for which

$$p_{Y_2}(y_2) = P\{Y_2 = y_2\} > 0$$

and the conditional probability $P\{A \mid Y_2 = y_2\}$ exist for such outcomes y_2 .

On the other hand, we have nothing at all to work with if Y_2 is continuous, since $P\{Y_2 = y_2\} = 0$ for all numbers y_2 (see Proposition 7.1 on p.116), since this results in $P\{Y_1 \in B_1 \mid Y_2 = y_2\}$ being **UNDEFINED** for all numbers y_2 !

To overcome this hurdle we will work with the conditional PMFs and PDFs

- $p_{Y_1|Y_2}(y_1 \mid y_2) = \frac{p_{Y_1, Y_2}(y_1, y_2)}{p_{Y_2}(y_2)}$, if Y_1 and Y_2 are discrete random variables,
- $f_{Y_1|Y_2}(y_1 \mid y_2) = \frac{f_{Y_1, Y_2}(y_1, y_2)}{f_{Y_2}(y_2)}$, if Y_1 and Y_2 are continuous random variables.

We close this remark by noticing that, in the case of discrete random variables, working with $Q\{Y_1 \in B_1\} = P\{Y_1 \in B_1 \mid Y_2 = y_2\}$ or with $p_{Y_1|Y_2}(y_1 \mid y_2)$ amounts to the same, because Q and $p_{Y_1|Y_2}$ satisfy

$$Q\{Y_1 \in B_1\} = \sum_{y_1 \in B_1} P\{Y_1 = y_1 \mid Y_2 = y_2\} = \sum_{y_1 \in B_1} p_{Y_1|Y_2}(y_1 \mid y_2). \quad \square$$

Compare the following remark to Remark 8.7 on p.165 for discrete random variables.

Remark 8.10. The following allows us to translate much that we have done with continuous random variables in connection with P to their analogues where we replace the (marginal) PDF $f_{Y_1}(y_1)$ with the conditional PDF $f_{Y_1|Y_2}(y_1 \mid y_2)$:

- Assume that $y_2 \in \mathbb{R}$ satisfies $f_{Y_2}(y_2) > 0$. Then the integrable function

$$f_{Y_1|Y_2}(\cdot \mid y_2) : y_1 \mapsto f_{Y_1|Y_2}(y_1 \mid y_2)$$

satisfies $\square f_{Y_1|Y_2}(y_1 \mid y_2) \geq 0$ for $-\infty < y_1 < \infty$ $\square \int_{-\infty}^{\infty} f_{Y_1|Y_2}(y_1 \mid y_2) dy_1 = 1$.

- According to Theorem 7.3 on p.118, $f_{Y_1|Y_2}(\cdot | y_2)$ is the PDF of the probability measure P_{y_2} on Ω , defined by

$$P_{y_2}\{a < Y_2 \leq b\} = \int_a^b f_{Y_1|Y_2}(y_1 | y_2) dy_1.$$

- Thus, all definitions, propositions and theorems in which an unspecified probability measure P is involved can be reformulated by replacing P with P_{y_2} .

This applies, among others, to the following concepts which were listed in Remark 8.7 on p.165 for discrete random variables:

- cumulative distribution function, • probability mass function
- probability density function • joint CDF • joint PMF • joint PDF
- expectation • variance • moments • moment generating function
- All that was said above extends to a random vector $\vec{U} = (U_1, \dots, U_k)$ in place of Y_1 . We only must replace $f_{Y_1, Y_2}(y_1, y_2)$ with $f_{\vec{U}, Y_2}(u_1, \dots, u_k, y_2)$, etc. \square

Definition 8.16 (Conditional expectation).

Let Y_1 and Y_2 be two random variables which are either jointly discrete or jointly continuous and $g : \mathbb{R} \rightarrow \mathbb{R}$. Let

$$(8.51) \quad E[g(Y_1) | Y_2 = y_2] := \sum_{y_1} g(y_1) p(y_1 | y_2) \quad (\text{discrete case}),$$

$$(8.52) \quad E[g(Y_1) | Y_2 = y_2] := \int_{-\infty}^{\infty} g(y_1) f(y_1 | y_2) dy_1 \quad (\text{continuous case}).$$

We call $E[g(Y_1) | Y_2 = y_2]$ the **conditional expectation** of $g(Y_1)$, given that $Y_2 = y_2$. \square

Remark 8.11. Note for the following that the function

$$\omega \mapsto E[g(Y_1) | Y_2 = Y_2(\omega)] = E[g(Y_1) | Y_2 = y_2] \Big|_{y_2=Y_2(\omega)}$$

defines a random variable on (Ω, P) . It is customary in many situations to suppress the argument ω and write

$$(8.53) \quad E[g(Y_1) | Y_2]$$

for this random variable. Clearly, if we write $Z(\omega)$ for $E[g(Y_1) | Y_2 = Y_2(\omega)]$, we can take its (unconditional) expectation

$$(8.54) \quad E[Z] = E[E[g(Y_1) | Y_2]].$$

In particular, if $g(y) = y$, we can take the expectation $E[E[Y_1 | Y_2]]$ of $E[Y_1 | Y_2]$. We will do so in the next theorem. \square

Theorem 8.17 (WMS Ch.05.11, Theorem 5.14).

Let Y_1 and Y_2 be either jointly continuous or jointly discrete random variables. Then

$$(8.55) \quad E[Y_1] = E[E[Y_1 | Y_2]].$$

See Remark 8.11 concerning the interpretation of the right-hand side.

PROOF: We give the proof for jointly continuous Y_1 and Y_2 . With the usual notation for joint PDF, marginal densities and conditional PDF we obtain

$$\begin{aligned} E[Y_1] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y_1 f_{Y_1, Y_2}(y_1, y_2) dy_1 dy_2 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y_1 f_{Y_1|Y_2}(y_1 | y_2) f_{Y_2}(y_2) dy_1 dy_2 \\ &= \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} y_1 f_{Y_1|Y_2}(y_1 | y_2) dy_1 \right) f_2(y_2) dy_2 \\ &= \int_{-\infty}^{\infty} E[Y_1 | Y_2 = y_2] f_{Y_2}(y_2) dy_2 = E[E[Y_1 | Y_2]]. \end{aligned}$$

The proof for the discrete case is done by doing summation instead of integration and replacing joint, marginal and conditional PDFs with the corresponding PMFs. ■

We define the conditional variance of Y_1 given $Y_2 = y_2$ by applying Definition 8.16 to the functions $g(y_1) = y_1$ and $g(y_1) = y_1^2$.

Definition 8.17 (Conditional variance).

Let Y_1 and Y_2 be two random variables which are either jointly discrete or jointly continuous. Let

$$(8.56) \quad \text{Var}[Y_1 | Y_2 = y_2] := E[Y_1^2 | Y_2 = y_2] - (E[Y_1 | Y_2 = y_2])^2.$$

We call $\text{Var}[Y_1 | Y_2 = y_2]$ the **conditional variance** of $g(Y_1)$, given that $Y_2 = y_2$. □

Theorem 8.18.

Let Y_1 and Y_2 be jointly discrete or jointly continuous random variables. Then

$$(8.57) \quad \text{Var}[Y_1 | Y_2] = E[(Y_1 - E[Y_1 | Y_2])^2 | Y_2],$$

$$(8.58) \quad \text{Var}[Y_1] = E[\text{Var}[Y_1 | Y_2]] + \text{Var}[E[Y_1 | Y_2]].$$

PROOF: We only give the proof of (8.58). Note that

$$(A) \quad \text{Var}[Y_1 | Y_2] = E[Y_1^2 | Y_2] - (E[Y_1 | Y_2])^2,$$

$$(B) \quad E[\text{Var}[Y_1 | Y_2]] = E[E[Y_1^2 | Y_2]] - E[(E[Y_1 | Y_2])^2].$$

By the definition of (unconditional) variance,

$$(C) \quad \text{Var}[E[Y_1 | Y_2]] = E[(E[Y_1 | Y_2])^2] - (E[E[Y_1 | Y_2]])^2.$$

Further,

$$\begin{aligned} \text{Var}[Y_1] &= E[Y_1^2] - (E[Y_1])^2 \\ &= E[E[Y_1^2 | Y_2]] - (E[E[Y_1 | Y_2]])^2 \\ &= E[E[Y_1^2 | Y_2]] - E[(E[Y_1 | Y_2])^2] + E[(E[Y_1 | Y_2])^2] - (E[E[Y_1 | Y_2]])^2 \\ &= E[E[Y_1^2 | Y_2] - (E[Y_1 | Y_2])^2] + \{E[(E[Y_1 | Y_2])^2] - (E[E[Y_1 | Y_2]])^2\} \\ &= E[\text{Var}[Y_1 | Y_2]] + \text{Var}[E[Y_1 | Y_2]]. \blacksquare \end{aligned}$$

8.7.3 Conditional Expectations as Optimal Mean Squared Distance Approximations

The presentation of the material presented here follows [1] Bickel and Doksum: Mathematical Statistics.

Introduction 8.3. One can measure the distance between two real-valued functions in several ways.

For example, one can define for $\varphi, \psi : A \rightarrow \mathbb{R}$,

$$\text{dist}_1(\varphi, \psi) := \max\{|\varphi(a) - \psi(a)| : a \in A\}.$$

In other words, one takes the maximum displacement over all arguments of φ and ψ . This “worst case scenario” as the advantage that it works for any kind of domain A , since all that is needed is that the function values are numeric.

However, it often makes more sense to consider the area between the curves defined by φ and ψ .

$$\text{dist}_2(\varphi, \psi) := \int_a^b |\varphi(x) - \psi(x)| dx.$$

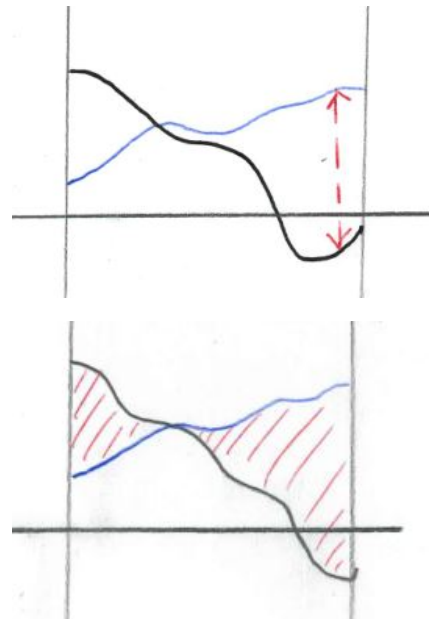
Doing so averages out all individual displacements $|\varphi(x) - \psi(x)|$ over all arguments and one obtains a measure of distance which is not distorted just one potential outlier.

There are mathematical reasons why one would rather work with the squared difference and consider

$$\text{dist}_3(\varphi, \psi) := \int_a^b |\varphi(x) - \psi(x)|^2 dx = \int_a^b (\varphi(x) - \psi(x))^2 dx.$$

Moreover, one can replace the ordinary integral $\int \dots dx$ with a weighted integral $\int \dots w(x) dx$ where $w(x) \geq 0$ for all x and define

$$\text{dist}_4(\varphi, \psi) := \int_a^b (\varphi(x) - \psi(x))^2 w(x) dx.$$



Here, bigger values $w(x)$ of the weight function w lead to a stronger contribution of $\varphi(x) - \psi(x)$ to the distance between φ and ψ

That last example shows us how the expectation of the difference of two functions of two continuous random variables can be viewed as a distance

$$\text{dist}(\varphi(Y_1), \psi(Y_2)) = E[(\varphi(Y_1) - \psi(Y_2))^2] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\varphi(y_1) - \psi(y_2))^2 f_{Y_1, Y_2}(y_1, y_2) dy_1 dy_2.$$

Since $E[(\varphi(Y_1) - \psi(Y_2))^2]$ also is defined for discrete random variables, we obtain for those a corresponding definition by replacing the joint PDF with the joint PMF and integration with summation:

$$\text{dist}(\varphi(Y_1), \psi(Y_2)) = E[(\varphi(Y_1) - \psi(Y_2))^2] = \sum_{y_1, y_2} (\varphi(y_1) - \psi(y_2))^2 p_{Y_1, Y_2}(y_1, y_2).$$

In either discrete or continuous case, we are particularly interested in the case $\varphi(y_1) = y_1$ and examine the distance

$$\text{dist}(Y_1, \psi(Y_2)) = E[(Y_1 - \psi(Y_2))^2]$$

for all possible functions $y_2 \mapsto \psi(y_2)$. It turns out that the minimum

$$\min\{\text{dist}(Y_1, \psi(Y_2)) : \text{all suitable functions } \psi\}$$

is attained by selecting $\psi : y_2 \mapsto E[Y_1 | Y_2 = y_2]$. \square

Lemma 8.1. *Let Y be a random variable on (Ω, P) that satisfies $E[Y^2] < \infty$. Then, $E[|Y|] < \infty$.*

PROOF:

$$\text{Let } A := |Y| < 1 \text{ and } Z := 1_A + |Y^2|, \text{ i.e., } Z(\omega) = \begin{cases} 1 + |Y^2|, & \text{if } |Y\omega| < 1, \\ |Y^2|, & \text{if } |Y\omega| \geq 1. \end{cases}$$

Since $|Y(\omega)| < 1$ for $\omega \in A$ and $Y(\omega)^2 \geq 1$ for $\omega \in A^c$, we obtain $|Y(\omega)| \leq Z(\omega)$ for all ω . Thus,

$$E[|Y|] \leq E[Z] \leq E[1] + E[Y^2].$$

The assertion follows. \blacksquare

Lemma 8.2. *Let Y be a random variable on (Ω, P) and $h : \mathbb{R} \rightarrow [0, \infty]$ defined by $a \mapsto E[(Y - a)^2]$. Then,*

- (a) *either $h(a) = \infty$ for all $a \in \mathbb{R}$,*
- (b) *or h attains a unique minimum at $a = E[Y]$.*

PROOF of (a): For fixed $a \in \mathbb{R}$, we define $F : \mathbb{R} \rightarrow \mathbb{R}$ by $F(y) := (y - a)^2 - ((1/2)y^2 - a^2)$. Then,

$$F'(y) = 2(y - a) - y = y - 2a \quad \text{and} \quad F''(y) = 1.$$

It follows that F attains a (unique) minimum at $y = 2a$. From $F(2a) = a^2 - (2a^2 - a^2) = 0$, we obtain that $F(y) \geq 0$ for all y . Thus, $(y - a)^2 \geq (1/2)y^2 - a^2$. This yields

$$(A) \quad \frac{1}{2}y^2 - a^2 \leq (y - a)^2 = y^2 - 2ay + a^2.$$

Next, we obtain from $(y - a)^2 \leq (y - a)^2 + (y + a)^2$ that

$$(B) \quad y^2 - 2ay + a^2 \leq (y^2 - 2ay + 2a^2) + (y^2 = 2ay + 2a^2) = 2y^2 + 2a^2.$$

Let $\omega \in \Omega$ and $y := Y(\omega)$. We combine **(A)** and **(B)** and obtain $\frac{1}{2}y^2 - a^2 \leq (y - a)^2 \leq 2y^2 + 2a^2$. Since ω is an arbitrary element of Ω , we have the following inequality of random variables:

$$(C) \quad \frac{1}{2}Y^2 - a^2 \leq (Y - a)^2 \leq 2Y^2 + 2a^2.$$

Taking expectations maintains inequalities. Since $E[(Y - a)^2] = h(a)$ and $E[Y^2] = h(0)$,

$$(D) \quad \frac{1}{2}h(0) - a^2 \leq h(a) \leq 2h(0) + 2a^2.$$

From this we see that either $[h(0) = \infty \text{ and in this case, } h(a) = \infty \text{ for all } a]$,
or $[h(0) < \infty \text{ and in this case, } h(a) < \infty \text{ for all } a]$.

PROOF of **(b)**: We assume for this part of the proof that $h(0) < \infty$, i.e., $E[(Y^2)] < \infty$. According to Lemma 8.1 we then also have $|E[Y]| < \infty$. Thus, we can write

$$(E) \quad \begin{aligned} h(a) &= E[(Y - a)^2] = E[(Y^2) - 2aE[Y] + a^2] \\ &= E[(Y^2) - (E[Y])^2 + (a^2 - 2aE[Y] + (E[Y])^2)] \\ &= Var[Y] + (a - E[Y])^2 \end{aligned}$$

It follows that h attains a unique minimum in height of $Var[Y]$ at $a = E[Y]$ and this concludes the proof of the lemma. ■

Theorem 8.19.

Assume that Y is a random variable and $\vec{X} = (X_1, \dots, X_k)$ is a random vector on (Ω, P) . Then, either $E[(Y - g \circ \vec{X})] = \infty$ for all real-valued functions $g : \mathbb{R}^k \rightarrow \mathbb{R}$ of k real arguments, or

$$E \left[(Y - E[Y | \vec{X}])^2 \right] \leq E \left[(Y - g \circ \vec{X})^2 \right],$$

for all such functions g . Further, this is a strict inequality if $E[Y | \vec{X}] \neq g \circ \vec{X}$.

Note that, as always, we consider equations and inequalities involving random variables to be true as long as they are satisfied on a set of probability 1.

PROOF: Let us fix $\vec{x} \in \mathbb{R}^k$ for which $E[Y | \vec{X} = \vec{x}]$ is defined.

- (a) In the case of discrete Y and \vec{X} this means that $p_{\vec{X}}(\vec{x}) > 0$ and then $B \mapsto \sum_{y \in B} p_{y|\vec{X}}(y | \vec{x})$ is a probability measure $P_{\vec{x}}$ on Ω for which we denote expectations by $E_{\vec{x}}[\dots]$. Further, for $\psi : \mathbb{R} \rightarrow \mathbb{R}$, $E[\psi(Y) | \vec{X} = \vec{x}] = E_{\vec{x}}[\psi(Y)]$
- (b) For continuous Y and \vec{X} this means that $f_{\vec{X}}(\vec{x}) > 0$. We have seen in Remark 8.10 on p.168 that $B \mapsto \int_B f_{y|\vec{X}}(y | \vec{x}) dy$ is a probability measure $P_{\vec{x}}$ on Ω for which we denote expectations by $E_{\vec{x}}[\dots]$. Further, for $\psi : \mathbb{R} \rightarrow \mathbb{R}$, $E[\psi(Y) | \vec{X} = \vec{x}] = E_{\vec{x}}[\psi(Y)]$
- (c) Thus, in both cases, all we have learned about ordinary expectations can be applied, for fixed \vec{x} , to the conditional expectations $E[\dots | \vec{X} = \vec{x}]$.

- (d) When we condition an expression on $\vec{X} = \vec{x}$, we can replace in that expression all occurrences of \vec{X} with \vec{x} .

It follows from (d) that

$$(A) \quad E \left[(Y - g(\vec{X}))^2 \mid \vec{X} = \vec{x} \right] = E \left[(Y - g(\vec{x}))^2 \mid \vec{X} = \vec{x} \right].$$

We can apply Lemma 8.2 with $E_{\vec{x}}(\dots)$ instead of $E(\dots)$ and the constant $g(\vec{x})$ instead of a and conclude that

$$(B) \quad E \left[(Y - g(\vec{x}))^2 \mid \vec{X} = \vec{x} \right] \geq E \left[(Y - E[Y \mid \vec{X}])^2 \mid \vec{X} = \vec{x} \right].$$

We apply both (A) and (B) and evaluate both sides of the resulting inequality for $\vec{x} = \vec{X}(\omega)$:

$$E \left[(Y - g(\vec{X}))^2 \mid \vec{X} = \vec{X}(\omega) \right] \geq E \left[(Y - E[Y \mid \vec{X}])^2 \mid \vec{X} = \vec{X}(\omega) \right].$$

As we have done before, we streamline this expression by replacing $\vec{X} = \vec{X}(\omega)$ with X :

$$E \left[(Y - g(\vec{X}))^2 \mid \vec{X} \right] \geq E \left[(Y - E[Y \mid \vec{X}])^2 \mid \vec{X} \right].$$

Taking expectations on both sides, we obtain

$$E \left[(Y - g(\vec{X}))^2 \right] \geq E \left[(Y - E[Y \mid \vec{X}])^2 \right].$$

We have shown the inequality that was asserted in the theorem.

We still must prove that this inequality is strict if $E[Y \mid \vec{X}] \neq g \circ \vec{X}$. To do so we apply the reasoning above to formula (E) of Lemma 8.2 and obtain

$$E \left[(Y - g(\vec{X}))^2 \right] = \text{Var}[Y] + E \left[(g(\vec{X}) - E[Y \mid \vec{X}])^2 \right].$$

Since $E \left[(g(\vec{X}) - E[Y \mid \vec{X}])^2 \right] > 0$ unless $P\{g(\vec{X}) \neq E[Y \mid \vec{X}]\} = 0$, the assertion at the end of the theorem follows. ■

The last theorem can be phrased as follows:

We interpret all functions $\phi(\vec{X})$ as all possible ways of creating random variables that only use the information available to \vec{X} where the quality of the approximation is measured by the **mean squared distance (MS Distance)** $E \left[(Y - g \circ \vec{X})^2 \right]$. Then

- the best MS Distance approximation of Y based on information provided by \vec{X} is $E[Y \mid \vec{X}]$

8.8 The Multinomial Probability Distribution

Introduction 8.4. In Definition 4.3 (p.77) of Chapter 4 (Combinatorial Analysis) we discussed multinomial coefficients

$$\binom{n}{n_1 n_2 \cdots n_k} = \frac{n!}{n_1! n_2! \cdots n_k!}$$

when counting the ways of classifying n items into k classes in such a way that n_1 items belong to class 1, n_2 items belong to class 2, ... n_k items belong to class k ($n_1 + \dots + n_k = n$). The multinomial probability distribution is based on those coefficients and generalizes the binomial distribution of Section 6.2 (Bernoulli Variables and the Binomial Distribution).

The binomial distribution is that of a random variable Y which counts the number of successes in n Bernoulli trials. (See Definition 6.5 on p.100 about Bernoulli trials.) To say this differently, Y counts the number of those Bernoulli trials which result in an outcome that falls into the “success class”.

The multinomial distribution will not be about a single random variable Y , but about a random vector $\vec{Y} = (Y_1, \dots, Y_k)$ of k random variables Y_j , which count the number of the n trials resulting in an outcome that falls into class j . What kind of trials are we talking about? We should expect those n random elements, let us call them X_1, \dots, X_n , to show some similarities to Bernoulli trials. Of course, there must be some significant differences. For example, each X_i will not have two outcomes (success or failure), but k outcomes corresponding to the k classes. \square

Definition 8.18 (Multinomial Sequence).

Let X_1, X_2, \dots be a finite or infinite sequence of random elements on a probability space (Ω, P) which take values in a set Ω' . We call this sequence a **multinomial sequence**, if the following are satisfied:

- (1) The sequence is iid.
- (2) There is some $k \in \mathbb{N}$ such that the outcome of each X_j is one of k distinct values $\omega'_1, \omega'_2, \dots, \omega'_k \in \Omega'$.

Since the X_j have identical distribution, there are probabilities p_1, p_2, \dots, p_k such that

- (3) $p_i := P\{X_j = \omega'_i\}$ is the same for all j and $p_1 + \dots + p_k = 1$.

If we consider a finite multinomial sequence X_1, X_2, \dots, X_n , we adopt the WMS notation and speak of a **multinomial experiment** of size n which consists of the **trials** X_j \square

Definition 8.19 (Multinomial distribution).

Assume that $\vec{Y} = (Y_1, Y_2, \dots, Y_k)$ is a vector of random variables which possesses the joint probability mass function

$$(8.59) \quad p_{\vec{Y}}(y_1, y_2, \dots, y_k) = \binom{n}{y_1, \dots, y_k} p_1^{y_1} p_2^{y_2} \cdots p_k^{y_k},$$

subject to the following conditions:

- $p_j \geq 0$ for $j = 1, 2, \dots, k$ and $\sum_{j=1}^k p_j = 1$.
- $y_i = 0, 1, 2, \dots, n$ for $i = 1, 2, \dots, k$ and $\sum_{i=1}^k y_i = n$.

Then we say that the random variables Y_i have a **multinomial distribution** with parameters n and $\vec{p} = (p_1, p_2, \dots, p_k)$. \square

Theorem 8.20.

Let $n \in \mathbb{N}$ and X_1, \dots, X_n be a multinomial sequence of size n . Let $p_j := P\{X_i = \omega'_j\}$. (That probability is the same for all i , since the X_i have identical distribution.)

Let $\vec{Y} = (Y_1, \dots, Y_k)$ be a vector of k random variables, such that each Y_j equals the number of the n trials resulting in an outcome that falls into class j . In other words,

$$(A) \quad Y_i(\omega) = y_i \Leftrightarrow X_j(\omega) = \omega'_j \text{ for exactly } y_i \text{ of the multinomial items } X_j.$$

Then \vec{Y} has a multinomial distribution with parameters n and $p_{\vec{Y}}(y_1, y_2, \dots, y_k)$.

PROOF: For fixed $\vec{y} = (y_1, \dots, y_k)$, the event $A := \{\vec{Y} = \vec{y}\}$ corresponds to all different ways that $\{1, 2, \dots, n\}$ can be partitioned into k subsets

$$(A) \quad \{1, 2, \dots, n\} = J_1 \uplus J_2 \uplus \dots \uplus J_k$$

such that each J_i contains y_i of those n indices. It follows from Theorem 4.6 on p.78 that

$$(B) \quad \text{there are } \binom{n}{y_1, y_2, \dots, y_k} \text{ different ways of creating such a partition.}$$

Thus, if we write

$$A(J_1, \dots, J_k) := \{X_{i_{m,1}} = \dots = X_{i_{m,y_m}} = \omega'_m \text{ for all } 1 \leq m \leq k\},$$

it follows that

$$(C) \quad P(A) = P\{\vec{Y} = \vec{y}\} = P\left(\uplus A(J_1, \dots, J_k)\right),$$

where this union is taken over all $\binom{n}{y_1, \dots, y_k}$ partitions J_1, \dots, J_k of $[1, n]_{\mathbb{Z}}$.

For a fixed $1 \leq m \leq k$, we write $J_m = \{i_{m,1} < i_{m,2} < \dots < i_{m,y_m}\}$. Since the X_j are independent,

$$P\{X_{i_{m,1}} = X_{i_{m,2}} = \dots = X_{i_{m,y_m}} = \omega'_m\} = P(\{X_{i_{m,1}} = \omega'_m\} \cap \dots \cap \{X_{i_{m,y_m}} = \omega'_m\}) = (p_m)^{y_m}$$

Since the X_j are independent not only for indices j belonging to J_m , but also across all J_m , it follows from the definition of $A(J_1, \dots, J_k)$ that

$$(D) \quad P(A(J_1, \dots, J_k)) = (p_1)^{y_1} (p_2)^{y_2} \dots (p_k)^{y_k}.$$

The right-hand side is independent of the particular partition J_1, \dots, J_k . We obtain from (B), (C) and (D) that

$$P\{\vec{Y} = \vec{y}\} = \binom{n}{y_1, \dots, y_k} (p_1)^{y_1} (p_2)^{y_2} \dots (p_k)^{y_k}.$$

Thus, \vec{Y} has the joint PMF that was specified in (8.59). We conclude that \vec{Y} has a multinomial distribution with parameters n and $p_{\vec{Y}}(y_1, y_2, \dots, y_k)$. ■

Example 8.6. Research by the marketing division of GreatWidgets Corp. has established that their customers' age is distributed as shown in the table to the right. A random sample of eight customers is taken. Assume that the proportions shown accurately reflect those of GreatWidgets Corp.

| Age | Proportion |
|------------------|------------|
| Group 1: 15 – 20 | 0.2 |
| Group 2: 21 – 30 | 0.2 |
| Group 3: 31 – 40 | 0.1 |
| Group 4: 41 – 50 | 0.2 |
| Group 5: > 50 | 0.3 |

what is the probability that the sample is composed as follows:

- Group 1: 1 person
- Group 2: 3 persons
- Group 4: 2 persons
- Group 5: 2 persons?

Solution:

We interpret the sample picks as the members X_1, \dots, X_8 of a multinomial sequence each of which has age group k as an outcome with probability p_k as indicated in the table.

Then the probability we are looking for is given by (8.59) on p.175

$$p_{\vec{y}}(y_1, y_2, y_3, y_4, y_5) = \frac{n!}{y_1! y_2! y_3! y_4! y_5!} p_1^{y_1} p_2^{y_2} p_3^{y_3} p_4^{y_4} p_5^{y_5}, \quad \square$$

In the context of this example we obtain

$$p(1, 3, 2, 0, 2) = \frac{8!}{1! 3! 2! 0! 2!} 0.2^1 0.2^2 0.1^0 0.2^2 0.3^0, = 0.009768. \quad \square$$

8.9 Order Statistics

A

The presentation of the material in this section is largely based on the 2015 Math 447 lecture notes of Prof. Xingye Qiao, Binghamton University

@@Author

Given are n random variables $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$. One can sort them, for any fixed $\omega \in \Omega$, in nondecreasing order. One obtains in this fashion a sequence, of size n , of numbers

$$Y_{(1)}(\omega) \leq Y_{(2)}(\omega) \leq Y_{(3)}(\omega) \leq \dots \leq Y_{(n)}(\omega).$$

Since these numbers depend on randomness ω , each $Y_{(j)}(\omega)$ represents an outcome of a random variable $Y_{(j)}$.

Example 8.7. Here are some examples.

(a) 70 students are randomly selected when exiting lecture hall and their age is measured in years. Those 70 ages, $A_1(\omega), \dots, A_{70}(\omega)$, are sorted in increasing order:

- $A_{(1)}(\omega)$ = height of the smallest person in the sample
- $A_{(2)}(\omega)$ = height of the second smallest person in the sample
- -----
- $A_{(j)}(\omega)$ = height of the j th smallest person in the sample
- -----
- $A_{(n)}(\omega)$ = height of the tallest person in the sample

Clearly, $A_{(1)}(\omega) \leq A_{(2)}(\omega) \leq A_{(3)}(\omega) \leq \dots \leq A_{(n)}(\omega)$.

Almost all of those ages will be one of 18, 19, ..., 25. Accordingly, it is not only possible that we encounter an index j that results in equality, $A_{(j)} = A_{(j+1)}$, but this will be the rule rather than the exception.

(b) Rather than considering the age of those 70 students, we now look at their height, measured in millimeters. Those 70 heights, $H_1(\omega), \dots, H_{70}(\omega)$, are sorted in increasing order.

Height can be considered a continuous random variable. Since the probability of two students having precisely the same height is zero, we may consider the outcomes $H_{(j)}$ distinct. Accordingly, we can replace “less or equal” with strict inequality and obtain

$$H_{(1)}(\omega) < H_{(2)}(\omega) < H_{(3)}(\omega) < \cdots < H_{(n)}(\omega). \quad \square$$

- We will deal in this section exclusively with continuous random variables.
- When considering a finite or infinite sequence Y_1, Y_2, Y_3, \dots of such random variables, we assume that they are iid (independent and identically distributed).

Definition 8.20 (Order statistics).

