

NEW YORK UNIVERSITY — COURANT INSTITUTE, MATH-UA 233

# Theory of Probability



*Maximilian Nitzschner*

12/13/2021

**Disclaimer:**

These are lecture notes for the course *Theory of Probability (MATH-UA 233)*, given at New York University in Fall 2021.

The primary textbook reference for this course is [1]. For some advanced topics and further reading, especially concerning more mathematical details, the book [2] may also be helpful.

These notes are preliminary and may contain typos. If you see any mistakes or think that the presentation is unclear and could be improved, please send an email to:  
[maximilian.nitzschner@cims.nyu.edu](mailto:maximilian.nitzschner@cims.nyu.edu). All comments and suggestions are appreciated.

# Contents

<b>0. Motivation</b>	<b>5</b>
<b>1. Outcomes, events and probability</b>	<b>6</b>
1.1. Sample spaces . . . . .	6
1.2. Elementary Combinatorics . . . . .	7
1.3. Events, $\sigma$ -algebras . . . . .	10
1.4. Probability . . . . .	14
<b>2. Conditional probability and stochastic independence</b>	<b>17</b>
2.1. Conditional probability . . . . .	17
2.2. The law of total probability and Bayes' theorem . . . . .	19
2.3. Stochastic independence . . . . .	20
<b>3. Discrete distributions</b>	<b>23</b>
<b>4. Continuous distributions</b>	<b>29</b>
<b>5. Random variables</b>	<b>36</b>
5.1. Definition of random variables . . . . .	36
5.2. Law and cumulative distribution of a real random variable . . . . .	37
5.3. Transformation of random variables . . . . .	42
<b>6. Expectation, variance and higher moments of random variables</b>	<b>45</b>
6.1. Expectation . . . . .	45
6.2. Variance . . . . .	49
<b>7. Joint distributions and independence of random variables</b>	<b>53</b>
7.1. Joint distributions of random variables . . . . .	53
7.2. Independence of random variables . . . . .	60
<b>8. Operations with random variables</b>	<b>64</b>
8.1. Extremes . . . . .	64
8.2. Sums of independent random variables . . . . .	65
<b>9. More on expectation</b>	<b>68</b>
9.1. Jensen's inequality . . . . .	68
9.2. Hölder's inequality . . . . .	69
<b>10. Covariance and correlation</b>	<b>70</b>

<b>11. Conditional distributions and conditional expectation</b>	<b>74</b>
11.1. Discrete conditional distributions . . . . .	74
11.2. Continuous conditional distributions . . . . .	76
11.3. Conditional expectation . . . . .	76
<b>12. Generating functions</b>	<b>81</b>
<b>13. Convergence in probability, almost sure convergence and the law of large numbers</b>	<b>84</b>
13.1. Convergence in probability and the weak law of large numbers . . . . .	84
13.2. Almost sure convergence and the strong law of large numbers . . . . .	85
13.3. Application: Monte-Carlo integration . . . . .	87
<b>14. The central limit theorem</b>	<b>89</b>
14.1. Convergence in distribution . . . . .	89
14.2. The central limit theorem . . . . .	90
<b>15. The Poisson Process</b>	<b>94</b>
<b>16. Markov chains: An overview</b>	<b>98</b>
<b>A. Appendix</b>	<b>102</b>
A.1. Multiple integrals . . . . .	102
A.2. Alternative Proof of the central limit theorem 14.4 . . . . .	108
A.3. More properties of convergence in distribution . . . . .	110

## 0. Motivation

The purpose of *probability theory* is to study and be able to make predictions about systems that involve *randomness*.

Consider as a very simple example throwing a single die, which is a process with a random outcome. A first objective will be to develop the mathematical description of characteristics of such a random experiment: This is the specification of a *stochastic model*. Loosely speaking:

*Probability theory is concerned with the description of random phenomena using stochastic models.* (0.1)

Here are some examples of phenomena calling for a probabilistic description:

- ▶ throwing a (fair) die or coin multiple times;
- ▶ describing the random movement of a particle in  $\mathbb{Z}^d$ ,  $d \geq 1$  (random walk): at every time step, the particle moves randomly to one of its neighboring sites, with “equal probability”;



Figure 0.1.: Left panel: Possible jumps for a particle at the origin; right panel: Position of the particle after 24 steps.

Some typical questions we could ask in this context are for instance:

- ▶ What is the average of the numbers of the die after a large number of throws?
- ▶ Where will the random particle be after a large number of steps?
- ▶ What is the approximate probability that the sum of the numbers coming up when throwing the die 1000 times exceeds 5000?
- ▶ Will the random particle ever come back to the origin?

In this course we develop techniques to answer some of these questions. Notably, we will see the *law of large numbers* and the *central limit theorem*, which address the first three questions above.

# 1. Outcomes, events and probability

(Reference: [1, Chapter 1-2], or [2, Sections 1.1-1.2, 2.1-2.3])

Our primary objective is to construct a mathematical model for a *random experiment*. Conceptually, this involves the specification of three quantities:

- a *set of outcomes* or *sample space*  $\Omega \neq \emptyset$ ; an element  $\omega \in \Omega$  should be interpreted as a possible realization / measurement of the random experiment.
- a *class of events*  $\mathcal{F} \subseteq \mathcal{P}(\Omega)$ , called  $\sigma$ -algebra; an event  $A \in \mathcal{F}$  is a subset of  $\Omega$ , and we aim at specifying its probability.
- a *probability measure*  $\mathbf{P}$ , which is a map from  $\mathcal{F}$  to  $[0, 1]$  that assigns a probability  $\mathbf{P}[A]$  to any given event  $A \in \mathcal{F}$ .

The triple  $(\Omega, \mathcal{F}, \mathbf{P})$  is called a *probability space*. In the following sections, we give precise definitions of these objects and present examples.

## 1.1. Sample spaces

**Definition 1.1.** A non-empty set  $\Omega$  consisting of the possible realizations of a random experiment is called *set of outcomes* or *sample space*. An element  $\omega \in \Omega$  is called an *outcome*.

*Example 1.2.* (i) Tossing a coin: The possible outcomes are heads and tails, which we denote by  $H$  and  $T$  respectively. In this case, we have

$$\Omega_1 = \{H, T\}. \quad (1.1)$$

(ii) Rolling a die: The outcomes are the integer numbers from 1 to 6, so

$$\Omega_2 = \{1, 2, 3, 4, 5, 6\}. \quad (1.2)$$

(iii) Tossing a coin *and* rolling a die: We define the sample space as the *Cartesian product* of  $\Omega_1$  and  $\Omega_2$ , namely

$$\begin{aligned} \Omega_3 &= \Omega_1 \times \Omega_2 \\ &= \{(\omega_1, \omega_2) ; \omega_1 \in \Omega_1, \omega_2 \in \Omega_2\} \\ &= \{(H, 1), (H, 2), \dots, (H, 6), (T, 1), (T, 2), \dots, (T, 6)\}. \end{aligned} \quad (1.3)$$

- (iv)  $n$ -fold coin toss (where  $n \in \mathbb{N} = \{1, 2, \dots\}$ ): Here, we need to record the outcome as an  $n$ -tuple

$$\begin{aligned}\Omega_4 &= \{(\underbrace{H, H, \dots, H, H}_{n \text{ elements}}), (H, H, \dots, H, T), (H, H, \dots, T, H), \dots, (T, T, \dots, T, T)\} \\ &= \{(\omega_1, \dots, \omega_n); \omega_i \in \{H, T\} \text{ for } 1 \leq i \leq n\} \\ &= \underbrace{\Omega_1 \times \dots \times \Omega_1}_{n \text{ times}} = \Omega_1^n.\end{aligned}\tag{1.4}$$

- (v) Tossing a coin infinitely many times: The natural choice for outcomes will be similar as in the previous example, but with sequences of infinite length rather than  $n$ -tuples. More precisely

$$\Omega_5 = \Omega_1^{\mathbb{N}} = \{(\omega_1, \omega_2, \dots); \omega_i \in \{H, T\} \text{ for } i \in \mathbb{N}\}.\tag{1.5}$$

- (vi) The number of customers in a shop during a given day:

$$\Omega_6 = \mathbb{N}_0 = \{0, 1, 2, \dots\}.\tag{1.6}$$

- (vii) The lifetime of a light bulb:

$$\Omega_7 = \mathbb{R}_0^+ = [0, \infty).\tag{1.7}$$

Let us point out that the sample spaces  $\Omega_1, \Omega_2, \Omega_3, \Omega_4$  and  $\Omega_6$  are countable<sup>1</sup>, whereas  $\Omega_5$  and  $\Omega_7$  are uncountable.

## 1.2. Elementary Combinatorics

In many elementary cases, the assumption that *all (finitely many) outcomes are equally likely* is justified (think of rolling a die or flipping a coin multiple times). In this situation, the probability space will be uniquely characterized by the number  $|\Omega| \in \mathbb{N}$ . We want to develop effective methods to count the number of outcomes of  $\Omega$  and events  $A \subseteq \Omega$ .

*Remark 1.3.* If  $N \in \mathbb{N}$  random experiments with finite sample spaces  $\Omega_1, \Omega_2, \dots, \Omega_N$  are performed successively, an appropriate choice for the sample space of the combined experiment is given by the *Cartesian product*

$$\begin{aligned}\Omega &= \prod_{j=1}^N \Omega_j := \Omega_1 \times \Omega_2 \times \dots \times \Omega_N \\ &= \{(\omega_1, \dots, \omega_N); \omega_j \in \Omega_j, 1 \leq j \leq N\}.\end{aligned}\tag{1.8}$$

The cardinality of  $\Omega$  is given by

$$|\Omega| = \prod_{j=1}^N |\Omega_j| = |\Omega_1| \cdot |\Omega_2| \cdot \dots \cdot |\Omega_N|.\tag{1.9}$$

We already saw this in Example 1.2, (iii).

<sup>1</sup>A set  $S$  is countable if it is empty or if there exists a surjective (onto) map  $\rho : \mathbb{N} \rightarrow S$ . This includes the case of finite  $S$ .

*Example 1.4.* Imagine a password contains four symbols, where the first two are (uppercase) roman letters, and the third and fourth are each a single digit (e.g. “TP21”). How many passwords can be formed with this set-up? Here we have:

$$\begin{aligned}\Omega_1 &= \Omega_2 = \{A, B, \dots, Z\}, \\ \Omega_3 &= \Omega_4 = \{0, 1, \dots, 9\}.\end{aligned}$$

A password is an element of  $\Omega = \Omega_1 \times \Omega_2 \times \Omega_3 \times \Omega_4$ . Therefore:

$$|\Omega| = \prod_{j=1}^4 |\Omega_j| = 26 \cdot 26 \cdot 10 \cdot 10 = 67600.$$

**Proposition 1.5.** *The number of choices of a sample of size  $r \in \mathbb{N}$  out of  $\{1, 2, \dots, n\}$  is given as follows:*

	<i>with repetitions</i>	<i>without repetitions</i>
<i>ordered</i>	$n^r$	$\frac{n!}{(n-r)!}$
<i>unordered</i>	$\binom{n+r-1}{r}$	$\binom{n}{r}$

For the case without repetitions, we additionally require  $r \leq n$ . In the table above we used the notations  $k! = k \cdot (k-1) \cdot \dots \cdot 1$  for  $k \in \mathbb{N}$  (and  $0! = 1$ ), as well as  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$  for  $0 \leq k \leq n$ .

*Proof.* ► Ordered samples, with repetitions: This is a special case of Remark 1.3. More precisely, we use

$$\Omega_1 = \{(\omega_1, \dots, \omega_r); \omega_j \in \{1, 2, \dots, n\}\} = \{1, 2, \dots, n\}^r, \quad (1.10)$$

with  $|\Omega_1| = n^r$ .

► Ordered samples, without repetitions: Here we use

$$\Omega_2 = \{(\omega_1, \dots, \omega_r); \omega_j \in \{1, 2, \dots, n\}, \omega_i \neq \omega_j \text{ for } i \neq j\}, \quad (1.11)$$

with  $|\Omega_2| = n \cdot (n-1) \cdot \dots \cdot (n-r+1)$ .

► Unordered samples, without repetitions: Here we use

$$\Omega_3 = \{\{\omega_1, \dots, \omega_r\}; \omega_j \in \{1, 2, \dots, n\}, \omega_i \neq \omega_j \text{ for } i \neq j\}. \quad (1.12)$$

Here  $r!|\Omega_3| = |\Omega_2| = \frac{n!}{(n-r)!}$  holds: This is because for  $r \in \{1, \dots, n\}$  different elements  $\omega_1, \dots, \omega_r$ , there are exactly  $r$  possibilities of reordering.



- Unordered samples, with repetitions: The sample space can be written as

$$\Omega_4 = \{(\omega_1, \dots, \omega_r) ; \omega_j \in \{1, 2, \dots, n\}, 1 \leq \omega_1 \leq \omega_2 \leq \dots \leq \omega_r \leq n\}. \quad (1.13)$$

We visualize an element of  $\Omega_4$  as follows: We separate the  $n$  numbers  $1, \dots, n$  by  $n - 1$  lines ( $|$ ), and for each instance of one of these numbers within the sequence  $(\omega_1, \dots, \omega_r)$ , we put a dot ( $\bullet$ ) in the respective bin.

*Example:* Let  $n = 6, r = 5$ . The element  $(1, 1, 3, 4, 6) \in \Omega_4$  corresponds to the string

$$\bullet \bullet || \bullet | \bullet || \bullet,$$

and the element  $(2, 2, 2, 5, 5) \in \Omega_4$  corresponds to the string

$$| \bullet \bullet \bullet ||| \bullet \bullet |.$$

The number of different strings corresponds therefore to the numbers of choices of a set of  $r$  elements (the dots) out of a set with  $n + r - 1$  elements (the strings consisting of dots and lines), which is exactly  $\binom{n+r-1}{r}$  by the previous step.

□

---

*End of Lecture 1*

*Example 1.6.* A committee of 12 persons consists of 3 representatives of group  $A$ , 4 of group  $B$  and 5 of group  $C$ . We want to choose a subcommittee of 5 persons, with

- one member of group  $A$ ,
- two members of group  $B$ ,
- two members of group  $C$ ?

Let  $S$  denote a set enumerating the different possible choices for the subcommittee. Note that we do not specify the order within the groups and obviously, there are no repetitions in the choice of the members. Thus we have

- $\binom{3}{1}$  choices for the member from group  $A$ ,
- $\binom{4}{2}$  choices for the member from group  $B$ ,
- $\binom{5}{2}$  choices for the member from group  $C$ ,

and thus

$$|E| = \binom{3}{1} \cdot \binom{4}{2} \cdot \binom{5}{2} = 180. \quad (1.14)$$

In Proposition 1.5, we essentially considered all possible ways of choosing samples of size  $r$  out of a set with  $n$  elements. Suppose now that  $n \in \mathbb{N}$  items are to be divided into  $k$  distinct groups of size  $n_k \in \mathbb{N}$  (so that  $n = n_1 + \dots + n_k$ ). How many such choices are possible?

- There are  $\binom{n}{n_1}$  choices for the first group,

- there are  $\binom{n-n_1}{n_2}$  choices for the second group,
- there are  $\binom{n-n_1-n_2}{n_3}$  choices for the third group, ...

and thus in total, there are

$$\begin{aligned}
 & \binom{n}{n_1} \cdot \binom{n-n_1}{n_2} \cdot \binom{n-n_1-n_2}{n_3} \cdot \dots \\
 &= \frac{n!}{n_1!(n-n_1)!} \cdot \frac{(n-n_1)!}{n_2!(n-n_1-n_2)!} \cdot \frac{(n-n_1-n_2)!}{n_3!(n-n_1-n_2-n_3)!} \cdot \dots \\
 &= \frac{n!}{\prod_{j=1}^k n_j!}.
 \end{aligned} \tag{1.15}$$

*Example 1.7.* Suppose a company has 10 employees. How many ways are there to assign tasks, if 5 employees are needed for task “A”, 3 are needed for task “B” and 2 are needed for task “C”? In this set-up, we have  $n = 10$  (the employees) and  $n_1 = 5$ ,  $n_2 = 3$ ,  $n_3 = 2$ , so:

$$\frac{10!}{5! \cdot 3! \cdot 2!} = 2520$$

possibilities. Incidentally, this is of course the same number as the number of reorderings of the “word” *AAAAABBBCC*.

The expression  $\frac{n!}{n_1! \dots n_k!}$  is called *multinomial coefficient*, and it generalizes the binomial coefficient  $\binom{n}{n_1} = \frac{n!}{n_1! n_2!}$ . Sometimes, the following abbreviation is used:

$$\binom{n}{n_1, n_2, \dots, n_k} = \frac{n!}{\prod_{j=1}^k n_j!}, \quad n_1, \dots, n_k \in \mathbb{N}_0, \sum_{j=1}^k n_j = n \in \mathbb{N}_0. \tag{1.16}$$

*Remark 1.8.* Using the multinomial coefficients, one can show the *multinomial theorem*: For  $x_1, \dots, x_k \in \mathbb{R}$  and  $n \in \mathbb{N}$ , one has

$$\left( \sum_{j=1}^k x_j \right)^n = \sum_{\substack{(n_1, \dots, n_k) \in \mathbb{N}_0^k \\ n_1 + \dots + n_k = n}} \binom{n}{n_1, n_2, \dots, n_k} \prod_{j=1}^k x_j^{n_j}. \tag{1.17}$$

This is a generalization of the well-known *binomial theorem*: For  $x, y \in \mathbb{R}$  and  $n \in \mathbb{N}$ , one has

$$(x + y)^n = \sum_{k=0}^n \binom{n}{k} x^k y^{n-k}. \tag{1.18}$$

### 1.3. Events, $\sigma$ -algebras

Suppose that we have fixed a sample space  $\Omega$ . In general we are interested in the occurrence of *events* that consist of a certain selection of outcomes. For instance consider rolling a die once (recall from Example 1.2, (ii) that

$$\Omega_2 = \{1, 2, 3, 4, 5, 6\}$$

is a reasonable choice for the sample space for this random experiment). The *event*

$$A = \text{“the upper face of the die shows an even number”} \quad (1.19)$$

can then be expressed as

$$A = \{2, 4, 6\} \subseteq \Omega_2. \quad (1.20)$$

$\leadsto$  **Naive definition:** An *event* is a subset  $A \subseteq \Omega$  of the sample space.

This works in the case where  $\Omega$  is countable (in particular, if  $\Omega$  is finite), but leads to an important complication when  $\Omega$  is uncountable (see Example 1.2, (v) and (vii)). It turns out that if we allow every subset  $A \subseteq \Omega$  for an uncountable  $\Omega$ , we cannot define a probability for  $A$  without running into problems. Fortunately, we can restrict our attention to smaller classes of subsets.

**Definition 1.9.** Let  $\Omega \neq \emptyset$ . The *power set*  $\mathcal{P}(\Omega)$  is the set of all subsets of  $\Omega$ , i.e.

$$\mathcal{P}(\Omega) = \{A; A \subseteq \Omega\}. \quad (1.21)$$

A  $\sigma$ -*algebra* on  $\Omega$  is a subset  $\mathcal{F} \subseteq \mathcal{P}(\Omega)$  that fulfills the following properties:

- (S1)  $\Omega \in \mathcal{F}$ .
- (S2) If  $A \in \mathcal{F}$ , then  $A^c = \Omega \setminus A \in \mathcal{F}$ .
- (S3) If for every  $j \in \mathbb{N}$ ,  $A_j \in \mathcal{F}$ , then  $\bigcup_{j=1}^{\infty} A_j = A_1 \cup A_2 \cup A_3 \cup \dots \in \mathcal{F}$ .

A set  $A \in \mathcal{F}$  is called an *event*. If  $\omega \in A$ , we say that the event  $A$  *occurs* (for the outcome  $\omega$ ). If  $\omega \notin A$ , we say that  $A$  *does not occur* (for the outcome  $\omega$ ).

*Remark 1.10.* (i) The power set  $\mathcal{P}(\Omega)$  itself is a  $\sigma$ -algebra (and we will usually use it if  $\Omega$  is countable, in particular if it is finite).

- (ii) The event  $\Omega$  always occurs in a random experiment, since  $\omega \in \Omega$  is always true. On the other hand, the event  $\emptyset = \Omega^c$  never occurs, since  $\omega \in \emptyset$  can never be true.
- (iii) In the previous definition, (S2) should be understood as follows: If  $A \in \mathcal{F}$  is an event, then  $A^c$ , which has the interpretation that  $A$  does not occur, should also be an event. Similarly (S3) means: If  $A_1, A_2, A_3, \dots$  are events, then  $\bigcup_{j=1}^{\infty} A_j$ , which has the interpretation that one of the  $A_j$  occurs, should also be an event.