Given n iid continuous random variables $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$, we sort them in increasing order. The resulting sequence of random variables, which we denote as $Y_{(j)}, j = 1, \dots, n$, then satisfies

$$(8.60) \quad Y_{(1)} \leq Y_{(2)} \leq Y_{(3)} \leq \cdots \leq Y_{(n)}.$$

We call $Y_{(j)}$ the **j th order statistic** of \vec{Y} .

See Example 8.7(b) why we may consider strictly increasing rather than nondecreasing. \square

Assumption 8.4.

Unless explicitly stated otherwise,

- $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$ denotes a list of n iid continuous random variables ($n \in \mathbb{N}$).
- $Y_1 \sim Y_2 \sim \cdots \sim Y_n$ implies $F_{Y_1} = F_{Y_2} = \cdots = F_{Y_n}$ and $f_{Y_1} = f_{Y_2} = \cdots = f_{Y_n}$.
We write $F(y) := F_{Y_j}(y)$ and $f(y) := f_{Y_j}(y)$ for the common CDF and PDF. \square

Remark 8.12. Note that

- The **first order statistic** or **smallest order statistic** is $Y_{(1)} = \min\{Y_1, \dots, Y_n\}$.
- The n th order statistic or **largest order statistic** is $Y_{(n)} = \max\{Y_1, \dots, Y_n\}$.
- A simple consequence of the definition of min and max are the following formulas:

$$(8.61) \quad Y_{(1)}(\omega) > y \Leftrightarrow \min(Y_1(\omega), \dots, Y_n(\omega)) > y \Leftrightarrow Y_j(\omega) > y \text{ for all } j \in [1, n]_{\mathbb{Z}},$$

$$(8.62) \quad Y_{(n)}(\omega) \leq y \Leftrightarrow \max(Y_1(\omega), \dots, Y_n(\omega)) \leq y \Leftrightarrow Y_j(\omega) \leq y \text{ for all } j \in [1, n]_{\mathbb{Z}}. \quad \square$$

Theorem 8.21 (CDF and PDF of the j th order statistic).

For $y \in \mathbb{R}$, the CDF of the k th order statistic ($k = 1, \dots, n$) satisfies the following:

$$(8.63) \quad F_{Y_{(1)}}(y) = 1 - [1 - F(y)]^n,$$

$$(8.64) \quad F_{Y_{(n)}}(y) = [F(y)]^n,$$

$$(8.65) \quad F_{Y_{(k)}}(y) = 1 - \sum_{j=0}^{k-1} \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j} = \sum_{j=k}^n \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j}.$$

For $y \in \mathbb{R}$, the PDF of the k th order statistic ($k = 1, \dots, n$) satisfies the following:

$$(8.66) \quad f_{Y_{(1)}(y)} = n [1 - F(y)]^{n-1} f(y),$$

$$(8.67) \quad f_{Y_{(n)}(y)} = n [F(y)]^{n-1} f(y),$$

$$(8.68) \quad f_{Y_{(k)}(y)} = \sum_{j=0}^{k-1} \binom{n}{j} f(y) \left(n [F(y)]^{n-1} - j [F(y)]^{j-1} \right).$$

$$(8.69) \quad f_{Y_{(k)}(y)} = n \binom{n-1}{k-1} f(y) \cdot [F(y)]^{k-1} \cdot [1 - F(y)]^{n-k}.$$

Note that the proofs are not given in the order of the seven formulas of the theorem.

PROOF of (8.64):

$$\begin{aligned} F_{Y_{(n)}(y)} &\stackrel{(8.62)}{=} P(\{Y_1 \leq y\} \cap \{Y_2 \leq y\} \cap \dots \cap \{Y_n \leq y\}) \\ &\stackrel{\text{indep}}{=} P\{Y_1 \leq y\} \cdot P\{Y_2 \leq y\} \cdot \dots \cdot P\{Y_n \leq y\} = [F(y)]^n. \end{aligned}$$

PROOF of (8.63):

$$\begin{aligned} P\{Y_{(1)} > y\} &\stackrel{(8.61)}{=} P(\{Y_1 > y\} \cap \{Y_2 > y\} \cap \dots \cap \{Y_n > y\}) \\ &\stackrel{\text{indep}}{=} P\{Y_1 > y\} \cdot P\{Y_2 > y\} \cdot \dots \cdot P\{Y_n > y\} = [1 - F(y)]^n. \end{aligned}$$

Thus, $F_{Y_{(1)}(y)} = 1 - P\{Y_{(1)} > y\} = 1 - [1 - F(y)]^n$.

PROOF of (8.66) and (8.67):

This follows from $\frac{d}{dy} (1 - [1 - F(y)]^n) = -n[1 - F(y)]^{n-1} (-f(y))$

and $\frac{d}{dy} ([F(y)]^n) = n[F(y)]^{n-1} f(y)$.

PROOF of (8.65):

This proof requires a lot more work than the proofs we have done so far. It will be done by constructing a binomial random variable.

- Since y is fixed, so is $p := F(y) = P\{Y_j \leq y\}$. (Identical for all j , since the Y_j are iid.)
- For $j = 1, \dots, n$, let $X_j(\omega) := \begin{cases} 1 & \text{if } Y_j(\omega) \leq y, \\ 0 & \text{else.} \end{cases}$ Let $U(\omega) := \sum_{j=1}^n X_j(\omega)$.
- We interpret $Y_j(\omega) \leq y$ as a success and $Y_j(\omega) > y$ as a failure. Then X_1, \dots, X_n form a 0–1 encoded Bernoulli sequence⁴⁹ and $U \sim \text{binom}(n, p)$, since U counts the number of successes.
- Observe that $Y_{(k)}(\omega) \leq y \Leftrightarrow Y_j(\omega) \leq y$ at least k times \Leftrightarrow there are at least k successes $\Leftrightarrow U(\omega) \geq k$. It does not matter whether or not there are more than k successes.
- Thus, $F_{Y_{(k)}(y)} = P\{Y_{(k)} \leq y\} = P\{U \geq k\} = \sum_{j=k}^n P\{U = j\} = 1 - \sum_{j=0}^{k-1} P\{U = j\}$.
- Since $U \sim \text{binom}(n, p)$ and $p = F(y)$, $F_{Y_{(k)}(y)} = 1 - \sum_{j=0}^{k-1} \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j}$.

⁴⁹See Definition 6.5 (Bernoulli items and variables) on p.100.

PROOF of (8.68):

This is done by differentiation. For each $j = 0, \dots, k-1$,

$$(A) \quad \begin{aligned} \frac{d}{dy} \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j} &= \binom{n}{j} \frac{d}{dy} ([F(y)]^j - F(y)^n) \\ &= \binom{n}{j} (j [F(y)]^{j-1} f(y) - n F(y)^{n-1} f(y)) \end{aligned}$$

$$\begin{aligned} \text{Thus, } f_{Y^{(k)}} &= \frac{d}{dy} \left[1 - \sum_{j=0}^{k-1} \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j} \right] \\ &= - \sum_{j=0}^{k-1} \frac{d}{dy} \binom{n}{j} ([F(y)]^j [1 - F(y)]^{n-j}) \\ &\stackrel{(A)}{=} \sum_{j=0}^{k-1} \binom{n}{j} f(y) (n [F(y)]^{n-1} - j [F(y)]^{j-1}). \end{aligned}$$

This finishes the proof of (8.68).

The proof of (8.69) is based on an entirely different approach. Before we do that proof, we first illustrate that approach by redoing those of (8.66) and (8.67). Those proofs are much simpler and are a good preparation for that of (8.69).

ALTERNATE PROOF of (8.66):

First we note the following for a continuous random variable U with density $f_U(u)$. Assume that $\delta > 0$ is very close to zero. Since we assumed for all our continuous random variables that they have continuous density, $f_U(\cdot) \approx \text{const} = f_U(u)$ on $]u, u + \delta[$.

(a) Thus, $P\{u < U \leq u + \delta\} = \int_u^{u+\delta} f_U(t) dt \approx f_U(u) \cdot \delta$.

(b) For the fixed y and some “really small” δy , we create three events:

□ L (for “left-hand side”) □ I (for “inside”) □ R (for “right-hand side”),

and a sequence of random elements X_1, \dots, X_n as follows.

□ $X_j(\omega) = L \Leftrightarrow Y_j(\omega) \leq y$. Then $P\{X_j = L\} = P\{Y_j \leq y\} = F(y)$.

□ $X_j(\omega) = I \Leftrightarrow y < Y_j(\omega) \leq y + \delta$, Then $P\{X_j = I\} = P\{y < Y_j \leq y + \delta\} \stackrel{(a)}{\approx} f_U(u) \cdot \delta$.

□ $X_j(\omega) = R \Leftrightarrow Y_j(\omega) > y + \delta$. Then $P\{X_j = R\} = P\{Y_j > y + \delta\} = 1 - F(y + \delta)$.

(c) By construction, the X_j form a multinomial sequence. Let $\vec{U} := (U_1, U_2, U_3)$, where

□ $U_1 := \#$ of indices j such that $X_j = L$,

□ $U_2 := \#$ of indices j such that $X_j = I$,

□ $U_3 := \#$ of indices j such that $X_j = R$.

(d) Then \vec{U} is multinomial with parameters $n, p_1 = F(y), p_2 = f(y)\delta, p_3 = 1 - F(y)$.

(e) Since we assume that $Y_{(j)}(\omega)$ is strictly increasing with j for all ω , it seems reasonable that, for “really small” δ , the following is true:

- If $Y_{(1)}(\omega) > y$, then $Y_{(j)}(\omega) > y + \delta$ for all $j > 1$.

- (f) Thus, $f_{Y_{(1)}}(y) \cdot \delta \stackrel{\text{(a)}}{\approx} P\{y < Y_{(1)} \leq y + \delta\}$
 $= P\{\text{exactly one of } Y_1, \dots, Y_n \in]y, y + \delta] \text{ and } Y_j > y + \delta \text{ for all other } j\}$
 $= P\{\text{none of the } X_j \text{ are } L \text{ and exactly one is } I \text{ and } n - 1 \text{ are } R\}$
 $= P\{U_1 = 0, U_2 = 1, U_3 = n - 1, \} \stackrel{\text{(d)}}{=} \binom{n}{0, 1, n - 1} [F(y)]^0 [f(y)\delta]^1 [1 - F(y + \delta)]^{n-1}.$
- (g) Since $\binom{n}{0, 1, n - 1} = \frac{n!}{0! \cdot 1! \cdot (n - 1)!} = n,$
 we obtain $f_{Y_{(1)}}(y) \cdot \delta \approx n [1 - F(y + \delta)]^{n-1} f(y)\delta.$

We divide both expressions by δ , then let $\delta \rightarrow 0$. Since $t \mapsto F(t)$ is continuous, $\lim_{\delta \rightarrow 0} F(y + \delta) = F(y)$.
 We conclude that the density of $Y_{(1)}$ is

$$f_{Y_{(1)}}(y) = n [1 - F(y)]^{n-1} f(y).$$

ALTERNATE PROOF of (8.67):

We can adapt the alternate proof for the density of $Y_{(1)}$ to obtain that of $Y_{(n)}$ as follows.

We keep all items through (e) and modify (f) and (g) as follows.

- (f') $f_{Y_{(n)}}(y) \cdot \delta \stackrel{\text{(a)}}{\approx} P\{y < Y_{(n)} \leq y + \delta\}$
 $= P\{\text{exactly one of } Y_1, \dots, Y_n \in]y, y + \delta] \text{ and } Y_j \leq y \text{ for all other } j\}$
 $= P\{\text{none of the } X_j \text{ are } R \text{ and exactly one is } I \text{ and } n - 1 \text{ are } L\}$
 $= P\{U_1 = n - 1, U_2 = 1, U_3 = 0, \} \stackrel{\text{(d)}}{=} \binom{n}{n - 1, 1, 0} [F(y)]^{n-1} [f(y)\delta]^1 [1 - F(y + \delta)]^0.$
- (g') Since $\binom{n}{n - 1, 1, 0} = \frac{n!}{(n - 1)! \cdot 1! \cdot 0!} = n,$
 we obtain $f_{Y_{(n)}}(y) \cdot \delta \approx n [F(y)]^{n-1} f(y)\delta.$

We divide both expressions by δ , then let $\delta \rightarrow 0$. We obtain the density of $Y_{(n)}$ as

$$f_{Y_{(n)}}(y) = n [F(y)]^{n-1} f(y).$$

PROOF of (8.69):

This time we adapt the alternate proof for the density of $Y_{(1)}$ to obtain that of $Y_{(k)}$ as follows.

We keep all items through (e) and modify (f) and (g) as follows.

- (f'') $f_{Y_{(k)}}(y) \cdot \delta \stackrel{\text{(a)}}{\approx} P\{y < Y_{(k)} \leq y + \delta\}$
 $= P\{\text{exactly one of } Y_1, \dots, Y_n \in]y, y + \delta] \text{ and } Y_j \leq y \text{ for } k - 1 \text{ indices } j$
 $\quad \text{and } Y_j > y \text{ for } n - k \text{ indices } j\}$
 $= P\{k - 1 \text{ of the } X_j \text{ are } L, n - k \text{ of the } X_j \text{ are } R \text{ and exactly one is } I\}$
 $= P\{U_1 = k - 1, U_2 = 1, U_3 = n - k\}$
 $\stackrel{\text{(d)}}{=} \binom{n}{k - 1, 1, n - k} [F(y)]^{k-1} [f(y)\delta]^1 [1 - F(y + \delta)]^{n-k}.$

(g'') Since $\binom{n}{k-1, 1, n-k} = \frac{n \cdot (n-1)!}{(k-1)! \cdot 1! \cdot (n-k)!} = n \cdot \binom{n-1}{k-1}$,
 we obtain $f_{Y_{(k)}}(y) \cdot \delta \approx n \cdot \binom{n-1}{k-1} [F(y)]^{k-1} f(y) \delta [1 - F(y + \delta)]^{n-k}$.

We divide both expressions by δ , then let $\delta \rightarrow 0$. Since $t \mapsto F(t)$ is continuous, $\lim_{\delta \rightarrow 0} F(y + \delta) = F(y)$.
 We conclude that the density of $Y_{(1)}$ is

$$f_{Y_{(k)}}(y) = n \binom{n-1}{k-1} [F(y)]^{k-1} f(y) [1 - F(y)]^{n-k}. \blacksquare$$

Remark 8.13. (8.65) yields (8.63) for $k = 1$ and (8.64) for $k = n$. This can be seen as follows:

Recall that

$$\begin{aligned} 1 &= (F(y) + [1 - F(y)])^n = \sum_{j=0}^n \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j} \\ \text{(A)} \quad &= \sum_{j=0}^{n-1} \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j} + \binom{n}{0} [F(y)]^0 [1 - F(y)]^n. \end{aligned}$$

If we evaluate (8.65) for $k = 1$ and $k = n$, we obtain

$$\begin{aligned} F_{Y_{(1)}}(y) &= 1 - \binom{n}{0} [F(y)]^0 [1 - F(y)]^n = 1 - 1 \cdot 1 \cdot [1 - F(y)]^n = [1 - F(y)]^n, \\ F_{Y_{(n)}}(y) &= 1 - \sum_{j=0}^{n-1} \binom{n}{j} [F(y)]^j [1 - F(y)]^{n-j} \stackrel{\text{(A)}}{=} \binom{n}{0} [F(y)]^0 [1 - F(y)]^n = [1 - F(y)]^n. \quad \square \end{aligned}$$

Remark 8.14. You may have noticed that there are two formulas for $f_{Y_{(k)}}(y)$.

(8.69) was shown by means of the “density approach” that utilized a limiting process $\delta \rightarrow 0$ in conjunction with the multinomial distribution. The proof was harder than that of (8.68). In return, (8.69) has computational advantages, since no more summation $\sum_{j=0}^{k-1}$ is required. \square

The next remark belongs thematically into Section 4.2 (Permutations) of Chapter 4. However, it has been placed here, since every order statistic

$$\vec{Y}_{(\bullet)} = (Y_{(1)}, \dots, Y_{(n)}).$$

is a (specific) permutation of $\vec{Y} = (Y_1, \dots, Y_n)$, and every other permutation

$$(Y_{i_1}, Y_{i_2}, \dots, Y_{i_n})$$

of $\vec{Y} = (Y_1, \dots, Y_n)$, possesses the same order statistic.

Remark 8.15. If we deal with a list $\vec{a} = (a_1, a_2, \dots, a_n)$ of distinct numbers, e.g.,

$$\text{(A)} \quad \vec{a} = (13.2, -3, 6.6, 2, -1.5),$$

then there is a uniquely determined permutation, $\vec{a}_{(\bullet)} = (a_{(1)}, a_{(2)}, \dots, a_{(n)})$ of \vec{a} , which has those a_j in increasing order. In other words,

$$a_{(1)} < a_{(2)} < \dots < a_{(n)}.$$

In the specific example **(A)**, we obtain

$$\vec{a}_{(\bullet)} = (-3, -1.5, 2, 6.6, 13.2).$$

If $\vec{b} = (b_1, b_2, \dots, b_n)$ is another list of distinct numbers, then

$$\vec{b}_{(\bullet)} = \vec{a}_{(\bullet)} \quad \Leftrightarrow \quad \vec{b} \text{ is a permutation of } \vec{a}.$$

Going back to our example, if

$$\begin{aligned} \vec{b} &= (13.2, 6.6, -1.5, -3, 2), \\ \vec{c} &= (13.2, -3, 6.6, 2, -1.51), \end{aligned}$$

then $\vec{b}_{(\bullet)} = \vec{a}_{(\bullet)}$, but $\vec{c}_{(\bullet)} \neq \vec{a}_{(\bullet)}$, since $\vec{a}_{(\bullet)}$ does not include the number -1.51 . \square

Theorem 8.22 (Joint PDF of the order statistic).

A: Let $\vec{y} \in \mathbb{R}^n$ satisfy

$$(8.70) \quad y_1 < y_2 < \dots < y_n.$$

For the vector $\vec{Y} = (Y_1, \dots, Y_n)$, let $\vec{Y}_{(\bullet)}$ be the vector of its associated order statistics, i.e.,

$$(8.71) \quad \vec{Y}_{(\bullet)} = (Y_{(1)}, \dots, Y_{(n)}).$$

Then its density function at \vec{y} is given by

$$(8.72) \quad f_{\vec{Y}_{(\bullet)}}(\vec{y}) = n! \cdot \prod_{j=1}^n f(y_j) = n! f(y_1) \cdots f(y_n).$$

B: If \vec{y} does not satisfy (8.70), then $f_{\vec{Y}_{(\bullet)}}(\vec{y}) = 0$.

FIRST PROOF:

Let Δ be a “small” n -dimensional cube with volume $Vol(\Delta)$ that is centered at \vec{y} . Study the proof of (8.65) of Theorem 8.21 on p.178. It explains (in the onedimensional case), why one can approximate

$$\begin{aligned} P\{\vec{Y} \in \Delta\} &\approx f_{\vec{Y}}(\vec{y}) \cdot Vol(\Delta), \\ P\{\vec{Y}_{(\bullet)} \in \Delta\} &\approx f_{\vec{Y}_{(\bullet)}}(\vec{y}) \cdot Vol(\Delta). \end{aligned}$$

A cube of sidelength 2ε has volume $Vol(\Delta) = (2\varepsilon)^n$. If we solve that equation for ε , we obtain

$$\varepsilon = \frac{Vol(\Delta)^{1/n}}{2}.$$

Since $y_1 < y_2 < \dots < y_n$, one can choose Δ and hence, $\varepsilon = Vol(\Delta)^{1/n}/2$, so small, that any two intervals $[y_i - \varepsilon, y_i + \varepsilon]$ and $[y_j - \varepsilon, y_j + \varepsilon]$ have empty intersection for $i \neq j$.

For the following, see Remark 8.15 on p.182. Note that

$$(A) \quad \begin{aligned} \vec{Y}_{(\bullet)}(\omega) \in \Delta &\Leftrightarrow y_k - \varepsilon \leq Y_{(k)}(\omega) \leq y_k + \varepsilon \text{ for all } k, \\ &\Leftrightarrow \text{for all } k, \exists j \text{ such that } y_k - \varepsilon \leq Y_j(\omega) \leq y_k + \varepsilon. \end{aligned}$$

We illustrate this point for $n = 3$, $Vol(\Delta) = 1/8$, $y_1 = 2.6$, $y_2 = 4.2$, $y_3 = 7.8$. $\varepsilon = (1/8^3)/2 = 0.25$. This is small enough for the intervals $y_j \pm 0.25$ to be disjoint.

There are $3! = 6$ different ways to have $\vec{Y}(\omega) \in \Delta$. They are:

- (1) $2.35 \leq Y_1(\omega) \leq 2.85$, $3.95 \leq Y_2(\omega) \leq 4.45$, $7.55 \leq Y_3(\omega) \leq 8.05$,
- (2) $2.35 \leq Y_1(\omega) \leq 2.85$, $3.95 \leq Y_3(\omega) \leq 4.45$, $7.55 \leq Y_2(\omega) \leq 8.05$,
- (3) $2.35 \leq Y_2(\omega) \leq 2.85$, $3.95 \leq Y_1(\omega) \leq 4.45$, $7.55 \leq Y_3(\omega) \leq 8.05$,
- (4) $2.35 \leq Y_2(\omega) \leq 2.85$, $3.95 \leq Y_3(\omega) \leq 4.45$, $7.55 \leq Y_1(\omega) \leq 8.05$,
- (5) $2.35 \leq Y_3(\omega) \leq 2.85$, $3.95 \leq Y_1(\omega) \leq 4.45$, $7.55 \leq Y_2(\omega) \leq 8.05$,
- (6) $2.35 \leq Y_3(\omega) \leq 2.85$, $3.95 \leq Y_2(\omega) \leq 4.45$, $7.55 \leq Y_1(\omega) \leq 8.05$,

Let us assume that $k = 2$, i.e., we consider the interval $[3.95, 4.45]$.

In (2) and (4), we choose $j = 3$ to obtain $Y_j \in [3.95, 4.45]$.

On the other hand, in (1) and (6), we choose $j = 2$ to obtain $Y_j \in [3.95, 4.45]$.

We refer you again to Remark 8.15 on p.182 to understand that (A) shows that

$$(B) \quad \begin{aligned} \vec{Y}_{(\bullet)}(\omega) \in \Delta &\Leftrightarrow \text{some permutation of } \vec{Y}(\omega) \in \Delta \\ &\Leftrightarrow \text{each permutation of } \vec{Y}(\omega) \in \Delta. \end{aligned}$$

- Since a list of n items has $n!$ permutations, there are $n!$ such (disjoint) events: There are $n!$ permutations (k_1, k_2, \dots, k_n) of $(1, 2, \dots, n)$ with corresponding event

$$\{y_1 - \varepsilon \leq Y_{k_1} \leq y_1 + \varepsilon\} \cap \{y_2 - \varepsilon \leq Y_{k_2} \leq y_2 + \varepsilon\} \cap \dots \cap \{y_n - \varepsilon \leq Y_{k_n} \leq y_n + \varepsilon\}.$$

- Since the Y_j are iid and $P\{y_i - \varepsilon \leq Y_{k_j} \leq y_i + \varepsilon\} \approx 2\varepsilon \cdot f(y_i)$ for each i and j , each such event has probability $\approx \prod_{j=1}^n f(y_j) \cdot (2\varepsilon)^n$.
- Thus, $f_{\vec{Y}_{(\bullet)}}(\vec{y}) \cdot Vol(\Delta) \approx n! \cdot \prod_{j=1}^n f(y_j) \cdot Vol(\Delta)$
- As $\Delta \rightarrow 0$, " \approx " becomes " $=$ " and then $f_{\vec{Y}_{(\bullet)}}(\vec{y}) = n! \cdot \prod_{j=1}^n f(y_j)$. ■

ALTERNATE PROOF:

- (a) We may assume that \vec{y} satisfies $y_1 < y_2 < \dots < y_n$, since $f_{\vec{Y}_{(\bullet)}}(\vec{y}) = 0$ otherwise.
 - For small enough dt_1, dt_2, dt_n , the intervals $[y_j, y_j + dt_j]$ are disjoint.
- (b) Thus, $[y_j \leq Y_{(j)}(\omega) \leq y_j + dt_j \text{ for all } j] \Leftrightarrow [\text{there is a permutation } i_1, i_2, \dots, i_n \text{ of the indices } 1, 2, \dots, n \text{ such that } y_j \leq Y_{i_j}(\omega) \leq y_j + dt_j \text{ for all } j]$

- (c) Thus, $[y_j \leq Y_{(j)}(\omega) \leq y_j + dt_j \text{ for all } j] \Leftrightarrow [\text{among the } X_i(\omega), \text{ exactly one is in } [y_1, y_1 + dt_1], \text{ exactly one is in } [y_2, y_2 + dt_2], \dots, \text{ exactly one is in } [y_n, y_n + dt_n]].$ (Thus, NONE are outside the union of those intervals.)
- (d) This can be interpreted as the counts of the outcomes of a multinomial sequence X_1, \dots, X_n , where $X_k(\omega)$ results in outcome $\#j$, if $y_j \leq Y_k \leq y_j + dt_j$.
- The probabilities $p_j = P\{X_k \text{ results in } \#j\}$ are, for small enough dt_j , equal to

$$p_j = P\{Y_i \in [y_j, y + dt_j]\} = \int_{y_j}^{y+dt_j} f(t) dt \approx f(t_j) dt_j.$$

- (e) From (b), (c), (d):

$$\begin{aligned} f_{\vec{Y}_{(\bullet)}}(\vec{y}) dt_1 \cdots dt_n &= P\{y_j \leq Y_{(j)}(\omega) \leq y_j + dt_j \text{ for all } j\} \\ &= P\{\text{there is a permutation } i_1, i_2, \dots, i_n \text{ of the indices } 1, 2, \dots, n \\ &\quad \text{such that } y_j \leq Y_{i_j} \leq y_j + dt_j \text{ for all } j\} \\ &= P\{\text{each } X_k \text{ has exactly one outcome } \#j \text{ for each } j = 1, \dots, n\} \\ &= \binom{n}{1, 1, \dots, 1} p_1^1 p_2^1 \cdots p_n^1 = \frac{n!}{1! \cdots 1!} \prod_j (f(t_j) dt_j). \end{aligned}$$

$$\text{Thus, } f_{\vec{Y}_{(\bullet)}}(\vec{y}) dt_1 \cdots dt_n = n! \prod_j f(t_j) (dt_1 \cdots dt_n).$$

- (f) We cancel $dt_1 \cdots dt_n$ on both sides and obtain $f_{\vec{Y}_{(\bullet)}}(\vec{y}) = n! \prod_j f(t_j)$. ■

Example 8.8. Find the formula for the joint density of $Y_{(1)}$ and $Y_{(n)}$.

Solution:

- (a) Note that, since the Y_j are continuous, “ $<$ ” and “ \leq ” can be interchanged and the same is true for “ $>$ ” and “ \geq ” when computing probabilities.
- (b) Also, applying $A = (A \cap B) \uplus A \cap B^c$ with $A = \{Y_{(n)} \leq y_n\}$ and $B = \{Y_{(1)} \leq y_1\}$ yields
- $$P\{Y_{(n)} \leq y_n\} = P\{Y_{(n)} \leq y_n, Y_{(1)} \leq y_1\} + P\{Y_{(n)} \leq y_n, Y_{(1)} > y_1\}.$$

We find the CDF as follows:

$$\begin{aligned} F_{Y_{(1)}, Y_{(n)}}(y_1, y_n) &\stackrel{\text{(b)}}{=} P\{Y_{(n)} \leq y_n\} - P\{Y_{(1)} > y_1, Y_{(n)} \leq y_n\} \\ &= P\{Y_j \leq y_n \text{ for all } j\} - P\{y_1 < Y_j \leq y_n \text{ for all } j\} \\ &= \prod_{j=1}^n P\{Y_j \leq y_n\} - \prod_{j=1}^n P\{y_1 < Y_j \leq y_n\} = [F(y_n)]^n - [F(y_n) - F(y_1)]^n. \end{aligned}$$

We used first independence, then identical distribution in the last line.

Differentiation of the above then gives us $f_{Y_{(1)}, Y_{(n)}}(y_1, y_n)$ as follows:

For convenience, we define $G(y_1, y_n) := F_{Y_{(1)}, Y_{(n)}}(y_1, y_n)$. Then,

$$G(y_1, y_n) = [F(y_n)]^n - [F(y_n) - F(y_1)]^n$$

Thus,

$$\frac{\partial G}{\partial y_1} = 0 - n[F(y_n) - F(y_1)]^{n-1} f(y_1) = n \cdot f(y_1) [F(y_n) - F(y_1)]^{n-1}$$

Thus,

$$\begin{aligned} f_{Y_{(1)}, Y_{(n)}}(y_1, y_n) &= \frac{\partial^2 G}{\partial y_1 \partial y_n} = n \cdot f(y_1) \cdot (n-1) [F(y_n) - F(y_1)]^{n-2} \cdot f(y_n) \\ &= n(n-1) \cdot f(y_1) f(y_n) \cdot [F(y_n) - F(y_1)]^{n-2} \end{aligned}$$

Alternate solution:

The PDF can be found by interpreting certain events related to finding the density as the outcomes of the following multinomial sequence, $\vec{X} = (X_1, \dots, X_n)$,

- (c) For a given j , the outcomes ω'_i and associated probabilities p_i for X_j are

$$\square \omega'_1: Y_j \text{ is close to } y_1 \Rightarrow p_1 = f(y_1) dy_1 \quad \square \omega'_2: Y_j \text{ is close to } y_n \Rightarrow p_2 = f(y_n) dy_n$$

$$\square \omega'_3: Y_j \text{ strictly inbetween } y_1 \text{ and } y_n \Rightarrow y_1 < Y_j < y_n \Rightarrow p_3 = F(y_n) - F(y_1).$$

Note that it is impossible that none of $\omega'_1, \omega'_2, \omega'_3$ happens and $Y_j < y_1$ or $Y_j > y_n$.

- (d) We denote by W_i the count of indices j such that $X_j = \omega'_i$.