- (iv) Let  $\Omega = \{1, 2, 3\}$  and consider the following subsets of  $\mathcal{P}(\Omega)$ :

$$\mathcal{F}_1 = \{\emptyset, \{1\}, \{2, 3\}, \{1, 2, 3\}\}$$

is a  $\sigma$ -algebra on  $\Omega$ . In this case the set  $\{2\}$  would not be an event (imagine an observer that cannot distinguish between the outcomes 2 and 3).

$$\mathcal{F}_2 = \{\emptyset, \{1\}, \{2, 3\}\}$$

is **not** a  $\sigma$ -algebra on  $\Omega$  (it violates (S1)).

$$\mathcal{F}_3 = \left\{ \emptyset, \{1\}, \{1, 2, 3\} \right\}$$

is **not** a  $\sigma$ -algebra on  $\Omega$  (it fulfills (S1) and (S3), but violates (S2)).

$$\mathcal{F}_4 = \left\{ \emptyset, \{1\}, \{2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\} \right\}$$

is **not** a  $\sigma$ -algebra on  $\Omega$  (it fulfills (S1) and (S2), but violates (S3)).

We draw some simple conclusions from Definition 1.9.

**Proposition 1.11.** *Let  $\Omega \neq \emptyset$  and  $\mathcal{F} \subseteq \mathcal{P}(\Omega)$  a  $\sigma$ -algebra.*

(i)  $\emptyset \in \mathcal{F}$ .

(ii) If for every  $j \in \mathbb{N}$ ,  $A_j \in \mathcal{F}$ , then  $\bigcap_{j=1}^{\infty} A_j \in \mathcal{F}$ .

(iii) If  $A, B \in \mathcal{F}$ , then  $A \cup B \in \mathcal{F}$ ,  $A \cap B \in \mathcal{F}$  and  $A \setminus B \in \mathcal{F}$ .

*Proof.* We first prove (i): Since  $\Omega \in \mathcal{F}$  by (S1) and  $\emptyset = \Omega^c = \Omega \setminus \Omega$ , we have that  $\emptyset \in \mathcal{F}$  by (S2).

We turn to (ii): By de Morgan's rules<sup>2</sup> we have that

$$\left( \bigcap_{j=1}^{\infty} A_j \right)^c = \bigcup_{j=1}^{\infty} \underbrace{A_j^c}_{\in \mathcal{F}, \text{ by (S2)}} \in \mathcal{F}, \text{ by (S3)}. \quad (1.22)$$

Therefore, we have again by (S2) that

$$\bigcap_{j=1}^{\infty} A_j = \left( \left( \bigcap_{j=1}^{\infty} A_j \right)^c \right)^c \in \mathcal{F}. \quad (1.23)$$

We now prove (iii): Set  $A_1 = A = \tilde{A}_1$ ,  $A_2 = B = \tilde{A}_2$  and  $A_j = \emptyset$ ,  $\tilde{A}_j = \Omega$  for  $j \geq 3$  (which are all in  $\mathcal{F}$ , using the assumption, (i) and (S1)). We then see that

$$A \cup B = A \cup B \cup \emptyset \cup \emptyset \cup \dots = \bigcup_{j=1}^{\infty} A_j \in \mathcal{F}, \quad (1.24)$$

$$A \cap B = A \cap B \cap \Omega \cap \Omega \cap \dots = \bigcap_{j=1}^{\infty} \tilde{A}_j \in \mathcal{F}, \quad (1.25)$$

<sup>2</sup>The *de Morgan rules* state that for any collection  $\{U_i\}_{i \in I}$  of subsets  $U_i \subseteq U$ , one has

$$\left( \bigcup_{i \in I} U_i \right)^c = \bigcap_{i \in I} U_i^c, \quad \left( \bigcap_{i \in I} U_i \right)^c = \bigcup_{i \in I} U_i^c$$

where we used (S2) and (ii), respectively. Finally, we have that

$$A \setminus B = A \cap \underbrace{B^c}_{\in \mathcal{F}, \text{ by (S2)}} \in \mathcal{F}. \quad (1.26)$$

□

We illustrate the set operations using again the example of rolling a single die.

*Example 1.12.* We use  $(\Omega, \mathcal{F}) = (\{1, 2, 3, 4, 5, 6\}, \mathcal{P}(\{1, 2, 3, 4, 5, 6\}))$  and consider the events

$A =$  “the upper face of the die shows an even number”  $= \{2, 4, 6\}$ ,

$B =$  “the upper face of the die shows a prime number”  $= \{2, 3, 5\}$ ,

$C =$  “the upper face of the die shows an odd number”  $= \{1, 3, 5\}$ .

From this, we obtain

$$\begin{aligned} B^c &= \{1, 4, 6\}, & A \cup B &= \{2, 3, 4, 5, 6\}, \\ A \cap B &= \{2\}, & A \cap C &= \emptyset. \end{aligned}$$

We see that the set  $B^c$  describes the event that  $B$  does not occur, the set  $A \cup B$  describes

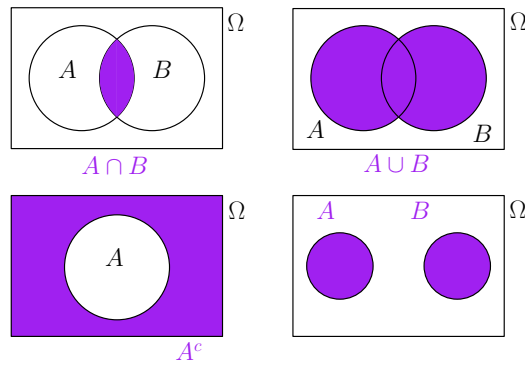


Figure 1.1.: Graphical representation of intersection, union and complement of sets (first three panels) and an example of two disjoint sets.

the event that  $A$  or  $B$  occurs<sup>3</sup> and  $A \cap B$  describes the event that  $A$  and  $B$  occur (both). The fact that  $A \cap C$  is the empty set corresponds to the fact that the events  $A$  and  $C$  are mutually exclusive.

---

*End of Lecture 2*

---

<sup>3</sup>As always in mathematics, the word “or” has a non-exclusive meaning: it includes the case where  $A$  and  $B$  occur both.

## 1.4. Probability

**Definition 1.13.** Let  $\Omega \neq \emptyset$  and  $\mathcal{F} \subseteq \mathcal{P}(\Omega)$  a  $\sigma$ -algebra on  $\Omega$ . A function  $\mathbf{P} : \mathcal{F} \rightarrow [0, 1]$  is called a *probability measure* (or simply a *probability*) if the following properties are fulfilled:

(P1)  $\mathbf{P}[\Omega] = 1$  (*normalization*).

(P2) If  $(A_j)_{j \in \mathbb{N}}$  is a sequence of events  $A_j \in \mathcal{F}$  that are *pairwise disjoint*, namely  $A_j \cap A_k = \emptyset$  for every  $j, k \in \mathbb{N}$  with  $j \neq k$ , then

$$\mathbf{P} \left[ \bigcup_{j=1}^{\infty} A_j \right] = \sum_{j=1}^{\infty} \mathbf{P}[A_j] \quad (\sigma\text{-additivity}). \quad (1.27)$$

The triple  $(\Omega, \mathcal{F}, \mathbf{P})$  is called a *probability space*.

*Example 1.14.* A very natural class of examples is given by considering

$$\emptyset \neq \Omega \text{ finite}, \quad \mathcal{F} = \mathcal{P}(\Omega), \quad (1.28)$$

and choosing the probability measure as follows:

$$\mathbf{P} : \mathcal{P}(\Omega) \rightarrow [0, 1], \quad \mathbf{P}[A] = \frac{|A|}{|\Omega|}, \quad (1.29)$$

where  $|\cdot|$  denotes the cardinality (i.e. the number of elements) of a set. The probability measure  $\mathbf{P}$  is the (discrete) *uniform distribution* on  $\Omega$ . The resulting probability space  $(\Omega, \mathcal{P}(\Omega), \mathbf{P})$  is sometimes called a *Laplace probability space*. It is characterized by the fact that

$$\mathbf{P}[\{\omega\}] = \frac{1}{|\Omega|}, \quad \text{for every } \omega \in \Omega, \quad (1.30)$$

meaning that every outcome has the same probability.

**Concrete example:** We roll a die twice and are interested in the probability that the number 6 shows up at least once. Assuming that the die is fair, we set consider the probability space  $(\Omega, \mathcal{F}, \mathbf{P})$  given by

$$\begin{aligned} \Omega &= \{1, 2, 3, 4, 5, 6\}^2 = \{(1, 1), (1, 2), \dots, (1, 6), (2, 1), \dots, (2, 6), \dots, (6, 6)\}, \\ \mathcal{F} &= \mathcal{P}(\Omega), \\ \mathbf{P}[A] &= \frac{|A|}{|\Omega|} = \frac{|A|}{36}, \quad \text{for all } A \in \mathcal{P}(\Omega), \end{aligned} \quad (1.31)$$

and the event in question is given by

$$\begin{aligned} B &= \text{“At least one 6 shows up”} \\ &= \{(1, 6), (2, 6), (3, 6), (4, 6), (5, 6), (6, 6), (6, 5), (6, 4), (6, 3), (6, 2), (6, 1)\}. \end{aligned} \quad (1.32)$$

We clearly have that

$$\mathbf{P}[\text{“At least one 6 shows up”}] = \mathbf{P}[B] = \frac{11}{36}. \quad (1.33)$$

Let us now give some elementary but important properties of probabilities.

**Proposition 1.15.** *Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $A, B, A_j \in \mathcal{F}$  for  $j \in \mathbb{N}$ . Then the following properties hold:*

- (i)  $\mathbf{P}[\emptyset] = 0$ ,
- (ii)  $\mathbf{P}[A^c] = 1 - \mathbf{P}[A]$ ,
- (iii) If  $A \subseteq B$ , then  $\mathbf{P}[A] \leq \mathbf{P}[B]$ ,
- (iv)  $\mathbf{P}[A \cup B] = \mathbf{P}[A] + \mathbf{P}[B] - \mathbf{P}[A \cap B]$ ,
- (v)  $\mathbf{P}\left[\bigcup_{j=1}^{\infty} A_j\right] \leq \sum_{j=1}^{\infty} \mathbf{P}[A_j]$ ,
- (vi) If  $A_1 \subseteq A_2 \subseteq \dots$  (we say that  $(A_j)_{j \in \mathbb{N}}$  is increasing), then  $\mathbf{P}\left[\bigcup_{j=1}^{\infty} A_j\right] = \lim_{n \rightarrow \infty} \mathbf{P}[A_n]$ ,
- (vii) If  $A_1 \supseteq A_2 \supseteq \dots$  (we say that  $(A_j)_{j \in \mathbb{N}}$  is decreasing), then  $\mathbf{P}\left[\bigcap_{j=1}^{\infty} A_j\right] = \lim_{n \rightarrow \infty} \mathbf{P}[A_n]$ .

*Proof.* We start with the proof of (i): Since  $\emptyset = \emptyset \cup \emptyset \cup \emptyset \cup \dots$  (and clearly  $\emptyset \cap \emptyset = \emptyset$ ), we see that

$$\mathbf{P}[\emptyset] \stackrel{(P2)}{=} \sum_{j=1}^{\infty} \mathbf{P}[\emptyset], \quad (1.34)$$

which can only be true if  $\mathbf{P}[\emptyset] = 0$ .

For (ii) note that  $A$  and  $A^c$  are disjoint and fulfill  $A \cup A^c = \Omega$ . We set  $A_1 = A$ ,  $A_2 = A^c$  and  $A_j = \emptyset$  for  $j \geq 3$ , so that

$$\begin{aligned} 1 &\stackrel{(P1)}{=} \mathbf{P}[\Omega] = \mathbf{P}[A \cup A^c \cup \emptyset \cup \emptyset \cup \dots] \\ &\stackrel{(P2)}{=} \mathbf{P}[A] + \mathbf{P}[A^c] + \underbrace{\mathbf{P}[\emptyset] + \mathbf{P}[\emptyset] + \dots}_{=0 \text{ by (i)}} \\ &= \mathbf{P}[A] + \mathbf{P}[A^c]. \end{aligned} \quad (1.35)$$

For the proof of (iii), consider  $\tilde{B} = B \setminus A (= \{\omega \in B; \omega \notin A\})$ , so that  $A \cup \tilde{B} = A \cup B = B$  and  $A \cap \tilde{B} = \emptyset$ . We find by the same argument as for (ii):

$$\mathbf{P}[B] = \mathbf{P}[A \cup \tilde{B} \cup \emptyset \cup \emptyset \cup \dots] = \mathbf{P}[A] + \underbrace{\mathbf{P}[\tilde{B}]}_{\geq 0} \geq \mathbf{P}[A]. \quad (1.36)$$

Note that this calculation shows the stronger statement

$$A \subseteq B \quad \Rightarrow \quad \mathbf{P}[B \setminus A] = \mathbf{P}[B] - \mathbf{P}[A]. \quad (1.37)$$

For (iv), we define  $D = B \setminus (A \cap B)$  and note that  $A \cap B \subseteq B$ ,  $A \cup B \stackrel{(\star)}{=} A \cup D$  and  $A \cap D = \emptyset$ . Argument for  $(\star)$ :

$$\begin{aligned} \omega \in A \cup B &\Leftrightarrow \omega \in A \text{ or } \omega \in B \\ &\Leftrightarrow \omega \in A \text{ or } \omega \in B \setminus (A \cap B) \Leftrightarrow \omega \in A \cup D. \end{aligned}$$

Thus, we see that

$$\begin{aligned} \mathbf{P}[A \cup B] &= \mathbf{P}[A \cup D \cup \emptyset \cup \emptyset \cup \dots] \stackrel{(P2)}{=} \mathbf{P}[A] + \underbrace{\mathbf{P}[D]}_{\stackrel{(1.37)}{=} \mathbf{P}[B] - \mathbf{P}[A \cap B]}} \\ &= \mathbf{P}[A] + \mathbf{P}[B] - \mathbf{P}[A \cap B]. \end{aligned} \quad (1.38)$$

Next, we prove (v). We define the sets

$$B_1 = A_1, \quad B_n = A_n \setminus \bigcup_{j=1}^{n-1} A_j, \quad n \geq 2. \quad (1.39)$$

The sets  $B_j$ ,  $j \in \mathbb{N}$ , are pairwise disjoint and fulfill  $\bigcup_{j=1}^{\infty} A_j = \bigcup_{j=1}^{\infty} B_j$  as well as  $B_j \subseteq A_j$  for every  $j \in \mathbb{N}$ . Therefore we have that

$$\mathbf{P}\left[\bigcup_{j=1}^{\infty} A_j\right] = \mathbf{P}\left[\bigcup_{j=1}^{\infty} B_j\right] \stackrel{(P2)}{=} \sum_{j=1}^{\infty} \mathbf{P}[B_j] \leq \sum_{j=1}^{\infty} \mathbf{P}[A_j]. \quad (1.40)$$

For (vi), we define again the same sets as in (1.39), but now note that  $(\star\star) \bigcup_{j=1}^n B_j = A_n$ ,  $n \in \mathbb{N}$ , and so

$$\mathbf{P}\left[\bigcup_{j=1}^{\infty} A_j\right] = \mathbf{P}\left[\bigcup_{j=1}^{\infty} B_j\right] \stackrel{(P2)}{=} \sum_{j=1}^{\infty} \mathbf{P}[B_j] = \lim_{n \rightarrow \infty} \sum_{j=1}^n \mathbf{P}[B_j] = \lim_{n \rightarrow \infty} \mathbf{P}[A_n], \quad (1.41)$$

having used (iv) and  $(\star\star)$  in the last step.

Finally, (vii) follows from (ii), (vi) and the fact that if  $A_1 \supseteq A_2 \supseteq \dots$ , then  $A_1^c \subseteq A_2^c \subseteq \dots$   $\square$

---

*End of Lecture 3*



## 2. Conditional probability and stochastic independence

(Reference: [1, Chapter 3], or [2, Sections 3.1, 3.3])

In this chapter we introduce the notion of *conditional probability*. Intuitively, the idea is that the existence of “partial knowledge” should influence how we determine the likelihood of a given outcome.

### 2.1. Conditional probability

Let us start with a very easy example.

*Example 2.1.* We throw two dice and ask for the probability that the sum of the numbers of both dice is smaller or equal to 7. We call this event  $A$ . Assuming that the dice are fair, this experiment is modelled by

$$(\Omega, \mathcal{F}, \mathbf{P}) = (\{1, 2, 3, 4, 5, 6\}^2, \mathcal{P}(\Omega), \mathbf{P}), \quad \mathbf{P}[\cdot] = \frac{|\cdot|}{36}. \quad (2.1)$$

Of course  $A$  is given by

$$A = \{(1, 1), (1, 2), (1, 3), (1, 4), (1, 5), (1, 6), (2, 1), (2, 2), (2, 3), (2, 4), (2, 5), (3, 1), (3, 2), (3, 3), (3, 4), (4, 1), (4, 2), (4, 3), (5, 1), (5, 2), (6, 1)\}. \quad (2.2)$$

So  $\mathbf{P}[A] = \frac{21}{36} = \frac{7}{12}$ . Now imagine we are given the information that one of the dice shows the number 6. We call this event  $B$ , i.e.

$$B = \{(1, 6), (2, 6), (3, 6), (4, 6), (5, 6), (6, 6), (6, 1), (6, 2), (6, 3), (6, 4), (6, 5)\}. \quad (2.3)$$

If we already *know* that  $B$  happens, how likely is  $A$ ? Clearly, the only outcomes of  $A$  that can still have occurred are

$$A \cap B = \{(1, 6), (6, 1)\}. \quad (2.4)$$

Thus, knowing that  $B$  occurred, we should now estimate the probability that  $A$  occurs by restricting the sample space  $\Omega$  to  $B$ , so:

$$\frac{\frac{|A \cap B|}{|\Omega|}}{\frac{|B|}{|\Omega|}} = \frac{|A \cap B|}{|B|} = \frac{2}{11}. \quad (2.5)$$

We elevate the term on the left-hand side to a definition.

**Definition 2.2.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Assume that the event  $B \in \mathcal{F}$  has a positive probability  $\mathbf{P}[B] > 0$ . We define the *conditional probability of  $A \in \mathcal{F}$  given  $B$*  by

$$\mathbf{P}[A|B] = \frac{\mathbf{P}[A \cap B]}{\mathbf{P}[B]}. \quad (2.6)$$

*Remark 2.3.* (i) If the events  $A$  and  $B$  are mutually exclusive ( $A \cap B = \emptyset$ ), then we always have  $\mathbf{P}[A|B] = 0$ , whenever the latter is defined.

(ii) One can rewrite the equation (2.6) as

$$\mathbf{P}[A \cap B] = \mathbf{P}[B] \cdot \mathbf{P}[A|B]. \quad (2.7)$$

This is sometimes called the *multiplication theorem*.

**Proposition 2.4.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Assume that the event  $B \in \mathcal{F}$  has a positive probability  $\mathbf{P}[B] > 0$ . Then  $\mathbf{P}[\cdot | B]$  defines a probability distribution on  $(\Omega, \mathcal{F})$  as well.

*Proof.* We need to verify that  $\mathbf{P}[\cdot | B]$  satisfies the axioms in Definition 1.13. First note that since  $\mathbf{P}[B] > 0$ , the expression  $\mathbf{P}[A|B]$  in (2.6) is well defined for every  $A \in \mathcal{F}$ . Moreover, since  $0 \leq \mathbf{P}[A \cap B] \leq \mathbf{P}[B]$  (using  $A \cap B \subseteq B$  and Proposition 1.15 (iii)), we see that indeed

$$\mathbf{P}[A|B] \in [0, 1] \quad \text{for every } A \in \mathcal{F}. \quad (2.8)$$

Furthermore, we have  $\Omega \cap B = B$ , and thus

$$\mathbf{P}[\Omega|B] = \frac{\mathbf{P}[\Omega \cap B]}{\mathbf{P}[B]} = \frac{\mathbf{P}[B]}{\mathbf{P}[B]} = 1, \quad (2.9)$$

verifying (P1). Finally, consider  $A_j \in \mathcal{F}$ ,  $j \in \mathbb{N}$ , pairwise disjoint. Then also the sets  $A_j \cap B$  are in  $\mathcal{F}$  (using Proposition 1.11 (iii)), and are pairwise disjoint. This implies

$$\begin{aligned} \mathbf{P} \left[ \bigcup_{j=1}^{\infty} A_j \middle| B \right] &= \frac{\mathbf{P} \left[ \left( \bigcup_{j=1}^{\infty} A_j \right) \cap B \right]}{\mathbf{P}[B]} \\ &= \frac{\mathbf{P} \left[ \bigcup_{j=1}^{\infty} (A_j \cap B) \right]}{\mathbf{P}[B]} \stackrel{(P2)}{=} \frac{\sum_{j=1}^{\infty} \mathbf{P}[A_j \cap B]}{\mathbf{P}[B]} \\ &= \sum_{j=1}^{\infty} \frac{\mathbf{P}[A_j \cap B]}{\mathbf{P}[B]} = \sum_{j=1}^{\infty} \mathbf{P}[A_j|B] \end{aligned} \quad (2.10)$$

concluding the proof of (P2). □

---

*End of Lecture 4*

## 2.2. The law of total probability and Bayes' theorem

In this section, we fix a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ .

The following result is known as the *law of total probability*.

**Theorem 2.5.** *Let  $B_1, \dots, B_n \in \mathcal{F}$  with  $\mathbf{P}[B_j] > 0$  for all  $1 \leq j \leq n$  and  $\bigcup_{j=1}^n B_j = \Omega$  with  $B_j \cap B_k = \emptyset$  for every  $j \neq k$ . Then we have for all  $A \in \mathcal{F}$  that*

$$\mathbf{P}[A] = \sum_{j=1}^n \mathbf{P}[B_j] \mathbf{P}[A|B_j]. \quad (2.11)$$

*Proof.* Note that

$$\begin{aligned} \mathbf{P}[A] &= \mathbf{P}[A \cap \Omega] = \mathbf{P}\left[A \cap \bigcup_{j=1}^n B_j\right] = \mathbf{P}\left[\bigcup_{j=1}^n (A \cap B_j)\right] \\ &\stackrel{\text{Prop. 1.15, (iv)}}{=} \sum_{j=1}^n \mathbf{P}[A \cap B_j] = \sum_{j=1}^n \mathbf{P}[B_j] \mathbf{P}[A|B_j]. \end{aligned} \quad (2.12)$$

□

We can right away combine the previous theorem with the definition of the conditional probability to obtain the following result, called *Bayes' theorem*:

**Theorem 2.6.** *Under the same assumptions as Theorem 2.5, we have for every  $1 \leq k \leq n$ , that*

$$\mathbf{P}[B_k|A] = \frac{\mathbf{P}[B_k] \mathbf{P}[A|B_k]}{\sum_{j=1}^n \mathbf{P}[B_j] \mathbf{P}[A|B_j]}. \quad (2.13)$$

*Proof.*

$$\mathbf{P}[B_k|A] \stackrel{(2.6)}{=} \frac{\mathbf{P}[B_k \cap A]}{\mathbf{P}[A]} \stackrel{(2.7), (2.11)}{=} \frac{\mathbf{P}[B_k] \mathbf{P}[A|B_k]}{\sum_{j=1}^n \mathbf{P}[B_j] \mathbf{P}[A|B_j]}. \quad (2.14)$$

□

*Example 2.7.* Biochemical tests for a certain marker / antigen / disease / ... within a population are never absolutely reliable. We consider a test with the following properties:

Let  $T$  denote the event “the test is positive”. Let  $M$  be the event “a given individual has the marker”. We assume that the test in question satisfies:

$$\begin{aligned} \mathbf{P}[T|M] &= 0.99 && (\text{sensitivity}), \\ \mathbf{P}[T^c|M^c] &= 0.99 && (\text{specificity}), \end{aligned} \quad (2.15)$$

and the marker we are looking for is such that for the population which is considered,

$$\mathbf{P}[M] = 0.01 \quad (\text{prevalence}). \quad (2.16)$$

What is the probability that someone who tests positive actually has the maker / antigen / disease?

Observe that  $\mathbf{P}[T|M^c] = 1 - 0.99 = 0.01$ . We use Bayes' theorem 2.6

$$\begin{aligned}\mathbf{P}[M|T] &= \frac{\mathbf{P}[T|M] \cdot \mathbf{P}[M]}{\mathbf{P}[T|M] \cdot \mathbf{P}[M] + \mathbf{P}[T|M^c] \cdot \mathbf{P}[M^c]} \\ &= \frac{0.99 \cdot 0.01}{0.99 \cdot 0.01 + 0.01 \cdot 0.99} \\ &= \frac{1}{2}.\end{aligned}\tag{2.17}$$

We see that because the trait under consideration is so rare, we find that half of those who test positive do not actually admit this trait, even though the test is fairly reliable!

*Remark 2.8.* Both the law of total probability and Bayes' theorem are valid if we have countably many pairwise disjoint  $B_1, B_2, \dots \in \mathcal{F}$  with  $\mathbf{P}[B_j] > 0$  for all  $j \in \mathbb{N}$  and  $\bigcup_{j=1}^{\infty} B_j = \Omega$ . In this case, (2.11) and (2.13) become

$$\mathbf{P}[A] = \sum_{j=1}^{\infty} \mathbf{P}[B_j] \mathbf{P}[A|B_j], \text{ and} \tag{2.18}$$

$$\mathbf{P}[B_k|A] = \frac{\mathbf{P}[B_k] \mathbf{P}[A|B_k]}{\sum_{j=1}^{\infty} \mathbf{P}[B_j] \mathbf{P}[A|B_j]}, \tag{2.19}$$

respectively.

### 2.3. Stochastic independence

We now introduce the notion of *stochastic independence* of events, which is one of the central concepts in probability theory and statistics. Again, we fix a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$  in this section.

**Heuristics:** The events  $A, B \in \mathcal{F}$  should be independent, if the occurrence of  $A$  has no influence on the occurrence of  $B$ , and vice versa. Specifically, if  $A$  happens, it should neither be more, nor less likely that  $B$  occurs, and vice versa, so

$$\begin{aligned}\mathbf{P}[A] &= \mathbf{P}[A|B] = \frac{\mathbf{P}[A \cap B]}{\mathbf{P}[B]}, \\ \mathbf{P}[B] &= \mathbf{P}[B|A] = \frac{\mathbf{P}[A \cap B]}{\mathbf{P}[A]},\end{aligned}\tag{2.20}$$

where we implicitly assumed that  $\mathbf{P}[A], \mathbf{P}[B] > 0$ . We turn this reasoning into a definition.

**Definition 2.9.** (i) The events  $A, B \in \mathcal{F}$  are called (*stochastically*) *independent*, if

$$\mathbf{P}[A \cap B] = \mathbf{P}[A] \cdot \mathbf{P}[B]. \tag{2.21}$$

- (ii) Let  $n \in \mathbb{N}$ ,  $n \geq 2$ . The events  $A_1, A_2, \dots, A_n \in \mathcal{F}$  are called *jointly (stochastically) independent*, if for every  $\{i_1, \dots, i_m\} \subseteq \{1, \dots, n\}$  with  $i_1, \dots, i_m$  pairwise distinct,

$$\mathbf{P}[A_{i_1} \cap \dots \cap A_{i_m}] = \mathbf{P}[A_{i_1}] \cdot \dots \cdot \mathbf{P}[A_{i_m}]. \quad (2.22)$$

*Remark 2.10.* (i) Note that both definitions include the case that a given event has probability zero. If we assume that  $\mathbf{P}[A] > 0$  and  $\mathbf{P}[B] > 0$ , then (2.20) and (2.21) are equivalent.

- (ii) The events  $\emptyset$  and  $\Omega$  are independent from any other given event. Intuitively, they contain “no additional information” on the probability.
- (iii) Equation (2.22) means that the occurrence of any subset of the events  $A_1, \dots, A_n$  does not give additional information on the occurrence of the others. For instance,

$$\mathbf{P}[A_1 | A_2 \cap \dots \cap A_n] = \frac{\mathbf{P}[A_1 \cap A_2 \cap \dots \cap A_n]}{\mathbf{P}[A_2 \cap \dots \cap A_n]} = \frac{\prod_{j=1}^n \mathbf{P}[A_j]}{\prod_{j=2}^n \mathbf{P}[A_j]} = \mathbf{P}[A_1], \quad (2.23)$$

provided that  $\mathbf{P}[A_2 \cap \dots \cap A_n] > 0$ .

- (iv) We stress that stochastic independence of two events  $A$  and  $B$  has nothing to do with them being disjoint as sets! In fact, If  $A$  and  $B$  are disjoint, then

$$0 = \mathbf{P}[\emptyset] = \mathbf{P}[A \cap B] = \mathbf{P}[A] \cdot \mathbf{P}[B], \quad (2.24)$$

so unless  $\mathbf{P}[A] = 0$  or  $\mathbf{P}[B] = 0$ , disjoint events  $A$  and  $B$  are not independent.

We illustrate the concept of independence with a number of examples.

*Example 2.11.* (i) We draw a card randomly from a standard card deck<sup>1</sup>, with

$$\Omega = \{(i, j) ; i \in \{\clubsuit, \spadesuit, \diamondsuit, \heartsuit\}, j \in \{1, 2, \dots, 13\}\}, \quad (2.25)$$

equipped with the discrete uniform distribution. Consider the events

$$\begin{aligned} A &= \{(\heartsuit, j) ; j \in \{1, \dots, 13\}\} = \text{drawing a } \heartsuit\text{-card,} \\ B &= \{(i, 1) ; i \in \{\clubsuit, \spadesuit, \diamondsuit, \heartsuit\}\} = \text{drawing an ace.} \end{aligned} \quad (2.26)$$

Clearly, we have that

$$A \cap B = \{(\heartsuit, 1)\}. \quad (2.27)$$

With this, we see that

$$\begin{aligned} \mathbf{P}[A] &= \frac{|A|}{|\Omega|} = \frac{13}{52} = \frac{1}{4}, & \mathbf{P}[B] &= \frac{|B|}{|\Omega|} = \frac{4}{52} = \frac{1}{13} \\ \mathbf{P}[A \cap B] &= \frac{|A \cap B|}{|\Omega|} = \frac{1}{52} = \mathbf{P}[A] \cdot \mathbf{P}[B]. \end{aligned} \quad (2.28)$$

So  $A$  and  $B$  are independent.

---

<sup>1</sup>With 52 *French-suited playing cards*.

(ii) Consider tossing a fair coin twice:

$$\Omega = \{(H, H), (H, T), (T, H), (T, T)\},$$

$A$  = “Heads” comes up in the first round

$$= \{(H, H), (H, T)\}, \quad \mathbf{P}[A] = \frac{1}{2},$$

$B$  = “Heads” comes up in the second round

$$= \{(H, H), (T, H)\}, \quad \mathbf{P}[B] = \frac{1}{2},$$

$$\mathbf{P}[A \cap B] = \frac{1}{4},$$

$C$  = “Heads” comes up exactly once

$$= \{(H, T), (T, H)\}, \quad \mathbf{P}[C] = \frac{1}{2},$$

$$\mathbf{P}[A \cap C] = \frac{1}{4} = \mathbf{P}[B \cap C],$$

However:  $\mathbf{P}[A \cap B \cap C] = \mathbf{P}[\emptyset] = 0 \neq \mathbf{P}[A] \cdot \mathbf{P}[B] \cdot \mathbf{P}[C]$ .

This shows that the events  $A$ ,  $B$  and  $C$  are *pairwise* independent (this means every two events out of  $\{A, B, C\}$  are independent), but not jointly independent.

We finish this section with the following result:

**Theorem 2.12.** *Let  $A_1, A_2, \dots, A_n \in \mathcal{F}$  be jointly independent. Then also  $B_1, B_2, \dots, B_n$  with  $B_i \in \{A_i, A_i^c\}$ , for  $1 \leq i \leq n$ , are jointly independent.*

*Proof.* We only show the case  $n = 2$  (the general case follows by induction).

$$\begin{aligned} (A_1 \cap A_2) \cup (A_1 \cap A_2^c) &= A_1 \quad (\text{disjoint union}) \\ \Rightarrow \underbrace{\mathbf{P}[A_1 \cap A_2] + \mathbf{P}[A_1 \cap A_2^c]}_{=\mathbf{P}[A_1] \cdot \mathbf{P}[A_2]} &= \mathbf{P}[A_1] \\ \Rightarrow \mathbf{P}[A_1 \cap A_2^c] &= \mathbf{P}[A_1] \cdot (1 - \mathbf{P}[A_2]) = \mathbf{P}[A_1] \cdot \mathbf{P}[A_2^c]. \end{aligned} \tag{2.29}$$

By changing the roles of  $A_1$  and  $A_2$ , we also have

$$\mathbf{P}[A_1^c \cap A_2] = \mathbf{P}[A_1^c] \cdot \mathbf{P}[A_2]. \tag{2.30}$$

We can finally use the same argument as in (2.29) (which implied the independence of  $A_1$  and  $A_2^c$  from the independence of  $A_1$  and  $A_2$ ) to infer the independence of  $A_1^c$  and  $A_2^c$  from the independence of  $A_1^c$  and  $A_2$ .  $\square$

---

*End of Lecture 5*

### 3. Discrete distributions

(Reference: [1, Section 4.2,4.6-4.8], or [2, Sections 2.1-2.5.1])

In the present chapter the most important discrete distributions are defined. We need to start with a short reminder on countable sets and sums over countable sets.

**Definition 3.1.** A set  $\Omega$  is called *countable* if it is empty or there is a surjective map  $\rho : \mathbb{N} \rightarrow \Omega$  (i.e. for every  $\omega \in \Omega$ , there exists  $j \in \mathbb{N}$  with  $\rho(j) = \omega$ ).<sup>1</sup>

For a countable set  $\Omega$  and a function  $f : \Omega \rightarrow [0, \infty)$ , we define the *sum of  $f$  over  $\Omega$*  by

$$\sum_{\omega \in \Omega} f(\omega) = \sup_{\substack{F \subseteq \Omega \\ F \text{ finite}}} \sum_{\omega \in F} f(\omega). \quad (3.1)$$

The expression on the right-hand side is an element of  $[0, \infty] = [0, \infty) \cup \{+\infty\}$ . In (3.1) and in the following, we use the convention that  $\sum_{\omega \in \emptyset} f(\omega) = 0$  for an empty sum. By slight abuse of notation, we also write  $\sum_{\omega \in A} f(\omega)$  instead of  $\sum_{\omega \in A} f|_A(\omega)$  for  $A \subseteq \Omega$  countable and infinite.

**Remark 3.2.** (i) The definition is consistent with finite sets: Indeed, if  $\Omega = \{\omega_1, \dots, \omega_N\}$  is a finite set, we immediately have that

$$\sum_{\omega \in \Omega} f(\omega) = \sum_{i=1}^N f(\omega_i). \quad (3.2)$$

(ii) Suppose that  $\Omega$  is countable, but infinite (such as  $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{Z}^2$ , ...). In this case, we can write  $\Omega = \{\omega_1, \omega_2, \dots\}$  by constructing a (not necessarily unique) bijective function  $\rho_* : \mathbb{N} \rightarrow \Omega$ ,  $\omega_j = \rho_*(j)$ , which we call *enumeration*. In this case, one can show that

$$\sum_{\omega \in \Omega} f(\omega) = \sum_{j=1}^{\infty} f(\omega_j) = \lim_{N \rightarrow \infty} \sum_{j=1}^N f(\omega_j), \quad (3.3)$$

and the limit on the right-hand side does not depend on the choice of the enumeration.

(iii) If  $A_j$ ,  $j \in \mathbb{N}$  are pairwise disjoint ( $A_j \cap A_k = \emptyset$  for  $j \neq k$ ) with  $\Omega = \bigcup_{j=1}^{\infty} A_j$ , then one can show the *rearrangement theorem*:

$$\sum_{\omega \in \Omega} f(\omega) = \sum_{j=1}^{\infty} \sum_{\omega \in A_j} f(\omega). \quad (3.4)$$

---

<sup>1</sup>In particular, finite sets are countable by this convention.

**Definition 3.3.** Let  $(\Omega, \mathcal{P}(\Omega), \mathbf{P})$  be a probability space with a countable sample space  $\Omega$ . The probability measure  $\mathbf{P}$  is then called a *discrete distribution*. The function

$$p : \Omega \rightarrow [0, 1], \quad p(\omega) = \mathbf{P}[\{\omega\}] \quad (3.5)$$

is the *probability mass function* of the distribution.

Obviously, if we are given a discrete distribution on  $(\Omega, \mathcal{P}(\Omega))$ , the probability mass function is uniquely determined. Conversely, every set  $(p(\omega))_{\omega \in \Omega}$  of non-negative numbers with  $\sum_{\omega \in \Omega} p(\omega) = 1$  determines a unique probability measure on  $(\Omega, \mathcal{P}(\Omega))$ , which is the statement of the following proposition.

**Proposition 3.4.** Let  $\Omega$  be countable and  $p : \Omega \rightarrow [0, 1]$  a map fulfilling

$$\sum_{\omega \in \Omega} p(\omega) = 1. \quad (3.6)$$

Then the map

$$\mathbf{P} : \mathcal{P}(\Omega) \rightarrow [0, 1], \quad A \mapsto \mathbf{P}[A] = \sum_{\omega \in A} p(\omega) \quad (3.7)$$

defines a probability measure on  $(\Omega, \mathcal{P}(\Omega))$ .

*Proof.* We first remark that since  $\Omega$  is countable and the real numbers  $(p(\omega))_{\omega \in \Omega}$  are non-negative, we have

$$\sum_{\omega \in \Omega} p(\omega) = \begin{cases} \sum_{i=1}^N p(\omega_i), & \Omega = \{\omega_1, \dots, \omega_N\} \text{ finite,} \\ \lim_{N \rightarrow \infty} \sum_{i=1}^N p(\omega_i), & \Omega = \{\omega_1, \omega_2, \dots\} \text{ infinite,} \end{cases} \quad (3.8)$$

and the value of the series does not depend on the choice of the enumeration (see Remark 3.2). Moreover, since  $A \subseteq \Omega$  is also countable, the expression for  $\mathbf{P}[A]$  in (3.7) is well-defined and in  $[0, 1]$ .

The condition (P1) is immediate by (3.6). For (P2), we consider  $A_j \in \mathcal{P}(\Omega)$  for  $j \in \mathbb{N}$  pairwise disjoint and use

$$\begin{aligned} \mathbf{P} \left[ \bigcup_{j=1}^{\infty} A_j \right] &= \sum_{\omega \in \bigcup_{j=1}^{\infty} A_j} p(\omega) \stackrel{(3.4)}{=} \sum_{j=1}^{\infty} \sum_{\omega \in A_j} p(\omega) \\ &= \sum_{j=1}^{\infty} \mathbf{P}[A_j]. \end{aligned} \quad (3.9)$$

□

---

*End of Lecture 6*

We will now present some of the most important discrete distributions.



**The discrete uniform distribution**  $\mathcal{U}(\Omega)$ 

$$\Omega \text{ finite, } \quad p(\omega) = \frac{1}{|\Omega|}. \quad (3.10)$$

This is just giving a name for the distribution considered already multiple times, see Example 1.14.

**The Bernoulli distribution**  $Ber(p)$ 

$$\Omega = \{0, 1\}, \quad p(1) = p \in [0, 1], \quad p(0) = 1 - p. \quad (3.11)$$

The Bernoulli distribution models random experiments in which a “success” occurs with probability  $p$ , and a “failure” occurs with probability  $1 - p$  (for instance, tossing a biased coin). Such experiments are also called *Bernoulli experiments*.