Then $\vec{W} = (W_1, W_2, W_3) \sim \text{multinomial}^{50}$ with joint PMF $p_{\vec{W}}(\vec{w})$ given by

$$p_{\vec{W}}(\vec{w}) = \binom{n}{w_1, w_2, w_3} p_1^{w_1} p_2^{w_2} p_3^{w_3}.$$

- Similar to what was done in the proofs of theorems 8.21 (CDF and PDF of the j th order statistic) and 8.22 (Joint PDF of the order statistic), we conclude from (c) and (d) that

$$\begin{aligned} \text{(e)} \quad f_{Y_{(1)}, Y_{(n)}}(y_1, y_n) dy_1 dy_n &= P\{Y_{(1)} \text{ is "dy}_1 \text{ close" to } y_1 \text{ and } Y_{(n)} \text{ is "dy}_n \text{ close" to } y_n\} \\ &= P\{\text{exactly one } Y_j \text{ is "dy}_1 \text{ close" to } y_1 \text{ and exactly one } Y_j \text{ is "dy}_n \text{ close" to } y_n \\ &\quad \text{and the other } Y_j \text{ (there are } n-2 \text{ left) are between } y_1 \text{ and } y_n\} \\ &= P\{W_1 = 1, W_2 = 1, W_3 = n-2\} = p_{\vec{W}}(1, 1, n-2) = \binom{n}{1, 1, n-2} p_1^1 p_2^1 p_3^{n-2}. \\ &= n(n-1) \cdot f(y_1) dy_1 \cdot f(y_n) dy_n \cdot [F(y_n) - F(y_1)]. \end{aligned}$$

$$\text{(f)} \quad \text{Thus, } f_{Y_{(1)}, Y_{(n)}}(y_1, y_n) dy_1 dy_n \stackrel{\text{(e)}}{=} n(n-1) \cdot f(y_1) \cdot f(y_n) \cdot [F(y_n) - F(y_1)]^{n-2} dy_1 dy_n.$$

- We cancel $dy_1 dy_n$ in that last equation and obtain

$$\text{(g)} \quad f_{Y_{(1)}, Y_{(n)}}(y_1, y_n) = n(n-1) \cdot f(y_1) \cdot f(y_n) \cdot [F(y_n) - F(y_1)]^{n-2}.$$

We have obtained the same result for the joint PDF of $Y_{(1)}$ and $Y_{(n)}$ as in the first solution. \square

Remark 8.16 (Sample median). Recall from Definition 7.4 (p th quantile) on p.119 that the median of a random variable U with CDF $F_U(\cdot)$ was its 0.5th quantile

$$\phi_{0.5} = \min\{u \in \mathbb{R} : F_U(u) \geq 0.5\}.$$

If U is continuous with a strictly increasing CDF, then $\phi_{0.5}$ is that unique value u , for which $F_U(u) = 0.5$. Thus, half of the area under the density $f_U(\cdot)$ is to the left of $\phi_{0.5}$ and the other half is to the right of $\phi_{0.5}$.

⁵⁰See Definition 8.19 (Multinomial distribution) on p.175.

Assume that $\vec{Y} = (Y_1, \dots, Y_n)$ describes the action of picking a sample of n real numbers. In other words, each Y_j is a random variable and each invocation $\vec{Y}(\omega)$ results in the specific sample $\vec{y} = (y_1, \dots, y_n)$, where

$$y_1 = Y_1(\omega), y_2 = Y_2(\omega), \dots, y_n = Y_n(\omega).$$

Further assume that the Y_j are continuous. Then we can assume that all sample picks Y_1, \dots, Y_n are distinct, so that the order statistic satisfies strict inequalities

$$(A) \quad Y_{(1)} < Y_{(2)} < \dots < Y_{(n)}.$$

The **sample median** of \vec{Y} is defined as follows.

- (a) If $n = 2k + 1$ is odd, then the sample median of \vec{Y} is the $(k + 1)$ th order statistic $Y_{(k+1)}$.
- (b) If $n = 2k$ is even, then the sample median of \vec{Y} is the (random) average $\frac{Y_k + Y_{k+1}}{2}$.

Two examples:

- (1) If $n = 7$, then the sample median is $Y_{(n+1)} = Y_{(4)}$. Three of the Y_j are to the left of $Y_{(4)}$ and the same number are to the right.
- (2) If $n = 8$, then the sample median of \vec{Y} is the average $\frac{Y_4 + Y_5}{2}$. Since we have strict inequalities in (A), four of the Y_j are to the left of the sample median and the same number are to the right.

The point to remember is that the sample median of an odd-sized sampling action is an order statistic, whereas that of an even sized one is not.

Example: Let us assume that the the sample picks of an odd sized sample $\vec{Y} = (Y_1, \dots, Y_{2n+1})$ are continuous and iid random variables. We can compute the PDF of the sample median as that of $Y_{(n+1)}$. This time we do so by associating a multinomial random vector with three outcomes: Either Y_j is near y_{n+1} or it is near one of the n values to the left or it is near one of the n values to the right. In that manner we obtain

$$f_{Y_{(n+1)}}(y) = \binom{2n+1}{n, 1, n} [F(y)]^n \cdot f(y) \cdot [1 - F(y)]^n. \quad \square$$

Remark 8.17. Here are two observations about n iid random variables Y_1, \dots, Y_n .

- (a) Assume that Y_{k_1}, \dots, Y_{k_n} is a permutation (ANY permutation!!) of Y_1, \dots, Y_n . Then the symmetry that results from iid implies that

$$P\{Y_1 < Y_2 < \dots < Y_n\} = P\{Y_{k_1} < Y_{k_2} < \dots < Y_{k_n}\}.$$

Since there are $n!$ permutations, each one of those probabilities equals $\frac{1}{n!}$.

- (b) Fix an arbitrary $k \in [1, n]_{\mathbb{Z}}$. Then

$$P\{Y = Y_{(1)}\} = P\{Y = Y_{(2)}\} = \dots = P\{Y = Y_{(n)}\}.$$

Since there are n such arrangements, each one of those probabilities equals $\frac{1}{n}$. \square

8.10 The Bivariate Normal Distribution (Optional)

Definition 8.21 (Bivariate normal distribution). ★

We say that two continuous random variables Y_1 and Y_2 have a **bivariate normal distribution**, or that they have a **joint normal distribution**, if their joint PDF is

$$(8.73) \quad f_{Y_1, Y_2}(y_1, y_2) = \frac{e^{-Q/2}}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}}, \quad -\infty < y_1 < \infty, \quad -\infty < y_2 < \infty,$$

$$\text{where } Q = \frac{1}{1-\rho^2} \left[\frac{(y_1 - \mu_1)^2}{\sigma_1^2} - 2\rho \frac{(y_1 - \mu_1)(y_2 - \mu_2)}{\sigma_1\sigma_2} + \frac{(y_2 - \mu_2)^2}{\sigma_2^2} \right].$$

We then also write $(Y_1, Y_2) \sim \mathcal{N}(\mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \rho)$. \square

Whereas we have marked this definition as optional, you should remember the following theorem.

Theorem 8.23.

If two random variables Y_1 and Y_2 are $\mathcal{N}(\mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \rho)$, then

(a) $Y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $Y_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$.

Thus, $E[Y_1] = \mu_1$, $\text{Var}[Y_1] = \sigma_1^2$, $E[Y_2] = \mu_2$, $\text{Var}[Y_2] = \sigma_2^2$.

(b) $\text{Cov}[Y_1, Y_2] = \sigma_1\sigma_2\rho$. Thus, ρ is the correlation coefficient of Y_1 and Y_2 .

PROOF (outline):

One proves (a) by showing that the marginal densities are

$$f_{Y_1}(y) = \frac{1}{\sigma_1\sqrt{2\pi}} e^{-(y-\mu_1)^2/(2\sigma_1^2)}, \quad f_{Y_2}(y) = \frac{1}{\sigma_2\sqrt{2\pi}} e^{-(y-\mu_2)^2/(2\sigma_2^2)}.$$

See (7.37) on p.133.

For the proof of (b), see Casella, Berger [3]. \blacksquare

Theorem 8.24.

If two jointly normal random variables Y_1 and Y_2 are uncorrelated, then they are independent.

PROOF: If $\rho = 0$, the joint PDF of Y_1 and Y_2 which was given in (8.73) is

$$f_{Y_1, Y_2}(y_1, y_2) = \frac{e^{-Q/2}}{2\pi\sigma_1\sigma_2},$$

where $Q = \frac{(y_1 - \mu_1)^2}{\sigma_1^2} - 0 + \frac{(y_2 - \mu_2)^2}{\sigma_2^2}$. Thus,

$$\begin{aligned} f_{Y_1, Y_2}(y_1, y_2) &= \frac{1}{(\sqrt{2\pi}\sigma_1)(\sqrt{2\pi}\sigma_2)} \exp \left\{ -\frac{(y_1 - \mu_1)^2}{2\sigma_1^2} - \frac{(y_2 - \mu_2)^2}{2\sigma_2^2} \right\} \\ &= \left(\frac{1}{(\sqrt{2\pi}\sigma_1)} \exp \left\{ -\frac{(y_1 - \mu_1)^2}{2\sigma_1^2} \right\} \right) \left(\frac{1}{(\sqrt{2\pi}\sigma_2)} \exp \left\{ -\frac{(y_2 - \mu_2)^2}{2\sigma_2^2} \right\} \right) \end{aligned}$$

It follows from Theorem 8.23(a) that $f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) f_{Y_2}(y_2)$. The independence of Y_1 and Y_2 follows from Theorem 8.4 on p.150. ■

Remark 8.18. The concept of joint normality can be extended from two random variables to an arbitrary number of random variables Y_1, \dots, Y_n . However, the definition of their joint PDF utilizes $n \times n$ matrices and their determinants. This requires some background in linear algebra and that is not a prerequisite for this course. □

9 Functions of Random Variables and their Distribution

This chapter essentially only contains enough material to serve as a reference and review “sheet”. You will not be able to properly understand the techniques noted here if you do not work through the many examples of the WMS text!

9.1 The Method of Distribution Functions

The Method of Distribution Functions is best explained by some examples.

Example 9.1. Find the CDF and PDF for $U := 2Y - 6$, where the density of the random variable Y is

$$(9.1) \quad f_Y(y) = \begin{cases} 8y, & \text{if } 0 \leq y \leq 1/2, \\ 0, & \text{else.} \end{cases}$$

Solution: Applying the distribution function method means the following:

- Find the CDF $F_U(u)$ of U □ Find the PDF $f_U(u)$ of U by differentiating $F_U(u)$
- Do this with help of the relation $U = 2Y - 6 \Leftrightarrow Y = \frac{U + 6}{2}$.

We obtain

$$F_U(u) = P\{U \leq u\} = P\{2Y - 6 \leq u\} = P\left\{Y \leq \frac{u + 6}{2}\right\} = F_Y\left(\frac{u + 6}{2}\right).$$

Note that

$$0 \leq y \leq \frac{1}{2} \Leftrightarrow 0 \leq \frac{u + 6}{2} \leq \frac{1}{2} \Leftrightarrow -6 \leq u \leq -5$$

Thus, $F_U(u) = 0$ for $u < -6$ and $F_U(u) = 1$ for $u > -5$.

For $-6 \leq u \leq -5$, i.e., $0 \leq y \leq \frac{1}{2}$, we must integrate:

$$P\left\{Y \leq \frac{u + 6}{2}\right\} = \int_0^{(u+6)/2} f_Y(y) dy = \int_0^{(u+6)/2} 8y dy = \frac{8}{2} \left(\frac{u + 6}{2}\right)^2.$$

We combine the cases $u < -6$; $-6 \leq u \leq -5$; $u > -5$ and obtain

$$F_U(u) = \begin{cases} 0, & \text{if } u < -6, \\ (u + 6)^2, & \text{if } -6 \leq u \leq -5, \\ 1, & \text{if } u > -5. \end{cases}$$

We differentiate this CDF and obtain the density function for U :

$$f_U(u) = \frac{dF_U(u)}{du} = \begin{cases} 2(u + 6), & \text{if } -6 \leq u \leq -5, \\ 0, & \text{else. } \square \end{cases}$$

Example 9.2 (WMS Ch.06.3, Example 6.3). The following is Example 6.3 of the WMS text. Its proof has been substantially rewritten.

Let (Y_1, Y_2) denote a random sample of size $n = 2$ from the uniform distribution on the interval $(0, 1)$. In other words, we assume that Y_1 and Y_2 are jointly continuous and have a joint PDF which is constant and not zero on the unit square.

The issue is to find the probability density function for $U := Y_1 + Y_2$.

Solution: It follows from the assumptions that Y_1 and Y_2 possess the same marginal PDF. The density function for each Y_i is

$$f(y) := f_{Y_1}(y) = f_{Y_2}(y) = \begin{cases} 1, & 0 \leq y \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

Since Y_1 and Y_2 are independent,

$$f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1)f_{Y_2}(y_2) = f(y_1)f(y_2) = \begin{cases} 1, & 0 \leq y_1 \leq 1, 0 \leq y_2 \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

Thus, $F_U(u) = P\{Y_1 + Y_2 \leq u\} = \iint_B f(y_1)f(y_2) dy_1 dy_2$, where, for a fixed u , the region of integration is

$$(A) \quad B := ([0, 1] \times [0, 1]) \cap \{(y_1, y_2) \in \mathbb{R}^2 : y_1 + y_2 \leq u\}.$$

We will separately treat the cases • $u \leq 0$ or $u \geq 2$ • $0 < u \leq 1$ • $1 < u < 2$.

Case 1: $u \leq 0$ or $u \geq 2$.

If $u \leq 0$, then $[0, 1] \times [0, 1]$ and $\{(y_1, y_2) \in \mathbb{R}^2 : y_1 + y_2 \leq u\}$ are disjoint. Thus, $B = \emptyset$ and $\iint_B \dots = 0$ and thus, $F_U(u) = 0$.

If $u \geq 2$, then $[0, 1] \times [0, 1] \subseteq \{(y_1, y_2) \in \mathbb{R}^2 : y_1 + y_2 \leq u\}$. Thus, $\iint_B \dots = \int_0^1 \int_0^1 \dots$ and thus, $F_U(u) = 1$.

Case 2: • $0 < u \leq 1$.

The graph of $y_1 + y_2 = u$ in the (y_1, y_2) plane is a straight line which intersects the vertical coordinate axis, $y_1 = 0$, at $y_2 = u$ and the horizontal coordinate axis, $y_2 = 0$, at $y_1 = u$. Thus, B is the triangle bounded by the coordinate axes and the line $y_1 + y_2 = u$. since it is half of a square with side length u , its area is $u^2/2$.

Of course, this also follows from the fact that $\iint_B \dots$ is achieved by first integrating, for $0 \leq y_1 \leq u$, over the vertical slice of B at y_1 and then integrating those integrals. Since the vertical slice of B at y_1 extends from $y_2 = 0$ to $y_1 + y_2 = u$, i.e., to $y_2 = u - y_1$

$$\begin{aligned} F_U(u) &= \iint_B 1 dy_1 dy_2 = \int_0^u \int_0^{u-y_1} 1 dy_2 dy_1 \\ &= \int_0^u (u - y_1) dy_1 = \left(uy_1 - \frac{y_1^2}{2} \right) \Big|_0^u = u^2 - \frac{u^2}{2} = \frac{u^2}{2}. \end{aligned}$$

Case 3: • $1 < u < 2$.

Let $\tilde{B} := ([0, 1] \times [0, 1]) \setminus \{(y_1, y_2) \in \mathbb{R}^2 : y_1 + y_2 \geq u\}$. Then

$$(B) \quad \tilde{B} = ([0, 1] \times [0, 1]) \cap \{(y_1, y_2) \in \mathbb{R}^2 : y_1 + y_2 \leq u\},$$

$$(C) \quad F_U(u) = 1 - P\{Y_1 + Y_2 \geq u\} = 1 - \iint_{\tilde{B}} 1 \, dy_1 \, dy_2$$

Now, the graph of $y_1 + y_2 = u$ in the (y_1, y_2) plane is a straight line which intersects the vertical line, $y_1 = 1$, at $y_2 = u - 1$ and the horizontal line, $y_2 = 0$, at $y_1 = u - 1$.

\tilde{B} is the right angle triangle bounded by the lines $y_1 = 1$, $y_2 = 1$ and $y_1 + y_2 = u$.

Its legs have length $1 - (u - 1) = 2 - u$. Thus, its area is half that of a square with side length $2 - u$. Thus, the area of \tilde{B} is $(2 - u)^2/2$. It follows from (C) that

$$F_U(u) = 1 - \text{area}(\tilde{B}) = 1 - \frac{4 - 4u + u^2}{2} = -1 + 2u - \frac{u^2}{2}.$$

This also could have been computed by iterated integration. In this case,

$$\begin{aligned} 1 - F_U(u) &= \iint_{\tilde{B}} 1 \, dy_1 \, dy_2 = \int_{u-1}^1 \int_{u-y_1}^1 1 \, dy_2 \, dy_1 \\ &= \int_{u-1}^1 (1 - u + y_1) \, dy_1 = \left((1 - u) + \frac{y_1^2}{2} \right) \Big|_{u-1}^1 \\ &= (1 - u)(2 - u) + \frac{1}{2} - \frac{(u-1)^2}{2} = 2 - 2u + \frac{u^2}{2}. \end{aligned}$$

We thus obtain, as before, $F_U(u) = 1 - (2 + 2u - u^2/2) = -1 + 2u - u^2/2$. \square

The problem of the next example is that of WMS Ch.6.4, Example 6.8. This instructor does not understand the reasoning given there and has provided a completely different proof. You find this example here rather than in the next section (section 9.2: The Method of Transformations in One Dimension), because it is solved with the techniques of this section.

Example 9.3. Let Y_1 and Y_2 be jointly continuous random variables with density function

$$f_{Y_1, Y_2}(y_1, y_2) = \begin{cases} e^{-(y_1 + y_2)}, & 0 \leq y_1, 0 \leq y_2, \\ 0, & \text{else.} \end{cases}$$

What are the CDF and PDF of $U := Y_1 + Y_2$?

Solution:

$$P\{U \leq u\} = P\{Y_1 + Y_2 \leq u\} = \iint_R e^{-y_1 - y_2} \, d\vec{y}$$

where $R =$ triangle with vertices $(0, u)$, $(0, 0)$, $(u, 0)$. Thus, for $u > 0$,

$$\begin{aligned} P\{U \leq u\} &= \int_0^u \left[\int_0^{u-y_1} e^{-y_1 - y_2} \, dy_2 \right] dy_1 = \int_0^u e^{-y_1} \left[-e^{-y_2} \Big|_0^{u-y_1} \right] dy_1 \\ &= \int_0^u e^{-y_1} [1 - e^{-(u-y_1)}] dy_1 = \int_0^u e^{-y_1} [1 - e^{y_1} e^{-u}] dy_1 \\ &= \int_0^u e^{-y_1} dy_1 - \int_0^u e^{-u} dy_1 = -e^{-y_1} \Big|_0^u - u e^{-u} \\ &= -(e^{-u} - 1) - u e^{-u} = 1 - (1 + u) e^{-u}. \end{aligned}$$

The derivative is (for $u > 0$)

$$\begin{aligned} f_U(u) &= \frac{d}{du}(1 - (1+u)e^{-u}) = -(1+u)'e^{-u} - (1+u)(e^{-u})' \\ &= -e^{-u} - (1+u)(-e^{-u}) = -e^{-u} + e^{-u} + ue^{-u} = ue^{-u}. \end{aligned}$$

Thus, the CDF is $F_U(u) = \begin{cases} 1 - (1+u)e^{-u}, & \text{if } u > 0, \\ 0, & \text{else} \end{cases}$

and the PDF is $f_U(u) = \begin{cases} ue^{-u}, & \text{if } u > 0, \\ 0, & \text{else.} \end{cases}$

The latter agrees with the WMS result. \square

Remark 9.1. In the following we use the arrow notation $\vec{y} = (y_1, \dots, y_n)$, $\vec{Y} = (Y_1, \dots, Y_n), \dots$

Summary of the Distribution Function Method

Goal: Find the PDF $f_U(u)$ for $U = g(\vec{Y})$, where $g : D \rightarrow \mathbb{R}$ has a domain $D \subseteq \mathbb{R}^n$ large enough to hold all arguments \vec{y} that are relevant for the problem.

- (1) Find the region $R = \{g \leq u\} = g^{-1}(] - \infty, u])$. (Thus, $R \subseteq \mathbb{R}^n$.)
- (2) Find the “boundary” $R^* = \{g = u\}$ of the region R .
- (3) Find the CDF $F_U(u) = P\{U \leq u\}$ by integrating $f(\vec{y})$ over the region R .
- (4) Find the the PDF $f_U(u) = \frac{dF_U(u)}{du}$ by differentiating $F_U(u)$.

Note for the above that, since g may not be invertible, g^{-1} denotes the preimage $g^{-1}(B) = \{\vec{y} : g(\vec{y}) \in B\}$, where $B \subseteq \mathbb{R}$. If, e.g., $B =] - \infty, u]$, then $R = g^{-1}(] - \infty, u])$, and (3) expresses

$$\begin{aligned} (9.2) \quad F_U(u) &= P\{U \leq u\} = P\{g(\vec{Y}) \leq u\} = P\{\omega : \vec{Y}(\omega) = \vec{y} \text{ such that } g(\vec{y}) \leq u\} \\ &= P\{Y \in R\} = \iint \cdots \int_R f_{\vec{Y}}(\vec{y}) d\vec{y}. \quad \square \end{aligned}$$

The next remark really should be considered another example for the distribution method. It has been marked as optional, so it will not be part of any exam or quiz. Nevertheless, you are strongly encouraged to work through its proof and increase your skills with respect to applying the distribution method.

Remark 9.2. ★ Let Y be a continuous random variable with PDF $f_Y(y)$ and let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a **symmetrical function** (also, **symmetric function**), i.e., $h(-y) = h(y)$ for all y . Also, assume that

- (1) $y \mapsto h(y)$ is differentiable (hence, continuous) everywhere.
- (2) $y \mapsto h(y)$ is injective for $y \geq 0$, i.e., $0 \leq y < y' \Rightarrow h(y) \neq h(y')$. (Thus, by symmetry, $h(y)$ also is injective for $y < 0$).

Continuous functions of a real variable are either strictly increasing or strictly decreasing on any subset of the domain where they are injective. (Draw a picture!) Thus, there are two possibilities.

- (1) h is strictly increasing on $[0, \infty[$ (and then, by symmetry, h is strictly decreasing on $[-\infty, 0]$). Also, h attains its global minimum at $y = 0$.
- (2) h is strictly decreasing on $[0, \infty[$ (and then, by symmetry, h is strictly increasing on $[-\infty, 0]$). Also, h attains its global maximum at $y = 0$.

In either case, there are no jumps for the continuous $h(\cdot)$. We will determine the CDF and PDF under the following assumptions: For any given $u \in \mathbb{R}$,

- (3) h is strictly increasing on $[0, \infty[$
- (4) $h(0) = 0$ and thus, $h(y) \geq 0$ for all y . Note that then $P\{U > 0\} = 1$ and $P\{U \leq 0\} = 0$.
- (6) Thus, if $u > 0$, then $U(\omega) \leq u \Leftrightarrow |Y(\omega)| \leq y = h^{-1}(u)$. Thus,

$$\begin{aligned} F_U(u) &= P\{U \leq u\} = P\{|Y| \leq h^{-1}(u)\} = P\{-h^{-1}(u) \leq Y \leq h^{-1}(u)\} \\ &= F_Y(h^{-1}(u)) - F_Y(-h^{-1}(u)) \quad \text{if } u > 0, \text{ i.e.,} \end{aligned}$$

$$(9.3) \quad F_U(u) = \begin{cases} 1, & \text{if } h(y) < u \text{ for all } y, \\ F_Y(h^{-1}(u)) - F_Y(-h^{-1}(u)), & \text{if there is } y = h^{-1}(u), \\ 0, & \text{if } u \leq 0. \end{cases}$$

We differentiate $\frac{d}{du}$ to obtain the density. We write $h^{-1}'(u) = \frac{dh^{-1}(u)}{du}$:

- $f_U(u) = h^{-1}'(u) f_Y(h^{-1}(u)) - (-1)h^{-1}'(u) f_Y(-h^{-1}(u))$

Thus,

$$(9.4) \quad f_U(u) = \begin{cases} h^{-1}'(u) [f_Y(h^{-1}(u)) + f_Y(-h^{-1}(u))] , & \text{if there is } y = h^{-1}(u), \\ 0, & \text{else. } \square \end{cases}$$

Example 9.4. As an example for that last remark, let us consider the function $h(y) = y^2$.⁵¹ h is strictly increasing on $[0, \infty[$ and its minimum is $h(0) = 0$. Thus, h satisfies the assumptions (3) and (4) of Remark 9.2. Since $\lim_{y \rightarrow \infty} y^2 = \infty$, the condition “if $h(y) < u$ for all y ” of (9.3) is never satisfied.

Further, the condition “if there is $y = h^{-1}(u)$ ” of (9.3) and (9.4) becomes “ $u \geq 0$ ”.

Thus, if $U = Y^2$, then $h^{-1}(u) = \sqrt{u}$ for $u \geq 0$ and $h^{-1}'(u) = 1/(2\sqrt{u})$. We obtain

$$f_U(u) = \begin{cases} \frac{1}{2\sqrt{u}} [f_Y(\sqrt{u}) + f_Y(-\sqrt{u})] , & \text{if } u > 0, \\ 0, & \text{else. } \square \end{cases}$$

Example 9.5. Assume that the random variable Y is $\mathcal{N}(0, 1)$, i.e., Y is standard normal. What is the distribution of $U := Y^2$?

For this example, let

$$(9.5) \quad \phi(y) := f_Y(y) = \frac{1}{\sqrt{2\pi}} e^{-y^2/2},$$

$$(9.6) \quad \Phi(y) := \int_{-\infty}^y \phi(t) dt.$$

In other words, ϕ is the PDF of Y and Φ is the CDF of Y .

Since $U \geq 0$, we have $f_U(u) = F_U(u) = 0$ for $u < 0$. Thus, we may assume that $u \geq 0$.

⁵¹That is WMS Example 6.4.

Then, $F_U(u) = P\{-\sqrt{u} \leq Y \leq \sqrt{u}\} = \Phi(\sqrt{u}) - \Phi(-\sqrt{u})$ and thus,

$$\begin{aligned} f_U(u) &= F'_U(u) = \frac{d}{du} [\Phi(\sqrt{u}) - \Phi(-\sqrt{u})] \\ &= \phi(\sqrt{u}) \frac{1}{2\sqrt{u}} + \phi(-\sqrt{u}) \frac{1}{2\sqrt{u}} = \phi(\sqrt{u}) \frac{1}{\sqrt{u}} = \frac{1}{\sqrt{2\pi}} e^{-(\sqrt{u})^2/2} \frac{1}{\sqrt{u}} \end{aligned}$$

Above, we used symmetry $\phi(-\sqrt{u}) = \phi(\sqrt{u})$ to obtain the equation before the last. Thus,

$$f_U(u) = \frac{1}{\sqrt{2\pi}} e^{-u/2} u^{-1/2} = \frac{u^{1/2-1}}{2^{1/2}\sqrt{\pi}} e^{-u/2}$$

One can show that $\Gamma(1/2) = \sqrt{\pi}$.⁵² We use that result without attempting to prove it and obtain, setting $\alpha := 1/2$ and $\beta := 2$,

$$f_U(u) = \frac{u^{1/2-1} e^{-u/2}}{2^{1/2}\Gamma(1/2)} = \frac{u^{\alpha-1} e^{-u/\beta}}{\beta^\alpha \Gamma(\alpha)}.$$

We finally remember that all this was done for $u \geq 0$ and that $f_U(u) = 0$ for $u < 0$.

$$f_U(u) = \begin{cases} \frac{u^{\alpha-1} e^{-u/\beta}}{\beta^\alpha \Gamma(\alpha)}, & \text{if } u \geq 0, \\ 0, & \text{else.} \end{cases}$$

It follows that the square of a $\mathcal{N}(0, 1)$ variable has a gamma(1/2, 2) distribution. Equivalently, it has a chi-square distribution with one degree of freedom. \square

Example 9.6. It is important that you recognize when there are significant shortcuts. It might be possible to obtain $F_U(u) = F_U(g^{-1}(y))$ without having to integrate the PDF. Here is an example.

Let the random variable Y be $\text{expon}(1)$. Find the CDF and PDF of $U := 2Y - 4$.

Solution:

(1) Here, $u = g(y) = 2y - 4$ has inverse $y = g^{-1}(u) = (u + 4)/2$.

(2) The CDF of Y is $F_Y(y) = \begin{cases} 1 - e^{-y}, & \text{if } y \geq 0, \\ 0, & \text{else.} \end{cases}$

(3) Thus, $F_U(u) = P\{U \leq u\} = P\{2Y - 4 \leq u\} = P\left\{Y \leq \frac{u+4}{2}\right\} = F_Y\left(\frac{u+4}{2}\right)$.

(4) From (2): $F_U(u) = \begin{cases} 1 - e^{-\frac{u+4}{2}}, & \text{if } \frac{u+4}{2} \geq 0, \\ 0, & \text{else.} \end{cases}$

(5) Thus, $F_U(u) = \begin{cases} 1 - e^{-\frac{u+4}{2}}, & \text{if } u \geq -4, \\ 0, & \text{else.} \end{cases}$

(6) We have obtained $F_U(u)$ without integrating a PDF.

(7) The density is $f_U(u) = F'_U(u) = \begin{cases} \frac{1}{2} e^{-\frac{u+4}{2}}, & \text{if } u \geq -4, \\ 0, & \text{else.} \end{cases}$ \square

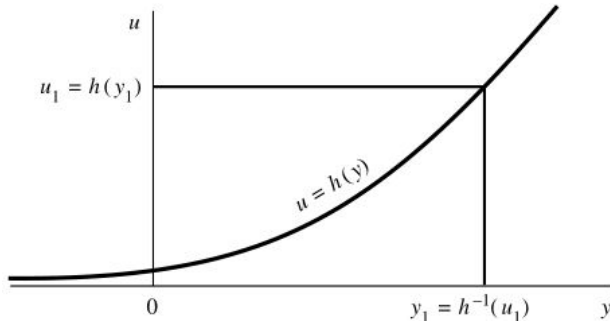
⁵²See, e.g., https://en.wikipedia.org/wiki/Gamma_function or Shilov, G. [9].

9.2 The Method of Transformations in One Dimension

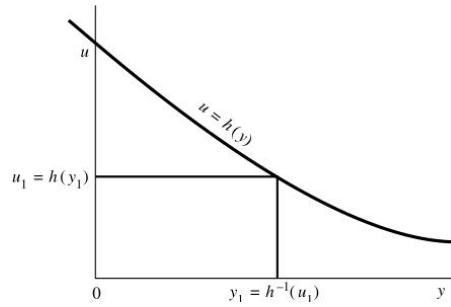
Introduction 9.1. We already encountered the method of transformations in Remark 9.2 on p.193. There we computed the CDF and PDF of the random variable $U = h(Y)$ for a continuous random variable Y and a symmetric and differentiable function $h(y)$ which was injective on the interval $B_1 = [0, \infty[$. (By symmetry, h also had those characteristics on $B_2 =]-\infty, 0[$.)

At the heart of the calculations was the fact that injectivity allowed us to compute, for a given u , a unique $y = h^{-1}(u)$ such that $h(y) = u$.

Since differentiable functions are continuous, injectivity on an interval B implies that h is either strictly increasing or strictly decreasing on B . See figures 9.1 and 9.2 below.



9.1 (Figure). **Strictly increasing function.**
Source: WMS Ch.6.4



9.2 (Figure). **Strictly decreasing function.**
Source: WMS Ch.6.4

Those figures illustrate the following.