**The Binomial distribution**  $Bin(n, p)$ 

$$\Omega = \{0, 1, \dots, n\}, \quad p(k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad p \in [0, 1], n \in \mathbb{N}. \quad (3.12)$$

Note that

$$\sum_{k=0}^n p(k) = \sum_{k=0}^n \binom{n}{k} p^k (1 - p)^{n-k} = (p + 1 - p)^n = 1. \quad (3.13)$$

The binomial distribution is extending the Bernoulli distribution in the following way: It models how many attempts out of  $n$  independent experiments with the same success parameter  $p \in [0, 1]$  are successful. To explain it, consider the auxiliary probability space

$$(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), \mathbf{Q}), \quad \mathbf{Q}[\{(\omega_1, \dots, \omega_n)\}] = \underbrace{p^{\sum_{j=1}^n \omega_j}}_{p^{\# \text{ of successes}}} \underbrace{(1 - p)^{n - \sum_{j=1}^n \omega_j}}_{(1-p)^{\# \text{ of failures}}}. \quad (3.14)$$

Here, the string  $(\omega_1, \dots, \omega_n) \in \{0, 1\}^n$  stands for the successes and failures of the experiment in the order observed, i.e.

$$\omega_j = \begin{cases} 1, & \text{if the } j\text{th experiment is a success,} \\ 0, & \text{if the } j\text{th experiment is a failure.} \end{cases} \quad (3.15)$$

For instance, the string  $(1, 0, 0, 1)$  means that the first and last of four experiments are successes, whereas the second and third experiments are failures. By the product structure in (3.14), the experiments are independent. Now consider the event (for  $0 \leq k \leq n$ )

$$E_k = \{(\omega_1, \dots, \omega_n) \in \{0, 1\}^n; \sum_{j=1}^n \omega_j = k\} = \text{“exactly } k \text{ successes”}. \quad (3.16)$$

We have

$$\mathbf{Q}[E_k] = |E_k| p^k (1 - p)^{n-k} = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (3.17)$$

*Example 3.5.* We throw a die 4 times and we are interested in the number of times that the number six shows up. This is modelled by the binomial distribution  $Bin(4, \frac{1}{6})$ . In the description (3.14) “0” stands for the occurrence of a number other than six (failure), whereas “1” stands for the occurrence of a six (success). In this example, we have

Probability	Outcomes in $\{0, 1\}^4$
$p(0) = (\frac{5}{6})^4$	(0, 0, 0, 0)
$p(1) = \binom{4}{1} \cdot \frac{1}{6} \cdot (\frac{5}{6})^3$	(1, 0, 0, 0), (0, 1, 0, 0), (0, 0, 1, 0), (0, 0, 0, 1)
$p(2) = \binom{4}{2} \cdot (\frac{1}{6})^2 \cdot (\frac{5}{6})^2$	(1, 1, 0, 0), (1, 0, 1, 0), (1, 0, 0, 1), (0, 1, 1, 0), (0, 1, 0, 1), (0, 0, 1, 1)
$p(3) = \binom{4}{3} \cdot (\frac{1}{6})^3 \cdot \frac{5}{6}$	(0, 1, 1, 1), (1, 0, 1, 1), (1, 1, 0, 1), (0, 0, 0, 1)
$p(4) = (\frac{1}{6})^4$	(1, 1, 1, 1)

We can also order the outcomes in the form of a “tree diagram” as follows:

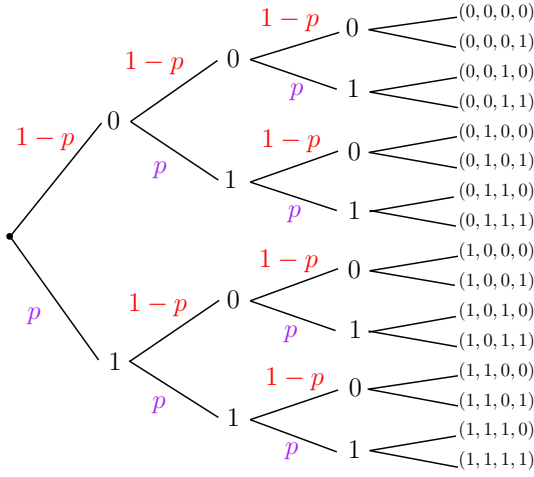


Figure 3.1.: Tree diagram of 4 successive independent Bernoulli experiments.

*Remark 3.6.* The above example shows that the same question “How likely is it that the number six comes up exactly twice when rolling a die four times?” is treated much more efficiently on the probability space

$$(\Omega = \{0, 1, 2, 3, 4\}, \mathcal{P}(\Omega), \mathbf{P}), \quad \mathbf{P} = Bin(4, \frac{1}{6}),$$

where we simply have

$$\mathbf{P}[\text{“2 sixes”}] = \mathbf{P}[\{2\}] = \binom{4}{2} \cdot \left(\frac{1}{6}\right)^2 \left(\frac{5}{6}\right)^2,$$

than on the probability space

$$(\tilde{\Omega} = \{0, 1\}^n, \mathcal{P}(\tilde{\Omega}), \mathbf{Q}), \quad \mathbf{Q}[\{(\omega_1, \dots, \omega_n)\}] = p^{\sum_{j=1}^n \omega_j} (1-p)^{n - \sum_{j=1}^n \omega_j},$$

where

$$\begin{aligned}\mathbf{Q}[\text{"2 sixes"}] &= \mathbf{Q}[\{(1, 1, 0, 0), (1, 0, 1, 0), (1, 0, 0, 1), (0, 1, 1, 0), (0, 1, 0, 1), (0, 0, 1, 1)\}] \\ &= \binom{4}{2} \cdot \left(\frac{1}{6}\right)^2 \left(\frac{5}{6}\right)^2.\end{aligned}$$

The information we need is already contained in the space  $(\Omega = \{0, 1, 2, 3, 4\}, \mathcal{P}(\Omega), \mathbf{P})$ . This concept of a *reduction in complexity* will motivate the study of random variables later.

### The Geometric distribution $Geo(p)$

$$\Omega = \mathbb{N}, \quad p(k) = (1-p)^{k-1}p, \quad p \in (0, 1). \quad (3.18)$$

Note that

$$\sum_{k=1}^{\infty} p(k) = \sum_{k=1}^{\infty} (1-p)^{k-1}p = \frac{p}{1-(1-p)} = 1. \quad (3.19)$$

The interpretation of the geometric distribution is the number of repetitions of a Bernoulli experiment (with success parameter  $p \in (0, 1)$ ) until the first success.

### The Hypergeometric distribution $H(N, M, n)$

$$\Omega = \{0, \dots, n\}, \quad p(k) = \frac{\binom{M}{k} \cdot \binom{N-M}{n-k}}{\binom{N}{n}}, \quad N, m, n \in \mathbb{N}, 0 \leq n, M \leq N. \quad (3.20)$$

The hypergeometric distribution should be understood as follows: Out of a set of  $N$  elements,  $M$  subelements have a certain favorable property. We choose uniformly at random an unordered sample of  $0 \leq n \leq N$  elements out of the large set without repetitions. Then  $p(k)$  denotes the probability that exactly  $0 \leq k \leq n$  have the the favorable property (this probability is always zero if  $M < k$ , which can happen if  $n > M$ ).

*Example 3.7.* In an urn there are 10 balls, three are green and seven are red. We draw (at once) four balls from the urn. Here

$$N = 10 \text{ (\# of balls in urn)}, \quad M = 3 \text{ (\# of green balls)}, \quad n = 4 \text{ (\# of balls drawn)}. \quad (3.21)$$

The probability that exactly two of the balls drawn are green is

$$p(2) = \frac{\binom{3}{2} \cdot \binom{7}{2}}{\binom{10}{4}}.$$

### The Poisson distribution $Pois(\lambda)$

$$\Omega = \mathbb{N}_0, \quad p(k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad \lambda > 0. \quad (3.22)$$

Note that

$$\sum_{k=0}^{\infty} p(k) = \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \cdot e^{\lambda} = 1. \quad (3.23)$$

The Poisson distribution is a natural distribution for modelling events that in principle can occur infinitely often (for instance the number of goals in a football game, or the number of raindrops falling in a given area during a given time, ...).

An important application is the following approximation of the Binomial distribution by the Poisson distribution:

**Proposition 3.8.** *Let  $\lambda > 0$  be fixed and  $p_n = \frac{\lambda}{n}$  for  $n \in \mathbb{N}$ . Then, for every  $k \in \mathbb{N}_0$ , it holds that*

$$\lim_{n \rightarrow \infty} \binom{n}{k} p_n^k (1 - p_n)^{n-k} = \frac{\lambda^k}{k!} e^{-\lambda}. \quad (3.24)$$

*Proof.*

$$\begin{aligned} \binom{n}{k} p_n^k (1 - p_n)^{n-k} &= \frac{n!}{k!(n-k)!} \frac{\lambda^k}{n^k} \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-k} \\ &= \frac{\lambda^k}{k!} \cdot \frac{n}{n} \cdot \frac{n-1}{n} \cdot \dots \cdot \frac{n-k+1}{n} \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-k} \\ &\rightarrow \frac{\lambda^k}{k!} \cdot 1 \cdot e^{-\lambda} \cdot 1, \quad \text{as } n \rightarrow \infty. \end{aligned} \quad (3.25)$$

□

---

*End of Lecture 7*

## 4. Continuous distributions

(Reference: [1, Sections 5.1, 5.3-5.6], or [2, Sections 1.2, 2.5.2, 2.6])

In this section, we introduce continuous distributions on  $\Omega \subseteq \mathbb{R}$ . Specifically, we want to be able to talk about probability spaces like  $(\mathbb{R}, \mathcal{F}, \mathbf{P})$  or  $([0, 1], \mathcal{F}, \mathbf{P})$  with appropriate choices of  $\mathcal{F}$  and  $\mathbf{P}$ . This requires some more details about  $\sigma$ -algebras, which we will present without proofs.

*Example 4.1.* Consider the arrival of a train with delay. We assume that its arrival is “uniformly distributed” between 1 PM and 2 PM. How can we model this? Suppose  $0 \hat{=} 1$  PM and  $1 \hat{=} 2$  PM. We split  $[0, 1]$  in half-open intervals of equal length  $\frac{1}{n}$ ,  $n \in \{2, 3, 4, \dots\}$  whose leftmost points are

$$\left\{ \frac{j}{n} ; 0 \leq j \leq n-1 \right\} = \left\{ 0, \frac{1}{n}, \frac{2}{n}, \dots, 1 - \frac{1}{n} \right\} \subseteq [0, 1]. \quad (4.1)$$

The probability for the train to arrive within one of these intervals  $\Delta_j = [\frac{j}{n}, \frac{j+1}{n})$ ,  $0 \leq j \leq n-1$  should be  $\frac{1}{n}$ . For  $0 \leq a < b < 1$ , we should have approximately

$$\mathbf{P}[a, b] \approx \sum_{a \leq \frac{j}{n} < b} \frac{1}{n} = \frac{1}{n} \sum_{a \leq \frac{j}{n} < b} 1 \longrightarrow \int_a^b \underbrace{1}_{=: f(x)} dx, \quad \text{as } n \rightarrow \infty. \quad (4.2)$$

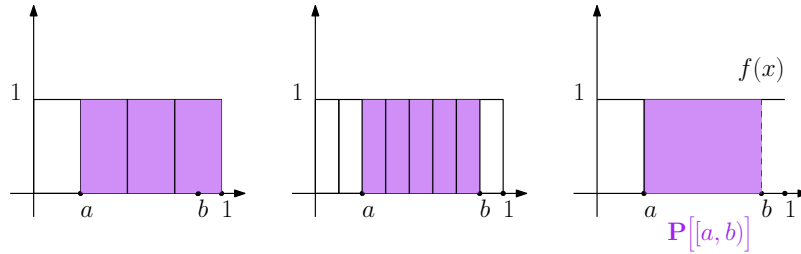


Figure 4.1.: The first and second panel represent the sum expression in (4.2) above for  $n = 4$  and  $n = 8$  respectively. The third panel represents the integral expression in the same equation.

Suppose now that the train has the highest chance to arrive around 2 PM. More specifically, let us assume that the probability to arrive within the interval  $\Delta_j = [\frac{j}{n}, \frac{j+1}{n})$  is  $\frac{2j}{n(n-1)}$  for  $0 \leq j \leq n-1$ . Note that

$$\sum_{j=0}^{n-1} \frac{2j}{n(n-1)} = \frac{2}{n(n-1)} \sum_{j=1}^{n-1} j = \frac{2}{n(n-1)} \cdot \frac{n(n-1)}{2} = 1. \quad (4.3)$$

This suggests that for  $0 \leq a < b < 1$ , we should have

$$\mathbf{P}[[a, b]] \approx \sum_{a \leq \frac{j}{n} < b} \frac{2j}{n(n-1)} = \frac{1}{n-1} \sum_{a \leq \frac{j}{n} < b} 2\frac{j}{n} \rightarrow \int_a^b \underbrace{2x}_{=:f(x)} dx, \quad \text{as } n \rightarrow \infty. \quad (4.4)$$

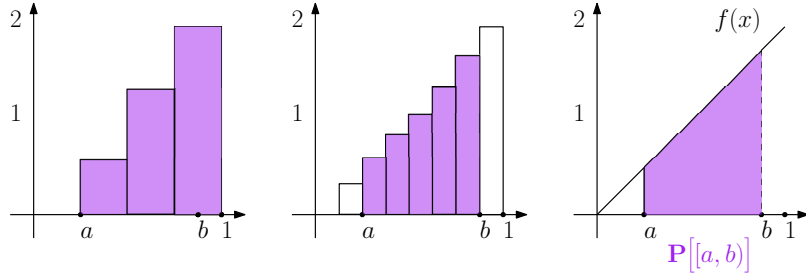


Figure 4.2.: The first and second panel represent the sum expression in (4.4) above for  $n = 4$  and  $n = 8$  respectively. The third panel represents the integral expression in the same equation.

The above considerations show that similarly to the probability mass function  $(p(\omega))_{\omega \in \Omega}$  for countable  $\Omega$ , and probabilities being defined as sums, we may want to define probabilities of intervals  $[a, b] \subseteq [0, 1)$  (or more generally  $[a, b] \subseteq \mathbb{R}$ ) as integrals over a function  $f$ . Let us turn this into a definition.

**Definition 4.2.** A piecewise continuous function<sup>1</sup>  $f : \mathbb{R} \rightarrow [0, \infty)$  is called *probability density function* if

$$\int_{-\infty}^{\infty} f(x) dx = \lim_{r, s \rightarrow \infty} \int_{-r}^s f(x) dx = 1. \quad (4.5)$$

We can then define for any interval  $[a, b] \subseteq \mathbb{R}$  with  $a < b$ :

$$\mathbf{P}[[a, b]] = \int_a^b f(x) dx \in [0, 1]. \quad (4.6)$$

The problem arises now if we want to define the probability  $\mathbf{P}[A]$  for sets  $A \subseteq \mathbb{R}$  which are *not* intervals. We would like our probability space to be  $([0, 1), \mathcal{P}([0, 1)), \mathbf{P})$  or  $(\mathbb{R}, \mathcal{P}(\mathbb{R}), \mathbf{P})$  with  $\mathbf{P}$  fulfilling (4.6). Something like this is however impossible, as the following deep result shows:

**Theorem 4.3.** *There is no probability measure  $\mathbf{P}$  on  $([0, 1), \mathcal{P}([0, 1)))$  such that  $\mathbf{P}[[a, b]] = b - a$  for every  $0 \leq a < b < 1$ .*

*Proof.* Measure Theory – Omitted. □

<sup>1</sup>This restriction is unnecessary, and can be weakened to “Lebesgue-measurable”. We will never encounter densities that are not piecewise continuous in this course!

The solution to this obstacle is to use a smaller  $\sigma$ -algebra than  $\mathcal{P}([0, 1])$  or  $\mathcal{P}(\mathbb{R})$ . We need the following preparations.

**Proposition 4.4.** *Let  $\Omega$  a non-empty set and  $\mathcal{E} \subseteq \mathcal{P}(\Omega)$ . The class defined by*

$$\sigma(\mathcal{E}) = \bigcap_{\substack{\mathcal{F} \subseteq \mathcal{P}(\Omega) \\ \mathcal{F} \text{ } \sigma\text{-algebra} \\ \mathcal{E} \subseteq \mathcal{F}}} \mathcal{F} \subseteq \mathcal{P}(\Omega) \quad (4.7)$$

*is a  $\sigma$ -algebra on  $\Omega$ , called the  $\sigma$ -algebra generated by  $\mathcal{E}$ .*

*Proof.* This follows from the fact that the intersection of (arbitrarily) many  $\sigma$ -algebras on  $\Omega$  is again a  $\sigma$ -algebra.  $\square$

The  $\sigma$ -algebra  $\sigma(\mathcal{E})$  is the smallest  $\sigma$ -algebra on  $\Omega$  that contains  $\mathcal{E}$ .

*Example 4.5.* (i) For any non-empty set  $\Omega$  we have  $\sigma(\{\Omega\}) = \{\emptyset, \Omega\}$ . Indeed, the right-hand side is a  $\sigma$ -algebra containing  $\{\Omega\}$ , and every  $\sigma$ -algebra must contain  $\{\emptyset, \Omega\}$ .

(ii) Consider  $\Omega = \{1, 2, 3, 4, 5, 6\}$ . What is  $\sigma(\{\{1\}\})$ ? Since this is a  $\sigma$ -algebra that contains  $\{1\}$  as an element, it must also contain  $\{1\}^c = \{2, 3, 4, 5, 6\}$  as an element. We then have

$$\sigma(\{\{1\}\}) = \{\emptyset, \{1\}, \{2, 3, 4, 5, 6\}, \Omega\},$$

since the expression on the right-hand side of the above equation is a  $\sigma$ -algebra itself.

The generated  $\sigma$ -algebra can also be understood as follows: Suppose specify the sets  $A \in \mathcal{E}$  to be observable events. Then  $\sigma(\mathcal{E})$  is precisely the “class of *all* events that can be observed within the model”. We use this idea to *define* a  $\sigma$ -algebra on  $\mathbb{R}$  (or on  $[a, b]$ ,  $[a, b]$ ,  $(a, b)$  and  $(a, b]$  for  $a < b$ ).

**Definition 4.6.** The *Borel  $\sigma$ -algebra*  $\mathcal{B}(\mathbb{R})$  on  $\mathbb{R}$  is defined by

$$\mathcal{B}(\mathbb{R}) = \sigma(\{[a, b]; a, b \in \mathbb{R}, a < b\}). \quad (4.8)$$

For  $a < b$  we also define the *Borel  $\sigma$ -algebras* on  $[a, b]$ ,  $[a, b)$ ,  $(a, b]$ ,  $(a, b)$  by

$$\begin{aligned} \mathcal{B}([a, b]) &= \{A \cap [a, b]; A \in \mathcal{B}(\mathbb{R})\}, \\ \mathcal{B}([a, b)) &= \{A \cap [a, b); A \in \mathcal{B}(\mathbb{R})\}, \\ \mathcal{B}((a, b]) &= \{A \cap (a, b]; A \in \mathcal{B}(\mathbb{R})\}, \\ \mathcal{B}((a, b)) &= \{A \cap (a, b); A \in \mathcal{B}(\mathbb{R})\}. \end{aligned} \quad (4.9)$$

By definition,  $\mathcal{B}([a, b])$  contains all intervals  $[a, b]$ ,  $a < b$ , so these sets are events. Are intervals like  $[a, b]$  or points also events?

*Example 4.7.* Let  $a, b \in \mathbb{R}$ ,  $a < b$ . Then  $\{a\}, [a, b], (a, b], (a, b) \in \mathcal{B}(\mathbb{R})$ . If  $A$  is a countable subset of  $\mathbb{R}$ , then  $A \in \mathcal{B}(\mathbb{R})$ .

Indeed, we have for instance that

$$\{a\} = \bigcap_{n=1}^{\infty} \underbrace{[a, a + \frac{1}{n})}_{\in \mathcal{B}(\mathbb{R})} \in \mathcal{B}(\mathbb{R}) \quad (4.10)$$

by Proposition 1.11, (ii). Then we have that

$$(a, b) = [a, b] \setminus \{a\} = \underbrace{[a, b]}_{\in \mathcal{B}(\mathbb{R})} \cap \underbrace{\{a\}^c}_{\in \mathcal{B}(\mathbb{R})} \in \mathcal{B}(\mathbb{R}), \quad (4.11)$$

again by Proposition 1.11, and the other claims about intervals follow similarly. If  $A \subseteq \mathbb{R}$  is countable, then

$$A = \{a_1, a_2, a_3, \dots\} = \bigcup_{n=1}^{\infty} \underbrace{\{a_n\}}_{\in \mathcal{B}(\mathbb{R})} \in \mathcal{B}(\mathbb{R}), \quad (4.12)$$

by (S3).

Every set that can be obtained from intervals or points by taking unions, intersections and complements is an element of  $\mathcal{B}(\mathbb{R})$ . Somewhat informally, we can say that “every set of our natural imagination is in  $\mathcal{B}(\mathbb{R})$ ”. However,  $\mathcal{B}(\mathbb{R}) \neq \mathcal{P}(\mathbb{R})$ , and this is what gives us a chance to define a probability measure  $\mathbf{P}$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ :

**Proposition 4.8.** *Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a probability density function. Then there exists a unique probability measure  $\mathbf{P}$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  with*

$$\mathbf{P}[[a, b]] = \int_a^b f(x) dx. \quad (4.13)$$

Such a probability measure  $\mathbf{P}$  is called a continuous distribution.