(1) If h is strictly increasing, then $h(y) \leq u_1 \Leftrightarrow y \leq h^{-1}(u_1)$. Thus,

$$(9.7) \quad \begin{aligned} P\{U \leq u\} &= P\{h(Y) \leq u\} = P\{h^{-1}[h(Y)] \leq h^{-1}(u)\} = P\{Y \leq h^{-1}(u)\}, \\ \text{i.e.,} \quad F_U(u) &= F_Y(h^{-1}(u)). \end{aligned}$$

We differentiate with respect to u and write $h^{-1'}(u)$ for $\frac{dh^{-1}(u)}{du}$. Then

$$f_U(u) = \frac{dF_U(u)}{du} = \frac{dF_Y(h^{-1}(u))}{du} = f_Y(h^{-1}(u)) \cdot h^{-1'}(u).$$

Since h is strictly increasing, $h^{-1'}(u) > 0$. Thus, $h^{-1'}(u) = |h^{-1'}(u)|$. Thus,

$$(9.8) \quad f_U(u) = f_Y(h^{-1}(u)) \cdot |h^{-1'}(u)|.$$

(2) If h is strictly decreasing, then $h(y) \leq u_1 \Leftrightarrow y \geq h^{-1}(u_1)$. Thus,

$$(9.9) \quad \begin{aligned} P\{U \leq u\} &= P\{h(Y) \leq u\} = P\{Y \geq h^{-1}(u)\} = 1 - P\{Y \leq h^{-1}(u)\}, \\ \text{i.e.,} \quad F_U(u) &= 1 - F_Y(h^{-1}(u)). \end{aligned}$$

We differentiate with respect to u . Then

$$f_U(u) = -\frac{dF_Y(h^{-1}(u))}{du} = f_Y(h^{-1}(u)) \cdot (-h^{-1'}(u)).$$

Since h is strictly decreasing, $h^{-1}'(u) < 0$. Thus, $-h^{-1}'(u) = |h^{-1}'(u)|$. Thus,

$$(9.10) \quad f_U(u) = f_Y(h^{-1}(u)) \cdot |h^{-1}'(u)|.$$

- (3) We compare (9.8) and (9.10) and see that they are equal. Thus, as long as h is either strictly increasing everywhere or strictly decreasing everywhere, (i.e., as long as f is invertible everywhere,)

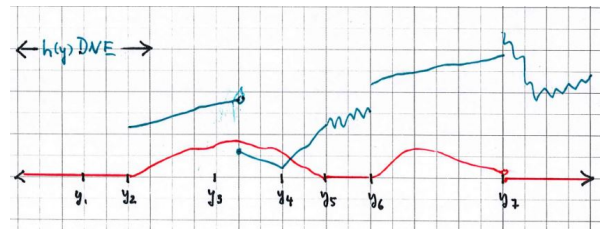
$$(9.11) \quad f_U(u) = f_Y(h^{-1}(u)) \cdot |h^{-1}'(u)| = f_Y(h^{-1}(u)) \cdot \left| \frac{d[h^{-1}(u)]}{du} \right|.$$

Since $\int_a^b f_Y(t) dt = \int_{[a,b] \cap \{\tilde{y}: f(\tilde{y}) \neq 0\}} f_Y(t) dt$ for any interval $[a, b]$, we only need to worry about the behavior of h for arguments belonging to

$$\text{suppt}(f_Y) := \{\tilde{y} : f(\tilde{y}) \neq 0\}.$$

It is customary to call $\text{suppt}(f_Y)$ the **support** of the density $f_Y(y)$.⁵³

- $\text{suppt}(f_Y) =]y_2, y_5[\cup]y_6, y_7[$. It does not matter what $h(y)$ does outside $\text{suppt}(f_Y)$.
- h must be injective on the support of f_Y .
- h changes direction at y_3 and y_4 , so the pieces $]y_2, y_3[$, $]y_3, y_4[$, $]y_4, y_5[$, must be treated separately. \square



The following theorem summarizes the observations of those introductory results:

Theorem 9.1.

Given are a continuous random variable Y with density $f_Y(y)$ and a differentiable function $h(y)$ which is either strictly increasing or strictly decreasing for all $y \in \text{suppt}(f_Y)$, i.e., for all y that satisfy $f_Y(y) > 0$. Then the PDF of $U := h(Y)$ is

$$(9.12) \quad f_U(u) = f_Y(h^{-1}(u)) \cdot |h^{-1}'(u)| = f_Y(h^{-1}(u)) \cdot \left| \frac{d[h^{-1}(u)]}{du} \right|.$$

PROOF: See the introduction 9.1. \blacksquare

Example 9.7 (Increasing function). Given is a random variable Y with the following PDF:

$$f_Y(y) = \begin{cases} 2y, & \text{if } 0 \leq y \leq 1, \\ 0, & \text{else.} \end{cases}$$

⁵³ ★ In general, one defines the support $\text{suppt}(h) := \{\tilde{x} : h(\tilde{x}) \neq 0\}$ for any real valued function $x \mapsto h(x)$.

Let $U := 4Y - 3$. Find the PDF for U by means of the transformation method.

Solution: We apply the transformation method with the strictly increasing function $u = h(y) = 4y - 3$. Then the inverse of h is $y = h^{-1}(u) = (u + 3)/4$, for all $u \in \mathbb{R}$.

- (1) We apply the transformation method with $u = h(y) = 4y - 3$ (strictly increasing).
- (2) Then the inverse of h is $y = h^{-1}(u) = (u + 3)/4$, for all $u \in \mathbb{R}$.
- (3) Further, $h^{-1}'(u) = 1/4$. Since $0 \leq (u + 3)/4 \leq 1 \Leftrightarrow -3 \leq u \leq 1$,

$$f_U(u) = \begin{cases} \frac{2(u+3)}{4} \cdot \frac{1}{4}, & \text{if } -3 \leq u \leq 1, \\ 0, & \text{else.} \end{cases} = \begin{cases} \frac{u+3}{8}, & \text{if } -3 \leq u \leq 1, \\ 0, & \text{else.} \quad \square \end{cases}$$

Example 9.8 (Decreasing function). Given is a random variable Y with the same PDF as in Example 9.7:

$$f_Y(y) = \begin{cases} 2y, & \text{if } 0 \leq y \leq 1, \\ 0, & \text{else.} \end{cases}$$

Let $U := -3Y + 2$. Find the PDF for U by means of the transformation method.

Solution: We apply the transformation method with the strictly decreasing function $u = h(y) = 2 - 3y$. Then the inverse of h is $y = h^{-1}(u) = (2 - u)/3$, for all $u \in \mathbb{R}$.

- (1) We apply the transformation method with $u = h(y) = 2 - 3y$ (strictly decreasing).
- (2) Then the inverse of h is $y = h^{-1}(u) = (2 - u)/3$, for all $u \in \mathbb{R}$.
- (3) Further, $h^{-1}'(u) = -1/3$. Since $0 \leq (2 - u)/3 \leq 1 \Leftrightarrow 0 \geq (u - 2) \geq -3 \Leftrightarrow -1 \leq u \leq 2$,

$$f_U(u) = \begin{cases} \frac{2(2-u)}{3} \cdot \left| \frac{-1}{3} \right|, & \text{if } -1 \leq u \leq 2, \\ 0, & \text{else.} \end{cases} = \begin{cases} \frac{4-2u}{9}, & \text{if } -1 \leq u \leq 2, \\ 0, & \text{else.} \quad \square \end{cases}$$

Example 9.9 (Distribution function method with two variables). Given are two jointly continuous random variables with uniform distribution on the triangle

$$B := \{(y_1, y_2) : 0 < y_2 < 1 - y_1 < 1\}.$$

Find the CDF of $U = Y_1 + Y_2$.

- (1) The joint PDF of (Y_1, Y_2) is $f_{Y_1, Y_2}(y_1, y_2) = \begin{cases} 2, & \text{if } 0 < y_2 < 1 - y_1 < 1, \\ 0, & \text{else.} \end{cases}$
- (2) $F_U(u) = P\{U \leq u\} = P\{Y_1 + Y_2 \leq u\} = \iint_{B \cap C} 2 \, d\vec{y}$, where $C = \{(y_1, y_2) : y_1 + y_2 \leq u\}$.
- (3) $(y_1, y_2) \in B \Rightarrow 0 < 1 - y_1 < 1 \Rightarrow 0 > y_1 - 1 > -1 \Rightarrow 0 < y_2 < 1$.
 $0 < y_2 < 1$ is obvious. Thus, $u \leq 0 \Rightarrow P\{U \leq u\} = 0$.
- (4) B is the triangle with vertices $(0, 0)$, $(0, 1)$ and $(1, 0)$.
For $u > 0$, C is the triangle with vertices $(0, 0)$, $(0, u)$ and $(u, 0)$

- (5) Thus, $0 < u < 1 \Rightarrow B \cap C = C \Rightarrow \iint_{B \cap C} 2 d\vec{y} = 2 \iint_C d\vec{y}$
- (6) Thus, from (5) & (2), $0 < u < 1 \Rightarrow B \cap C = C \Rightarrow F_U(u) = 2 \iint_C d\vec{y}$.
 $\iint_C \cdots d\vec{y}$ is done by integrating, for each fixed $0 < y_1 < u$, over that part of the vertical line $\{y_2 : y_2 = y_1\}$ that is within C . That is the segment $0 < y_2 < u - y_1$.
- (7) Thus, $0 < u < 1 \Rightarrow F_U(u) = 2 \int_0^u \int_0^{u-y_1} dy_2 dy_1$
 $= 2 \int_0^u (u - y_1 - 0) dy_1 = 2u^2 - 2 \frac{y_1^2}{2} \Big|_0^u = u^2$.
- (8) From (4), $u \geq 1 \Rightarrow B \cap C = B = \text{suppt}(f_U) \Rightarrow F_U(u) = 1$.
- (9) Thus, from (3) & (7) & (8), $F_U(u) = \begin{cases} 0, & \text{if } u \leq 0, \\ u^2, & \text{if } 0 < u < 1, \\ 1, & \text{if } u \geq 1. \end{cases}$
- Differentiation yields $f_U(u) = \begin{cases} 2u, & \text{if } 0 < u < 1, \\ 0, & \text{if } u \leq 0 \text{ or } u \geq 1. \end{cases} \square$

Remark 9.3. In the following we use the arrow notation $\vec{y} = (y_1, \dots, y_n)$, $\vec{Y} = (Y_1, \dots, Y_n)$, ...

Summary of the Transformation Method

Goal: Find the PDF $f_U(u)$ for $U = h(Y)$, where

- $h : R \rightarrow \mathbb{R}$ has a domain $R \subseteq \mathbb{R}$ large enough to hold all arguments y that are relevant for the problem. That requires that R contains the support of the PDF f_Y (the set where f_Y is not zero).
- h is invertible on R . In other words, h is injective on R : If $y \in R$ and $u = h(y)$, then there is no $\tilde{y} \in R$ such that $\tilde{y} \neq y$ and $h(\tilde{y}) = u$.
- Thus h has an inverse $u \mapsto h^{-1}(u)$ which maps any u that is a function value $u = h(y)$ back to y . Do not confuse this genuine inverse function of $h(\cdot)$ with the preimage function $B \mapsto h^{-1}(B) = \{y \in Y : h(y) \in B\}$! That one maps **sets to sets**!
- We require that h is either strictly increasing or strictly decreasing for those $y \in R$ where $f_Y(y) > 0$. This assumption guarantees that h is injective and its inverse $u \mapsto h^{-1}(u)$ exists on the support of f_Y .

To find the PDF $f_U(u)$ for $U = h(Y)$, proceed as follows:

- (1) Find the inverse function, $y = h^{-1}(u)$, for those u that correspond to y with $f_Y(y) \neq 0$.
- (2) Find the derivative $\frac{dh^{-1}}{du} = \frac{dh^{-1}(u)}{du} = h^{-1}'(u)$.
- (3) Finally, compute $f_U(u)$ as follows: $f_U(u) = f_Y(h^{-1}(u)) \left| \frac{dh^{-1}(u)}{du} \right|$. \square

Remark 9.4. The transformation method still works if h is not either strictly increasing or decreasing on $\text{suppt}(g)$, as long as h is injective and R can be subdivided by intervals on which h is either strictly increasing or strictly decreasing.

As an example, consider $u := h(y) := \begin{cases} y, & \text{if } y \leq 0, \\ e^{-y}, & \text{if } y > 0. \end{cases}$

- On $]-\infty, 0]$, h is strictly increasing with inverse $y = h^{-1}(u) = u$. This inverse has derivative $h^{-1}'(u) = 1 > 0$.
- On $]0, \infty[$, h is strictly decreasing with inverse $y = h^{-1}(u) = -\ln(u)$. This inverse has derivative $h^{-1}'(u) = -1/u < 0$.
- Obviously if $y \leq 0$, then $y \leq 0 \Leftrightarrow u \leq 0$. Moreover, $y > 0 \Leftrightarrow 0 < u = e^{-y} < 1$.
- Thus, $f_U(u) = \begin{cases} f_Y(u) \cdot |1| = f_Y(u), & \text{if } u \leq 0, \\ f_Y(e^{-u}) \cdot |-1/u| = \frac{f_Y(-\ln(u))}{u}, & \text{if } 0 < u < 1, \\ 0, & \text{else. } \square \end{cases}$

9.3 The Method of Transformations in Multiple Dimension

Introduction 9.2. In Chapter 9.2 (The Method of Transformations in Multiple Dimension), we looked for ways to compute the density $f_U(u)$ of the transform $U = h(Y)$ of a continuous random variable Y by means of a function h which maps real numbers y to real numbers $u = h(y)$. Theorem 9.1 on p.197 provided us with an explicit formula for the PDF $f_U(u)$ of the transformed random variable $U = h(Y)$:

$$(9.13) \quad f_U(u) = f_Y(h^{-1}(u)) \cdot |h^{-1}'(u)| = f_Y(h^{-1}(u)) \cdot \left| \frac{d[h^{-1}(u)]}{du} \right|.$$

- (1) Since $|h^{-1}'(u)|$ appears in that formula, $h^{-1}(u)$ must exist and be differentiable.
- (2) That in turn requires that h is differentiable, in particular continuous.
- (3) Moreover, neither $h'(y)$ nor $h^{-1}'(u)$ can be zero, since $h'(y) \cdot h^{-1}'(u) = 1$.

Existence of $h^{-1}(u)$ requires h to be injective on the support of the PDF f_Y :

- (4) If u_0 is the function value $u_0 = h(y)$ of some argument y that satisfies $f_Y(y) > 0$,
 - then there is no other argument \tilde{y} that also satisfies $u_0 = h(\tilde{y})$ and $f_Y(\tilde{y}) > 0$.

Since h is continuous, (4) is satisfied if h is either strictly increasing or strictly decreasing for all y in the support of h , so we replaced (4) with that simpler assumption.

We now look for an n -dimensional analogue. If you have attended a linear algebra course, you are knowledgeable about $n \times n$ matrices and their determinants. If your background about those subjects is limited to a course in multivariable calculus, then assume that $n = 2$ or $n = 3$. We study

- random vectors $\vec{Y} = (Y_1, \dots, Y_n)$, where each coordinate Y_j is a random variable.
- functions $\vec{u} = \vec{h}(\vec{y})$ that map n -dimensional arguments \vec{y} to n -dimensional function values \vec{y} , have continuous partial derivatives $\frac{\partial h_i}{\partial y_j}$ for $i, j \in [1, n]_{\mathbb{Z}}$ and that satisfy a multidimensional analogue of (4):
- (5) If the vector \vec{u}_0 is a function value $\vec{u}_0 = \vec{h}(\vec{y})$ of some argument \vec{y} that satisfies $f_{\vec{Y}}(\vec{y}) > 0$, (here, $f_{\vec{Y}}(\vec{y})$) is the PDF of the jointly continuous random variables Y_1, \dots, Y_n ,
 - then there is no other argument \vec{y} that also satisfies $\vec{u}_0 = \vec{h}(\vec{y})$ and $f_{\vec{Y}}(\vec{y}) > 0$.

These two conditions guarantee the invertibility of the function $\vec{y} \mapsto \vec{u} = \vec{h}(\vec{y})$: This inverse function $\vec{h}^{-1}(\cdot)$ is defined by the relation

$$\vec{u} = \vec{h}(\vec{y}) \Leftrightarrow \vec{y} = \vec{h}^{-1}(\vec{u}).$$

Since the function values $\vec{y} = \vec{h}^{-1}(\vec{u})$ belong to \mathbb{R}^n , $\vec{h}^{-1}(\cdot)$ consists of n coordinate functions $h_1^{-1}(\cdot), h_2^{-1}(\cdot), \dots, h_n^{-1}(\cdot)$. They are defined by the equations

$$(9.14) \quad h_1^{-1}(\vec{u}) = y_1, \quad h_2^{-1}(\vec{u}) = y_2, \quad \dots, \quad h_n^{-1}(\vec{u}) = y_n.$$

In the onedimensional case, the existence of continuous $\frac{dh}{du}$ which satisfies $\left| \frac{dh}{du} \right| \neq 0$ implies that of a continuous and non-zero derivative $\frac{dh^{-1}}{dy}$. Further,

$$(9.15) \quad \frac{dh^{-1}}{dy} = 1 / \frac{dh}{du}.$$

In the n -dimensional case, we must replace the condition $\left| \frac{dh}{du} \right| \neq 0$ with the condition

$$(5) \quad J^{-1} := \det \begin{bmatrix} \frac{\partial h_1}{\partial y_1} & \frac{\partial h_1}{\partial y_2} & \dots & \frac{\partial h_1}{\partial y_n} \\ \frac{\partial h_2}{\partial y_1} & \frac{\partial h_2}{\partial y_2} & \dots & \frac{\partial h_2}{\partial y_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial h_n}{\partial y_1} & \frac{\partial h_n}{\partial y_2} & \dots & \frac{\partial h_n}{\partial y_n} \end{bmatrix} \neq 0.$$

The choice of the symbol J^{-1} for this determinant will become clear in a moment. The assumptions (5) and (6) are sufficient for the existence of all partial derivatives $\frac{\partial h_i^{-1}}{\partial u_j}$ and their continuity. They form an $n \times n$ matrix and one can show that its determinant, which we denote by J , also does not vanish. In other words,

$$(9.16) \quad J = \det \begin{bmatrix} \frac{\partial h_1^{-1}}{\partial u_1} & \frac{\partial h_1^{-1}}{\partial u_2} & \dots & \frac{\partial h_1^{-1}}{\partial u_n} \\ \frac{\partial h_2^{-1}}{\partial u_1} & \frac{\partial h_2^{-1}}{\partial u_2} & \dots & \frac{\partial h_2^{-1}}{\partial u_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial h_n^{-1}}{\partial u_1} & \frac{\partial h_n^{-1}}{\partial u_2} & \dots & \frac{\partial h_n^{-1}}{\partial u_n} \end{bmatrix} \neq 0.$$

Moreover, the determinants J^{-1} and J satisfy the analogue of (9.15):

$$(9.17) \quad J^{-1} = \frac{1}{J}. \quad \square$$

Before we examine how this material about the matrices of the partial derivatives and their determinants can be used to compute the joint PDF of the random vector $\vec{U}(\omega) = \vec{h}(\vec{Y}(\omega))$ and before state our findings as a formal theorem, we illustrate the above with the following example.

Example 9.10 (The joint PDF of two independent, exponential random variables – Part 1). In this twodimensional example, the function $\vec{h} = (h_1, h_2)$ is defined as follows:

$$(9.18) \quad u_1 := h_1(y_1, y_2) := 2y_1 + y_2,$$

$$(9.19) \quad u_2 := h_2(y_1, y_2) := y_1 - 2y_2.$$

(1) We show that this function can be inverted by solving these equations for $\vec{y} = (y_1, y_2)$.

- $u_1 - 2u_2 \stackrel{(9.18)}{=} y_2 + 4y_2 = 5y_2 \Rightarrow y_2 = u_1/5 - 2u_2/5.$
- Thus, $y_1 \stackrel{(9.19)}{=} u_2 + 2y_2 = u_2 + (1/5)[2u_1 - 4u_2] = (2u_1)/5 + u_2/5.$

We have found the inverse function $\vec{h}^{-1} = (h_1^{-1}, h_2^{-1})$ to be

$$(9.20) \quad h_1^{-1}(u_1, u_2) = y_1 = \frac{1}{5}(2u_1 + u_2),$$

$$(9.21) \quad h_2^{-1}(u_1, u_2) = y_2 = \frac{1}{5}(u_1 - 2u_2).$$

We will continue in Example 9.11 on p.204. \square

In the introduction, we informally discussed the following result from multivariable calculus which we are rephrasing here in the language of joint PDFs of continuous random variables and which is at the heart of this section. It is so lengthy that we spread it over several boxes. As mentioned before, assume that $n \leq 3$ if you do not have sufficient knowledge of linear algebra.

Theorem 9.2.

- Let $\vec{Y} = (Y_1, \dots, Y_n)$ be a vector of random variables with joint PDF $f_{\vec{Y}}(\vec{y})$ and let R be a “nice” subset of \mathbb{R}^n which is so big that it hosts all outcomes $\vec{Y}(\omega)$ of \vec{Y} .
- Let the function $\vec{h} : R \rightarrow \mathbb{R}^n$; $\vec{y} \mapsto \vec{u} = \vec{h}(\vec{y})$ satisfy the following.
- \square \vec{h} has continuous partial derivatives $\frac{\partial h_i}{\partial y_j}$ for all $1 \leq i, j \leq n$.
- \square If the vector \vec{u} is a function value $\vec{u} = \vec{h}(\vec{y})$ of some argument \vec{y} that satisfies $f_{\vec{Y}}(\vec{y}) > 0$, then there is no other argument \vec{y} that satisfies all those conditions.

Then \vec{h} has an inverse $\vec{h}^{-1} = (h_1^{-1}, h_2^{-1}, \dots, h_n^{-1})$ which is defined by the relation

$$\vec{u} = \vec{h}(\vec{y}) \Leftrightarrow \vec{y} = \vec{h}^{-1}(\vec{u}).$$

We can write this for the coordinate functions $h_i(\cdot)$ and $h_j^{-1}(\cdot)$ as follows:

$$(9.22) \quad u_1 = h_1(\vec{y}), \dots, u_n = h_n(\vec{y}) \quad \text{and} \quad y_1 = h_1^{-1}(u), \dots, y_n = h_n^{-1}(u).$$

Also, all partial derivatives $\frac{\partial h_i^{-1}}{\partial u_j}$ exist and are continuous for $1 \leq i, j \leq n$.

$$(9.23) \quad \text{Let } \frac{d\vec{h}}{d\vec{y}} := \begin{bmatrix} \frac{\partial h_1}{\partial y_1} & \frac{\partial h_1}{\partial y_2} & \dots & \frac{\partial h_1}{\partial y_n} \\ \frac{\partial h_2}{\partial y_1} & \frac{\partial h_2}{\partial y_2} & \dots & \frac{\partial h_2}{\partial y_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial h_n}{\partial y_1} & \frac{\partial h_n}{\partial y_2} & \dots & \frac{\partial h_n}{\partial y_n} \end{bmatrix}, \quad \frac{dh^{-1}}{d\vec{u}} := \begin{bmatrix} \frac{\partial h_1^{-1}}{\partial u_1} & \frac{\partial h_1^{-1}}{\partial u_2} & \dots & \frac{\partial h_1^{-1}}{\partial u_n} \\ \frac{\partial h_2^{-1}}{\partial u_1} & \frac{\partial h_2^{-1}}{\partial u_2} & \dots & \frac{\partial h_2^{-1}}{\partial u_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial h_n^{-1}}{\partial u_1} & \frac{\partial h_n^{-1}}{\partial u_2} & \dots & \frac{\partial h_n^{-1}}{\partial u_n} \end{bmatrix}.$$

$$(9.24) \quad \text{Let } J^{-1} := J^{-1}(\vec{y}) := \det\left(\frac{d\vec{h}}{d\vec{y}}\right), \quad J := J(\vec{u}) := \det\left(\frac{dh^{-1}}{d\vec{u}}\right).$$

- We add another assumption: $J^{-1}(\vec{y}) \neq 0$ for all y that satisfy $f_{\vec{y}}(\vec{y}) > 0$.

$$(9.25) \quad \text{Then } J(h(\vec{y})) \neq 0 \quad \text{and} \quad J(h(\vec{y})) = 1/J^{-1}(\vec{y}).$$

Further, the density of the transform $\vec{U} = h(\vec{Y})$ is computed as

$$(9.26) \quad f_{\vec{U}}(\vec{u}) = f_{\vec{Y}}(h^{-1}(\vec{u})) \cdot |J(\vec{u})|.$$

PROOF: Beyond the scope of this course. It needs knowledge not only of linear algebra, but also of the so called implicit function theorem. ■

Before we give some examples to illustrate this theorem, we make a remark about some of the notation introduced there and then give a name to the determinant J^{-1} of the matrix $\frac{d\vec{h}}{d\vec{y}}$ of the partial derivatives of h .

Remark 9.5. In the onedimensional case ($n = 1$), the situation is as follows.

- \mathbb{R}^n is the set \mathbb{R} of real numbers, • $\vec{u} = \vec{h}(\vec{y})$ becomes $u = h(y)$ for real numbers y and u ,
- the 1×1 “matrix” of “partial” derivatives is $h'(y) = \frac{dh}{dy}$.

Considering that last point, it seems natural to write $\frac{d\vec{h}}{d\vec{y}}$ for the $n \times n$ matrix of partial derivatives $\frac{\partial h_i}{\partial y_j}$ and this author chose to do so. However, you will find either different notation⁵⁴ or, like in the WMS text, no dedicated symbols at all. That works well enough with 2×2 matrices. □

Definition 9.1 (Jacobian and Jacobian matrix).

⁵⁴For example, Williamson, Richard E. and Trotter, Hale [12] uses the notation $\vec{h}'(\vec{y})$, the multidimensional analogue of $h'(y)$.

The matrix $\frac{d\vec{h}}{d\vec{y}}$ of the partial derivatives of the function $\vec{y} \mapsto \vec{h}(\vec{y})$ is called the **Jacobian matrix** of $\vec{h}(\cdot)$. We refer to its determinant, $J^{-1}(\vec{y}) = \det\left(\frac{d\vec{h}}{d\vec{y}}\right)$, as the **Jacobian**, sometimes also the **Jacobian determinant**, of $\vec{h}(\cdot)$. \square

Notation 9.1 (Jacobian).

- Stewart writes $\frac{\partial(u_1, \dots, u_n)}{\partial(y_1, \dots, y_n)} := \det\left(\frac{d\vec{h}^{-1}}{d\vec{u}}\right)$ and $\frac{\partial(y_1, \dots, y_n)}{\partial(u_1, \dots, u_n)} := \det\left(\frac{d\vec{h}^{-1}}{d\vec{u}}\right)$
 - Thus, the expression $J = J(\vec{u}) = \det\left(\frac{d\vec{h}^{-1}}{d\vec{u}}\right)$, which appears in
 (9.26) $f_{\vec{Y}}(\vec{u}) = f_{\vec{Y}}(h^{-1}(\vec{u})) \cdot |J(\vec{u})|$, is the Jacobian of $h^{-1}(\vec{u})$ and not of $h(\vec{y})$.
 - This author follows the great majority of books on multivariable calculus in defining the the Jacobian as the determinant of $\frac{d\vec{h}}{d\vec{y}}$.
 - Be aware that WMS chooses instead to call $J = \det\frac{d\vec{h}^{-1}}{d\vec{u}}$ the Jacobian.
 - The reason seems to be that most books on probability and statistics agree on using the letter J for $\det\frac{d\vec{h}^{-1}}{d\vec{u}}$ (without giving a name to that determinant) and WMS does not want to use the somewhat lengthy “the reciprocal of the Jacobian” in its frequent references to J .
- \square

Example 9.11 (The joint PDF of two independent, exponential random variables – Part 2). In Example 9.10 on p.202, we defined $\vec{u} = \vec{h}(\vec{y})$ as follows:

$$u_1 = h_1(y_1, y_2) = 2y_1 + y_2, \quad u_2 = h_2(y_1, y_2) = y_1 - 2y_2.$$

We computed its inverse $\vec{u} = \vec{h}^{-1}(\vec{u}) =$ and obtained

$$y_1 = h_1^{-1}(u_1, u_2) = \frac{1}{5}(2u_1 + u_2), \quad y_2 = h_2^{-1}(u_1, u_2) = \frac{1}{5}(u_1 - 2u_2).$$

Observe that both \vec{h} and \vec{h}^{-1} are defined for all points in \mathbb{R}^2 .

The partial derivatives of \vec{h} are

$$\frac{\partial h_1}{\partial y_1} = 2, \quad \frac{\partial h_1}{\partial y_2} = 1, \quad \frac{\partial h_2}{\partial y_1} = 1, \quad \frac{\partial h_2}{\partial y_2} = -2.$$

Those of \vec{h}^{-1} are

$$\frac{\partial h_1^{-1}}{\partial u_1} = \frac{2}{5}, \quad \frac{\partial h_1^{-1}}{\partial u_2} = \frac{1}{5}, \quad \frac{\partial h_2^{-1}}{\partial u_1} = \frac{1}{5}, \quad \frac{\partial h_2^{-1}}{\partial u_2} = \frac{-2}{5}.$$

Further,

$$\frac{d\vec{h}}{d\vec{y}} = \begin{bmatrix} 2 & 1 \\ 1 & -2 \end{bmatrix}, \quad \frac{d\vec{h}^{-1}}{d\vec{u}} = \begin{bmatrix} \frac{2}{5} & \frac{1}{5} \\ \frac{1}{5} & -\frac{2}{5} \end{bmatrix},$$

Since the determinant of a 2×2 matrix $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$, is $ad - bc$, we obtain

$$J^{-1} = (2)(-2) - (1)(1) = -5, \quad J = \begin{pmatrix} 2 \\ \frac{1}{5} \end{pmatrix} \begin{pmatrix} -2 \\ \frac{1}{5} \end{pmatrix} - \begin{pmatrix} 1 \\ \frac{1}{5} \end{pmatrix} \begin{pmatrix} 1 \\ \frac{1}{5} \end{pmatrix} = \frac{-4 - 1}{25} = \frac{-1}{5},$$

Observe that $J = \frac{1}{J^{-1}}$, validates what was stated in (9.25) on p.203.

We will continue in Example 9.12. \square

Example 9.12 (The joint PDF of two independent, exponential random variables – Part 3). In Example 9.10 on p.202, we defined $\vec{u} = \vec{h}(\vec{y})$ as follows:

$$(9.27) \quad u_1 = h_1(y_1, y_2) = 2y_1 + y_2, \quad u_2 = h_2(y_1, y_2) = y_1 - 2y_2.$$

In its continuation, Example 9.11 above, we obtained $J = \text{const} = \frac{-1}{5}$ for the reciprocal of the Jacobian of \vec{h} .

We are ready to specify the random variables that we wish to transform by means of $\vec{h}(\cdot)$.

- Assume that Y_1 and Y_2 are independent expon(2) random variables.
- Let $U_1 := h_1(\vec{Y}) = 2Y_1 + Y_2$, $U_2 := h_2(\vec{Y}) = Y_1 - 2Y_2$.
- Apply Theorem 9.2 on p.202 to compute the joint density $f_{\vec{U}}(u_1, u_2)$ of $\vec{U} = \vec{h}(\vec{Y})$.