*Proof.* Measure Theory – Omitted. □

---

*End of Lecture 8*

Let us record the following simple observations:

**Lemma 4.9.** *Let  $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathbf{P})$  be a probability space with a continuous distribution  $\mathbf{P}$  and let  $a, b \in \mathbb{R}$  with  $a < b$ . Then*

$$\mathbf{P}[\{a\}] = 0, \quad (4.14)$$

$$\mathbf{P}[[a, b]] = \mathbf{P}[(a, b]] = \mathbf{P}[(a, b)) = \int_a^b f(x) dx, \quad (4.15)$$

$$\mathbf{P}[(a, b]] = \mathbf{P}[(a, b)) = \int_a^b f(x) dx, \quad (4.16)$$

$$\mathbf{P}[[a, \infty)) = \mathbf{P}[(a, \infty)) = \int_a^{\infty} f(x) dx \quad (4.17)$$



*Proof.* Let  $\varepsilon > 0$ . Then  $\{a\} \subseteq [a, a + \varepsilon)$ , so

$$\mathbf{P}[\{a\}] \leq \mathbf{P}[a, a + \varepsilon] = \int_a^{a+\varepsilon} f(x) dx. \quad (4.18)$$

which tends to zero as  $\varepsilon \rightarrow 0$ .<sup>2</sup> Thus we have (4.14). The other claims follow easily.  $\square$

Let us give some of the most important continuous distributions.

**The uniform distribution**  $\mathcal{U}([a, b])$

$$f(x) = \frac{1}{b-a} \mathbb{1}_{[a,b]}(x) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & x \notin [a, b], \end{cases} \quad a < b. \quad (4.19)$$

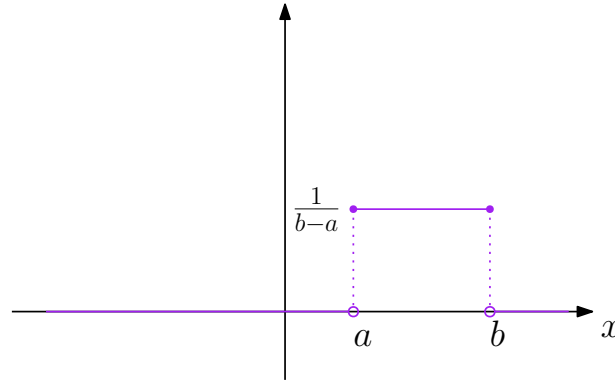


Figure 4.3.: Probability density function for  $\mathcal{U}([0, 1])$ .

**The exponential distribution**  $\mathcal{E}(\lambda)$

$$f(x) = \lambda e^{-\lambda x} \mathbb{1}_{[0, \infty)}(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0 \end{cases} \quad \lambda > 0. \quad (4.20)$$

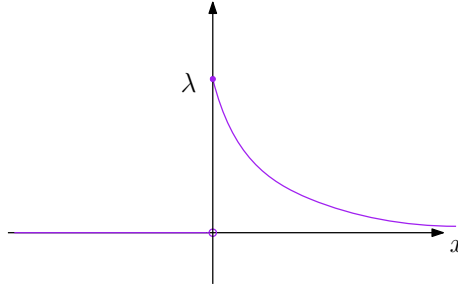
Note that this is indeed a probability density function, since

$$\int_{-\infty}^{\infty} f(x) dx = \int_0^{\infty} \lambda e^{-\lambda x} dx = \left[ -e^{-\lambda x} \right]_0^{\infty} = 1. \quad (4.21)$$

Typical application: The lifetime of radioactive isotopes is  $\mathcal{E}(\lambda)$ -distributed. An important property of the exponential distribution is the fact that it is “memoryless”. Indeed, one has

$$\mathbf{P}[(s+t, \infty) | (s, \infty)] = \mathbf{P}[(t, \infty)], \quad (4.22)$$

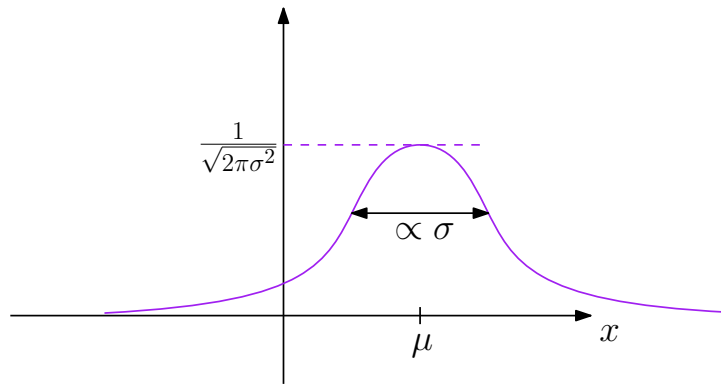
<sup>2</sup>This is immediately clear if  $f$  is continuous in  $a$ . However, we do not exclude a case where  $\lim_{\varepsilon \downarrow 0} f(a + \varepsilon) = \infty$ , as long as  $f$  is still integrable (for instance  $f(x) = \frac{1}{2\sqrt{x}} \mathbb{1}_{(0,1]}(x)$  at  $a = 0$ ). In this case, at the point  $a$  one has to argue that  $\int_a^{a+\varepsilon} f(x) dx + \int_{a+\varepsilon}^{a+1} f(x) dx = c \leq 1$ , and since  $\int_a^{a+1} f(x) dx = \lim_{\varepsilon \downarrow 0} \int_{a+\varepsilon}^{a+1} f(x) dx = c$ , we must have  $\int_a^{a+\varepsilon} f(x) dx \rightarrow 0$ , since  $f \geq 0$ .

Figure 4.4.: Probability density function for  $\mathcal{E}(\lambda)$ .

for  $s, t > 0$ . This is mimicking a similar property of the geometric distribution in the discrete case.

**The normal distribution**  $\mathcal{N}(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, \quad \mu \in \mathbb{R}, \sigma > 0. \quad (4.23)$$

Figure 4.5.: Probability density function for  $\mathcal{N}(\mu, \sigma^2)$ .

Let us explain why this is indeed a probability density function:

$$\int_{-\infty}^{\infty} f(x) dx \stackrel{y=\frac{x-\mu}{\sigma}}{=} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy =: I. \quad (4.24)$$

Now we see that

$$\begin{aligned} \left( \int_{-\infty}^{\infty} e^{-\frac{1}{2}y^2} dy \right)^2 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(x^2+y^2)} dx dy = \int_0^{2\pi} \int_0^{\infty} r e^{-\frac{1}{2}r^2} dr d\varphi \\ &= 2\pi \left[ -e^{-\frac{1}{2}r^2} \right]_0^{\infty} = 2\pi, \end{aligned} \quad (4.25)$$

therefore  $I = 1$ . We will use the notation

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}, \quad \Phi(x) = \int_{-\infty}^x \varphi(t) dt \quad (4.26)$$

for the probability density and distribution function of a *standard normal distribution*  $\mathcal{N}(0, 1)$ .  
The function  $f$  in (4.23) or  $\varphi$  in (4.26) is often also called a *Gaussian (bell) curve*.

---

*End of Lecture 9*

## 5. Random variables

(Reference: [1, Section 4.1-4.2, 4.10, 5.1], or [2, Sections 1.3])

### 5.1. Definition of random variables

Let us consider rolling two dice. We are interested in the sum of the outcomes of the two dice. If the dice are fair, the probability space modelling this problem is given by (1.31), namely

$$\begin{aligned}\Omega &= \{1, 2, 3, 4, 5, 6\}^2 = \{(1, 1), (1, 2), \dots, (6, 6)\}, \\ \mathcal{F} &= \mathcal{P}(\Omega), \\ \mathbf{P} &= \mathcal{U}(\Omega), \quad \text{i.e. } \mathbf{P}[A] = \frac{|A|}{|\Omega|} = \frac{|A|}{36}, \quad \text{for all } A \in \mathcal{P}(\Omega).\end{aligned}\tag{5.1}$$

The sum of the outcomes is given by the number

$$S(\omega) = \omega_1 + \omega_2,\tag{5.2}$$

and it can attain values in  $\{2, 3, \dots, 12\}$ . We can calculate the probability that the sum of the outcomes is 3 as follows:

$$\mathbf{P}[\{\omega \in \Omega; S(\omega) = 3\}] = \mathbf{P}[\{(1, 2), (2, 1)\}] = \frac{2}{36} = \frac{1}{18}.\tag{5.3}$$

More generally, if we ask for the probability that the sum of the outcomes is  $k \in \{2, \dots, 12\}$ :

$$\mathbf{P}[\{\omega \in \Omega; S(\omega) = k\}] = \mathbf{P}[S^{-1}(\{k\})].\tag{5.4}$$

We observe that for the function  $S : \Omega \rightarrow \mathbb{R}$  and every  $A \in \mathcal{B}(\mathbb{R})$ , the set

$$\{S \in A\} = \{\omega \in \Omega; S(\omega) \in A\} = S^{-1}(A) \in \mathcal{P}(\Omega)\tag{5.5}$$

is an event. The function  $S$  is called a random variable. Let us give a general definition.

**Definition 5.1.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. A function

$$X : \Omega \rightarrow \mathbb{R}\tag{5.6}$$

is called a (*real*) random variable if

$$X^{-1}(A) \in \mathcal{F} \quad \text{for all } A \in \mathcal{B}(\mathbb{R}).\tag{5.7}$$

We also say that the function  $X$  is  $\mathcal{F} - \mathcal{B}(\mathbb{R})$ -measurable. In this case we also write  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ .

**Remark 5.2.** (i) It can be argued that a function  $X : \Omega \rightarrow \mathbb{R}$  is a real random variable if and only if

$$\{X \leq a\} = \{\omega \in \Omega; X(\omega) \leq a\} \in \mathcal{F} \quad \text{for all } a \in \mathbb{R}. \quad (5.8)$$

- (ii) If  $\Omega$  is countable and  $\mathcal{F} = \mathcal{P}(\Omega)$ , then every function  $X : \Omega \rightarrow \mathbb{R}$  is a random variable.
- (iii) We can also view  $X$  as a map from  $\Omega$  to  $\Omega_X = \{X(\omega); \omega \in \Omega\}$ . This is in particular useful, if  $\Omega_X$  is itself countable. In this case, we say that  $X$  is a *discrete random variable*.
- (iv) Note that the definition of a random variable does not depend on the probability measure, but only on the  $\sigma$ -algebra  $\mathcal{F}$  used.

## 5.2. Law and cumulative distribution of a real random variable

**Definition 5.3.** Let  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be a real random variable. The *law of  $X$  under  $\mathbf{P}$*  is the probability measure  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$

$$\mathbf{P}_X[B] = \mathbf{P}[X^{-1}(B)], \quad B \in \mathcal{B}(\mathbb{R}). \quad (5.9)$$

**Remark 5.4.** (i) It is easy to see that  $\mathbf{P}_X$  is indeed a probability measure, so  $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathbf{P}_X)$  is a probability space.

- (ii) In the case where  $\Omega_X$  is countable, we can also consider  $\mathbf{P}_X$  as a probability measure on  $(\Omega_X, \mathcal{P}(\Omega_X))$ . In this case, (5.9) becomes

$$\mathbf{P}_X[B] = \mathbf{P}[X^{-1}(B)], \quad B \in \mathcal{P}(\Omega_X). \quad (5.10)$$

- (iii) We write as abbreviation  $X \sim \mathbf{P}_X$ . For instance,  $X \sim \text{Bin}(n, p)$  means that  $\mathbf{P}_X = \text{Bin}(n, p)$ .
- (iv) If the law  $\mathbf{P}_X$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  is a continuous distribution, i.e. it is defined in terms of a probability density  $f_X$  as in (4.13), we say that  $X$  is a *continuous random variable*.

We now introduce another notion related to the law of a random variable.

**Definition 5.5.** Let  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be a real random variable. For  $x \in \mathbb{R}$ , we set

$$F_X(x) = \mathbf{P}_X[(-\infty, x]] = \mathbf{P}[X \leq x]. \quad (5.11)$$

The function  $F_X$  is called the *cumulative distribution function* of (the law of)  $X$ .

The most important cases for us will be random variables with a discrete or continuous distribution.

- If  $X$  is a discrete real random variable, then we can write

$$F_X(x) = \sum_{y \leq x} \mathbf{P}_X[\{y\}] = \sum_{y \leq x} \mathbf{P}[X = y]. \quad (5.12)$$

► If  $X$  is a continuous random variable, then we have

$$F_X(x) = \int_{-\infty}^x f_X(t) dt, \quad (5.13)$$

where  $\mathbf{P}_X[[a, b]] = \int_a^b f_X(t) dt$ , i.e.  $f_X$  is the probability density of the law of  $X$ . Since  $f_X$  is assumed to be piecewise continuous, we have (by the fundamental theorem of Calculus) the important identity

$$f_X(x) = F'_X(x), \quad \text{at all points of continuity of } f_X. \quad (5.14)$$

*Example 5.6.* (i) Let  $X \sim \text{Bin}(2, \frac{1}{2})$ . In this case, we have

$$\begin{aligned} F_X(x) &= \begin{cases} 0, & x < 0, \\ \mathbf{P}[X = 0], & x \in [0, 1), \\ \mathbf{P}[X = 0] + \mathbf{P}[X = 1], & x \in [1, 2), \\ 1, & x \geq 2, \end{cases} \\ &= \begin{cases} 0, & x < 0, \\ \frac{1}{4}, & x \in [0, 1), \\ \frac{3}{4}, & x \in [1, 2), \\ 1, & x \geq 2. \end{cases} \end{aligned} \quad (5.15)$$

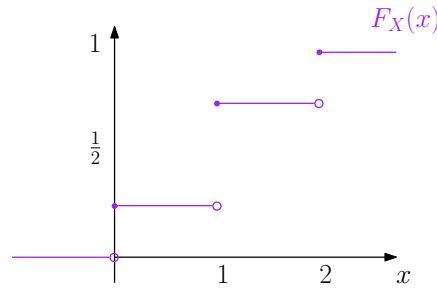
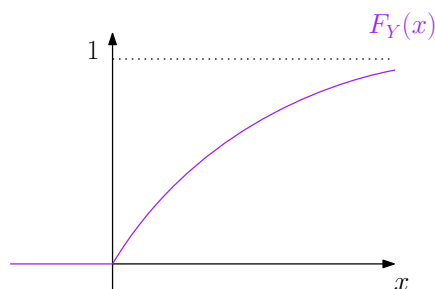


Figure 5.1.: Cumulative distribution function of  $X \sim \text{Bin}(2, \frac{1}{2})$

(ii) Let  $Y \sim \mathcal{E}(\lambda)$ ,  $\lambda > 0$ . Then the cumulative distribution function of  $Y$  is given by

$$F_Y(x) = \int_{-\infty}^x \lambda e^{-\lambda t} \mathbb{1}_{[0, \infty)}(t) dt = \begin{cases} 0, & x < 0, \\ 1 - e^{-\lambda x}, & x \geq 0. \end{cases} \quad (5.16)$$

Figure 5.2.: Cumulative distribution function of  $Y \sim \mathcal{E}(\lambda)$ 

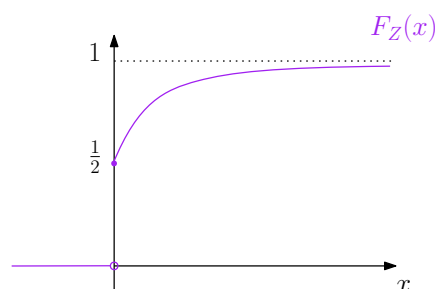
- (iii) In the previous examples (i) and (ii), the random variable  $X$  is discrete and  $Y$  is continuous. Let us stress that this is not a dichotomy. For instance, let  $Y \sim \mathcal{N}(0, 1)$ . Now let

$$Z = Y \cdot \mathbb{1}_{\{Y \geq 0\}}. \quad (5.17)$$

Here we have

$$F_Z(x) = \begin{cases} 0, & x < 0, \\ \Phi(x) = \int_{-\infty}^x \varphi(t) dt, & x \geq 0. \end{cases} \quad (5.18)$$

Note that  $Z$  is neither continuous ( $\mathbf{P}[Z = 0] = \frac{1}{2}$ ), nor discrete (it can attain all values in  $[0, \infty)$ ).

Figure 5.3.: Cumulative distribution function of  $Z$  as defined in (5.17).

---

*End of Lecture 10*

We will now collect some general facts about cumulative distribution functions.

**Lemma 5.7.** *Let  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be a real random variable. Its cumulative distribution function  $F = F_X$  satisfies the following properties:*

- (i)  $F(x) \in [0, 1]$  for all  $x \in \mathbb{R}$ .
- (ii)  $F$  is non-decreasing.

(iii)  $F$  is right continuous, i.e.

$$\lim_{\varepsilon \downarrow 0} F(x + \varepsilon) = F(x). \quad (5.19)$$

(iv)  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$ .

*Proof.* Claim (i) follows since  $F(x) = \mathbf{P}[X \in (-\infty, x]] \in [0, 1]$ .

For claim (ii), we use the fact that for  $x \leq x'$ , we have  $(-\infty, x] \subseteq (-\infty, x']$  and so

$$F(x) = \mathbf{P}_X [(-\infty, x]] \leq \mathbf{P}_X [(-\infty, x']] = F(x'). \quad (5.20)$$

For claim (iii) we apply Proposition 1.15, (vii), to the probability measure  $\mathbf{P}_X$  and the sets  $A_n = (-\infty, x_n]$ , where  $x_n \downarrow x$  for some  $x \in \mathbb{R}$  (here we mean  $x_n \geq x_{n+1}$  and  $x_n \geq x$  for every  $n \in \mathbb{N}$  and  $x_n \rightarrow x$ ). Then

$$F_X(x_n) = \mathbf{P}_X [(-\infty, x_n]] \rightarrow \mathbf{P}_X [(-\infty, x]] = F_X(x), \quad \text{as } n \rightarrow \infty, \quad (5.21)$$

since  $\bigcap_{n=1}^{\infty} (-\infty, x_n] = (-\infty, x]$ .

For (iv), we consider a sequence of real numbers  $(a_n)_{n \in \mathbb{N}}$  with  $a_n \uparrow \infty$  as  $n \rightarrow \infty$ . Then

$$F_X(a_n) = \mathbf{P}_X [(-\infty, a_n]] \rightarrow \mathbf{P}_X [\mathbb{R}] = 1, \quad \text{as } n \rightarrow \infty, \quad (5.22)$$

since  $\mathbb{R} = \bigcup_{n=1}^{\infty} (-\infty, a_n]$  and  $(-\infty, a_n] \subseteq (-\infty, a_{n+1}]$  for all  $n \in \mathbb{N}$ , upon using Proposition 1.15, (vi). The other claim follows similarly, again by using Proposition 1.15, (vii).  $\square$

In fact, the above properties characterize distribution functions in the following sense:

**Theorem 5.8.** *Let  $F : \mathbb{R} \rightarrow \mathbb{R}$  satisfying the properties (i) – (iv) of Lemma 5.7. Then there exists a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$  and a random variable  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  such that  $F = F_X$ . The law  $\mathbf{P}_X$  of  $X$  is uniquely determined by  $F$ .*

*Proof.* We define  $X : ((0, 1), \mathcal{B}((0, 1))) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  as follows:

$$X(\omega) = \sup\{y \in \mathbb{R} ; F(y) < \omega\}. \quad (5.23)$$

Note that

$$\{\omega \in (0, 1) ; X(\omega) \leq x\} = \{\omega \in (0, 1) ; \omega \leq F(x)\}, \quad x \in \mathbb{R}. \quad (5.24)$$

Indeed, if  $\omega \leq F(x)$ , then  $x \notin \{y \in \mathbb{R} ; F(y) < \omega\}$ , which implies  $x \geq X(\omega)$ .

On the other hand, if  $\omega \in (0, 1)$  with  $F(x) < \omega$ , since  $F$  is right continuous, there is  $\varepsilon > 0$  with  $F(x + \varepsilon) < \omega$ . Therefore  $X(\omega) \geq x + \varepsilon > x$ . This means that  $F(x) < \omega$  implies that  $X(\omega) > x$ . Consequently,  $x \geq X(\omega)$  implies  $\omega \leq F(x)$ .