Solution:

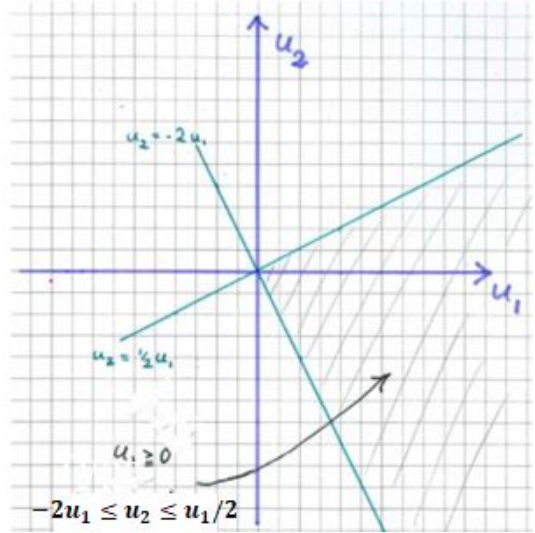
$$(a) \quad f_{\vec{Y}}(\vec{y}) = f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) \cdot f_{Y_2}(y_2) = \begin{cases} \frac{1}{4} e^{-(y_1+y_2)/2}, & \text{if } y_1, y_2 > 0, \\ 0, & \text{else.} \end{cases}$$

$$(b) \quad \text{We recall that } y_1 = \frac{1}{5}(2u_1 + u_2) \text{ and } y_2 = \frac{1}{5}(u_1 - 2u_2). \text{ Thus,}$$

$$\begin{aligned} f_{\vec{U}}(\vec{u}) &= f_{U_1, U_2}(u_1, u_2) = \frac{1}{4} \exp \left\{ - \left(\frac{1}{5}(2u_1 + u_2) + \frac{1}{5}(u_1 - 2u_2) \right) / 2 \right\} \cdot \left| -\frac{1}{5} \right| \\ &= \frac{1}{20} \exp \left\{ \frac{-1}{10} (2u_1 + u_2 + u_1 - 2u_2) \right\} = \frac{1}{20} \exp \left\{ \frac{3u_1 - u_2}{-10} \right\} = \frac{1}{20} \exp \left\{ \frac{u_2 - 3u_1}{10} \right\}. \end{aligned}$$

- **BUT ONLY IF** $y_1 = h_1^{-1}(\vec{u}) \geq 0$ **AND** $y_2 = h_2^{-1}(\vec{u}) \geq 0$! What are those vectors \vec{u} ?

- (c) $y_1 \geq 0$ and $y_2 \geq 0 \Leftrightarrow 2u_1 + u_2 \geq 0$
and $u_1 - 2u_2 \geq 0$
- (d) $y_1 \geq 0$ and $y_2 \geq 0 \stackrel{(9.27)}{\Rightarrow} u_1 = 2y_1 + y_2 \geq 0$.
- (e) From (c): $2u_1 + u_2 \geq 0 \Rightarrow u_2 \geq -2u_1$
- (f) From (c): $u_1 - 2u_2 \geq 0 \Rightarrow u_1 \geq 2u_2$
 $\Rightarrow u_2 \leq \frac{u_1}{2}$
- (g) From (d), (e), (f): $h_1^{-1}(\vec{u}) \geq 0$ and $h_2^{-1}(\vec{u}) \geq 0 \Leftrightarrow u_1 \geq 0$ and $-2u_1 \leq u_2 \leq \frac{u_1}{2}$.
- The figure to the right shows that those are the points enclosed by the quadrant which is obtained when rotating the first quadrant clockwise, by an angle of 60°



where $h_1^{-1}(u_1, u_2) > 0$ and $h_2^{-1}(u_1, u_2) > 0$

(h) Thus, if we denote this quadrant by R ,

$$f_{\vec{U}}(\vec{u}) = \begin{cases} \frac{1}{20} e^{(u_2-3u_1)/10}, & \text{if } \vec{u} \in R, \\ 0, & \text{else.} \end{cases}$$

At this point we know how to integrate with respect to the PDF of $\vec{U} = \vec{h}(\vec{Y})$. We can replace the integral $d\vec{u}$ over the region R by an iterated integral $du_2 du_1$ as follows.

For a fixed $u_1 > 0$, the integration bounds for u_2 are $-2u_1 \leq u_2 \leq \frac{u_2}{2}$. (See (g)). Thus,

$$\iint_{\mathbb{R}^2} \dots f_{\vec{U}}(\vec{U}) d\vec{u} = \iint_R \dots \frac{1}{20} e^{(u_2-3u_1)/10} d\vec{u} = \int_0^\infty \int_{-2u_1}^{u_2/2} \dots \frac{1}{20} e^{(u_2-3u_1)/10} du_2 du_1$$

For example, if $w = g(\vec{U}) = g(u_1, u_2)$ is a real-valued function of $(u_1, u_2) \in \mathbb{R}^2$, then

$$E[g(\vec{U})] = \int_0^\infty \int_{-2u_1}^{u_2/2} g(\vec{u}) \frac{1}{20} e^{(u_2-3u_1)/10} du_2 du_1 \quad \square$$

9.4 The Method of moment-generating Functions

Assumption 9.1. Unless stated otherwise, we will assume in this entire section that

- (a) $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$ denotes a list of n random variables ($n \in \mathbb{N}$).
 - Either all Y_j are discrete, or they all are continuous random variables.
- (b) $h : D \rightarrow \mathbb{R}; \vec{y} \mapsto u = h(\vec{y}) = h(y_1, \dots, y_n)$
is a function with domain $D \subseteq \mathbb{R}^n$ (this covers $\mathbb{R} = \mathbb{R}^1$ for $n = 1$), such that
 - there is no issue with the existence of the PMF or PDF of $U := h(\vec{Y})$.
 - All MGFs, $m_{Y_j}(t) = E[e^{tY_j}]$ and $m_U(t) = E[e^{tU}]$ exist if $|t|$ is small enough, i.e., there is some $\delta > 0$ such that those MGFs exist for $-\delta < t < \delta$.
- (c) Those assumptions also hold for differently named (vectors of) random variables and functions, e.g. $V = g(\vec{Y}) = g(\tilde{Y}_1, \dots, \tilde{Y}_k)$. \square

Introduction 9.3. The moment–generating function method for finding the probability distribution of a function of random variables Y_1, Y_2, \dots, Y_n is based on Proposition 6.5 on p.111 (Section 6.5: Moments, Central Moments and Moment Generating Functions). It was stated without proof and asserts that the following is true under the conditions stated in Assumption 9.1:

Assume that two random variables Y and \tilde{Y} possess identical k th moments about the origin for all $k = 1, 2, \dots$. In other words, assume that

$$E[Y^1] = E[\tilde{Y}^1], E[Y^2] = E[\tilde{Y}^2], E[Y^3] = E[\tilde{Y}^3], \dots$$

Then $P_Y = P_{\tilde{Y}}$, i.e., Y and \tilde{Y} have the same distribution. \square

We have the following uniqueness theorem.

Theorem 9.3 (The MGF determines the distribution).

Given are two random variables Y and \tilde{Y} . If their moment–generating functions $m_Y(t)$ and $m_{\tilde{Y}}(t)$ exist and coincide in a small interval that is centered at $t = 0$,

- *Then $P_Y = P_{\tilde{Y}}$, i.e., Y and \tilde{Y} have the same probability distribution.*

PROOF:

Theorems 6.18 on p.111 and 7.9 on p.127 allow us to conclude that

$$E[Y^k] = \left. \frac{d^k}{dt^k} m_Y(t) \right|_{t=0} = \left. \frac{d^k}{dt^k} m_{\tilde{Y}}(t) \right|_{t=0} = E[\tilde{Y}^k] \text{ for all } k \in \mathbb{N}.$$

It follows from Proposition 6.5 on p.111 that $P_Y = P_{\tilde{Y}}$ \blacksquare

Remark 9.6.

To find the distribution of $U = h(\vec{Y}) = h(Y_1, Y_2, \dots, Y_n)$ by means of the MGF method, proceed as follows:

- Compute the MGF $m_U(t) = E[e^{tU}]$ of U
- Does this MGF match that of a random variable V with a known distribution? You may want to consult a list of MGFs like the one in Appendix 2 of [11] Wackerly, Mendenhall, Scheaffer, R.L.
- Then you are done, since Theorem 9.3 (The MGF determines the distribution) guarantees that $P_U = P_V$.

Of course, the devil is in the details. In most cases, you will not succeed in finding that matching MGF, unless one or both of the following are satisfied:

- U is a linear function $U = a_1 Y_1 + \dots + a_n Y_n$, with constant $a_j \in \mathbb{R}$.
- The random variables Y_1, \dots, Y_n are independent and $h(\vec{y}) = h_1(y_1) \cdot h_2(y_2) \cdot \dots \cdot h_n(y_n)$, for suitable functions $h_j(y)$.

We will examine some very important and general cases that illustrate all this. \square

Example 9.13 (WMS Ch.06.5, Example 6.10). Suppose that Y is a normally distributed random variable with mean μ and variance σ^2 . Show that

$$Z := \frac{Y - \mu}{\sigma}$$

has a standard normal distribution, i.e., $Z \sim \mathcal{N}(0, 1)$.

Solution:

- (a) According to Proposition 7.6 on p.133, $m_Y(t) = e^{\mu t + (\sigma^2 t^2)/2}$.
- (b) Any random variable W is independent from any constant (real number) a .
- (c) Thus, according to Theorem 8.10 on p.157, the random variables $h_1(W) = e^{tW}$ and $h_2(a) = e^{-at}$ are independent, and $E[e^{tW} \cdot e^{-at}] = E[e^{tW}] \cdot e^{-at}$.
- (d) Thus if $U = Y - \mu$, then $m_U(t) = E[e^{tY - t\mu}] = E[e^{tY} e^{-t\mu}] = E[e^{tY}] \cdot e^{-t\mu}$.
- Thus, $m_U(t) = m_Y(t) e^{-t\mu} \stackrel{\text{(a)}}{=} e^{\mu t + (\sigma^2 t^2)/2} \cdot e^{-t\mu} = e^{(\sigma^2 t^2)/2}$.
 - Since $Z = U/\sigma$, $m_Z(t) = m_U(t/\sigma) = e^{(\sigma^2 (t/\sigma)^2)/2} = e^{t^2/2}$.
- (e) We use Proposition 7.6 once more and see that $t \mapsto e^{t^2/2}$ is the MGF of a standard normal random variable. Thus, $Z \sim \mathcal{N}(0, 1)$. \square

Example 9.14 (WMS Ch.06.5, Example 6.11). Let Z be a normally distributed random variable with mean 0 and variance 1. Use the method of moment-generating functions to find the probability distribution of Z^2 .

Solution:

The moment-generating function for Z^2 is

$$\begin{aligned} \text{(A)} \quad m_{Z^2}(t) &= E(e^{tZ^2}) = \int_{-\infty}^{\infty} e^{tz^2} f(z) dz = \int_{-\infty}^{\infty} e^{tz^2} \frac{e^{-z^2/2}}{\sqrt{2\pi}} dz \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-(z^2/2)(1-2t)} dz = \int_{-\infty}^{\infty} \psi(z) dz, \end{aligned}$$

where

$$\begin{aligned} \psi(z) &= \exp \left[- \left(\frac{z^2}{2} \right) (1 - 2t) \right] / \sqrt{2\pi} \\ &= \exp \left[- \left(\frac{z^2}{2} \right) / (1 - 2t)^{-1} \right] / \left(\sqrt{2\pi} (1 - 2t)^{-1/2} \cdot \frac{1}{(1 - 2t)^{-1/2}} \right). \end{aligned}$$

We define $\sigma := (1 - 2t)^{-1/2}$ and obtain

$$\psi(z) = \exp \left[- \left(\frac{z^2}{2} \right) / \sigma^2 \right] / \left(\sqrt{2\pi} \sigma \cdot \frac{1}{\sigma} \right) = e^{-z^2/(2\sigma^2)} \cdot \frac{\sigma}{\sqrt{2\pi} \sigma} = \sigma \varphi(z),$$

where $\varphi(z)$ is the density of a $\mathcal{N}(0, \sigma)$ random variable. Thus, $\int_{-\infty}^{\infty} \varphi(z) dz = 1$. It follows from (A) and $\psi(z) = \sigma \varphi(z)$ and $\sigma := (1 - 2t)^{-1/2}$ that

$$m_{Z^2}(t) = \int_{-\infty}^{\infty} \psi(z) dz = \int_{-\infty}^{\infty} (1 - 2t)^{-1/2} \varphi(z) dz = \frac{1}{(1 - 2t)^{1/2}} \int_{-\infty}^{\infty} \varphi(z) dz = \frac{1}{(1 - 2t)^{1/2}}.$$

According to Proposition 7.8 on p.136, $t \mapsto \frac{1}{(1 - 2t)^{1/2}}$ is the MGF of a random variable which follows a gamma(1/2, 2) distribution which is, by definition 7.11 on p.137, also known as a χ^2 distribution with one degree of freedom. We obtained this result previously in Example 9.5 on p.194 by the method of distribution functions. \square

Theorem 9.4 (MGF of a sum of functions of independent variables).

Given are n independent random variables Y_1, Y_2, \dots, Y_n with MGFs $m_{Y_1}(t), m_{Y_2}(t), \dots, m_{Y_n}(t)$. and n real-valued functions $h_1(y_1), \dots, h_n(y_n)$ of real numbers y_1, \dots, y_n .

Let $U := h_1(Y_1) + h_2(Y_2) + \dots + h_n(Y_n)$. Then (under the conditions of Assumption 9.1 on 206)

$$(9.28) \quad m_U(t) = m_{h_1(Y_1) + \dots + h_n(Y_n)} = \prod_{j=1}^n m_{h_j(Y_j)}(t).$$

PROOF:

For each $j = 1, \dots, n$, let $g_j(y) := e^{th_j(y)}$. Consider a fixed t . Since functions of independent random variables are independent random variables, the random variables $V_j := g_j(Y_j) = e^{th_j(Y_j)}$ are independent. We apply Theorem 8.10 on p.157 and obtain

$$\begin{aligned} m_U(t) &= E[e^{t(V_1 + V_2 + \dots + V_n)}] \\ &= E[e^{tV_1}] \dots E[e^{tV_n}] = E[e^{th_1(Y_1)}] \dots E[e^{th_n(Y_n)}] \\ &= m_{h_1(Y_1)}(t) \cdot m_{h_2(Y_2)}(t) \cdot \dots \cdot m_{h_n(Y_n)}(t). \quad \blacksquare \end{aligned}$$

Corollary 9.1 (WMS Ch.06.5, Theorem 6.2).

Let Y_1, Y_2, \dots, Y_n be independent random variables with moment-generating functions $m_{Y_1}(t), m_{Y_2}(t), \dots, m_{Y_n}(t)$, respectively. Then

$$(9.29) \quad m_{Y_1 + \dots + Y_n}(t) = \prod_{j=1}^n m_{Y_j}(t) = m_{Y_1}(t) \cdot m_{Y_2}(t) \cdot \dots \cdot m_{Y_n}(t).$$

PROOF:

This follows from applying Theorem 9.4 to the functions $h_j(y_j) = y_j$. \blacksquare

Next, we generalize But its great importance gives it the status of a theorem.

Theorem 9.5 (Linear combinations of uncorrelated normal variables are normal).

Given are n uncorrelated, $\mathcal{N}(\mu_j, \sigma_j^2)$ random variables Y_j , ($j = 1, \dots, n$). In other words, each Y_j is normal with expectation μ_j and standard deviation σ_j . Let $a_1, \dots, a_n \in \mathbb{R}$. Then

$$(9.30) \quad \sum_{j=1}^n a_j Y_j \sim \mathcal{N} \left(\sum_{j=1}^n a_j \mu_j, \sum_{j=1}^n a_j^2 \sigma_j^2 \right).$$

Thus, the linear combination of uncorrelated normal random variables is normal with expectation and variance being the linear combinations of the individual expectations and variances.

PROOF:

First off, we recall that one of the special properties of normal random variables is that they are uncorrelated if and only if they are independent. Thus we can use everything that applies to independent random variables.

Consider a fixed t and define

$$U := \sum_{j=1}^n a_j Y_j.$$

We apply Theorem 9.4 (MGF of a sum of functions of independent variables) on p.209 with the functions $h_j(y_j) = a_j y_j$ and obtain

$$\begin{aligned} m_U(t) &= \prod_{j=1}^n m_{a_j Y_j}(t) = \prod_{j=1}^n m_{Y_j}(a_j t) \\ &= \prod_{j=1}^n \exp \left\{ (\sigma_j^2/2)(a_j t)^2 + \mu_j(a_j t) \right\} \end{aligned}$$

Here we used that a $\mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2)$ variable has MGF $e^{\tilde{\sigma}^2 t^2/2 + \tilde{\mu} t}$. See Proposition 7.6 on p.133. Thus,

$$\begin{aligned} m_U(t) &= \exp \left\{ \sum_{j=1}^n (\sigma_j^2/2)(a_j t)^2 + \mu_j(a_j t) \right\} \\ &= \exp \left\{ \left(\sum_{j=1}^n (\sigma_j^2 a_j^2/2) t^2 \right) + \left(\sum_{j=1}^n (\mu_j a_j) t \right) \right\} \\ &= \exp \left\{ \left(\sum_{j=1}^n (a_j^2 \sigma_j^2) \right) / 2 \cdot t^2 + \left(\sum_{j=1}^n (a_j \mu_j) \right) \cdot t \right\} \end{aligned}$$

By Proposition 7.6, the last expression is the MGF of a $\mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2)$ variable with

$$\tilde{\mu} = \sum_{j=1}^n (a_j \mu_j), \quad \tilde{\sigma}^2 = \sum_{j=1}^n (a_j^2 \sigma_j^2).$$

Since distributions of random variables are determined by their MGFs,

$$U \sim \mathcal{N} \left(\sum_{j=1}^n a_j \mu_j, \sum_{j=1}^n a_j^2 \sigma_j^2 \right). \quad \blacksquare$$

Remark 9.7. It is a consequence of Theorem 9.5 that the sum of two independent random variables also is normal. The following counterexample shows that we cannot drop the assumption of independence. It is cited in many books on probability and can be found, e.g., in [7] Pishro-Nik, Hossein: Introduction to Probability, Statistics, and Random Processes.

Assume that U and V are independent random variables with distributions

- $U \sim \mathcal{N}(0, 1)$,
- $V \sim \text{binom}(n = 1, p = 0.5)$.

Let

$$W(\omega) := \begin{cases} U(\omega), & \text{if } V(\omega) = 1, \\ -U(\omega), & \text{if } V(\omega) = 0. \end{cases}$$

- (a) Show that $W \sim \mathcal{N}(0, 1)$.
 (b) Let $Y := U + W$. Show that Y is not a continuous random variable.

It follows from (b) that Y is not normal, since normal random variables are continuous. \square

Solution to (a): Note that the PDF of U is symmetric, i.e., $f_U(u) = f_U(-u)$ for all $u \in \mathbb{R}$. Thus, for all u ,

$$P\{U \leq u\} = \int_{-\infty}^u f_U(t) dt = \int_{-u}^{\infty} f_U(t) dt = P\{U \geq -u\} = P\{-U \leq u\}.$$

It follows that U and $-U$ have the same distribution and thus, $-U \sim \mathcal{N}(0, 1)$.⁵⁵

Now, we show that $W \sim \mathcal{N}(0, 1)$. Let $w \in \mathbb{R}$. Then,

$$\begin{aligned} P\{W \leq w\} &= P\{W \leq w, V = 0\} + P\{W \leq w, V = 1\} \\ &= P\{W \leq w \mid V = 0\} P\{V = 0\} + P\{W \leq w \mid V = 1\} P\{V = 1\} \\ &= \frac{1}{2} P\{-U \leq w \mid V = 0\} + \frac{1}{2} P\{U \leq w \mid V = 1\} \end{aligned}$$

We use the independence of U and V followed by $U \sim -U$ and obtain

$$P\{W \leq w\} = \frac{1}{2} (P\{-U \leq w\} + P\{U \leq w\}) = \frac{1}{2} (P\{U \leq w\} + P\{U \leq w\}) = P\{U \leq w\}.$$

Thus, $W \sim U$. Since U is standard normal, so is W . We have proven (a).

Solution to (b): It follows from the definition of W and $Y := U + W$, that

$$Y(\omega) := \begin{cases} 2U(\omega), & \text{if } V(\omega) = 1, \\ 0, & \text{if } V(\omega) = 0. \end{cases}$$

It follows that the CDF F_Y of Y has a jump

- $F_Y(0) - F_Y(0-) = P\{Y = 0\} = 1/2$

at $y = 0$. Thus, Y is not a continuous random variable and we have shown (b). \blacksquare

Theorem 9.6.

Given are n independent, gamma(α_j, β) random variables Y_j , ($j = 1, \dots, n$). In other words, each Y_j is gamma with the same scale parameter β . Then

$$(9.31) \quad \sum_{j=1}^n Y_j \sim \text{gamma} \left(\sum_{j=1}^n \alpha_j, \beta \right).$$

Thus, the sum of independent gamma random variables with the same scale parameter β is gamma with the shape parameter being the sum of the shape parameters, and scale parameter β .

⁵⁵This result should not come as a surprise since, for $n = 1$ and $a_1 = -1$, Theorem 9.5 on p.209 states the following: If $Y_1 \sim \mathcal{N}(\mu, \sigma^2)$, then $-Y_1 \sim \mathcal{N}(-\mu, \sigma^2)$. Note though, that the proof given here shows that U and $-U$ have the same distribution whenever U has a symmetric PDF. Also note that $U \sim -U$ holds if U is discrete with a symmetric PMF, i.e., $p_U(u) = P\{U = u\} = P\{U = -u\} = p_U(-u)$, for all u .

PROOF:

Consider a fixed t and define

$$U := \sum_{j=1}^n Y_j.$$

We apply Theorem 9.4 (MGF of a sum of functions of independent variables) on p.209 and recall that the MGF of a gamma($\tilde{\alpha}$, $\tilde{\beta}$) variable \tilde{Y} is, according to Proposition 7.8 on p.136, $m_{\tilde{Y}} = (1 - \tilde{\beta}t)^{-\tilde{\alpha}}$. We obtain

$$\begin{aligned} m_U(t) &= \prod_{j=1}^n m_{a_j Y_j}(t) = \prod_{j=1}^n m_{Y_j}(a_j t) \\ &= \prod_{j=1}^n \frac{1}{(1 - \beta t)^{\alpha_j}} = \frac{1}{(1 - \beta t)^{\sum_{j=1}^n \alpha_j}}. \end{aligned}$$

Since distributions of random variables are determined by their MGFs,

$$U \sim \text{gamma} \left(\sum_{j=1}^n \alpha_j, \beta \right). \blacksquare$$

Corollary 9.2.

Let Y_1, Y_2, \dots, Y_n be independent χ^2 variables such that each Y_j has ν_j degrees of freedom. Then

$$(9.32) \quad m_{Y_1 + \dots + Y_n}(t) \sim \chi^2 \left(\sum_{j=1}^n \nu_j \text{ df} \right).$$

PROOF:

This follows immediately from Theorem 9.6, Since χ^2 variables with ν_j df are gamma($\nu_j/2, 2$). \blacksquare

10 Limit Theorems

Introduction 10.1. In this section we will discuss the ways in which a sequence Y_n of random variables can have a random variable Y as its limit. Before we go there, let us quickly review convergence of a sequence $(y_n)_n$ of real numbers and of a sequence of functions $f_n : A \rightarrow \mathbb{R}$, with all members f_n defined on a subset A of \mathbb{R}^k , where $k = 1, 2, \dots$. Note that $k = 1$ covers the situation where the arguments are real numbers. Some examples of number sequences:

- If $y_n = \frac{3 - 2n}{5 + n^2 - 6n}$, then $\lim_{n \rightarrow \infty} y_n = \frac{3}{5}$, and the sequence converges to $\frac{3}{5}$.
- If $y_n = (-1)^n$, then $\lim_{n \rightarrow \infty} y_n$ does not exist.
- If $y_n = \sum_{j=1}^n n$, then $\lim_{n \rightarrow \infty} y_n = \infty$. Recall that convergence only happens if the limit is a real number. Thus, $(y_n)_n$ does not “converge to ∞ ”. Rather, this sequence diverges. ⁵⁶

For the following examples of function sequences, let us agree that, if $f_n, f : A \rightarrow \mathbb{R}$, where $A \subseteq \mathbb{R}$, then “pointwise convergence” ⁵⁷ of the functions f_n to the function f simply means that

$$(10.1) \quad \lim_{n \rightarrow \infty} f_n(a) = f(a) \quad \text{for all } a \in A.$$

- Let $f_n, f, g, h : [0, 1] \rightarrow \mathbb{R}$ be the functions

$$(10.2) \quad \square f_n(x) := x^n \quad \square f(x) := \begin{cases} 0, & \text{if } 0 \leq x < 1, \\ 1, & \text{if } x = 1, \end{cases} \quad \square g(x) := 0, \quad \square h(x) := x.$$

The situation with respect to pointwise convergence is as follows:

- f is the pointwise limit of the sequence f_n .
- Even though g is the pointwise limit of the sequence f_n on $[0, 1[$, it is not the pointwise limit on $[0, 1]$, since $\lim_{n \rightarrow \infty} f_n(x) = g(x) = 0$, for $0 \leq x < 1$, but $\lim_{n \rightarrow \infty} f_n(1) = 1$, whereas $g(1) = 0$.
- h is not the pointwise limit of the sequence f_n (except on $\{0, 1\}$).

Did you notice that no use was made of the fact that the domain $[0, 1]$ of those functions is a set of numbers?

- Assume instead that Ω is some arbitrary, nonempty set (not necessarily a probability space). Further assume that there are functions $f_n, f : \Omega \rightarrow \mathbb{R}$. We still have the notion of pointwise convergence of the functions f_n to the function f : (10.1) becomes

$$(10.3) \quad \lim_{n \rightarrow \infty} f_n(\omega) = f(\omega) \quad \text{for all } \omega \in \Omega$$

and one certainly can examine whether or not the above is true for any kind of domain, i.e., for any nonempty set Ω .

We will not discuss vector-valued sequences. However, for completeness sake, we give the following example.

⁵⁶There is no such thing as divergence to $\pm\infty$. Thus, you must say that (y_n) diverges, **not** that (y_n) diverges to ∞ .

⁵⁷The formal definition of pointwise limits will be given in Section 10.1 (Four Kinds of Limits for Sequences of Random Variables).

- If $\vec{y}_n = ((-1)^n, \cos(2/n))$, then $\lim_{n \rightarrow \infty} \vec{y}_n$ does not exist, since the limit of a vector-valued sequence is, by definition, the vector of the limits of the coordinates. The second coordinate sequence, $y_n = \cos(2/n)$, converges to the number 1. Since the first coordinate sequence, $y_n = (-1)^n$, does not have a limit, neither does $(\vec{y}_n)_n$. Thus this sequence does not converge.

After these preliminary remarks, let us consider sequences of random variables. We recall that all random variables Y are functions

$$Y : (\Omega, \mathfrak{F}, P) \rightarrow \mathbb{R} \quad \omega \mapsto Y(\omega).$$

They take their arguments ω in a probability space $(\Omega, \mathfrak{F}, P)$ and map them to numeric outcomes $y = Y(\omega)$.

- The σ -algebra is of no significance in this chapter, so we keep ignoring it and simply consider the probability space (Ω, P) .
- On the other hand, the arguments ω play an essential role and we will often replace “ Y ” with “ $\omega \mapsto Y(\omega)$ ” to remind the reader that we are dealing with functions of ω .
- If $(Y_n)_n$ is a sequence of random variables $(\Omega, P) \rightarrow \mathbb{R}$. Then each $\omega \in \Omega$ comes with its own sequence $(Y_n(\omega))_n$ of real numbers.
- One obvious question to ask about those sequences $Y_n(\omega)$ of real numbers is this one:
 - Does $\lim_{n \rightarrow \infty} Y_n(\omega)$ exist and will it be a real number (rather than $\pm\infty$) for all $\omega \in \Omega$?
 - If so, then the assignment $\omega \mapsto Y(\omega) := \lim_{n \rightarrow \infty} Y_n(\omega)$ defines a real-valued function $Y : (\Omega, P) \rightarrow \mathbb{R}$, i.e., another random variable. What are its properties?
- Not quite so obvious: □ Does the presence of the probability measure P on Ω give additional insight about the convergence behavior of the functions $\omega \mapsto Y_n(\omega)$?
- In contrast to the deterministic case where the only mode of convergence available to us is pointwise convergence,⁵⁸ we will see in Section 10.1 (Four Kinds of Limits for Sequences of Random Variables) that the presence of a probability P allows us to consider additional modes of convergence:
 - convergence almost surely,
 - convergence in probability measure,
 - convergence in distribution. □

10.1 Four Kinds of Limits for Sequences of Random Variables

The following definition is a central place for all the different convergence modes of sequences of random variables that are of interest to us. We will examine each one in detail.

Definition 10.1 (Convergence of Random Variables).

⁵⁸This is not entirely true: If Ω is a subset of \mathbb{R} or of \mathbb{R}^k , then there is the notion of **uniform convergence**, $f_n(\cdot) \rightarrow f(\cdot)$. We will not be concerned with uniform convergence in this course.

Let Y_n ($n \in \mathbb{N}$) and Y be random variables on a probability space (Ω, P) . We define

$$(10.4) \quad Y_n \xrightarrow{\text{pw}} Y \text{ or } \text{pw} - \lim_{n \rightarrow \infty} Y_n = Y, \quad \text{if } \lim_{n \rightarrow \infty} Y_n(\omega) = Y(\omega), \text{ for all } \omega \in \Omega,$$

$$(10.5) \quad Y_n \xrightarrow{\text{a.s.}} Y \text{ or } \text{a.s.} - \lim_{n \rightarrow \infty} Y_n = Y, \quad \text{if } P\{\omega \in \Omega : \lim_{n \rightarrow \infty} Y_n(\omega) = Y(\omega)\} = 1,$$

$$(10.6) \quad Y_n \xrightarrow{P} Y \text{ or } P - \lim_{n \rightarrow \infty} Y_n = Y, \quad \text{if } \forall \varepsilon > 0 \lim_{n \rightarrow \infty} P\{\omega \in \Omega : |Y_n(\omega) - Y(\omega)| > \varepsilon\} = 0,$$

$$(10.7) \quad Y_n \xrightarrow{D} Y, \text{ if } \lim_{n \rightarrow \infty} F_{Y_n}(y) = F_Y(y), \forall y \in \mathbb{R} \text{ where the CDF } F_Y \text{ of } Y \text{ is continuous.}$$

We also say:

If $Y_n \xrightarrow{\text{pw}} Y$, Y is the **pointwise limit** of the Y_n , or: Y_n **converges pointwise** to Y .

If $Y_n \xrightarrow{\text{a.s.}} Y$, Y is the **almost sure limit** of the Y_n , or: Y_n **converges almost surely** to Y .

If $Y_n \xrightarrow{P} Y$, Y is the **limit in probability**; of the Y_n , or: Y_n **converges in probability** to Y .

If $Y_n \xrightarrow{D} Y$, Y is the **limit in distribution** of the Y_n , or: Y_n **converges in distribution** to Y .

Example 10.1. Consider $\Omega := [0, 1]$ as a probability space (Ω, P) by defining

$$P([a, b]) := b - a, \text{ for } 0 \leq a < b \leq 1.$$

In other words, P is the uniform distribution on $[0, 1]$.

We rename the functions f_n, f, g, h of (10.2) in the introduction to Y_n, Y, U, V , since doing so will make it less confusing to examine the convergence behavior of the sequence. This particularly applies to converges in probability and in distribution. Accordingly, we define

$$Y_n(\omega) := \omega^n, \quad U(\omega) = 0, \quad V(\omega) := \omega, \quad (\text{for } 0 \leq \omega \leq 1) \quad Y(\omega) := \begin{cases} 0, & \text{if } 0 \leq \omega < 1, \\ 1, & \text{if } \omega = 1. \end{cases}$$

Part I: Pointwise and a.s convergence

Pointwise convergence behavior of the Y_n corresponds to that of (10.2):

- Y is the pointwise limit of the sequence Y_n ,
- U is the pointwise limit of the Y_n on $[0, 1[$ only, but not on Ω ,
- V is not the pointwise limit of the Y_n (except for $\omega = 0$) or $\omega = 1$).

With respect to almost sure convergence, we see that

- $Y_n \xrightarrow{\text{a.s.}} Y$, since $\{\lim_{n \rightarrow \infty} Y_n = Y\} = [0, 1] = \Omega$, and $P(\Omega) = 1$.
- $Y_n \xrightarrow{\text{a.s.}} U$, since $\{\lim_{n \rightarrow \infty} Y_n \neq U\} = \{1\}$, and $P(\{1\}) = 0$.
- $(Y_n)_n$ does not converge to V a.s., since $P\{\lim_{n \rightarrow \infty} Y_n = V\} = P\{0, 1\} = 0 \neq 1$.

Part II: Convergence in probability

Next, we examine convergence in probability. We will see that a sequence of random variables can have more than one P -limit by showing the following: The sequence $\omega \mapsto Y_n(\omega) = \omega^n$ has both $\omega \mapsto U(\omega) = 0$ and $\omega \mapsto Y(\omega) = 1$ if $\omega = 1$ and 0 else as P -limits.

By definition of $P\text{-}\lim_{n \rightarrow \infty} Y_n = \tilde{Y}$, we must prove that, for any fixed, but arbitrary $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} P\{|Y_n - \tilde{Y}| > \varepsilon\} = 0. \quad \text{See (10.6).}$$

Since this probability decreases as ε increases and we must show that it approaches 0 as $n \rightarrow \infty$, we only need to worry about the very small ε . Thus, we may assume that $0 < \varepsilon < 1$.

We observe that, for $Y_n(\omega) = \omega^n$ and $0 < \varepsilon < 1$,

$$\begin{aligned} \text{(A)} \quad & [|Y_n(\omega)| \geq \varepsilon \Leftrightarrow \omega^n \geq \varepsilon \Leftrightarrow \omega \geq \varepsilon^{1/n}] \\ & \Rightarrow [P\{|Y_n| \geq \varepsilon\} = P([\varepsilon^{1/n}, 1]) = 1 - \varepsilon^{1/n}]. \end{aligned}$$

$$\text{(B)} \quad 0 < \varepsilon < 1 \Rightarrow \lim_{n \rightarrow \infty} \varepsilon^{1/n} = 1 \Rightarrow \lim_{n \rightarrow \infty} (1 - \varepsilon^{1/n}) = 0.$$

Part II (1): We now prove that $P\text{-}\lim_{n \rightarrow \infty} Y_n = Y$:

$$\begin{aligned} \text{(a)} \quad & [|Y_n(\omega) - Y(\omega)| \geq \varepsilon \Leftrightarrow |Y_n(\omega)| \geq \varepsilon \text{ and } \omega \neq 1] \\ & \Rightarrow [P\{|Y_n - Y| \geq \varepsilon\} \leq P\{|Y_n| \geq \varepsilon\} \stackrel{\text{(A)}}{=} 1 - \varepsilon^{1/n} \stackrel{\text{(B)}}{\rightarrow} 0, \text{ as } n \rightarrow \infty.]. \end{aligned}$$

Thus, $\lim_{n \rightarrow \infty} P\{|Y_n - Y| \geq \varepsilon\} = 0$.

Part II (2): We now prove that $P\text{-}\lim_{n \rightarrow \infty} Y_n = U$:

- We could repeat the proof for the P -convergence of Y_n to Y with very minor modifications and the reader is encouraged to do so. Instead, we will use that result to show that $P\text{-}\lim_{n \rightarrow \infty} Y_n = U$
- Since the outcome $\{1\}$ has probability zero and $Y(\omega) = U(\omega)$ for $\omega \neq 1$,

$$\begin{aligned} P\{|Y_n - Y| \geq \varepsilon\} &= P\{|Y_n - Y| \geq \varepsilon \text{ and } \omega \neq 1\} \\ &= P\{|Y_n - U| \geq \varepsilon \text{ and } \omega \neq 1\} = P\{|Y_n - U| \geq \varepsilon\}. \end{aligned}$$

- Since $\lim_{n \rightarrow \infty} P\{|Y_n - Y| \geq \varepsilon\} = 0$,

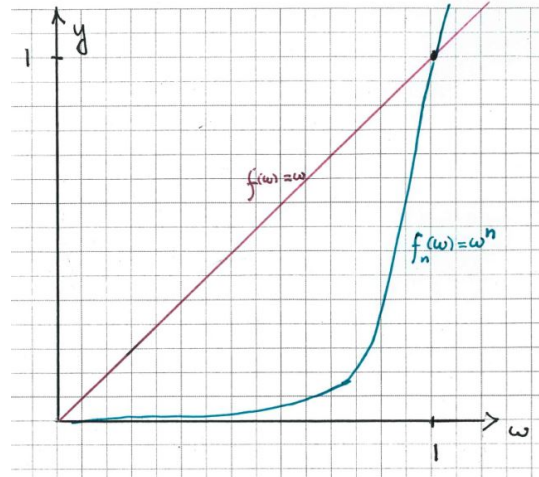
$$\lim_{n \rightarrow \infty} P\{|Y_n - U| \geq \varepsilon\} = \lim_{n \rightarrow \infty} P\{|Y_n - Y| \geq \varepsilon\} = 0.$$

Thus, $P\text{-}\lim_{n \rightarrow \infty} Y_n = U$.

Part II (3): Next, we show that it is not true that $(Y_n)_n$ converges in probability to V .

We argue by picture rather than giving an exact proof, since that would require some very tedious of terms containing $\ln(k)$.

- The picture makes it very clear that $\varepsilon = 1/10 \Rightarrow \omega - \omega^n > \varepsilon$ for $\frac{49}{100} \leq \omega \leq \frac{51}{100}$ and $n \geq 100$.
 Thus, $P\{|Y_n - V| \geq \varepsilon\} \geq \varepsilon \cdot \left(\frac{51}{100} - \frac{49}{100}\right) = \frac{2}{1000}$.
 Thus, $\lim_{n \rightarrow \infty} P\{|Y_n - V| \geq \varepsilon\} = 0$ is not true.
- Since $\lim_{n \rightarrow \infty} P\{|Y_n - V| \geq \varepsilon\} = 0$ must hold for ALL ε and we showed that this is not so for $\varepsilon = \frac{1}{10}$, it follows that $(Y_n)_n$ does not converge in probability to V .



Part III: Convergence in distribution

We will show that Y_n does not converge to V in distribution as follows.

- Recall that $P[a, b] = b - a$ for all $0 \leq a < b \leq 1$. Let $y \in \mathbb{R}$.
- Since $V(\omega) = \omega$, $F_V(y) = P\{V \leq y\} = P\{\omega \in \Omega : V(\omega) \leq y\} = P]0, y] = y$.
- Since $Y_n(\omega) = \omega^n$, $F_{Y_n}(y) = P\{Y_n \leq y\} = P\{\omega \in \Omega : \omega^n \leq y\} = P]0, y^{1/n}] = y^{1/n}$.
- Thus, for $0 < y < 1$, $F_V(y) = y$, whereas, $\lim_{n \rightarrow \infty} F_{Y_n}(y) = 0 \neq F_V(y)$.
- Since all those y are points of continuity for F_V , it follows that $(Y_n)_n$ does not converge in distribution to V .

On the other hand, the theorem that follows now shows that $(Y_n)_n$ converges in distribution to Y and U , since we have shown convergence in probability to those random variables. \square

Theorem 10.1 (Relationship between the modes of convergence).

Let Y and Y_1, Y_2, \dots be random variables on a probability space (Ω, P) . Then,

$$(10.8) \quad Y_n \xrightarrow{pw} Y \Rightarrow Y_n \xrightarrow{a.s.} Y \Rightarrow Y_n \xrightarrow{P} Y \Rightarrow Y_n \xrightarrow{D} Y.$$

PROOF:

I: It is obvious that $Y_n \xrightarrow{pw} Y \Rightarrow Y_n \xrightarrow{a.s.} Y$ for the following reason:

- For each $n \in \mathbb{N}$, let $A_n := \{\omega \in \Omega : \lim_{n \rightarrow \infty} Y_n(\omega) \neq Y(\omega)\}$.
- Then, for $k \in \mathbb{N}$, $Y_n \xrightarrow{pw} Y \Rightarrow A_k = \emptyset \Rightarrow P(A_k) = 0 \Rightarrow \lim_{n \rightarrow \infty} P(A_n) = 0 \Rightarrow Y_n \xrightarrow{a.s.} Y$.

II: The proofs that $Y_n \xrightarrow{a.s.} Y \Rightarrow Y_n \xrightarrow{P} Y$ and $Y_n \xrightarrow{P} Y \Rightarrow Y_n \xrightarrow{D} Y$ are outside the scope of this course. Fairly accessible proofs for those who can work with sets like

$$\bigcap_{n \geq 1} \left(\bigcup_{j \geq n} \{\omega \in \Omega : |Y_j(\omega) - Y(\omega)| \geq \varepsilon\} \right)$$

and are familiar with the exact definition of convergence of sequences ⁵⁹ can be found at this [Wikipedia](#) link. ■

There are many theorems concerning the convergence of random variables. We only mention here the following two which will be used later in this chapter.

Theorem 10.2 (Slutsky's Theorem). ★

Let Y_1, Y_2, \dots and U_1, U_2, \dots be two sequences of random variables. Let Y be another random variable and c a constant such that

- $Y_n \xrightarrow{D} Y$ (convergence in distribution)
- $U_n \xrightarrow{P} c$ (convergence in probability)

Then,

$$(10.9) \quad Y_n + U_n \xrightarrow{D} Y + c,$$

$$(10.10) \quad Y_n \cdot U_n \xrightarrow{D} cY,$$

$$(10.11) \quad \frac{Y_n}{U_n} \xrightarrow{D} \frac{Y_n}{c}, \quad \text{assuming that } c \neq 0.$$

PROOF: Omitted. See, e.g., [1] Bickel and Doksum: Mathematical Statistics.

Theorem 10.3 (Convergence is maintained under continuous transformations). ★

Let Y_1, Y_2, \dots and Y be random variables on some probability space (Ω, P) . Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be continuous. Then,

$$Y_n \xrightarrow{a.s.} Y \Rightarrow f \circ Y_n \xrightarrow{a.s.} f \circ Y.$$

$$Y_n \xrightarrow{P} Y \Rightarrow f \circ Y_n \xrightarrow{P} f \circ Y.$$

$$Y_n \xrightarrow{D} Y \Rightarrow f \circ Y_n \xrightarrow{D} f \circ Y.$$

PROOF: Omitted. ⁶⁰ ■

Example 10.2 (Convergence in probability but not a.s.). ★

Consider the “sliding hump” example. ⁶¹ As our probability space we choose $\Omega := [0, 1]$, the unit interval in \mathbb{R} , with the probability measure defined by $P([a, b]) := b - a$.

(a) We partition Ω into the two intervals $I_1 = [0, 1/2]$ and $I_2 =]1/2, 1]$.

- For $n = 1, 2$, let $Y_n(\omega) := \begin{cases} 1, & \text{if } \omega \in I_n, \\ 0, & \text{else.} \end{cases}$

⁵⁹ x_n converges to $x \Leftrightarrow$ for all $\varepsilon > 0$ one can find $N \in \mathbb{N}$ such that $|x_n - x| < \varepsilon$ whenever $n \geq N$.

⁶⁰A proof can be found at this [Convergence of random variables](#) (Mann–Wald theorem, general transformation theorem) [Wikipedia](#) link.

⁶¹See this [StackExchange](#) link.

- (b) We partition Ω into the three intervals $I_3 = [0, 1/3]$, $I_4 =]1/3, 2/3]$, and $I_5 =]2/3, 1]$, then into $I_6 = [0, 1/4]$, $I_7 =]1/4, 2/4]$, $I_8 =]2/4, 3/4]$, and $I_9 =]3/4, 1]$, and so on
- We define random variables Y_n as in (a): For $n \in \mathbb{N}$, let $Y_n(\omega) := \begin{cases} 1, & \text{if } \omega \in I_n, \\ 0, & \text{else.} \end{cases}$
- (c) Then the sequence Y_n converges in probability to the (deterministic) random variable $\omega \mapsto Y(\omega) := 0$. A proof is given directly after this example.
- (d) But this sequence of random variables does not converge almost surely. In fact, there is no $0 \leq \omega \leq 1$ for which $\lim_{n \rightarrow \infty} Y_n(\omega)$ exist:
- Fix $\omega \in [0, 1]$. By construction, there are indices $n_1 = n_1(\omega) < n_2 = n_2(\omega) < n_3 = n_3(\omega) < \dots$, such that $\omega \in I_{n_k}$ and I_{n_k} has length $1/k$. (Thus, $P(I_{n_k}) = 1/k$.)
- (e) Let $\omega' \in [0, 1]$; $\omega' \neq \omega$. The subsequences $n_k(\omega)$ and $n_k(\omega')$ will differ for all k so large that $\frac{1}{k} < \frac{|\omega - \omega'|}{2}$, i.e., $\frac{2}{k} < |\omega - \omega'|$, since $\omega \in I_{n_k(\omega)}$ and $\omega' \in I_{n_k(\omega')} \Rightarrow I_{n_k(\omega)} \cap I_{n_k(\omega')} = \emptyset$. (Draw a picture!)
- (f) It follows for such big k , that $Y_{n_k(\omega)}(\omega) = 1$ and $Y_{n_k(\omega)}(\omega') = 0$. On the other hand, $Y_{n_k(\omega')}(\omega) = 0$ and $Y_{n_k(\omega')}(\omega') = 1$. Thus, the full sequences $Y_n(\omega)$ does not have a limit, since it would have to be 1 along the subsequence $n_k(\omega)$ and 0 along the subsequence $n_k(\omega')$.
- (g) ω is arbitrary in $\Omega = [0, 1]$. This shows that there is no $\omega \in \Omega$ for which $\lim_{n \rightarrow \infty} Y_n(\omega)$ exists. \square

PROOF that (Y_n) converges in probability:

If we write $|I_n|$ for the length of the interval I_n , then

- (h) $\square |I_n| = 1 \Leftrightarrow n = 1 \quad \square |I_n| = 1/2 \Leftrightarrow n = 2, 3 \quad \square |I_n| = 1/3 \Leftrightarrow n = 4, 5, 6$.
Thus, if $s_1 = 1$, $s_2 = s_1 + 2$, $s_3 = s_2 + 3, \dots, s_k = s_{k-1} + k = \sum_{j=1}^k j = \frac{k \cdot (k+1)}{2}, \dots$,
- (i) then $I_n = 1/k \Leftrightarrow n = s_{k-1} + 1, s_{k-1} + 2, \dots, s_{k-1} + k \Leftrightarrow s_{k-1} < n \leq s_k$.
- (j) It should be clear that $[n \rightarrow \infty] [k \rightarrow \infty]$ For a proof: $\square \Leftarrow$ follows from $n \geq k$.
 \square For the other direction, we observe that $n \stackrel{(i)}{\leq} 2s_k = 2k(k+1) < 2(k+1)^2$, i.e., $\sqrt{n/2} - 1 < k$. Thus, $[n \rightarrow \infty] \Rightarrow [k \rightarrow \infty]$ and \Rightarrow follows.
- (k) Since $Y_n(\omega) := \begin{cases} 1, & \text{if } \omega \in I_n, \\ 0, & \text{else} \end{cases}$ for $n \in \mathbb{N}$, we obtain $P\{|Y_n| \geq \varepsilon\} = 0$ for $\varepsilon \leq 1$ and, with n_k as defined in (k), $P\{|Y_{n_k}| \geq \varepsilon\} = \frac{1}{k}$ for $0 < \varepsilon \leq 1$. Thus, $P\{|Y_{n_k}| \geq \varepsilon\} \leq \frac{1}{k}$ for $\varepsilon > 0$.
- (l) Fix $\varepsilon > 0$ and $k \in \mathbb{N}$. $|I_n|$ and hence, $P\{|Y_n| > \varepsilon\}$ is nonincreasing with n . Thus, $n \geq n_k \Rightarrow P\{|Y_n| > \varepsilon\} \leq P\{|Y_{n_k}| > \varepsilon\} = \frac{1}{k}$. Since $[n \rightarrow \infty] \stackrel{(j)}{\Rightarrow} [k \rightarrow \infty]$, it follows that $\lim_{n \rightarrow \infty} P\{|Y_n| > \varepsilon\} = 0$ and this shows that $Y_n \xrightarrow{P} 0$. \blacksquare

\square

10.2 Two Laws of Large Numbers

Our knowledge of convergence in probability and almost surely enables us to understand the weak law and the strong law of large numbers. Recall that the “id” part of any iid sequence (Y_n) implies that $E[Y_1] = E[Y_2] = \dots$ and $Var[Y_1] = Var[Y_2] = \dots$.

Theorem 10.4 (Weak Law of Large Numbers).

Let Y_1, Y_2, \dots be an iid sequence of random variables on a probability space (Ω, P) with finite variance: $\sigma^2 := var[Y_n] < \infty$. Let $\mu := E[Y_n]$. Then,

$$(10.12) \quad \frac{Y_1 + Y_2 + \dots + Y_n}{n} \text{ converges in probability to } \mu, \text{ i.e.,}$$

$$[\varepsilon > 0] \Rightarrow \left[\lim_{n \rightarrow \infty} P \left\{ \left| \frac{1}{n} \sum_{j=1}^n Y_j - \mu \right| > \varepsilon \right\} = 0 \right]$$

PROOF: Let

$$\omega \mapsto \bar{Y}_n(\omega) := \frac{Y_1(\omega) + Y_2(\omega) + \dots + Y_n(\omega)}{n} = \frac{1}{n} \sum_{j=1}^n Y_j(\omega).$$

We have seen in Example 8.5 (Variance of the sample mean) on p.163, that

$$(A) \quad \mu_{\bar{Y}_n} = E[\bar{Y}_n] = \mu, \quad \text{and} \quad \sigma_{\bar{Y}_n}^2 = Var[\bar{Y}_n] = \frac{\sigma^2}{n}.$$

We apply Tchebysheff’s inequality 7.53 on p.140 with $k = \varepsilon\sqrt{n}/\sigma$ and obtain from (A), that

$$P \{ |\bar{Y}_n - \mu| > \varepsilon \} \leq \frac{1}{(n\varepsilon^2/\sigma^2)} = \frac{\sigma^2}{n\varepsilon^2} \rightarrow 0, \text{ as } n \rightarrow \infty$$

This proves that $P\text{-}\lim_{n \rightarrow \infty} \bar{Y}_n = \mu$. ■

Remark 10.1. We have previously encountered the random variable \bar{Y}_n under the name \bar{Y} , as the sample mean of a sample of size n . See Example 8.5 (Variance of the sample mean) on p.163.

It is considered bad form to use a subscript for the sample mean. We chose to do so in this section about the laws of large numbers anyway, since we are not dealing with this sample mean in the context of samples of a fixed size, but we are examining what happens as this size approaches infinity. □

Remark 10.2. We have learned in Theorem 10.1 (Relationship between the modes of convergence) on p.217, that almost sure convergence implies convergence in probability. One can interpret this in the following manner:

- It is harder to establish almost sure convergence, since it is a more powerful tool for proving that some mathematical property is true.
- Accordingly, it would be wonderful if one could strengthen a theorem that proves convergence in probability for some sequence of random variables, to show that this convergence actually happens almost surely.

- It turns out that this is possible for the weak law of large numbers (Theorem 10.4 on p.220. It is called the **weak** law of large numbers because there also is a **strong** law of large numbers which replaces the conclusion $P\text{-}\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n Y_j = \mu$ with $\text{a.s.}\text{-}\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n Y_j = \mu$. We will study that next. \square

Our knowledge of convergence in probability and almost surely enables us to understand the weak law and the strong law of large numbers. Recall that the “id” part of any iid sequence (Y_n) implies that $E[Y_1] = E[Y_2] = \dots$ and $\text{Var}[Y_1] = \text{Var}[Y_2] = \dots$.

Theorem 10.5 (Strong Law of Large Numbers).

Let Y_1, Y_2, \dots be an iid sequence of random variables on a probability space (Ω, P) .
Let $\mu := E[Y_n]$. Then,

$$(10.13) \quad \frac{Y_1 + Y_2 + \dots + Y_n}{n} \text{ converges almost surely to } \mu, \text{ i.e.,}$$

$$P \left\{ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n Y_j \neq \mu \right\} = 0.$$

PROOF:

Outside the scope of these lecture notes. \blacksquare

Example 10.3 (Infinite Monkey Theorem). A monkey has been granted eternal life. It is continually hitting at random the keys of a wordprocessor that will never break down.

The keyboard has a customized layout that makes it equally likely for each key, at any given key stroke, to be selected by the monkey. (For example, there is no CAPS key. Rather, there are separate keys for “a” and “A”, “b” and “B”,)

What is the probability that, in this infinite sequence of letters, there is a contiguous block that constitutes the collected work of William Shakespeare? We expect a flawless result: No typos, correct punctuation, CAPS exactly when required,?

Solution:

- There are K different keys that are being hit, at each stroke, with equal probability.
- Only one of them is correct at any given time and the others are failures.
- Thus, the sequence X_1, X_2, \dots of key strokes is a sequence of independent random items with constant success probability $p_j = p = 1/K$.
- We consider the indices $1, 2, 3, \dots$ as points in time, so X_{753} is the key that was hit at time $j = 753$.
- The author does not know how many letters Shakespeares collected work (“S-C-W”) consists of, but this certainly is a finite number. Let us denote it by N .

Let $Y_1 := 1$, if X_1, X_2, \dots, X_N form S-C-W. Let $Y_1 := 0$, else.

Let $Y_2 := 1$, if $X_{N+1}, X_{N+2}, \dots, X_{2N}$ form S-C-W. Let $Y_2 := 0$, else.

Let $Y_j := 1$, if $X_{(j-1)N+1}, X_{(j-1)N+2}, \dots, X_{jN}$ form S-C-W. Let $Y_j := 0$, else.

- If $i \neq j$, then Y_i and Y_j depend on “disjoint” chunks $(X_{(i-1)N+1}, X_{(i-1)N+2}, \dots, X_{iN})$ and $(X_{(j-1)N+1}, X_{(j-1)N+2}, \dots, X_{jN})$ of the independent X_k . Thus, Y_i and Y_j are independent.
- Also, both are $\text{binom}(1, (1/K)^N)$ (Bernoulli trials).
- Thus, $(Y_n)_n$ is an iid sequence with expectations $\mu = (1/K)^N$.
- By the strong law of large numbers, there is an event $A \subseteq \Omega$ such that $P(A) = 1$ and

$$\omega \in A \Rightarrow \lim_{n \rightarrow \infty} \sum_{j=1}^n Y_j(\omega) / n = \mu = \left(\frac{1}{K}\right)^N > 0.$$

- Since we divide the sum by n , the limit is zero if only finitely many $Y_j(\omega)$ are 1. Thus,

$$\omega \in A \Rightarrow Y_j(\omega) = 1, \text{ infinitely often!}$$

- Since $P(A) = 1$ and Y_i denotes the completion of the n th collection of Shakespeare’s works:
- With probability 1, the monkey will produce an infinite number of Shakespeare’s entire collection! \square

10.3 Sampling Distributions

Introduction 10.2. Back in Chapter 5.2 (Random Sampling and Urn Models With and Without Replacement), we gave Definition 5.2 (Sampling as a Random element) on p.88 of a sampling action.

- A sampling action of size n was nothing but a vector $\vec{X} = (X_1, X_2, \dots, X_n)$ of random elements. What makes it a sampling action is the interpretation of $\omega \mapsto X_j(\omega)$ as the j th pick of an item from a population of interest and the intent to use the outcomes $x_j = X_j(\omega)$ for inferences about that population.

These sample picks may happen with or without replacement. Sampling with replacement is desirable from a mathematical point of view, since we may consider the sample picks as having identical distribution. Thus,

$$F_{X_1}(x) = F_{X_2}(x) = \dots = F_{X_n}(x) \quad (x \in \mathbb{R});$$

This in turn implies that, if the sample picks are real-valued functions of ω i.e., they are random variables, they all have the same expectation and variance and

Moreover, nothing is assumed about the independence of the sample picks. To have it would be extremely desirable from a mathematical perspective. For example, if the X_j are jointly continuous random variables, knowledge of the marginal densities yields the joint density, because,

$$f_{\vec{X}}(\vec{x}) = f_{X_1}(x_1) \cdot f_{X_2}(x_2) \cdot \dots \cdot f_{X_n}(x_n) \quad (\vec{x} \in \mathbb{R}^n).$$

Unfortunately, identical distribution and independence are simplifications of the real world. This is even true when one considers n rolls of a die. ⁶² The surface on which the die is rolled is not perfectly even, so that negates identical distribution. If several people take turns, then the different ways in which they throw the die creates a dependency. Of course, it is very likely that those differences, if we are able to detect them, are so minuscule that they can be ignored.

But there are many examples where those deviations are so large that we cannot work under the iid assumption. This need not necessarily occur in a real world application. It can also be part of the probabilistic models we create: Whenever we assume that we sample without replacement from a finite population, the probabilistic makeup of the items remaining in that population changes with every item we happen to pick for our sample.

Consider sampling at random from an urn that initially contains R red and $N - R$ black balls. If X_j is red, then there will be less of a probability of X_{j+1} being red, than if X_j was black. Hence, the X_j are neither independent, nor identically distributed.

However, those sample picks constitute a simple random sample action according to Definition 5.3 (Simple Random Sample) on p.89:

- A sampling action $\vec{X} = (X_1, X_2, \dots, X_n)$ of size n from a population of size $N \geq n$ is called a simple random sampling action (SRS action), if it is done without replacement and if each one of the potential outcomes $\vec{x} = \vec{X}(\omega)$ has equal chance of being selected.

If the sample size of an SRS action is large, but small when compared to the size of the population, then treating it as iid will result in insignificant computational differences. ⁶³ This observation is one of the reasons that even the more restrictive definition of an SRS action is of a generality we are not looking for in this chapter. We follow [5] Hogg, McKean, Craig: Introduction to Mathematical Statistics.

A typical statistical problem can be described as follows: We have a random variable Y that we know about, but we do not know its distribution, given by its CDF $F_Y(y)$.

Our insufficient knowledge of Y can manifest itself in two different ways:

- (I) We know the type of distribution, but not all of its parameters. For example, we may know that Y is normal with $\sigma^2 = 3.65$, but its mean μ is unknown.
- (II) We do not even know the type of distribution: Does Y follow a Poisson distribution or is it normal or exponential or?

We deal in this section with problem (I). \square

Example 10.4. Some more problem (I) examples are the following:

- (a) $Y \sim \text{binom}(64, p)$, with unknown success probability p . We write $p_Y(y; p)$ for the PMF to make explicit the role of the unknown parameter, p .
- (b) $Y \sim \mathcal{N}(\mu, \sigma^2)$, where both μ and σ^2 are unknown. We write $f_Y(y; \mu, \sigma)$ for the PDF to make explicit the role of the unknown parameters, μ and σ .
- (c) $Y \sim \text{expon}(\beta)$, with unknown β . We write $f_Y(y; \beta)$ for the PDF.
- (d) $Y \sim \text{gamma}(\alpha, 3)$, with unknown α . We write $f_Y(y; \alpha)$ for the PDF. \square

⁶²Interpret X_j as the j th pick from the population of all rolls of that die.

⁶³We mentioned this in Remark 5.2 on p.88.

Remark 10.3. The examples just given suggest now to handle the general case. Since the random variable Y is given and we know its distribution except for one or several parameters, we know its PMF $p_Y(y)$ in the discrete case or PDF $f_Y(y)$ in the continuous case. It is customary to write θ or $\vec{\theta}$ for the unknown parameter or **parameters of the distribution** and to write Θ for the **parameter space**, i.e., the set of all parameters we consider for the problem.⁶⁴

Thus, in Example 10.3(a), $\Theta = [0, 1]$. In Example 10.3(b), $\Theta =]-\infty, \infty[\times]0, \infty[$.

Problem (I) can now be formulated as follows:

- Given is a random variable Y of which we know its distribution except for one or several parameters.
 - We know the PMF $p_Y(y; \theta)$ if Y is discrete. □ We know the PDF $f_Y(y; \theta)$ if Y is continuous.
- How can we find a good, possibly optimal, procedure to estimate θ from the sample? that we have drawn or intend to draw from the population?

It seems obvious enough, that this estimate must be a function

$$\theta = T(\vec{y}) = T(y_1, \dots, y_n) = T(\vec{Y}(\omega)) = T(Y_1(\omega), \dots, Y_n(\omega)).$$

In the context of a sampling action, we refer to the specific list of numbers, $\vec{y} = (y_1, \dots, y_n)$, as the values or **realizations**, of the sampling action. □

We had stated in the introduction that we will restrict the scope of the sampling actions in this section to the iid case.

Definition 10.2 (Random sampling action from a distribution).

Let Y be a random variable on a probability space (Ω, P) . We call a vector $\vec{Y} = (Y_1, \dots, Y_n)$ a **random sampling action from the distribution of Y** , or also, a **random sampling action on Y** , if

- each Y_j has the same distribution as Y
- the random variables Y_1, \dots, Y_n are iid. □

That definition allows us to restate the essence of Remark 10.3 as follows: We expect a procedure to estimate the parameter θ of a PMF $p_Y(y; \theta)$ or PDF $f_Y(y; \theta)$ to be a random variable $\omega \mapsto T(\vec{Y}(\omega))$.

There is a special name for transforms $\vec{y} \mapsto T(\vec{y})$ of a random sampling action on Y .

Definition 10.3 (Statistic).

⁶⁴It is unfortunate that this standard notation for parameters to be estimated is at odds with the other standard which uses the CAPS version of a letter to denote a random item and the corresponding small letter to denote an outcome of this random element. (For example, $y = Y(\omega)$).

Let Y be a random variable on a probability space (Ω, P) and $\vec{Y} = (Y_1, \dots, Y_n)$ a random sampling action on Y . Let

$$T : \mathbb{R}^n \mapsto \mathbb{R}; \quad \vec{y} \mapsto T(\vec{y})$$

be some function that can be applied to the sampling action \vec{Y} . We call the random variable

$$\omega \mapsto T(\vec{Y}(\omega))$$

a **statistic** of that sampling action. We call the distribution of that random variable,

$$(10.14) \quad B \mapsto P_{T \circ \vec{Y}}(B) = P\{T(\vec{Y}) \in (B)\} = P\{\omega \in \Omega : T(\vec{Y}(\omega)) \in B\}$$

its **sampling distribution**. Once the sampling action has been performed and the corresponding realization $\vec{y} = \vec{Y}(\omega)$ has been obtained, we call $t = T(\vec{Y}(\omega))$ the realization of the statistic. \square

Theorem 10.6.