We then equip the space  $(0, 1), \mathcal{B}((0, 1))$  with the uniform distribution  $\mathbf{P} = \mathcal{U}((0, 1))$ .<sup>1</sup> Then the law of  $X$  has cumulative distribution function given by  $F$ . Indeed:

$$F_X(x) = \mathbf{P}[X \leq x] \stackrel{(5.24)}{=} \mathbf{P}[(0, F(x)]] = F(x). \quad (5.25)$$

The proof of the second part (uniqueness of  $\mathbf{P}_X$ ) is omitted.  $\square$

Let us briefly explain the role of the discontinuity points of a cumulative distribution function  $F$ . If we look at (i), (iii) in Example 5.6, we see that a jump of the cumulative distribution function at a point  $x \in \mathbb{R}$  correspond to the probability  $\mathbf{P}_X[\{x\}] = \mathbf{P}[X = x]$ . More formally:

**Lemma 5.9.** *Let  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be a real random variable with cumulative distribution function  $F = F_X$ . Then, for every  $x \in \mathbb{R}$ , we have*

$$\mathbf{P}[X = x] = F(x) - F(x-), \quad (5.26)$$

where  $F(x-)$  (the left limit of  $F$  at  $x$ ) is defined as

$$F(x-) = \lim_{\varepsilon \downarrow 0} F(x - \varepsilon). \quad (5.27)$$

*Proof.* Note that since  $F$  is non-decreasing, the limit in (5.27) is well defined and equal to  $\lim_{n \rightarrow \infty} F(x - \frac{1}{n})$ . The claim (5.26) then follows from Proposition 1.15, (vii), the fact that

$$\{x\} = \bigcap_{n=1}^{\infty} (x - \frac{1}{n}, x] \quad (5.28)$$

and

$$F(x) - F(x - \frac{1}{n}) = \mathbf{P}[(x - \frac{1}{n}, x]]. \quad (5.29)$$

$\square$

---

#### End of Lecture 11

To summarize the previous results, we see that a cumulative distribution function  $F$  uniquely determines a probability measure  $\mathbf{P}$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  and vice versa. Let us also stress the fact that the cumulative distribution function is really associated to the *law* of a random variable, and not the random variable itself. This motivates the following definition.

**Definition 5.10.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  two random variables. We say that  $X$  and  $Y$  are *equal in distribution* or *identically distributed* if

$$\mathbf{P}_X = \mathbf{P}_Y \quad (\Leftrightarrow F_X(x) = F_Y(x), \text{ for all } x \in \mathbb{R}). \quad (5.30)$$

This is denoted as  $X \stackrel{d}{=} Y$ .

---

<sup>1</sup>Here,  $\mathcal{U}((0, 1))$  stands for the uniform distribution on the open interval  $(0, 1)$ , viewed as a probability measure on  $((0, 1), \mathcal{B}((0, 1)))$ . Since continuous distributions give zero mass to points, this is essentially the same distribution as the uniform distribution  $\mathcal{U}([0, 1])$  on  $[0, 1]$ , viewed as a probability measure on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ .

*Example 5.11.* (i) Consider again throwing two dice. We use the probability space (5.1). The result of the first and second die are given by the random variables

$$\begin{aligned} X : \{1, \dots, 6\}^2 &\rightarrow \{1, \dots, 6\}, & X(\omega_1, \omega_2) &= \omega_1, \\ Y : \{1, \dots, 6\}^2 &\rightarrow \{1, \dots, 6\}, & Y(\omega_1, \omega_2) &= \omega_2. \end{aligned} \quad (5.31)$$

We see that  $X \sim \mathcal{U}(\{1, \dots, 6\})$  and  $Y \sim \mathcal{U}(\{1, \dots, 6\})$ , so  $X \stackrel{d}{=} Y$ . Note that of course  $X \neq Y$ , since for instance  $X(1, 2) = 1 \neq 2 = Y(1, 2)$ .

(ii) Let  $X$  be any continuous random variable and  $f_X$  the probability density of its law. Assume that  $f_X$  is an even function, i.e.

$$f_X(x) = f_X(-x), \quad \text{for all } x \in \mathbb{R}. \quad (5.32)$$

Then  $X \stackrel{d}{=} -X$ . Indeed, we have

$$\begin{aligned} F_{-X}(x) &= \mathbf{P}[-X \leq x] = \mathbf{P}[X \geq -x] = 1 - \mathbf{P}[X < -x] \\ &= 1 - \int_{-\infty}^{-x} f_X(t) dt \\ &= 1 + \int_{\infty}^x f_X(-t) dt \\ &= 1 - \int_x^{\infty} f_X(t) dt \\ &= \mathbf{P}[X \leq x] = F_X(x). \end{aligned} \quad (5.33)$$

Here we repeatedly used the results of Lemma 4.9. A more concrete example: If  $X \sim \mathcal{N}(0, \sigma^2)$ ,  $\sigma > 0$ , we have that  $-X \sim \mathcal{N}(0, \sigma^2)$  as well.

The example above already gives a hint that random variables are a useful tool for algebraic manipulations. We will see more of this in the next section.

### 5.3. Transformation of random variables

*Example 5.12.* We measure the temperature of a liquid in  $^{\circ}\text{C}$  and want to transform it into  $^{\circ}\text{F}$ .

$^{\circ}\text{C}$  : random variable  $X$ ,

$^{\circ}\text{F}$  : random variable  $Y$ .

We can use the known formula

$$Y = \frac{9}{5} \cdot X + 32, \quad \text{more generally } Y = a \cdot X + b, \quad (5.34)$$

for  $a \neq 0, b \in \mathbb{R}$ . We assume that  $X$  is a continuous random variable with probability density (distribution) function  $f_X$  ( $F_X$ ), for instance  $X \sim \mathcal{N}(\mu, \sigma^2)$ . What is the distribution of  $Y$ , i.e. how do  $f_Y$  or  $F_Y$  look like? For simplicity, we also assume  $a > 0$ .

$$\begin{aligned} F_Y(y) &= \mathbf{P}_Y[(-\infty, y]] = \mathbf{P}[Y \leq y] = \mathbf{P}[aX + b \leq y] = \mathbf{P}\left[X \leq \frac{y-b}{a}\right] = F_X\left(\frac{y-b}{a}\right) \\ \Rightarrow f_Y(y) &= \frac{d}{dy} F_X\left(\frac{y-b}{a}\right) = \frac{1}{a} \cdot f_X\left(\frac{y-b}{a}\right). \end{aligned} \quad (5.35)$$

If  $Y = a \cdot X + b$  with general  $a \neq 0$ , we have

$$f_Y(y) = \frac{1}{|a|} \cdot f_X\left(\frac{y-b}{a}\right). \quad (5.36)$$

In the special case where  $X \sim \mathcal{N}(\mu, \sigma^2)$ , we have

$$\begin{aligned} f_X(x) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \\ \Rightarrow f_Y(y) &= \frac{1}{\sqrt{2\pi}\sigma|a|} e^{-\frac{1}{2\sigma^2 a^2}(y-b-a\mu)^2} \\ \Rightarrow Y = aX + b &\sim \mathcal{N}(a\mu + b, a^2\sigma^2). \end{aligned} \quad (5.37)$$

Formula (5.36) is the linear transformation rule for continuous random variables. We now show a general transformation rule for continuous random variables.

**Theorem 5.13.** *Let  $X$  be a continuous random variable and  $f_X$  the probability density function of its law. Suppose that  $g : \mathbb{R} \rightarrow \mathbb{R}$  is strictly increasing or strictly decreasing and differentiable.<sup>2</sup> Then*

$$Y = g(X) \quad (5.38)$$

*is also a continuous random variable and has density*

$$f_Y(y) = \begin{cases} f_X(g^{-1}(y)) \cdot \left| \frac{d}{dy} g^{-1}(y) \right|, & y = g(x) \text{ with } f_X(x) > 0, \\ 0, & \text{else.} \end{cases} \quad (5.39)$$

*Proof.* Let  $g$  be strictly increasing. Then, for  $[a, b] \subseteq \{g(x); f_X(x) > 0\}$ , we have

$$\begin{aligned} \mathbf{P}_Y[[a, b]] &= \mathbf{P}[g(X) \in [a, b]] = \mathbf{P}[X \in [g^{-1}(a), g^{-1}(b)]] \\ &= \int_{g^{-1}(a)}^{g^{-1}(b)} f_X(x) dx = \int_a^b \underbrace{f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y)}_{=f_Y(y)} dy. \end{aligned} \quad (5.40)$$

The proof for  $g$  strictly decreasing is similar. □

<sup>2</sup>To be very precise, we also need to make sure that  $Y = g(X)$  is still a random variable (i.e. that it is  $\mathcal{F} - \mathcal{B}(\mathbb{R})$ -measurable. This follows from the fact that  $g$  is differentiable and thus  $\mathcal{B}(\mathbb{R}) - \mathcal{B}(\mathbb{R})$ -measurable (in fact continuity is sufficient) and the simple fact that the composition of measurable functions is measurable.

*Example 5.14.* Let  $X \sim \mathcal{U}([0, 1])$  and consider  $Y = \exp(X) = e^X$ . The function  $\exp$  satisfies the requirements of the previous theorem, and its inverse is  $\log$ . Moreover,  $f_X(x) \neq 0$  if and only if  $x \in [0, 1]$ . We have

$$f_Y(y) = \begin{cases} \frac{d}{dy} \log(y) = \frac{1}{y}, & y \in [1, e], \\ 0, & y \notin [1, e]. \end{cases} \quad (5.41)$$

In the next example, we use the same method as in Theorem 5.13 to introduce the  $\chi^2$  distribution (with one degree of freedom).

*Example 5.15.* Let  $X \sim \mathcal{N}(0, 1)$  and  $Y = X^2$ . We want to calculate  $f_Y(y)$ . Unfortunately the function  $g : \mathbb{R} \rightarrow \mathbb{R}$ ,  $x \mapsto x^2$  is not strictly increasing / decreasing, but we can still use the same idea as in the proof of the transformation rule. Indeed, we see that  $g^{-1}(y) = \sqrt{y}$  and  $\frac{d}{dy} g^{-1}(y) = \frac{1}{2\sqrt{y}}$  for  $y > 0$ . Then, for  $0 \leq a < b < \infty$ , we find

$$\begin{aligned} \mathbf{P}[Y \in [a, b]] &= \mathbf{P}[X \in [\sqrt{a}, \sqrt{b}]] + \mathbf{P}[X \in (-\sqrt{b}, -\sqrt{a}]] \\ &= 2\mathbf{P}[X \in [\sqrt{a}, \sqrt{b}]] \\ &= 2 \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{y}{2}} \frac{1}{2\sqrt{y}} dy. \end{aligned} \quad (5.42)$$

We used the symmetry of  $X \sim \mathcal{N}(0, 1)$  and the fact that  $\mathbf{P}_X$  does not give mass to points. It follows that

$$f_Y(y) = \frac{1}{\sqrt{2\pi}} y^{-\frac{1}{2}} e^{-\frac{y}{2}} \mathbb{1}_{[0, \infty)}(y). \quad (5.43)$$

We say that a random variable  $Y$  with a law given by the density is  $\chi^2$ -distributed (with one degree of freedom).

## 6. Expectation, variance and higher moments of random variables

(Reference: [1, Section 4.3-4.5, 4.9, 5.2], or [2, Sections 4,1,4.3])

### 6.1. Expectation

We now introduce the notion of the *expectation* or *expected value* of a real random variable. The idea is to somehow quantify the “typical” or “average” value of a random variable  $X$ . Let us motivate the definition with an example.

*Example 6.1.* We consider a game where we throw a die once, and get the following rewards:

- ▶ \$ 1 if the die shows 1 or 2,
- ▶ \$ 2 if the die shows 3 or 4,
- ▶ \$ 4 if the die shows 5, and
- ▶ \$ 8 if the die shows 6.

What would be a “fair price” for playing this game? We would like to stakes to be such that we do not lose money on average by playing. To describe the (random) return in one round, we can define the random variable  $X : \{1, 2, 3, 4, 5, 6\} \rightarrow \{1, 2, 4, 8\}$ , by

$$X(\omega) = \mathbb{1}_{\{1,2\}}(\omega) + 2 \cdot \mathbb{1}_{\{3,4\}}(\omega) + 4 \cdot \mathbb{1}_{\{5\}}(\omega) + 8 \cdot \mathbb{1}_{\{6\}}(\omega). \quad (6.1)$$

Assume that we play  $n \in \mathbb{N}$  times, and  $n_1, n_2, \dots, n_6$  denotes the number of 1, 2, ..., 6 that show up in  $n$  rounds. Our return (in \$) after  $n$  steps is

$$1 \cdot n_1 + 1 \cdot n_2 + 2 \cdot n_3 + 2 \cdot n_4 + 4 \cdot n_5 + 8 \cdot n_6. \quad (6.2)$$

So, the average return in one round is

$$1 \cdot \frac{n_1}{n} + 1 \cdot \frac{n_2}{n} + 2 \cdot \frac{n_3}{n} + 2 \cdot \frac{n_4}{n} + 4 \cdot \frac{n_5}{n} + 8 \cdot \frac{n_6}{n}. \quad (6.3)$$

The idea is now that for large  $n$ , the relative fractions  $\frac{n_i}{n}$  should be close to  $\frac{1}{6}$ . This gives us the value

$$\begin{aligned} \mathbf{E}[X] &= 1 \cdot \frac{1}{6} + 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 8 \cdot \frac{1}{6} \\ &= 1 \cdot \mathbf{P}[X = 1] + 2 \cdot \mathbf{P}[X = 2] + 4 \cdot \mathbf{P}[X = 4] + 8 \cdot \mathbf{P}[X = 6] \\ &= 18 \cdot \frac{1}{6} = 3. \end{aligned} \quad (6.4)$$

This somehow suggests that the “fair” price to play the game is \$ 3.

---

End of Lecture 12

This example motivates the definition of the expectation.

**Definition 6.2.** (i) Let  $X$  be a discrete real random variable with values in  $\Omega_X (\subseteq \mathbb{R})$  and let  $p_X$  be the probability mass function of its law  $\mathbf{P}_X$ . We define the *expectation of  $X$*  as

$$\mathbf{E}[X] = \sum_{\omega \in \Omega_X} \omega \cdot p_X(\omega), \quad (6.5)$$

if  $\sum_{\omega \in \Omega_X} |\omega| p_X(\omega) < \infty$ .

(ii) Let  $X$  be a continuous real random variable and let  $f_X$  be the probability density function of its law  $\mathbf{P}_X$ . We define the *expectation of  $X$*  as

$$\mathbf{E}[X] = \int_{-\infty}^{\infty} x \cdot f_X(x) dx, \quad (6.6)$$

if  $\int_{-\infty}^{\infty} |x| f_X(x) dx < \infty$ .

**Remark 6.3.** (i) For real random variables  $X$  that are neither discrete nor continuous (such as we saw in Example 5.6, (iii)), we typically cannot easily define the expectation by a formula as above. We refer to Section [2, Sections 4,1] for the case of general  $X$ .

(ii) The expectation only depends on the law  $\mathbf{P}_X$  of  $X$ . In other words: If  $X \stackrel{d}{=} Y$  and the expectation of  $X$  exists, then the expectation of  $Y$  exists as well and  $\mathbf{E}[X] = \mathbf{E}[Y]$ .

Let us give a couple of examples.

**Example 6.4.** (i) Let  $X = c \in \mathbb{R}$ . Then

$$\mathbf{E}[X] = c, \quad (6.7)$$

since  $X$  is a discrete random variable and  $\Omega_X = \{c\}$ ,  $\mathbf{P}[X = c] = 1$ .

(ii) Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $A \in \mathcal{F}$ . Then  $\mathbb{1}_A$ , defined by

$$\mathbb{1}_A(\omega) = \begin{cases} 1, & \omega \in A, \\ 0, & \omega \notin A, \end{cases} \quad (6.8)$$

is a random variable and

$$\mathbf{E}[\mathbb{1}_A] = \mathbf{P}[A]. \quad (6.9)$$

Indeed,  $\mathbb{1}_A^{-1}(B) \in \{\emptyset, A, A^c, \Omega\}$  for every  $B \in \mathcal{B}(\mathbb{R})$ , and  $\Omega_{\mathbb{1}_A} = \{0, 1\}$ . Therefore, we have

$$\mathbf{E}[\mathbb{1}_A] = 0 \cdot \mathbf{P}[\mathbb{1}_A = 0] + 1 \cdot \underbrace{\mathbf{P}[\mathbb{1}_A = 1]}_{=\mathbf{P}[A]}. \quad (6.10)$$

## 6. Expectation, variance and higher moments of random variables

- (iii)  $X \sim \text{Pois}(\lambda)$ , where  $\lambda > 0$ . The corresponding probability mass function is  $p_X(k) = \frac{\lambda^k}{k!} e^{-\lambda}$  (for  $k \in \mathbb{N}_0$ ), and

$$\mathbf{E}[X] = \sum_{k=0}^{\infty} k p_X(k) = \lambda \sum_{k=0}^{\infty} k \cdot \frac{\lambda^{k-1}}{k!} e^{-\lambda} = \lambda e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} = \lambda. \quad (6.11)$$

- (iv)  $X \sim \mathcal{U}([a, b])$  for  $a < b$ . The corresponding probability density function is  $f_X(x) = \frac{1}{b-a} \mathbb{1}_{[a,b]}(x)$ . We have

$$\mathbf{E}[X] = \int_a^b x \cdot \frac{1}{b-a} dx = \frac{1}{b-a} \left[ \frac{x^2}{2} \right]_a^b = \frac{a+b}{2}. \quad (6.12)$$

- (v)  $X \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ . Here the probability density function is  $f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ .

$$\begin{aligned} \mathbf{E}[X] &= \int_{-\infty}^{\infty} x \cdot \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \int_{-\infty}^{\infty} y \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{y^2}{2\sigma^2}} dy \\ &\quad + \mu \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{y^2}{2\sigma^2}} dy}_{= (*)} = 0 + \mu = \mu, \end{aligned} \quad (6.13)$$

where we used that the expression  $(*)$  is again a probability density function of  $\mathcal{N}(0, \sigma^2)$ , and therefore its integral is one.

- (vi)  $X \sim \text{Geo}(p)$ ,  $p \in (0, 1)$  has expectation

$$\mathbf{E}[X] = \frac{1}{p}. \quad (6.14)$$

- (vii)  $X \sim \mathcal{E}(\lambda)$ ,  $\lambda > 0$  has expectation

$$\mathbf{E}[X] = \frac{1}{\lambda}. \quad (6.15)$$

- (viii)  $X \sim \text{Bin}(n, p)$ ,  $n \in \mathbb{N}$ ,  $p \in (0, 1)$  has expectation

$$\mathbf{E}[X] = np. \quad (6.16)$$

- (ix) Consider the probability distribution on  $\mathbb{N}$  characterized by the probability mass function  $p(k) = \frac{6}{\pi^2} \frac{1}{k^2}$ .<sup>1</sup> Let  $X \sim \mathbf{P}$ . Then the expectation of  $X$  does not exist. Indeed:

$$\sum_{k=1}^{\infty} k \cdot p_X(k) = \sum_{k=1}^{\infty} \frac{6}{\pi^2} \frac{1}{k} = \infty. \quad (6.17)$$

---

<sup>1</sup>The prefactor is chosen since  $\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6}$ . This can be shown using Fourier series.

The claims in (vi) and (vii) are Exercises. Claim (viii) will be shown very easily later, after introducing the notion of independent random variables.