Let Y be a random variable on a probability space (Ω, P) and $\vec{Y} = (Y_1, \dots, Y_n)$ a random sampling action on Y . Let $T_1, T_2, \dots, T_k : \mathbb{R}^n \mapsto \mathbb{R}$ be statistics for that sample action. Let

$$T^* : \mathbb{R}^k \mapsto \mathbb{R}; \quad (t_1, \dots, t_k) \mapsto T^*(t_1, \dots, t_k).$$

Then, setting $\vec{t} = (t_1, \dots, t_k)$ and $\vec{T} = (T_1, \dots, T_k)$, the composition

$$T^* \circ \vec{T} \circ \vec{Y} : \omega \mapsto T^*(\vec{T}[\vec{Y}(\omega)]) = T^*(T_1[\vec{Y}(\omega)], \dots, T_k[\vec{Y}(\omega)])$$

also is a statistic of \vec{Y} .

PROOF:

Left as an exercise which is very easy for the student who has had exposure to functions $\mathbb{R}^n \rightarrow \mathbb{R}^k$ with dimensions n and/or k that can exceed the value 3. \blacksquare

The last theorem can be stated succinctly and without mathematical symbols as follows:

A function of a function of the data is a function of the data.

Here is an example of a statistic which is so important that it deserves its own definition. It also is used to illustrate Theorem 10.6.

Definition 10.4 (Sample variance).

Let $\vec{Y} = (Y_1, \dots, Y_n)$ be a random sample action on a random variable Y .

The **sample variance** is defined as the random variable

$$(10.15) \quad \omega \mapsto S^2(\omega) := \frac{1}{n-1} \sum_{j=1}^n (Y_j(\omega) - \bar{Y}(\omega))^2.$$

We further call $\omega \mapsto S(\omega) := \sqrt{S^2(\omega)}$ the **sample standard deviation**.

We will often write s^2 and s for the realizations $S^2(\omega)$ and $S(\omega)$ that result from creating the sample. \square

Example 10.5. For the following examples assume that $\vec{Y} = (Y_1, \dots, Y_n)$ is a random sample action on a random variable Y .

- (a) In Example 8.5 (Variance of the sample mean) on p.163, we considered the sample mean

$$\omega \mapsto \bar{Y}(\omega) = \frac{1}{n} \sum_{j=1}^n Y_j(\omega). \quad \bar{Y} \text{ is a statistic: The transform is } T(\vec{Y}) = \frac{1}{n} \sum_{j=1}^n Y_j.$$

We also mentioned that this statistic is an obvious choice for estimating the parameter $\mu = E[Y]$ of the underlying random variable Y .

- (b) Sample variance S^2 and sample standard deviation S which were defined above are statistics. This can be shown with the help of Theorem 10.6 on p.225 as follows. Let

$$t_1 = T_1(\vec{y}) = y_1, t_2 = T_2(\vec{y}) = y_2, \dots, t_n = T_n(\vec{y}) = y_n, t_{n+1} = T_{n+1}(\vec{y}) = \bar{y}.$$

$$T^*(t_1, \dots, t_n, t_{n+1}) = \frac{1}{n-1} \sum_{j=1}^n (t_j - t_{n+1})^2$$

Then $S^2 = T^*(T_1(\vec{Y}), \dots, T_n(\vec{Y}), T_{n+1}(\vec{Y}))$. By Theorem 10.6, S^2 is a statistic for the random sampling action \vec{Y} . We apply this theorem again to the function $T^{**} : t^* \mapsto \sqrt{t^*}$ and obtain that the standard deviation S is a statistic, since $S = T^{**}(S^2)$.

- (c) The j th order statistic, $Y_{(j)}$ is indeed a statistic, since knowledge of all values of a list y_1, \dots, y_n of real numbers uniquely determines which one is the j th largest value in that list.
- (d) The **sample range**, $R = Y_{(n)} - Y_{(1)}$, is a statistic, since it is a function (the difference) of the two statistics $Y_{(n)}$ and $Y_{(1)}$. \square

Example 10.6 (WMS Ch.07.1, Example 7.1). Example 7.1 of the WMS text discusses in quite big detail the sampling distribution of the statistic \bar{Y} for a sample of three independent rolls of a balanced die. You are strongly encouraged to study it. \square

Theorem 10.7 (WMS Ch.07.2, Theorem 7.1). ()

Let Y_1, Y_2, \dots, Y_n be a random sampling action of size n from a normal distribution with mean μ and variance σ^2 , i.e., we sample on a random variable $Y \sim \mathcal{N}(\mu, \sigma^2)$. Then the sample mean \bar{Y} follows a normal distribution with mean μ and variance σ^2/n .

PROOF: That is an immediate consequence of Theorem 9.5 (Linear combinations of uncorrelated normal variables are normal) on p.209. ■

Theorem 10.8 (WMS Ch.07.2, Theorem 7.2).

Let $\vec{Y} = (Y_1, \dots, Y_n)$ be a random sampling action on $Y \sim \mathcal{N}(\mu, \sigma^2)$. Let $Z_j = (Y_j - \mu)/\sigma$ for $j = 1, 2, \dots, n$. Then $\vec{Z} = (Z_1, \dots, Z_n)$ is a random sampling action on a standard normal variable. (In particular, the Z_j are iid.) Further,

$$(10.16) \quad \sum_{j=1}^n Z_j^2 = \sum_{j=1}^n \left(\frac{Y_j - \mu}{\sigma} \right)^2$$

follows a χ^2 distribution with n degrees of freedom.

PROOF: It follows from Theorem 9.5 (Linear combinations of uncorrelated normal variables are normal) on p.209 that the linear combination $Z_j = (Y_j - \mu)/\sigma$ is standard normal. It follows from Theorem 9.4 (MGF of a sum of functions of independent variables) on p.209 that the Z_j are iid. It follows from Theorem 9.6 on p.211 that $\sum_{j=1}^n Z_j^2 \sim \chi^2(\text{df} = n)$. ■

The following is Example Example 6.13 of the WMS text.

Proposition 10.1. ★

Let Y_1 and Y_2 be independent standard normal random variables. Then $Y_1 + Y_2$ and $Y_1 - Y_2$ are independent and normally distributed, both with mean 0 and variance 2.

PROOF: See WMS Ch.06.6, Example 6.13. ■

Theorem 10.9 (Independence of sample mean and sample variance in normal populations).

Let $\vec{Y} = (Y_1, \dots, Y_n)$ be a random sample action on a $\mathcal{N}(\mu, \sigma^2)$ random variable Y . Then, $\vec{Z} = (Z_1, \dots, Z_n)$ is a random sample action on a standard normal variable Z . Further,

$$(a) \quad \frac{(n-1)S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{j=1}^n (Y_j - \bar{Y})^2 \sim \chi^2(\text{df} = (n-1))$$

(b) \bar{Y} and S^2 are independent random variables.

PROOF: ★ See the proof of WMS Ch.07.2, Theorem 7.3 for the case $n = 2$. ■

- The sample mean \bar{Y} was a natural choice to estimate the mean $\mu = E[Y]$ of a random variable X .
- It seems just as natural to use the sample variance S^2 to estimate $\sigma^2 = \text{Var}[Y]$. We will see that, if Y follows a normal distribution, this choice turns out to be mathematically sound.

The t distribution which we define next is a means towards that end.

Definition 10.5 (Student’s t -distribution ⁶⁵).

Let Z and W be independent random variables such that Z is standard normal and W is χ^2 with ν df. Let

$$(10.17) \quad T = \frac{Z}{\sqrt{W/\nu}}$$

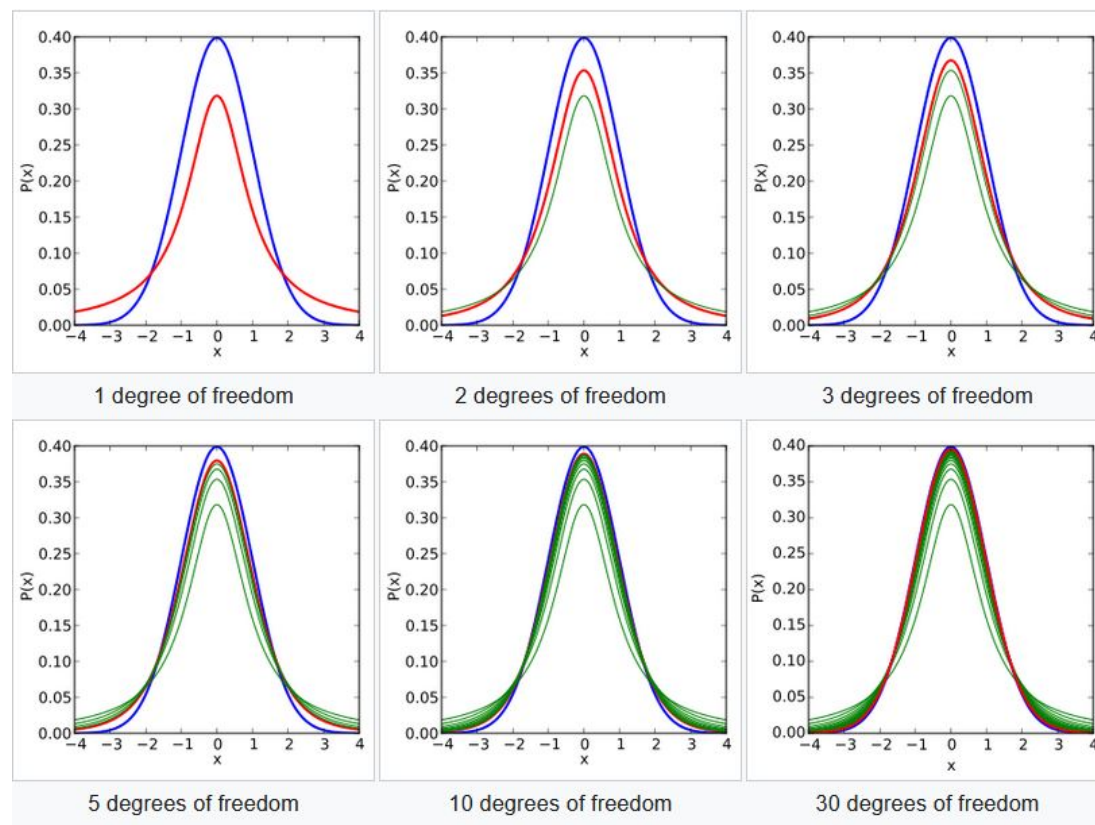
Then we refer to the distribution P_T of T as a **t -distribution** or **Student’s t -distribution** with ν df. We also write that as $T \sim t(\nu)$ or $T \sim t(\text{df} = \nu)$. \square

Remark 10.4.

- One can prove that $E[T] = 0$ for any ν , and $Var[T] = \frac{\nu}{\nu - 2}$ for $\nu > 2$.

The density of the t -distribution looks very similar to that of a normal density. Both have a symmetrical, bell shaped graph. But note the following:

- Since it does not depend on ν , $E[T] = 0$ is not a parameter of the t -distribution.
- Since $\frac{\nu}{\nu - 2} > 1$, the tails are fatter than those of a $\mathcal{N}(0, 1)$ variable. See Figure 10.1. \square



10.1 (Figure). **densities of the standard normal and t distribution.** Source: [Wikipedia](#).

⁶⁵Named after the English statistician William S. Gosset (1876 – 1937). Georg Ferdinand Ludwig Philipp Cantor (1845 – 1918), Gosset was Head Brewer of the Guinness Brewery in Dublin, Ireland and published his papers under the pseudonym "Student".

Theorem 10.10.

Let $Y \sim \mathcal{N}(\mu, \sigma^2)$ and $\vec{Y} = (Y_1, \dots, Y_n)$ be a random sample action on Y . Let

$$(10.18) \quad T := \frac{\bar{Y} - \mu}{S/\sqrt{n}}.$$

Then T follows a t -distribution with $(n - 1)$ df.

PROOF: Let

$$(A) \quad Z := \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \quad \text{and} \quad W := \frac{(n-1)S^2}{\sigma^2}.$$

We have seen that $Z \sim \mathcal{N}(0, 1)$ and $W \sim \chi^2(\text{df} = n - 1)$. Since \bar{Y} and S^2 are independent by Theorem 10.9 on p.227, Z as a function of \bar{Y} only and W as a function of S^2 only also are independent. Thus,

$$T = \frac{Z}{\sqrt{W/(n-1)}} \stackrel{(A)}{=} \frac{(\bar{Y} - \mu)/(\sigma/\sqrt{n})}{[\sqrt{(n-1)S^2/\sigma^2}]/(n-1)} = \left(\frac{\bar{Y} - \mu}{S\sqrt{n}} \right)$$

has at distribution with $(n - 1)$ df. ■

Example 10.7 (WMS Ch.07.2, Example 7.6). Example 7.6 of the WMS text discusses a practical example of the Student's t -distribution that discusses how to estimate the unknown variance of a normal random variable from a sample. You are strongly encouraged to study it. □

The next and last distribution tied to random sampling on a normal variable that we give in this section allows us to compare the variances of two random sampling actions on normal random variables that represent two independent populations. This is used in the so called analysis of variance (ANOVA) to decide whether the means of several independent normal populations all coincide or whether at least two of them are different.

Definition 10.6 (F -distribution).

Given are two independent random variables $W_1 \sim \chi^2(\text{df} = \nu_1)$ and $W_2 \sim \chi^2(\text{df} = \nu_2)$, with ν_1 and ν_2 df, respectively. Then we say that

$$F = \frac{W_1/\nu_1}{W_2/\nu_2}$$

follows an **F distribution** with ν_1 **numerator degrees of freedom** and ν_2 **denominator degrees of freedom**. □

Remark 10.5. ★ One can show that

- $\nu_2 > 2 \Rightarrow E[F] = \frac{\nu_2}{\nu_2 - 2}$,
- $\nu_2 > 4 \Rightarrow \text{Var}[F] = \frac{2\nu_2^2(\nu_1 + \nu_2 - 2)}{\nu_1(\nu_2 - 2)^2(\nu_2 - 4)}$. □

Theorem 10.11.

Consider two random sampling actions of sizes n_1 and n_2 on random variables $Y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $Y_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ from two independent populations, with sample variances S_1^2 and S_2^2 . Let

$$(10.19) \quad F := \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2}.$$

Then F follows an F distribution with $(n_1 - 1)$ numerator df and $(n_2 - 1)$ denominator df.

PROOF: Let

$$W_1 := \frac{(n_1 - 1)S_1^2}{\sigma_1^2}, \quad W_2 := \frac{(n_2 - 1)S_2^2}{\sigma_2^2}.$$

Since the random sampling actions are independent, so are their sample variances S_1^2 and S_2^2 , and so are the transforms W_1 of S_1^2 and W_2 of S_2^2 . By Definition 10.6 of an F distribution,

$$\frac{W_1/\nu_1}{W_2/\nu_2} = \frac{[(n_1 - 1)S_1^2/\sigma_1^2]/[(n_1 - 1)]}{[(n_2 - 1)S_2^2/\sigma_2^2]/[(n_2 - 1)]} = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2}$$

follows an F distribution with $(n_1 - 1)$ numerator df and $(n_2 - 1)$ denominator df. ■

Example 10.8 (WMS Ch.07.2, Example 7.7). Example 7.6 of the WMS text discusses another practical example of the Student's F distribution. You are strongly encouraged to study it. □

10.4 The Central Limit Theorem

Introduction 10.3. In section 10.3 (Sampling Distributions) we were able to determine the sampling distributions of some very important statistics that can be computed from the realization of a random sample action \vec{Y} on some random variable Y . But there was very restrictive assumption on that underlying random variable

- Y had to follow a normal distribution.

We will find a solution for determining the sampling distribution of the sample mean, $\bar{Y} = \frac{1}{n} \sum_{j=1}^n Y_j$, even if Y is not normal.

- It is an **asymptotic solution**, i.e., its comes in form of a $U = \lim_{n \rightarrow \infty} U_n$ theorem.
- Here, U_n is a statistic $T_n \circ \vec{Y}$, which we can compute from (the realization of) \vec{Y} and $\bar{Y}_n := \frac{1}{n} \sum_{j=1}^n Y_j$, a very natural approximation of \bar{Y} , can also be computed from U_n
- n denotes the sample size. Thus, the sample must be sufficiently large to allow us to ignore the discrepancy between U_n and U .

We have learned that there are four different kinds of limits which occur in connection with a sequence of random variables. We will discuss in this chapter the central limit theorem. It allows us to show the existence of the least desirable of those four limits, the limit in distribution. But that is not as bad as it sounds for the following reason.

- For large enough n , the CDF of U_n is close to that of U . Since the CDF determines the probabilities of all important events B , we can approximate $P\{U_n \in B\} \approx P\{U \in B\}$, □

We will state and prove the limit theorem which was mentioned in the introduction above, after the following important theorem that relates convergence in distribution, $Y_n \xrightarrow{D} Y$, to (pointwise) convergence, $m_{Y_n}(t) \rightarrow m_Y(t)$ of the associated MGFs.

Theorem 10.12 (Lévy–Cramér continuity theorem). ★

Let Y_1, Y_2, \dots be a sequence of random variables (iid is not assumed) with associated CDFs F_{Y_1}, F_{Y_2}, \dots and MGFs $m_{Y_1}(t), m_{Y_2}(t), \dots$.

Let Y be a random variable with associated CDF F_Y and MGF $m_Y(t)$. Then,

$$(10.20) \quad \begin{aligned} & [m_{Y_n}(t) \rightarrow m_Y(t) \text{ as } n \rightarrow \infty, \text{ for all } t \in \mathbb{R}] \\ \Rightarrow & [F_{Y_n}(y) \rightarrow F_Y(y) \text{ as } n \rightarrow \infty, \text{ for all } y \text{ where } F_Y(\cdot) \text{ is continuous.}] \end{aligned}$$

PROOF: Outside the scope of this course. ■

Theorem 10.13 (Central Limit Theorem).

Central Limit Theorem:

Let $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$ be a vector of iid random variables with common expectation $E[Y_j] = \mu$ and finite variance $Var[Y_j] = \sigma^2$. Let Z be a standard normal variable and

$$U_n := \frac{\sum_{j=1}^n Y_j - n\mu}{\sigma \cdot \sqrt{n}} = \frac{\bar{Y}_n - \mu}{\sigma/\sqrt{n}}, \quad \text{where } n \in \mathbb{N}, \bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i.$$

Then, U_n converges to Z in distribution as $n \rightarrow \infty$. In other words,

$$\lim_{n \rightarrow \infty} P\{U_n \leq u\} = P\{Z \leq u\} = \int_{-\infty}^u \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \quad \text{for all } u.$$

PROOF:

- (1) Let $\tilde{Y}_n := Y_n - \mu$. The \tilde{Y}_n are iid, with $E[\tilde{Y}_j] = 0$, $Var[\tilde{Y}_j] = \sigma^2$ and MGF $m(t) := m_{\tilde{Y}_n}(t)$. By Corollary 9.1 on p.209, $m_{\tilde{Y}_1 + \dots + \tilde{Y}_n}(t) = [m(t)]^n$. Thus.

$$(2) \quad m_{U_n}(t) = E \left[\exp \left\{ \sum_{j=1}^n \tilde{Y}_j \cdot \frac{t}{\sigma\sqrt{n}} \right\} \right] = m_{\tilde{Y}_1 + \dots + \tilde{Y}_n} \left(\frac{t}{\sigma\sqrt{n}} \right) = \left[m \left(\frac{t}{\sigma\sqrt{n}} \right) \right]^n.$$

- (3) According to Theorem 10.12 (Lévy–Cramér continuity theorem), it suffices to show that

$$\lim_{n \rightarrow \infty} m_{U_n}(t) = m_Z(t) = e^{-t^2/2}.$$

Equivalently, since $x \mapsto e^x$ is continuous, it suffices to show that

$$(4) \quad \lim_{n \rightarrow \infty} \ln m_{U_n}(t) = \frac{t^2}{2}.$$

- (5) Let $h := \frac{t}{\sigma\sqrt{n}}$. Then $n = \frac{t^2}{\sigma^2 h^2}$. Thus, by (2),

$$\ln m_{U_n}(t) = n \ln m(h) = \frac{t^2}{\sigma^2 h^2} \ln m(h) = \frac{t^2}{\sigma^2} \left(\frac{\ln m(h)}{h^2} \right).$$

Thus,

$$(6) \quad \lim_{n \rightarrow \infty} \ln m_{U_n}(t) = \frac{t^2}{\sigma^2} \lim_{h \rightarrow 0} \frac{\ln m(h)}{h^2}.$$

Since $m(0) = e^0 = 1$, the right-hand limit is of the form $0/0$. We use L'Hôpital's rule ⁶⁶ twice in a row and obtain, since $m(t) = m_{\tilde{Y}_n}(t)$ and hence, $m''(0) = E[\tilde{Y}_n^2]$,

$$(7) \quad \begin{aligned} \lim_{h \rightarrow 0} \frac{\ln m(h)}{h^2} &= \lim_{h \rightarrow 0} \frac{[1/m(h)] m'(h)}{2h} = \lim_{h \rightarrow 0} \frac{m'(h)}{2hm(h)} \\ &= \lim_{h \rightarrow 0} \frac{m''(h)}{2m(h) + 2hm'(h)} = \frac{m''(0)}{2m(0) + 0} = \frac{m''_{\tilde{Y}_n}(0)}{2} = \frac{E[\tilde{Y}_n^2]}{2}. \end{aligned}$$

(8) Since $\tilde{Y}_n = Y_n - \mu$ and hence, $E[\tilde{Y}_n^2] = E[(Y_n - \mu)^2] = \text{Var}[Y_n] = \sigma^2$, (7) implies

$$\lim_{h \rightarrow 0} \frac{\ln m(h)}{h^2} = \frac{\sigma^2}{2}.$$

$$\text{Thus, by (6), } \lim_{n \rightarrow \infty} \ln m_{U_n}(t) = \frac{t^2}{\sigma^2} \cdot \frac{\sigma^2}{2} = \frac{t^2}{2}.$$

We have shown (4) and this finishes the proof. ■

Remark 10.6. In statistical applications the CLT often is employed as follows: Carefully designed statistical techniques have resulted in the estimate $\mu = \mu_0$ for μ , the unknown mean of the population of interest. But this has been quite some time ago. Today there is reason to believe that this value is now outdated and one wants to obtain supporting evidence for that claim.

- We make $\mu = \mu_0$ our working hypothesis.
- An SRS \vec{Y} of size n is taken and $c_0 := \frac{\bar{y} - \mu_0}{\sigma/\sqrt{n}}$ is computed from the sample mean realization $\bar{y} = \sum_{j=1}^n y_j$ which one obtains from the realization $\vec{y} = \vec{Y}(\omega)$ of the sample.
- If $\vec{Y}(\omega)$ is far away from μ_0 , then $\alpha_0 := P\left\{\left|\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}}\right| > c_0\right\}$ will be very small.

For example, assume that $c_0 = 3$, i.e., $|\bar{y} - \mu_0| = 3 \cdot (\sigma/\sqrt{n})$. The r.v. $\omega \rightarrow \bar{Y}(\omega)$ satisfies

$$E[\bar{Y}] = E[Y] = \mu = \mu_0 \quad \text{and} \quad \text{Var}[\bar{Y}] = \frac{\text{Var}[Y]}{n} = \frac{\sigma^2}{n}, \quad \text{i.e.,} \quad \frac{\sigma}{\sqrt{n}} = \text{SD}(\bar{Y}),$$

Thus, $c_0 = 3$ signifies that this r.v. is three SDs away from its mean. According to the CLT, $\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}}$ is approximately standard normal and we can employ the the 68%–95%–99.7% rule for the normal distribution (the empirical rule). It tells us that only about 0.3% of the probability is outside the ± 3 SD range:

$$\alpha_0 \approx 1 - 0.997 = 0.003.$$

That is the probability that a \bar{Y} belonging to a random sample like ours (with the same sample size) is 3 SDs or more away from μ_0 .

⁶⁶in the form $\lim_{x \rightarrow 0} \frac{f(x)}{g(x)} = \lim_{x \rightarrow 0} \frac{f'(x)}{g'(x)}$

- So it was just the luck of the draw that let us obtain a sample that only has a chance of one in 333 of being picked. **Or is there another explanation?**

How about this? $\alpha_0 = 0.05$ was obtained contingent on the hypothesis that μ still equals μ_0 . Let us change our point of view and assume that there was nothing unusual about our sample.

- We reject the hypothesis $\mu = \mu_0$, since the data obtained from the sample suggest that $|\bar{Y} - \mu| < |\bar{Y} - \mu_0|$ and that necessitates $\mu \neq \mu_0$.
- In the extreme, we could dispense with any effort to find a well founded estimate of μ . Instead, we act as if our particular sample serves that purpose and replace μ_0 with $\mu_1 := \bar{y}$.

In the extreme, we could dispense with any effort to find a well founded estimate of μ . Instead, we act as if our particular sample serves that purpose and replace μ_0 with $\mu_1 := \bar{y}$. But of course, that generally is not a good idea and one should follow the established process to obtain a new estimate of μ . \square

Remark 10.7. This is a continuation of the previous example.

- The procedure outlined there to decide whether or not to reject the hypothesis $\mu = \mu_0$ involved the computation of the expression $c_0 := \frac{\bar{y} - \mu_0}{\sigma/\sqrt{n}}$.
- However, knowledge of the population variance $\sigma^2 = \text{Var}[Y_j]$ of a sample pick Y_j from that population is the exception rather than the rule and σ^2 must be estimated from the sample. The obvious way of doing so is use of the sample variance realization $s^2 = S^2(\omega)$.
- We have the following problem. The CLT asserts that, for large enough n , $\omega \mapsto \frac{\bar{Y}(\omega) - \mu_0}{\sigma/\sqrt{n}}$ is approximately standard normal. We used that fact to compute $P\left\{\left|\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}}\right| > c_0\right\}$ and we based the decision to reject or not reject the hypothesis $\mu = \mu_0$ on that number.
- But what happens if we replace σ with $S(\omega)$? If the random variable $\omega \mapsto \frac{\bar{Y}(\omega) - \mu_0}{S(\omega)/\sqrt{n}}$ also is approximately standard normal for large n , then our problem is solved. \square

To show that the CLT indeed remains in force if σ^2 is replaced by S^2 , we must collect some material.

Theorem 10.14 (Student t converges to normal distribution).

Let T_1, T_2, \dots be a sequence of random variables such that $T_j \sim t(df = j)$. Then T_j converges in distribution to a standard normal variable.

PROOF: Omitted. ⁶⁷ Note though that the graphs of the t -PDFs shown in Remark 10.4 on p.228 visually support the assertion of this theorem.

Lemma 10.1. ★ Let $\vec{y} := (y_1, \dots, y_n) \in \mathbb{R}^n$, ($n \in \mathbb{N}$), and $\bar{y} := \sum_{j=1}^{\infty} y_j$ the arithmetic mean of \vec{y} . Then,

$$(a) \quad \sum_{j=1}^n (y_j - c)^2 = \sum_{j=1}^n (y_j - \bar{y})^2 + \sum_{j=1}^n (\bar{y} - c)^2,$$

⁶⁷A proof can be found at this [StackExchange](#) link.

(b) \bar{y} minimizes the expression $\sum_{j=1}^n (y_j - c)^2$, where $c \in \mathbb{R}$:

$$\sum_{j=1}^n (y_j - c)^2 \geq \sum_{j=1}^n (y_j - \bar{y})^2 \quad \text{for all } c \in \mathbb{R},$$

PROOF: To show (a), we observe that

$$\begin{aligned} (10.21) \quad \sum_{j=1}^n (y_j - \bar{y})(\bar{y} - c) &= \bar{y} \sum_{j=1}^n y_j + \bar{y} \cdot c \sum_{j=1}^n 1 - c \sum_{j=1}^n y_j - \bar{y} \cdot \bar{y} \sum_{j=1}^n 1 \\ &= \bar{y}(n\bar{y}) + (\bar{y}c)n - c(n\bar{y}) - (\bar{y}^2)n = 0. \end{aligned}$$

Hence,

$$\begin{aligned} \sum_{j=1}^n (y_j - c)^2 &= \sum_{j=1}^n (y_j - \bar{y} + \bar{y} - c)^2 \\ &= \sum_{j=1}^n (y_j - \bar{y})^2 + 2 \sum_{j=1}^n (y_j - \bar{y})(\bar{y} - c) + \sum_{j=1}^n (\bar{y} - c)^2 \\ &\stackrel{(10.21)}{=} \sum_{j=1}^n (y_j - \bar{y})^2 + \sum_{j=1}^n (\bar{y} - c)^2. \end{aligned}$$

This proves (a). Clearly, the last expression is minimal when the right-hand summation term vanishes, i.e., when $\bar{y} = c$. This proves (b). ■

Corollary 10.1. ★

The sample variance $S^2 = \frac{1}{n-1} \sum_{j=1}^n (Y_j - \bar{Y})^2$ of any sample $\vec{Y} := (Y_1, \dots, Y_n)$, ($n \in \mathbb{N}$), satisfies

$$(n-1)S^2 = \sum_{j=1}^n Y_j^2 - n\bar{Y}^2.$$

PROOF: We apply formula (a) of Lemma 10.1 with $c = 0$ and obtain

$$\sum_{j=1}^n Y_j^2 = \sum_{j=1}^n (Y_j - \bar{Y})^2 + \sum_{j=1}^n \bar{Y}^2 = \sum_{j=1}^n (Y_j - \bar{Y})^2 + n \cdot \bar{Y}^2.$$

Thus,

$$(n-1)S^2 = \sum_{j=1}^n (Y_j - \bar{Y})^2 = \sum_{j=1}^n Y_j^2 - n\bar{Y}^2. \quad \blacksquare$$

Theorem 10.15 (Sample variance converges to population variance).

Let $\vec{Y} := (Y_1, \dots, Y_n) \in \mathbb{R}^n$, ($n \in \mathbb{N}$), be a random sampling action from the distribution of a random variable Y with finite variance $\sigma^2 < \infty$.

Then the sample variance $S_n^2 = \frac{1}{n-1} \sum_{j=1}^n (Y_j - \bar{Y})^2$ converges a.s. (hence, also in probability and in distribution) to σ^2 .