**Theorem 6.5.** Let  $g : \mathbb{R} \rightarrow \mathbb{R}$ .

(i) If  $X$  is a discrete real random variable with probability mass function  $(p_X(\omega))_{\omega \in \Omega_X}$ , then

$$\mathbf{E}[g(X)] = \sum_{\omega \in \Omega_X} g(\omega)p_X(\omega), \quad \text{if } \sum_{\omega \in \Omega_X} |g(\omega)|p_X(\omega) < \infty. \quad (6.18)$$

(ii) If  $X$  is a continuous real random variable with probability density function  $f_X$  (and  $g$  piecewise continuous<sup>2</sup>), then

$$\mathbf{E}[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x)dx, \quad \text{if } \int_{-\infty}^{\infty} |g(x)|f_X(x)dx < \infty. \quad (6.19)$$

*Proof.* For (i), let  $Y = g(X)$ , then  $\Omega_Y = \{y_1, y_2, \dots\}$ . Let  $A_i = g^{-1}(\{y_i\})$  for  $i \in \mathbb{N}$ . Clearly,  $\Omega_X = \bigcup_{i=1}^{\infty} A_i$ , and the  $A_j$  are pairwise disjoint. We have

$$\begin{aligned} \mathbf{E}[g(X)] &= \mathbf{E}[Y] = \sum_{i=1}^{\infty} y_i \cdot \mathbf{P}_Y[\{y_i\}] = \sum_{i=1}^{\infty} y_i \sum_{\omega_j \in A_i} p_X(\omega_j) \\ &= \sum_{i=1}^{\infty} \sum_{\omega_j \in A_i} g(\omega_j) \cdot p_X(\omega_j) \\ &= \sum_{\omega \in \Omega_X} g(\omega)p_X(\omega). \end{aligned} \quad (6.20)$$

For (ii), the proof of the general case is more complicated and relies on Measure Theory. For the special case where  $g$  is strictly increasing / strictly decreasing and differentiable, we can however use Theorem 5.13 and see that for  $Y = g(X)$

$$\mathbf{E}[g(X)] = \int_{-\infty}^{\infty} y f_Y(y) dy = \int_{-\infty}^{\infty} y f_X(g^{-1}(y)) \cdot \left| \frac{d}{dy} g^{-1}(y) \right| dy = \int_{-\infty}^{\infty} g(x) f_X(x) dx. \quad (6.21)$$

□

**Corollary 6.6.** Let  $X$  be a real random variable with finite expectation  $\mathbf{E}[X]$ . Then, for  $a, b \in \mathbb{R}$ ,

$$\mathbf{E}[aX + b] = a\mathbf{E}[X] + b. \quad (6.22)$$

*Proof.* We only prove this statement for  $X$  continuous or discrete. Without loss of generality, assume that  $X$  is continuous (the discrete case proceeds analogously). Consider  $g(x) = ax + b$ . Then

$$\mathbf{E}[g(X)] = \int_{-\infty}^{\infty} (ax + b)f_X(x)dx = a \underbrace{\int_{-\infty}^{\infty} xf_X(x)dx}_{=\mathbf{E}[X]} + b \underbrace{\int_{-\infty}^{\infty} f_X(x)dx}_{=1}. \quad (6.23)$$

A similar calculation shows that the integral  $\int_{-\infty}^{\infty} |g(x)|f_X(x)dx$  is finite. □

<sup>2</sup>This can be weakened to requiring that  $g : (\mathbb{R}, \mathcal{B}(\mathbb{R})) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  measurable.



**Theorem 6.7.** Let  $X, Y$  be two real random variables, both with finite expectation. Then

$$\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y]. \quad (6.24)$$

We will show this later after introducing the joint distribution of random variables in the next section.

## 6.2. Variance

**Definition 6.8.** Let  $X$  be a continuous or discrete real random variable.

(i) For  $k \in \mathbb{N}$ , the  $k$ -th moment of  $X$  is defined by

$$\mu_k = \mathbf{E}[X^k], \quad \text{if } \mathbf{E}[|X|^k] < \infty. \quad (6.25)$$

Note that if  $\mathbf{E}[|X|^k] < \infty$  for some  $k \in \mathbb{N}$ , then  $\mathbf{E}[|X|^\ell] < \infty$  for all  $1 \leq \ell \leq k$ , since  $|X|^\ell \leq 1 + |X|^k$ .

(ii) Assume that  $X$  has a finite second moment,  $\mathbf{E}[X^2] < \infty$ . We define the variance of  $X$  as

$$\text{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2]. \quad (6.26)$$

We also define the standard deviation of  $X$  by

$$\sigma(X) = \sqrt{\text{Var}[X]}. \quad (6.27)$$

The variance is a measure how much the distribution of  $X$  typically spreads around its expectation. If it is large, the distribution is well spread-out. If it is small, the the distribution is concentrated around the expectation. Both the notion of  $k$ th moment and variance only depend on the law  $\mathbf{P}_X$  of  $X$ , similar as for the expectation.

---

*End of Lecture 13*

**Proposition 6.9.** Let  $X$  be a continuous or discrete real random variable with  $\mathbf{E}[X^2] < \infty$  and  $a, b \in \mathbb{R}$ . Then

(i)

$$\text{Var}[X] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2 = \mu_2 - \mu_1^2. \quad (6.28)$$

(ii)

$$\text{Var}[aX + b] = a^2 \text{Var}[X]. \quad (6.29)$$

*Proof.* We first prove (i):

$$\begin{aligned} \text{Var}[X] &= \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2 - 2X\mathbf{E}[X] + \mathbf{E}[X]^2] \\ &= \mathbf{E}[X^2] - 2\mathbf{E}[X]^2 + \mathbf{E}[X]^2 = \mathbf{E}[X^2] - \mathbf{E}[X]^2. \end{aligned} \quad (6.30)$$

For (ii), we calculate

$$\begin{aligned} \text{Var}[aX + b] &= \mathbf{E}[(aX + b - a\mathbf{E}[X] - b)^2] \\ &= \mathbf{E}[a^2(X - \mathbf{E}[X])^2] = a^2 \text{Var}[X]. \end{aligned} \quad (6.31)$$

□

Let us give some examples.

*Example 6.10.* (i) Let  $X \sim \text{Ber}(p)$ ,  $p \in (0, 1)$ . We have

$$\begin{aligned}\mathbf{E}[X] &= 0 \cdot (1 - p) + 1 \cdot p = p, \\ \mathbf{E}[X^2] &= 0^2 \cdot (1 - p) + 1^2 \cdot p = p, \\ \text{Var}[X] &= p - p^2 = p(1 - p).\end{aligned}\tag{6.32}$$

(ii) Let  $X \sim \mathcal{U}(\{1, \dots, 6\})$ . We have

$$\begin{aligned}\mathbf{E}[X] &= \sum_{k=1}^6 k \cdot \mathbf{P}[X = k] = 3.5, \\ \mathbf{E}[X^2] &= \sum_{k=1}^6 k^2 \cdot \mathbf{P}[X = k] = \frac{91}{6}, \\ \Rightarrow \quad \text{Var}[X] &= \frac{91}{6} - \frac{49}{4} = \frac{70}{24} = \frac{35}{12} \approx 2.92.\end{aligned}\tag{6.33}$$

(iii) Let  $X \sim \mathcal{N}(0, 1)$ . We already saw that  $\mathbf{E}[X] = 0$  (see (6.13)). We now evaluate

$$\begin{aligned}\text{Var}[X] &= \mathbf{E}\left[(X - \underbrace{\mathbf{E}[X]}_{=0})^2\right] = \int_{-\infty}^{\infty} x^2 \varphi(x) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-\frac{1}{2}x^2} dx \\ &= \frac{1}{\sqrt{2\pi}} \left[ -xe^{-\frac{1}{2}x^2} \right]_{-\infty}^{\infty} + \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}x^2} dx \\ &= 0 + 1 = 1.\end{aligned}\tag{6.34}$$

Now let  $Y = \sigma X + \mu$  for  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ . Then

$$\begin{aligned}Y &\sim \mathcal{N}(\mu, \sigma^2) \quad (\text{by (5.37)}) \\ \text{Var}[Y] &= \sigma^2 \text{Var}[X] = \sigma^2 \quad (\text{by (6.29)}).\end{aligned}\tag{6.35}$$

We have established that

$$Z \sim \mathcal{N}(\mu, \sigma^2) \quad \Rightarrow \quad \mathbf{E}[Z] = \mu, \quad \text{Var}[Z] = \sigma^2.\tag{6.36}$$

In other words: The standard deviation of  $\mathcal{N}(\mu, \sigma^2)$  is exactly  $\sigma$ .

(iv)  $X \sim \text{Pois}(\lambda)$ . Recall that  $\mathbf{E}[X] = \lambda$ . Then

$$\begin{aligned}\mathbf{E}[X(X-1)] &= \sum_{k=0}^{\infty} k(k-1) \frac{\lambda^k}{k!} e^{-\lambda} \\ &= \lambda^2 \sum_{k=2}^{\infty} \frac{\lambda^{k-2}}{(k-2)!} e^{-\lambda} = \lambda^2 \\ \Rightarrow \quad \mathbf{E}[X^2] &= \mathbf{E}[X(X-1)] + \mathbf{E}[X] = \lambda^2 + \lambda, \\ \Rightarrow \quad \text{Var}[X] &= \mathbf{E}[X^2] - \mathbf{E}[X]^2 = \lambda^2 + \lambda - \lambda^2 = \lambda.\end{aligned}\tag{6.37}$$

## 6. Expectation, variance and higher moments of random variables

(v)  $X \sim \text{Bin}(n, p)$ ,  $n \in \mathbb{N}$ ,  $p \in (0, 1)$  has variance

$$\text{Var}[X] = np(1 - p). \quad (6.38)$$

We now argue why the variance is indeed a useful quantification how the distribution is spread out. We study the expression  $\mathbf{P}[|X - \mathbf{E}[X]| \geq \varepsilon]$  for  $\varepsilon > 0$ . The following inequality is called the *Markov inequality*.

**Theorem 6.11.** *Let  $X$  be a non-negative random variable with finite expectation. Then, for  $\varepsilon > 0$ ,*

$$\mathbf{P}[X \geq \varepsilon] \leq \frac{\mathbf{E}[X]}{\varepsilon}. \quad (6.39)$$

*Proof.* We prove the case where  $X$  is discrete. The case of continuous  $X$  is similar. Let  $\Omega_X = \{\omega_1, \omega_2, \dots\}$ . Then

$$\begin{aligned} \mathbf{P}[X \geq \varepsilon] &= \sum_{\substack{i=1 \\ \omega_i \geq \varepsilon}}^{\infty} \mathbf{P}[X = \omega_i] \\ &\leq \sum_{\substack{i=1 \\ \omega_i \geq \varepsilon}}^{\infty} \frac{\omega_i}{\varepsilon} \mathbf{P}[X = \omega_i] \\ &\leq \frac{1}{\varepsilon} \sum_{i=1}^{\infty} \omega_i \mathbf{P}[X = \omega_i] = \frac{1}{\varepsilon} \mathbf{E}[X]. \end{aligned} \quad (6.40)$$

□

From the Markov inequality, we have the following result, called the *Chebyshev inequality*.

**Theorem 6.12.** *Let  $\mathbf{E}[X^2] < \infty$ . For any  $a \in \mathbb{R}$  and  $\varepsilon > 0$ , one has*

$$\mathbf{P}[|X - a| \geq \varepsilon] \leq \frac{\mathbf{E}[(X - a)^2]}{\varepsilon^2}. \quad (6.41)$$

*In particular, for  $a = \mathbf{E}[X]$  one has*

$$\mathbf{P}[|X - \mathbf{E}[X]| \geq \varepsilon] \leq \frac{\text{Var}[X]}{\varepsilon^2}. \quad (6.42)$$

*Proof.* We apply the Markov inequality (6.39) to the non-negative random variable  $(X - a)^2$ . It follows that

$$\begin{aligned} \mathbf{P}[|X - a| \geq \varepsilon] &= \mathbf{P}[(X - a)^2 \geq \varepsilon^2] \\ &\stackrel{(6.39)}{\leq} \frac{\mathbf{E}[(X - a)^2]}{\varepsilon^2}. \end{aligned} \quad (6.43)$$

□

Chebyshev's inequality gives a bound how likely it is that the result of a random number deviates by a certain amount from its expectation. In particular, we have the following “ $k\sigma$ -rules”:

**Corollary 6.13.** *Let  $\mathbf{E}[X^2] < \infty$ .*

(i) *If  $\sigma = \sigma(X) = \sqrt{\text{Var}(X)} > 0$ , we have*

$$\mathbf{P}[|X - \mathbf{E}[X]| \geq k\sigma] \leq \frac{1}{k^2}, \quad (6.44)$$

*for  $k > 0$ . In particular:*

$$\begin{aligned} \mathbf{P}[|X - \mathbf{E}[X]| \geq 2\sigma] &\leq \frac{1}{4}, \\ \mathbf{P}[|X - \mathbf{E}[X]| \geq 3\sigma] &\leq \frac{1}{9}. \end{aligned} \quad (6.45)$$

(ii) *If  $\text{Var}[X] = 0$ , then*

$$\mathbf{P}[X = \mathbf{E}[X]] = 1. \quad (6.46)$$

*Remark 6.14.* The Chebyshev inequality is very general, but only gives very rough bounds. Consider for instance  $X \sim \mathcal{N}(\mu, \sigma^2)$  for  $\mu \in \mathbb{R}, \sigma > 0$ . Then

$$\begin{aligned} \mathbf{P}[|X - \mu| < k\sigma] &= \mathbf{P}\left[\underbrace{\left|\frac{X - \mu}{\sigma}\right|}_{\sim \mathcal{N}(0,1)} < k\right] \\ &= \mathbf{P}\left[-k < \frac{X - \mu}{\sigma} \leq k\right] = \Phi(k) - \Phi(-k) \\ &= 2\Phi(k) - 1. \end{aligned} \quad (6.47)$$

From this it follows that

$$\mathbf{P}[|X - \mu| \geq k\sigma] = 1 - (2\Phi(k) - 1) = 2 - 2\Phi(k). \quad (6.48)$$

For instance, we have

$$\begin{aligned} \mathbf{P}[|X - \mu| \geq \sigma] &\approx 2 - 2 \cdot 0.84 = 0.32, \\ \mathbf{P}[|X - \mu| \geq 2\sigma] &\approx 2 - 2 \cdot 0.98 = 0.04, \\ \mathbf{P}[|X - \mu| \geq 3\sigma] &\approx 2 - 2 \cdot 0.9987 = 0.0026. \end{aligned} \quad (6.49)$$

The last equality in (6.47) follows from the symmetry of the density  $\varphi$ .

$$\begin{aligned} \Phi(x) &= \int_{-\infty}^x \varphi(t) dt = 1 - \int_x^{\infty} \varphi(t) dt \\ &= 1 - \int_{-\infty}^{-x} \varphi(t) dt \\ &= 1 - \Phi(-x). \end{aligned} \quad (6.50)$$

## 7. Joint distributions and independence of random variables

(Reference: [1, Sections 6.1-6.2], or [2, Sections 3.3,3.4])

### 7.1. Joint distributions of random variables

In several situations, it is necessary to study multiple random variables at once and the way they are related. For instance, consider

$$\begin{aligned} X &= \text{length of a randomly chosen fish of species A,} \\ Y &= \text{age of a randomly chosen fish of species A.} \end{aligned} \tag{7.1}$$

Of course we expect that there should be a certain dependence between the values of  $X$  and  $Y$ . In this context, we want to study the joint distribution of  $X$  and  $Y$ , that is the random vector (!)  $(X, Y) \in \mathbb{R}^2$ .

Let us start with a simple example. Consider tossing a coin three times.

$$\begin{aligned} \Omega &= \{0, 1\}^3 = \{(\omega_1, \omega_2, \omega_3) ; \omega_i \in \{0, 1\}\}, \\ X(\omega) &= \omega_1, \quad \Omega_X = \{0, 1\}, \\ Y(\omega) &= \sum_{i=1}^3 \omega_i, \quad \Omega_Y = \{0, 1, 2, 3\}. \end{aligned} \tag{7.2}$$

The random vector  $(X, Y)$  then takes values in  $\Omega_{(X,Y)} = \Omega_X \times \Omega_Y \subseteq \mathbb{R}^2$ . Since  $\Omega_X \times \Omega_Y$  is a finite set, we can consider  $Z = (X, Y)$  as a random variable with values in  $\Omega_X \times \Omega_Y$  and study its distribution on

$$\begin{aligned} &(\Omega_X \times \Omega_Y, \mathcal{P}(\Omega_X \times \Omega_Y)), \text{ where } \Omega_X \times \Omega_Y = \{(x, y) ; x \in \Omega_X, y \in \Omega_Y\}. \\ \mathbf{P}_{(X,Y)}[A] &= \mathbf{P}[(X, Y)^{-1}(A)], \quad A \in \mathcal{P}(\Omega_X \times \Omega_Y). \end{aligned} \tag{7.3}$$

We equip  $(\Omega, \mathcal{P}(\Omega))$  with the uniform distribution  $\mathbf{P} = \mathcal{U}(\Omega)$ . We say that  $\mathbf{P}_{(X,Y)}$  is the *joint distribution* of  $X$  and  $Y$ . As an example:

$$\begin{aligned} A &= \{(0, 1), (1, 2)\}, \quad (X, Y)^{-1}(A) = \{(0, 1, 0), (0, 0, 1), (1, 1, 0), (1, 0, 1)\} \\ \Rightarrow \quad \mathbf{P}_{(X,Y)}[A] &= \frac{|(X, Y)^{-1}(A)|}{|\Omega|} = \frac{4}{8} = \frac{1}{2}. \end{aligned} \tag{7.4}$$

We say that

$$p_{X,Y}(x_i, y_j) = \mathbf{P}_{(X,Y)}[\{(x_i, y_j)\}], \quad (x_i, y_j) \in \Omega_X \times \Omega_Y \quad (7.5)$$

is the *joint probability mass function of  $X$  and  $Y$* . Let us compute all of the values in the example above:

$x_i \setminus y_j$	0	1	2	3
0	1/8	2/8	1/8	0
1	0	1/8	2/8	1/8

Can we obtain the probability mass function of  $X$  or  $Y$  from  $p_{(X,Y)}$ ? Let more generally be  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be two discrete random variables and let  $\Omega_X = \{x_1, x_2, \dots\}$  and  $\Omega_Y = \{y_1, y_2, \dots\}$ .

$$\begin{aligned} p_X(x_i) &= \mathbf{P}[\{X = x_i\}] = \sum_{j=1}^{\infty} \mathbf{P}[X = x_i, Y = y_j] \\ &= \sum_{j=1}^{\infty} p_{(X,Y)}(x_i, y_j). \end{aligned} \quad (7.6)$$

Similarly, we have

$$p_Y(y_j) = \sum_{i=1}^{\infty} p_{(X,Y)}(x_i, y_j). \quad (7.7)$$

The probability distributions  $\mathbf{P}_X$  and  $\mathbf{P}_Y$ , which are given in terms of  $p_X$  and  $p_Y$  are called the *marginal distributions* of  $(X, Y)$ . Turning back to the example above, can fill in the table

$x_i \setminus y_j$	0	1	2	3	$p_X(x_i)$
0	1/8	2/8	1/8	0	1/2
1	0	1/8	2/8	1/8	1/2
$p_Y(y_j)$	1/8	3/8	3/8	1/8	1

We want to make this discussion more general by

- considering more than two random variables,
- dropping the assumption that the random variables are discrete.

To this end, we want to define an equivalent notion of measurability as in Definition 5.1 for a vector  $(X_1, \dots, X_n)$  consisting of  $n$  random variables  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  defined on the same probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ . As previously, we want to speak about the law of the random vector  $(X_1, \dots, X_n)$ ,  $\mathbf{P}_{(X_1, \dots, X_n)}[B]$  for certain subsets of  $\mathbb{R}^n$ , similar as in Definition 5.3. To do this, we need to introduce a  $\sigma$ -algebra on  $\mathbb{R}^n$ .

---

*End of Lecture 14*

**Definition 7.1.** Consider the system of subsets of  $\mathbb{R}^n$  given by

$$\mathcal{E} = \{[a_1, b_1) \times [a_2, b_2) \times \dots \times [a_n, b_n); a_i < b_i, a_i, b_i \in \mathbb{R}\}. \quad (7.8)$$

The Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R}^n)$  on  $\mathbb{R}^n$  is defined by

$$\mathcal{B}(\mathbb{R}^n) = \sigma(\mathcal{E}). \quad (7.9)$$

Again, every set of our imagination is contained in  $\mathcal{B}(\mathbb{R}^n)$ , such as every countable set of points, lines, hyperplanes, cubes, cylinders, balls, ...

The following lemma ensures that we can consider a vector of random variable as a random variable with values in  $\mathbb{R}^n$ .