PROOF: ★ Let $U_n := \frac{n-1}{n} S_n^2$ and $\bar{Y}_n := \bar{Y} = \frac{1}{n} \sum_{j=1}^n Y_j$. By Corollary 10.1,

$$(A) \quad U_n = \frac{1}{n} \sum_{j=1}^n Y_j^2 - \bar{Y}_n^2.$$

Since the sample picks Y_j are iid, so are their squares. Note that

$$E[Y_j^2] = \text{Var}[Y_j] + (E[Y_j])^2 = \sigma^2 + \mu^2$$

We apply the Strong Law of Large Numbers to the iid sequences Y_j^2 and Y_j and obtain

$$(B) \quad \text{a.s.-} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n Y_j^2 = \sigma^2 + \mu^2, \quad \text{a.s.-} \lim_{n \rightarrow \infty} \bar{Y}_n = \mu.$$

Next, we apply Theorem 10.15 (Sample variance converges to population variance) on p.234 to the continuous function $x \mapsto x^2$. It follows from $\text{a.s.-} \lim_{n \rightarrow \infty} \bar{Y}_n = \mu$ obtained in (B), that

$$(C) \quad \text{a.s.-} \lim_{n \rightarrow \infty} \bar{Y}_n^2 = \mu^2.$$

It now follows from the definition of U_n and from (A) and (B) and (C), that

$$\text{a.s.-} \lim_{n \rightarrow \infty} S_n^2 = \text{a.s.-} \lim_{n \rightarrow \infty} \frac{n-1}{n} S_n^2 = \text{a.s.-} \lim_{n \rightarrow \infty} U_n = (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

It follows from Theorem 10.1 (Relationship between the modes of convergence) on p.217 that convergence $S_n^2 \rightarrow \sigma^2$ also takes place in probability and in distribution. ■

We now are able to provide a version of the CLT which allows us to work with $\omega \mapsto \frac{\bar{Y}(\omega) - \mu_0}{S(\omega)/\sqrt{n}}$ instead of $\omega \mapsto \frac{\bar{Y}(\omega) - \mu_0}{\sigma/\sqrt{n}}$ and solves the issue brought up in Remark 10.7 on p.233.

Theorem 10.16 (CLT – Sample variance version).

Let $\vec{Y} = (Y_1, Y_2, \dots, Y_n)$ be a vector of iid random variables with common expectation $E[Y_j] = \mu$ and finite variance $\text{Var}[Y_j] = \sigma^2$. Let Z be a standard normal variable. For $n \in \mathbb{N}$, let

$$\bar{Y}_n := \frac{1}{n} \sum_{i=1}^n Y_i, \quad S_n^2 := \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2, \quad S_n := \sqrt{S_n^2}, \quad W_n := \frac{\bar{Y}_n - \mu}{S_n/\sqrt{n}}.$$

(Thus, \bar{Y}_n and S_n are sample mean and sample standard deviation of the RSA \vec{Y}).

Then W_n converges to Z in distribution as $n \rightarrow \infty$.

PROOF: ★ ⁶⁸ Let $U_n := \frac{\bar{Y}_n - \mu}{\sigma/\sqrt{n}}$.

According to the standard version of the CLT (Theorem 10.13 on p.231) $U_n \xrightarrow{D} Z$ and, according to Theorem 10.15 (Sample variance converges to population variance) on p.234, $S_n^2 \xrightarrow{D} \sigma^2$.

By Theorem 10.3 (Convergence is maintained under continuous transformations) on p.218,

$$\sigma U_n \xrightarrow{D} \sigma Z \quad \text{and} \quad S_n = \sqrt{S_n^2} \xrightarrow{D} \sqrt{\sigma^2} = \sigma.$$

Since the limit σ of S_n is constant, we can apply Slutsky's theorem (Theorem 10.2 on p.218) and obtain

$$W_n = \frac{\sigma U_n}{S} \xrightarrow{D} \frac{\sigma Z}{\sigma} = Z. \quad \blacksquare$$

Remark 10.8. Note that it follows from Theorem 10.10 on p.229 that, in the special case that the sample picks Y_j are $\mathcal{N}(\mu, \sigma^2)$,

$$W_n = \frac{\bar{Y}_n - \mu}{S_n/\sqrt{n}} \sim t(\text{df} = n - 1).$$

For that reason, one would rather approximate W_n with a $t(\text{df} = n - 1)$ distribution than a standard normal distribution, if the following was true:

- (1) The population is known to approximately follow a normal distribution.
- (2) The sample size is rather small (rule of thumb: $n < 40$). For such small n , the distribution of W_n may be too far away from $\mathcal{N}(0, 1)$, the limit for $n \rightarrow \infty$. \square

Example 10.9 (WMS Ch.07.3, Example 7.8). Example 7.8 of the WMS text discusses a practical example of the use of the CLT (SAT scores). You are strongly encouraged to study it. \square

Example 10.10 (WMS Ch.07.3, Example 7.9). Example 7.9 of the WMS text discusses another practical example of the use of the CLT (checkout counter service times). You are strongly encouraged to study it. \square

Example 10.11 (WMS Ch.07.4, Example 7.10). Example 7.10 of the WMS text also discusses an application of the CLT The approximation of a binomial distribution with a normal distribution. You are strongly encouraged to study it. \square

Example 10.12 (WMS Ch.07.4, Example 7.11). Example 7.11 of the WMS text also discusses the so called **continuity correction** that should be done when one approximates a binomial distribution with a normal distribution. You are strongly encouraged to study that example. \square

⁶⁸Adapted from [stats stackexchange](#) link.

11 Sample Problems for Exams

11.1 Practice Midterm 1 for Math 447 - Chris Haines

Here are some commented excerpts of a practice exam for the first midterm. It was written by Prof. Christopher Haines and forwarded to me by Prof. Adam Weisblat, both at Binghamton University (October 2023).

Exercise 11.1. Practice Midterm 1 (C. Haines) – # 01

SKIPPED

Answer: N/A ■

Exercise 11.2. Practice Midterm 1 (C. Haines) – # 02

The Lakers and Heat are playing in the NBA Finals. The series is a best-of-seven (first team to win four games clinches the series). The Lakers will win each game with probability $3/4$.

- Given that the Heat won game one, what is the probability the Lakers go on to win the series?
- Given that the Heat win at least two games in the series, what is the probability the Lakers go on to win the series?

Solution:

We denote a sequence of games as $\vec{x} = (x_1, x_2, \dots, x_n)$, where $n \leq 7$ and $x_j = H$ if the Heat win game j and $x_j = L$ if the Lakers win game j . Note that $n < 7$ is possible, for example, if $\vec{x} = (H, H, H, H)$. (The series is finished.)

Solution to (a):

- Let $A := \{ \text{The Lakers win the series} \}$
- Let $B := \{ \text{The Heat win game \#1} \}$
-

Assume that $\vec{x} \in A \cap B$. Then $x_1 = H$ and

- either $x_2 = x_3 = x_4 = x_5 = L \Rightarrow$ one choice
- or one of x_2, \dots, x_4 is H and the other three and x_5 are $L \Rightarrow \binom{4}{1} = 4$ choices
- or two of x_2, \dots, x_5 are H and the other three and x_6 are $L \Rightarrow \binom{5}{2} = 10$ choices
- Thus, $P(A \cap B) = 1 \cdot \frac{1}{4} \cdot \left(\frac{3}{4}\right)^4 + 4 \cdot \left(\frac{1}{4}\right)^2 \cdot \left(\frac{3}{4}\right)^4 + 10 \cdot \left(\frac{1}{4}\right)^3 \cdot \left(\frac{3}{4}\right)^4$

We obtain $P(A | B) = P(A \cap B) / P(B) = 1701 / 2048$. ■

Solution to (b): Note that my solution differs from that given in the original (see course materials page!)

- Let $A := \{ \text{The Lakers win the series} \}$,
- $B := \{ \text{The Heat win at least 2 games} \}$,
- $B_2 := \{ \text{The Heat win precisely 2 games} \}$.
- $B_3 := \{ \text{The Heat win precisely 3 games} \}$,
- Then $A \cap B = A \cap (B_2 \uplus B_3)$ (Heat cannot win more than 3 if Lakers win the series).

To compute $P(A \cap B) = P(A \cap B_2) + P(B_3 \cap B_3)$, we note that

- either $\vec{x} \in A \cap B_2 \Leftrightarrow$ exactly two of x_1, \dots, x_5 are H and $x_6 = L \Rightarrow \binom{5}{2} = 10$ choices
- or $\vec{x} \in A \cap B_3$, i.e., exactly 3 of x_1, \dots, x_6 are H and $x_7 = L \Rightarrow \binom{6}{3} = \frac{6 \cdot 5 \cdot 4}{3!} = 20$ choices
- Thus, $P(A \cap B) = 10 \cdot \left(\frac{1}{4}\right)^2 \cdot \left(\frac{3}{4}\right)^4 + 20 \cdot \left(\frac{1}{4}\right)^3 \cdot \left(\frac{3}{4}\right)^4$

Next, we compute $P(B^c)$.

- Let $B_0 := \{ \text{The Heat win precisely 0 games} \}$. Then $\vec{x} \in B_0 \Leftrightarrow x_1 = x_2 = x_3 = x_4 = L \Rightarrow 1$ choice
- Let $B_1 := \{ \text{The Heat win precisely 1 game} \}$. Then $\vec{x} \in B_1 \Leftrightarrow$ exactly one of x_1, \dots, x_4 is H and $x_5 = L \Rightarrow 4$ choices
- Further, $P(B^c) = P(B_0) + P(B_1) = \left(\frac{3}{4}\right)^4 + 4 \cdot \frac{1}{4} \left(\frac{3}{4}\right)^4 = 2 \left(\frac{3}{4}\right)^4$.

Thus,

$$P(A | B) = \frac{P(A \cap B)}{1 - P(B^c)} = \frac{10 \cdot \left(\frac{1}{4}\right)^2 \cdot \left(\frac{3}{4}\right)^4 + 20 \cdot \left(\frac{1}{4}\right)^3 \cdot \left(\frac{3}{4}\right)^4}{1 - 2 \left(\frac{3}{4}\right)^4} \blacksquare$$

12 Other Appendices

12.1 Greek Letters

The following section lists all greek letters that are commonly used in mathematical texts. You do not see the entire alphabet here because there are some letters (especially upper case) which look just like our latin alphabet letters. For example: $A = \text{Alpha}$ $B = \text{Beta}$. On the other hand there are some lower case letters, namely epsilon, theta, sigma and phi which come in two separate forms. This is not a mistake in the following tables!

| | | | |
|-----------------------|-------------------|--------------------|----------------|
| α alpha | θ theta | ξ xi | ϕ phi |
| β beta | ϑ theta | π pi | φ phi |
| γ gamma | ι iota | ρ rho | χ chi |
| δ delta | κ kappa | ϱ rho | ψ psi |
| ϵ epsilon | \varkappa kappa | σ sigma | ω omega |
| ε epsilon | λ lambda | ς sigma | |
| ζ zeta | μ mu | τ tau | |
| η eta | ν nu | υ upsilon | |

| | | | |
|----------------|------------------|--------------------|----------------|
| Γ Gamma | Λ Lambda | Σ Sigma | Ψ Psi |
| Δ Delta | Ξ Xi | Υ Upsilon | Ω Omega |
| Θ Theta | Π Pi | Φ Phi | |

12.2 Notation

This appendix on notation has been provided because future additions to this document may use notation which has not been covered in class. It only covers a small portion but provides brief explanations for what is covered.

For a complete list check the list of symbols and the index at the end of this document.

Notation 12.1. a) If two subsets A and B of a space Ω are disjoint, i.e., $A \cap B = \emptyset$, then we often write $A \uplus B$ rather than $A \cup B$ or $A + B$. The complement $\Omega \setminus A$ of A is denoted A^c .

b) $\mathbb{R}_{>0}$ or \mathbb{R}^+ denotes the interval $]0, +\infty[$, $\mathbb{R}_{\geq 0}$ or \mathbb{R}_+ denotes the interval $[0, +\infty[$,

c) The set $\mathbb{N} = \{1, 2, 3, \dots\}$ of all natural numbers excludes the number zero. We write \mathbb{N}_0 or \mathbb{Z}_+ or $\mathbb{Z}_{\geq 0}$ for $\mathbb{N} \uplus \{0\}$. $\mathbb{Z}_{\geq 0}$ is the B/G notation. It is very unusual but also very intuitive. \square

References

- [1] Peter J. Bickel and Kjell A. Doksum. Mathematical Statistics. Holden-Day, San Francisco, 1st edition, 1977.
- [2] Thomas Björk. Arbitrage Theory in Continuous Time. Oxford University Press, 2nd edition, 2004.
- [3] George Casella and Roger L. Berger. Statistical Inference. Cengage, 2nd edition, 2001.
- [4] Saeed Ghahramani. Fundamentals of Probability:With Stochastic Processes. Chapman and Hall, 4th edition, 2018.
- [5] Robert Hogg, Joseph McKean, and Allen Craig. Introduction to Mathematical Statistics. Pearson, 8th edition.
- [6] Vladislav Kargin. Lecture Notes for the Introduction to Probability Course. May 24, 2022 edition, 2022.
- [7] Hossein Pishro-Nik. Introduction to Probability, Statistics, and Random Processes. Kappa Research, LLC, 2014.
- [8] Sheldon M. Ross. A First Course in Probability. Macmillan, New York, 3rd edition, 198.
- [9] Georgi Evgen'evitch Shilov. Elementary Real and Complex Analysis. Dover, Mineola, 1st edition, 1996.
- [10] Steve E. Shreve. Stochastic Calculus for Finance I: The Binomial Asset Pricing Model. Springer, 1st edition, 2003.
- [11] D. Wackerly, W. Mendenhall, and R.L. Scheaffer. Mathematical Statistics with Applications. Thomson Brooks/Cole, 7th edition, 2008.
- [12] Richard E. Williamson and Hale F. Trotter. Multivariable Mathematics. Prentice Hall, 3rd edition, 1995.

List of Symbols

- $A_n \downarrow A$ – nonincreasing set seq. , 31
 $A_n \uparrow A$ – nondecreasing set seq. , 31
 $F_Y(y)$ – CDF of random var. Y , 115
 $[a, b[,]a, b]$ – half-open intervals , 24
 $[a, b]$ – closed interval , 24
 C_k^n – nbr of combinations , 75
 P_r^n – permutation , 73
 $\binom{n}{r}$ – nbr of combinations , 75
 \Rightarrow – implication , 18
 \emptyset – empty set, 16
 $\exists!$ – exists unique , 23
 \exists – exists , 23
 \forall – for all , 23
 $\mathfrak{P}(\Omega), 2^\Omega$ – power set , 21
 $\pm\infty$ – \pm infinity , 24
 $|x|$ – absolute value , 25
 $]a, b[_\mathbb{Q}$ – interval of rational #s , 25
 $]a, b[_\mathbb{Z}$ – interval of integers , 25
 $]a, b[$ – open interval , 24
 $x \in X$ – element of a set, 15
 $x \notin X$ – not an element of a set, 15
 $x_n \downarrow x$ – nonincreasing seq. , 31
 $x_n \uparrow x$ – nondecreasing seq. , 31
 A^c – complement of A , 19
 \mathbb{N}_0 – nonnegative integers, 24
 \mathbb{R}^+ – positive real numbers, 24
 $\mathbb{R}_{>0}$ – positive real numbers, 24
 $\mathbb{R}_{\geq 0}$ – nonnegative real numbers, 24
 $\mathbb{R}_{\neq 0}$ – non-zero real numbers, 24
 \mathbb{R}_+ – nonnegative real numbers, 24
 $\mathbb{Z}_{\geq 0}$ – nonnegative integers, 24
 \mathbb{Z}_+ – nonnegative integers, 24

 $(x_i)_{i \in I}$ – family , 32
 1_A – indicator function of A , 64
 $2^\Omega, \mathfrak{P}(\Omega)$ – power set , 21
 $\binom{n}{n_1 n_2 \dots n_k}$ – multinom. coeff. , 77
 $\binom{n}{k}$ – binomial coeff. , 77
 μ'_k – k th moment , 110
 μ_k – k th central moment , 111, 127
 μ'_k – k th moment , 127
 ϕ_p – p th quantile , 119
 ρ – correlation coeff. , 160
 σ_Y – standard dev, discr. r.v. , 98
 σ_Y^2 – variance, cont. r.v. , 127
 σ_Y^2 – variance, discr. r.v. , 98
 $\text{binom}(n, p)$, 101
 θ – distribution parameter , 224
 Θ – parameter space , 224
 $\text{Cov}[Y_1, Y_2]$ – covariance , 159
 $E(Y)$ – expected value , 122
 $E[g(Y_1) \mid Y_2 = y_2]$ – conditional expectation , 169
 $E[Y]$ – expected value , 93
 $m(t)$ – MGF , 111
 R – sample range , 226
 S – sample standard deviation , 226
 s – sample standard deviation , 226
 S^2 – sample variance , 226
 s^2 – sample variance , 226
 $SD(Y)$ – standard dev, discr. r.v. , 98
 $\text{Var}[Y_1 \mid Y_2 = y_2]$ – conditional variance , 170
 $\text{Var}[Y]$ – variance, cont. r.v. , 127
 $\text{Var}[Y]$ – variance, discr. r.v. , 98
 $Y_n \xrightarrow{\text{a.s.}} Y$ – almost sure limit , 215
 $Y_n \xrightarrow{D} Y$ – limit in distrib. , 215
 $Y_n \xrightarrow{\text{pw}} Y$ – pointwise limit , 215
 $Y_n \xrightarrow{P} Y$ – limit in probab. , 215
 $\Gamma(\alpha)$ – gamma function , 135
 \Leftrightarrow – if and only if, 16
 \mathbb{N}, \mathbb{N}_0 , 239
 $\mathbb{R}^+, \mathbb{R}_{>0}$, 239
 $\mathbb{R}_+, \mathbb{R}_{\geq 0}$, 239
 $\mathbb{R}_{>0}, \mathbb{R}^+$, 239
 $\mathbb{R}_{\geq 0}, \mathbb{R}_+$, 239
 $\mathbb{Z}_+, \mathbb{Z}_{\geq 0}$, 239
 \mathfrak{B} – Borel σ -algebra of \mathbb{R} , 52
 \mathfrak{B}^n – Borel σ -algebra of \mathbb{R}^n , 52
 $\mathcal{N}(\mu, \sigma^2)$ – normal with μ, σ^2 , 133
 $\mathcal{N}(\mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \rho)$ – bivariate normal , 188
 $\sigma\{\mathcal{A}\}$ – σ -algebra generated by \mathcal{A} , 51
 $\text{suppt}(f_Y)$ – support of the PDF f_Y , 197
 $|X|$ – size of a set , 22
 $\{\}$ – empty set, 16
 $A \uplus B$ – disjoint union , 239
 $A \cap B$ – A intersection B , 17
 $A \setminus B$ – A minus B , 18
 $A \subset B$ – Do not use, 16
 $A \subseteq B$ – A is subset of B , 16

$A \subsetneq B$ – A is strict subset of B , 16
 $A \Delta B$ – symmetric difference of A and B , 18
 $A \uplus B$ – A disjoint union B , 17
 A^c – complement, 239
 $B \supset A$ – Do not use, 16
 $B \supsetneq A$ – B is strict superset of A , 16
 $B(\alpha, \beta)$, 139
 $f : X \rightarrow Y$ – function, 27
 $f^{-1}(B)$ – indirect image, preimage, 61
 $f_{Y_1|Y_2}(y_1 | y_2)$ – conditional PDF, 148
 $P(A | B)$ – conditional probab, 55, 147
 $(\Omega, \mathfrak{F}, P)$ – probability space, 44
 (S, \mathcal{S}, P) – sample space, 44
 $\chi^2(\nu)$ – chi-square with ν df, 137
 \mapsto – maps to, 27
 $\sigma\{\mathcal{A}\}$ – σ -algebra generated by \mathcal{A} , 51
 $A \cup B$ – A union B , 17
 $A \supseteq B$ – A is superset of B , 16
 $f|_A$ – restriction of f , 29
 $F_{Y_1, Y_2}(y_1, y_2)$ – joint CDF, 143
 P – measure, 44
 $p_{Y_1, Y_2}(y_1, y_2)$ – joint PMF, 144
 $X_1 \times X_2 \cdots \times X_n$ – cartesian product, 35
 $Y_{(j)}$ – j th order statistic, 178
 $\text{beta}(\alpha, \beta)$ – beta with α, β , 139
 $\text{chi-square}(\nu)$ – chi-square with ν df, 137
 $\text{expon}(\beta)$ – exponential with β , 137
 $\text{gamma}(\alpha, \beta)$ – gamma with α, β , 135
 $\text{geom}(p)$, 102
 $\text{poisson}(\lambda)$, 108
 $\text{uniform}(\theta_1, \theta_2)$ – uniform distrib, 128

Index

- $\chi^2(\nu)$ (chi-square distribution), 137
- σ -algebra, 43
 - Borel σ -algebra, 52
- σ -algebra generated by collection of sets, 51
- σ -field, 43
- 0–1 encoded Bernoulli trial, 100
- 68%–95%–99.7% rule, 132

- absolute value, 25
- absolutely convergent series, 37
- almost sure convergence, 215
- almost sure limit, 215
- argument, 27
- assignment operator, 27
- asymptotic solution, 230

- Bayes formula, 84
- Bernoulli random element, 100
- Bernoulli sequence, 100
- Bernoulli trial, 100
 - 0–1 encoded, 100
 - failure probability, 100
 - success probability, 100
- Bernoulli variable, 100
- beta probability distribution, 139
- beta(α, β), 139
- bijective, 28
- binom(n, p) distribution, 101
- binomial coefficients, 77
- binomial distribution, 101
- binomial theorem, 77
- bivariate cumulative distribution function, 143
- bivariate normal distribution, 188
- bivariate probability mass function, 144
- Borel σ -algebra, 52
- Borel set, 52

- carrier, 44
- carrier set, 44
- cartesian product, 35
- CDF, 115
 - conditional, 147
 - joint, 143
- central moment of a random variable, 111
- characteristic function, 64

- chi-square distribution, 137
- chi-square with ν df (chi-square distribution), 137
- chi-square(ν) (chi-square distribution), 137
- closed interval, 24
- codomain, 27
- coefficient
 - binomial, 77
 - multinomial, 77
- collection, 17
 - indexed, 17
- combination, 75
- combinatorics, 71
- complement, 19
- conditional CDF, 147
- conditional distribution function, 147
- conditional expectation, 169
- conditional PDF, 148
- conditional PMF, 147
- conditional probability, 55
- conditional probability density function, 148
- conditional probability mass function, 147
- conditional variance, 170
- continuous random variable, 116
- continuous uniform probability distribution, 128, 154
- convergence
 - almost surely, 215
 - in distribution, 215
 - in probability, 215
 - pointwise, 215
 - uniform, 214
- convergence in distribution, 215
- convergence in probability, 215
- correction factor, 108
- correlation
 - negative, 158, 160
 - positive, 158, 160
 - zero, 158, 160
- correlation coefficient, 160
- countable set, 31
- covariance, 159
- cumulative distribution function, 115
 - bivariate, 143

- joint, 143
- De Morgan's Law, 21
- decreasing, 31
- degrees of freedom, 137
 - chi-square distribution, 137
 - denominator, 229
 - numerator, 229
- denominator degrees of freedom, 229
- density function
 - marginal, 146
- dependent random elements, 149
- dependent random variables, 149
- determinant
 - Jacobian, 204
- deterministic sample, 88
- deterministic sampling action, 87
- df = degrees of freedom, 137
- discrete probability space, 49
- discrete random variable, 67
- discrete random vector, 67
- disjoint, 17
- distribution, 66
 - binomial, 101
 - marginal, 146
 - multinomial, 175
 - parameter, 224
 - uniform, 128, 154
- distribution function, 115
 - conditional, 147
 - joint, 143
- domain, 27
- dummy variable (setbuilder), 16
- element of a set, 15
- empirical probability, 5
- empirical rule, 132
- empty set, 16
- equiprobability, 47
- estimator, 164
 - unbiased, 164
- event, 5
 - independence, 57, 58
 - mutually exclusive, 44
- events generated by random elements, 148
- exclusive events, 44
- expectation
 - conditional, 169
 - expectation - continuous r.v., 122
 - expectation - discrete r.v., 93
 - expected value, 156
 - expected value - continuous r.v., 122
 - expected value - discrete r.v., 93
 - experiment
 - multinomial, 175
 - $\text{expon}(\beta)$ (exponential distribution), 137
 - exponential distribution, 137
 - extension of a function, 29
- F distribution, 229
- failure probability, 100
- family, 32
- finite sequence, 29
- first quartile, 119
- function, 27
 - argument, 27
 - assignment operator, 27
 - codomain, 27
 - domain, 27
 - extension, 29
 - function value, 27
 - inverse, 28
 - linear, 161
 - maps to operator, 27
 - range, 27
 - restriction, 29
 - support, 197
 - symmetric, 193
 - symmetrical, 193
- function value, 27
- gamma distribution, 135
- gamma function, 135
- $\text{gamma}(\alpha, \beta)$, 135
- $\text{geom}(p)$ distribution, 102
- geometric distribution, 102
- graph, 27
- greek letters, 239
- half-open interval, 24
- histogram
 - left skewed, 135
 - right skewed, 135
- hypergeometric distribution, 106

- identity, 69
- identity function, 69
- iid sequence, 100
- increasing, 30, 31
- independent and identically distributed, 100
- independent events, 57, 58
- independent random elements, 149
- independent random variables, 149
- index set, 32
- indexed collection, 17
- indexed family, 32
- indicator function, 64
- infinite sequence, 29
- injective, 28
- integer, 23
- interval
 - closed, 24
 - half-open, 24
 - open, 24
- inverse function, 28
- irrational number, 23

- Jacobian, 204
- Jacobian determinant, 204
- Jacobian matrix, 204
- joint CDF, 143
- joint cumulative distribution function, 143
- joint distribution function, 143, 144
- joint normal distribution, 188
- joint PDF, 145
- joint PMF, 144
- joint probability density function, 145
- joint probability mass function, 144
- jointly continuous random variables, 145

- largest order statistic, 178
- left skewed, 135
- left tailed, 135
- limit
 - almost sure, 215
 - in probability, 215
 - pointwise, 215
- limit in probability, 215
- linear function, 161

- maps to operator, 27
- marginal density function, 146
- marginal distribution, 146
- marginal PDF, 146
- marginal PMF, 146
- marginal probability mass function, 146
- Markov inequality, 140
- maximum, 25
- mean, 156
- mean - continuous r.v., 122
- mean - discrete r.v., 93
- mean squared distance, 174
- measurable, 65
- median, 119
 - sample median, 187
- member of a set, 15
- member of the family, 32
- memoryless property, 138
- MGF (moment-generating function), 111
- moment about its mean, 111
- moment about the origin, 110
- moment of a random variable, 110
- moment-generating function, 111
- MS distance, 174
- multinomial coefficients, 77
- multinomial distribution, 175
- multinomial experiment, 175
- multinomial sequence, 175
- multiplicative law of probability, 56
- mutually disjoint, 6, 8–10, 17
- mutually exclusive, 44

- natural number, 23
- negative binomial distribution, 105
- negative correlation, 158, 160
- nondecreasing, 30, 31
- nonincreasing, 31
- normal distribution
 - bivariate, 188
 - joint, 188
- normal probability distribution, 133
- numerator degrees of freedom, 229

- open interval, 24
- or
 - exclusive, 23
 - inclusive, 23
- order statistic, 178
 - largest, 178
 - smallest, 178

- outcome, 5, 11
 - sample space, 11
- parameter of a distribution, 224
- parameter space, 224
- partition, 21, 33
- partitioning, 21, 33
- PDF
 - conditional, 148
 - joint, 145
 - marginal, 146
- PDF (probability density function), 117
- percentile, 119
- permutation, 73
- PMF
 - conditional, 147
 - joint, 144
 - marginal, 146
- PMF (probability mass function), 91
- pointwise convergence, 215
- pointwise limit, 215
- Poisson probability distribution, 108
- poisson(λ), 108
- positive correlation, 158, 160
- power set, 21
- preimage, 61
- probability, 44
 - conditional, 55
 - empirical, 5
- probability density function, 50, 117
 - conditional, 148
 - joint, 145
- probability distribution, 66
- probability function, 91
- probability mass function, 91
 - conditional, 147
 - joint, 144
 - marginal, 146
- probability measure, 11, 44
- probability space, 11, 44
 - discrete, 49
- proof by cases, 20
- quantile, 119
- quartile
 - first, 119
 - third, 119
- r.v. = random variable, 67
- random action, 7
- random element, 68
 - dependence, 149
 - events generated by, 148
 - independence, 149
- random item, 68
- random sampling action
 - from a distribution, 224
 - on a random variable, 224
- random variable, 67
 - central moment, 111
 - continuous, 116
 - expectation, 122
 - expected value, 122
 - mean, 122
 - dependence, 149
 - discrete, 67
 - expectation, 93
 - expected value, 93
 - mean, 93
 - variance, 98
 - distribution function, 115
 - independence, 149
 - moment, 110
 - moment about its mean, 111
 - moment about the origin, 110
 - moment-generating function, 111
 - standard deviation, 98
 - standard normal, 133
 - uncorrelated, 158, 160
 - uniform, 128
- random variables
 - jointly continuous, 145
- random vector, 67
 - discrete, 67
- range, 27
 - sample, 226
- rational number, 23
- real number, 23
- realization, 11, 86, 88
- realizations, 224
- rearrangement
 - sequence, 37
 - series, 37
- restriction of a function, 29

- right continuous function, 116
- right skewed, 135
- right tailed, 135
- rv = random variable, 67
- sample, 86, 88
 - deterministic, 88
 - realization, 86, 88
 - realizations, 224
- sample mean, 164
- sample point, 11, 44
- sample range, 226
- sample space, 11, 44
- sample standard deviation, 226
- sample variance, 226
- sampling action, 11, 86, 88
- sampling distribution, 225
- sampling procedure, 86, 88
- sampling process, 86, 88
- scale parameter, 135
- sequence, 29
 - finite, 29
 - finite subsequence, 30
 - infinite, 29
 - multinomial, 175
 - start index, 29
 - subsequence, 30
- series
 - absolutely convergent, 37
- set, 15
 - countable, 31
 - difference, 18
 - difference set, 18
 - disjoint, 17
 - intersection, 17, 33
 - mutually disjoint, 17
 - proper subset, 16
 - proper superset, 16
 - setbuilder notation, 15
 - size, 22
 - strict subset, 16
 - strict superset, 16
 - subset, 16
 - superset, 16
 - symmetric difference, 18
 - uncountable, 31
 - union, 17, 33
 - shape parameter, 135
 - sigma-algebra, 43
 - sigma-field, 43
 - simple random sample, 89
 - simple random sampling action, 89
 - singleton = singleton set, 8
 - size, 22
 - smallest order statistic, 178
 - SRS, 89
 - SRS action, 89
 - standard deviation, 98
 - sample, 226
 - standard normal, 133
 - start index, 29
 - statistic, 225
 - strictly decreasing, 31
 - strictly increasing, 30, 31
 - Student's t -distribution, 228
 - subsequence, 30
 - finite, 30
 - success probability, 100
 - support, 197
 - surjective, 28
 - symmetric function, 193
 - symmetrical function, 193
 - t -distribution, 228
 - Tchebysheff inequalities, 140
 - third quartile, 119
 - triangle inequality, 25
 - unbiased estimator, 164
 - uncorrelated random variables, 158, 160
 - uncountable set, 31
 - uniform convergence, 214
 - uniform probability distribution, 128
 - uniform random variable, 128
 - uniform random vector, 154
 - uniform(θ_1, θ_2), 128
 - universal set, 18
 - urn model with replacement, 90
 - urn model without replacement, 90
 - variance
 - conditional, 170
 - sample, 226
 - variance - discrete r.v., 98
 - vector, 35

zero correlation, [158](#), [160](#)