**Lemma 7.2.** Let  $(\Omega, \mathcal{F})$  be a measurable space and  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$  be maps. Then

$$X_1, \dots, X_n \text{ are } \mathcal{F}\text{-}\mathcal{B}(\mathbb{R}) \text{ measurable} \Leftrightarrow (X_1, \dots, X_n) : \Omega \rightarrow \mathbb{R}^n \text{ is } \mathcal{F}\text{-}\mathcal{B}(\mathbb{R}^n) \text{ measurable.} \quad (7.10)$$

The latter means that

$$(X_1, \dots, X_n)^{-1}(A) \in \mathcal{F} \quad \text{for all } A \in \mathcal{B}(\mathbb{R}^n). \quad (7.11)$$

**Definition 7.3.** Let  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be  $n$  real random variables.

(i) The joint law of  $X_1, \dots, X_n$  under  $\mathbf{P}$  is the probability measure on  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$

$$\mathbf{P}_{(X_1, \dots, X_n)}[B] = \mathbf{P}[(X_1, \dots, X_n)^{-1}(B)], \quad B \in \mathcal{B}(\mathbb{R}^n). \quad (7.12)$$

(ii) For  $x_1, \dots, x_n \in \mathbb{R}$ , we set

$$\begin{aligned} F_{X_1, \dots, X_n}(x_1, \dots, x_n) &= \mathbf{P}_{(X_1, \dots, X_n)}[(-\infty, x_1] \times \dots \times (-\infty, x_n)] \\ &= \mathbf{P}[X_1 \leq x_1, \dots, X_n \leq x_n]. \end{aligned} \quad (7.13)$$

The function  $F_{X_1, \dots, X_n}$  is called the *joint cumulative distribution function* of  $X_1, \dots, X_n$  / the *cumulative distribution function of the law of  $(X_1, \dots, X_n)$* .<sup>1</sup>

Note that by Lemma 7.2, the event  $(X_1, \dots, X_n)^{-1}(B)$  under the probability in (7.12) is in  $\mathcal{F}$ . We are now ready to define multivariate continuous distributions. For more details on the multiple integrals appearing below we refer to the Appendix A.1.

**Definition 7.4.** A (measurable) function  $f : \mathbb{R}^n \rightarrow [0, \infty)$  is called a *probability density function* if

$$\iint \dots \int_{\mathbb{R}^n} f(x_1, \dots, x_n) d^n x = 1. \quad (7.14)$$

We define  $\mathbf{P}$  to be the (unique) probability measure on  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$  that fulfills

$$\mathbf{P}[[a_1, b_1) \times [a_2, b_2) \times \dots \times [a_n, b_n)] = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_n \dots dx_2 dx_1. \quad (7.15)$$

Such a probability measure is called a (*multivariate*) *continuous distribution*.

<sup>1</sup>This wording reflects that there are two ways to view how  $X_1, \dots, X_n$  “interact”: We may view them as  $n$  different  $\mathbb{R}$ -valued functions, or we look at the single,  $\mathbb{R}^n$ -valued function  $(X_1, \dots, X_n)$ . In the former interpretation, we may think about  $\mathbf{P}_{(X_1, \dots, X_n)}$  as the joint law of  $X_1, \dots, X_n$ , in the latter we can say that  $\mathbf{P}_{(X_1, \dots, X_n)}$  is the law of the vector-valued random variable  $(X_1, \dots, X_n)$ .

*Example 7.5.* The function

$$f(x, y) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x^2 + y^2)\right) \quad (7.16)$$

is a probability density function on  $\mathbb{R}^2$ , the density of a *multivariate standard normal distribution*  $\mathcal{N}(0, I_{2 \times 2})$ . We will study such multivariate normal distributions later.

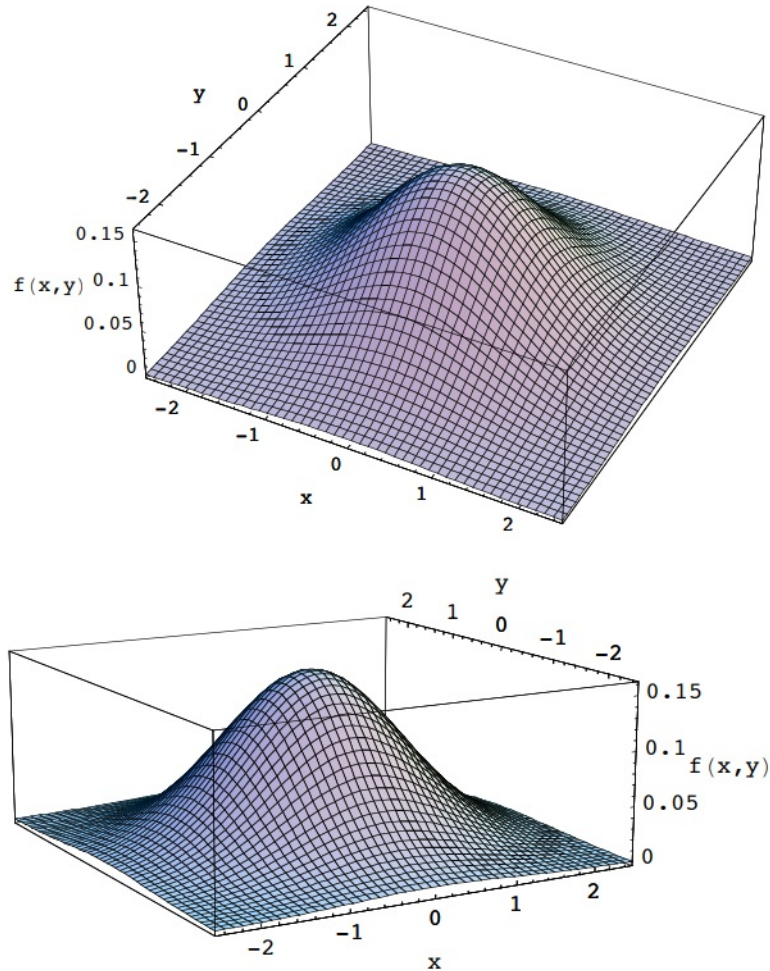


Figure 7.1.: Plot of the density  $f$  in (7.16)

*Remark 7.6.* (i) One can show that the measure  $\mathbf{P}$  in the above definition exists (using measure theory), and for any  $A \in \mathcal{B}(\mathbb{R}^n)$ , one has

$$\mathbf{P}[A] = \iint \dots \int_A f(x_1, \dots, x_n) d^n x, \quad (7.17)$$



The definition in (7.17) is understood as a Lebesgue-integral in general (not treated in this course). Typically we are only interested in sets  $A$  over which this integral can be performed as an iterated Riemann integral.

- (ii) Equality in distributions of random vectors  $(X_1, \dots, X_n)$  and  $(Y_1, \dots, Y_n)$  is defined as in Definition 5.10. Similarly as in the one-dimensional case, the law  $\mathbf{P}_{(X_1, \dots, X_n)}$  is uniquely determined by  $F_{X_1, \dots, X_n}$ .

Let us now restrict our attention to the case of either discrete or continuous real random vectors.

**Definition 7.7.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  random variables.

- (i) If all random variables are discrete with values in the countable sets  $\Omega_{X_1}, \dots, \Omega_{X_n} \subseteq \mathbb{R}$ , then  $(X_1, \dots, X_n)$  only takes values in the countable set  $\Omega_{X_1} \times \dots \times \Omega_{X_n}$ . We say that  $(X_1, \dots, X_n)$  is a *discrete (real) random vector*. The law  $\mathbf{P}_{(X_1, \dots, X_n)}$  on  $(\Omega_{X_1} \times \dots \times \Omega_{X_n}, \mathcal{P}(\Omega_{X_1} \times \dots \times \Omega_{X_n}))$  is characterized by

$$\begin{aligned} p_{X_1, \dots, X_n}(x_1, \dots, x_n) &= \mathbf{P}_{(X_1, \dots, X_n)}[\{(x_1, \dots, x_n)\}] \\ &= \mathbf{P}[X_1 = x_1, \dots, X_n = x_n], \quad x_i \in \Omega_{X_i}, 1 \leq i \leq n. \end{aligned} \quad (7.18)$$

Here,  $(p_{X_1, \dots, X_n}(x_1, \dots, x_n))_{(x_1, \dots, x_n) \in \Omega_{X_1} \times \dots \times \Omega_{X_n}}$  is the *joint probability mass function* of  $X_1, \dots, X_n$ . For  $I \subseteq \{1, \dots, n\}$  one has

$$p_{(X_i; i \in I)}(x_i; i \in I) = \mathbf{P}[X_i = x_i; i \in I] = \sum_{\substack{x_j \in \Omega_{X_j} \\ j \notin I}} p_{X_1, \dots, X_n}(x_1, \dots, x_n) \quad (7.19)$$

for the marginal joint probability mass function of  $(X_i)_{i \in I}$ .

- (ii) If the joint law  $\mathbf{P}_{(X_1, \dots, X_n)}$  of  $X_1, \dots, X_n$  is a continuous distribution on  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ , i.e. it is defined in terms of a multivariate probability density  $f_{X_1, \dots, X_n}$  as in (7.17), then  $(X_1, \dots, X_n)$  is a *continuous (real) random vector*, or  $X_1, \dots, X_n$  are *jointly continuous*. The function  $f_{X_1, \dots, X_n}$  is called the *joint probability density function* of  $X_1, \dots, X_n$ . For  $I \subseteq \{1, \dots, n\}$  one has

$$f_{(X_i; i \in I)}(x_i; i \in I) = \iint \dots \int_{\mathbb{R}^{n-|I|}} f_{X_1, \dots, X_n}(x_1, \dots, x_n) \prod_{j \notin I} (dx_j) \quad (7.20)$$

for the marginal joint probability density function of  $(X_i)_{i \in I}$ .

**Remark 7.8.** Note that it contrary to the case where all  $X_1, \dots, X_n$  are discrete, it is *not* the case that  $(X_1, \dots, X_n)$  is a continuous random vector, if all  $X_1, \dots, X_n$  are all continuous. As a counterexample, take  $X = Y \sim \mathcal{N}(0, 1)$ . Then the vector  $(X, Y) = (X, X)$  is not a continuous random vector, since

$$\mathbf{P}[(X, Y) \in \Delta] = \mathbf{P}_{(X, Y)}[\Delta] = 1, \quad \text{where } \Delta = \{(t, t); t \in \mathbb{R}\}. \quad (7.21)$$

This cannot be true if  $\mathbf{P}_{(X, Y)}$  was given in terms of a multivariate probability density  $f$ , since  $\int_{\Delta} f(x, y) dx dy = 0$ .

We want to practice a bit how to work with joint cumulative distribution functions / probability density functions / probability mass functions.

**Proposition 7.9.** *Let  $(X_1, \dots, X_n)$  be a continuous random vector. We have*

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f_{X_1, \dots, X_n}(t_1, \dots, t_n) dt_n \dots dt_1. \quad (7.22)$$

*Proof.* First, we see that

$$\mathbf{P}_{(X_1, \dots, X_n)}[A_1 \times A_2 \times \dots \times \underbrace{\{a\}}_{\text{position } j} \times \dots \times A_n] = 0, \quad (7.23)$$

for all  $A_1, \dots, A_{j-1}, A_{j+1}, \dots, A_n \in \mathcal{B}(\mathbb{R}), a \in \mathbb{R}$ .

This follows similarly as (4.14) in the one-dimensional case. Therefore:

$$\mathbf{P}_{(X_1, \dots, X_n)}[(a_1, x_1] \times (a_2, x_2] \times \dots \times (a_n, x_n]] = \int_{a_1}^{x_1} \int_{a_2}^{x_2} \dots \int_{a_n}^{x_n} f(t_1, \dots, t_n) dt_n \dots dt_2 dt_1. \quad (7.24)$$

We see that  $(a_1, x_1] \times \dots \times (-k, x_n] \subseteq (a_1, x_1] \times \dots \times (-k, x_n]$  and

$$\bigcup_{k=1}^{\infty} (a_1, x_1] \times \dots \times (-k-1, x_n] = (a_1, x_1] \times \dots \times (-\infty, x_n]. \quad (7.25)$$

By Proposition 1.15, (vi), it follows that

$$\begin{aligned} \mathbf{P}_{(X_1, \dots, X_n)}[(a_1, x_1] \times (a_2, x_2] \times \dots \times (-\infty, x_n]] \\ = \lim_{k \rightarrow \infty} \int_{a_1}^{x_1} \int_{a_2}^{x_2} \dots \int_{-k}^{x_n} f(t_1, \dots, t_n) dt_n \dots dt_2 dt_1 \\ = \int_{a_1}^{x_1} \int_{a_2}^{x_2} \dots \int_{-\infty}^{x_n} f(t_1, \dots, t_n) dt_n \dots dt_2 dt_1. \end{aligned} \quad (7.26)$$

Repeating this procedure for the other integration variables yields the claim.  $\square$

*Example 7.10.* (i) Let  $X, Y$  be distributed with the joint density

$$f_{X,Y}(x, y) = \begin{cases} 2x + 2y - 4xy, & x, y \in [0, 1], \\ 0, & x, y \notin [0, 1]. \end{cases} \quad (7.27)$$

We have that

$$\begin{aligned} f_X(x) &= \int_0^1 (2x + 2y - 4xy) dy = 2x + 1 - 2x = 1, & x \in [0, 1], \\ f_Y(y) &= \int_0^1 (2x + 2y - 4xy) dx = 1, & y \in [0, 1]. \end{aligned} \quad (7.28)$$

Therefore  $X \sim \mathcal{U}([0, 1])$ ,  $Y \sim \mathcal{U}([0, 1])$ . Let us calculate  $\mathbf{P}[X^2 \leq Y]$ :

$$\begin{aligned} \mathbf{P}[X^2 \leq Y] &= \int_0^1 \int_{x^2}^1 (2x + 2y - 4xy) dy dx \\ &= \frac{19}{30}. \end{aligned} \quad (7.29)$$

(ii) We now consider  $X, Y$  with the joint density

$$f_{X,Y}(x, y) = \mathbb{1}_{[0,1]^2}(x, y). \quad (7.30)$$

This is the *uniform distribution on  $[0, 1]^2$* , denoted by  $\mathcal{U}([0, 1]^2)$ . We see that

$$\begin{aligned} f_X(x) &= \int_0^1 1 \cdot dy = 1, & x \in [0, 1], \\ f_Y(y) &= \int_0^1 1 \cdot dx = 1, & y \in [0, 1]. \end{aligned} \quad (7.31)$$

So again,  $X \sim \mathcal{U}([0, 1])$  and  $Y \sim \mathcal{U}([0, 1])$ . We calculate again  $\mathbf{P}[X^2 \leq Y]$ :

$$\begin{aligned} \mathbf{P}[X^2 \leq Y] &= \int_0^1 \int_{x^2}^1 1 \cdot dy dx = \int_0^1 (1 - x^2) dx \\ &= \frac{2}{3}. \end{aligned} \quad (7.32)$$

This shows that even if  $X$  and  $Y$  have the same marginal densities, their joint distribution can be different from case to case.

**Theorem 7.11.** *Let  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ .*

(i) *If  $X_1, \dots, X_n$  are discrete real random variables with joint probability mass function given by  $(p_{X_1, \dots, X_n}(x_1, \dots, x_n))_{(x_1, \dots, x_n) \in \Omega_{X_1} \times \dots \times \Omega_{X_n}}$ , then*

$$\begin{aligned} \mathbf{E}[g(X_1, \dots, X_n)] &= \sum_{(x_1, \dots, x_n) \in \Omega_{X_1} \times \dots \times \Omega_{X_n}} g(x_1, \dots, x_n) p_{X_1, \dots, X_n}(x_1, \dots, x_n), \\ &\text{if } \sum_{(x_1, \dots, x_n) \in \Omega_{X_1} \times \dots \times \Omega_{X_n}} |g(x_1, \dots, x_n)| p_{X_1, \dots, X_n}(x_1, \dots, x_n) < \infty. \end{aligned} \quad (7.33)$$

(ii) *If  $(X_1, \dots, X_n)$  is a continuous random vector and  $X_1, \dots, X_n$  have the joint probability density function  $f_{X_1, \dots, X_n}$  (and  $g$  is measurable), then*

$$\begin{aligned} \mathbf{E}[g(X_1, \dots, X_n)] &= \int \int \dots \int_{\mathbb{R}^n} g(x_1, \dots, x_n) f_{X_1, \dots, X_n}(x_1, \dots, x_n) d^n x, \\ &\text{if } \int \int \dots \int_{\mathbb{R}^n} |g(x_1, \dots, x_n)| f_{X_1, \dots, X_n}(x_1, \dots, x_n) d^n x < \infty. \end{aligned} \quad (7.34)$$

*Proof.* This is analogous to Theorem 6.5.  $\square$

We can now prove the additivity of the expectation claimed in the previous section:

*Proof of Theorem 6.7.* We first assume that  $X$  and  $Y$  are both discrete. Then, by applying (7.33) with the function  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $g(x, y) = x + y$ :

$$\begin{aligned} \mathbf{E}[X + Y] &= \sum_{x \in \Omega_X, y \in \Omega_Y} (x + y) p_{X,Y}(x, y) \\ &= \sum_{x \in \Omega_X} x \underbrace{\sum_{y \in \Omega_Y} p_{X,Y}(x, y)}_{=p_X(x)} + \sum_{y \in \Omega_Y} y \underbrace{\sum_{x \in \Omega_X} p_{X,Y}(x, y)}_{=p_Y(y)} \\ &= \mathbf{E}[X] + \mathbf{E}[Y]. \end{aligned} \quad (7.35)$$

Now let  $X$  and  $Y$  be jointly continuous. Then

$$\begin{aligned} \mathbf{E}[X + Y] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + y) f_{X,Y}(x, y) dx dy \\ &= \int_{-\infty}^{\infty} x \underbrace{\left( \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \right)}_{=f_X(x)} dx + \int_{-\infty}^{\infty} y \underbrace{\left( \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx \right)}_{=f_Y(y)} dy \\ &= \mathbf{E}[X] + \mathbf{E}[Y]. \end{aligned} \quad (7.36)$$

We omit the proof of all other cases.  $\square$

---

*End of Lecture 15*

## 7.2. Independence of random variables

Consider two random variables  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . If  $A, B \in \mathcal{B}(\mathbb{R})$ , we can consider the two events

$$X^{-1}(A) \quad \text{and} \quad Y^{-1}(B). \quad (7.37)$$

We already have a notion of independence of these events: If we cannot obtain any information about  $X^{-1}(A)$  from knowing whether  $Y^{-1}(B)$  has occurred, then the two events are independent. If this is the case for *any* choice of the sets  $A$  and  $B$ , we will not be able to obtain any information from  $X$  about  $Y$  or vice versa, since  $\{X^{-1}(A); A \in \mathcal{B}(\mathbb{R})\}$  contains all possible events that we can observe by knowing  $X$ <sup>2</sup>, and  $\{Y^{-1}(B); B \in \mathcal{B}(\mathbb{R})\}$  contains all possible events we can observe by knowing  $Y$ . We come to the general definition of stochastic independence of random variables.

---

<sup>2</sup>We will not elaborate on this further, but this is called the *generated  $\sigma$ -algebra of  $X$* , written  $\sigma(X)$ . For instance, if  $X : \Omega = \{1, 2, 3, 4, 5, 6\} \rightarrow \mathbb{R}$ ,  $X(\omega) = \mathbb{1}_{\{1,3,5\}}(\omega)$  (i.e.  $X(\omega) = 1$  if  $\omega$  is odd and  $X(\omega) = 0$  if  $\omega$  is even). Then  $\sigma(X) = \{\emptyset, \{1, 3, 5\}, \{2, 4, 6\}, \Omega\}$ . This means that the “information” contained in knowing the value of  $X(\omega)$  is exactly whether or not  $\omega$  is even, namely whether  $\omega \in \{1, 3, 5\}$  or  $\omega \in \{2, 4, 6\}$ .

**Definition 7.12.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space.

- (i) The random variables  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  are *independent*, if

$$\mathbf{P}[X_1 \in A_1, X_2 \in A_2, \dots, X_n \in A_n] = \prod_{i=1}^n \mathbf{P}[X_i \in A_i], \quad \text{for } A_1, \dots, A_n \in \mathcal{B}(\mathbb{R}). \quad (7.38)$$

- (ii) The random variables  $(X_i; i \in \mathcal{I})$ , where  $X_i : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  and  $\mathcal{I}$  is an arbitrary set, are *independent*, if for all finite sets  $\{i_1, \dots, i_n\} \subseteq \mathcal{I}$  (with  $i_1, \dots, i_n$  pairwise distinct):

$$\mathbf{P}[X_{i_1} \in A_1, X_{i_2} \in A_2, \dots, X_{i_n} \in A_n] = \prod_{j=1}^n \mathbf{P}[X_{i_j} \in A_j], \quad \text{for } A_1, \dots, A_n \in \mathcal{B}(\mathbb{R}). \quad (7.39)$$

Of course (i) is a special case of (ii) with  $\mathcal{I} = \{1, \dots, n\}$ . Let us discuss some special cases:

**Proposition 7.13.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  random variables.

- (i)  $X_1, \dots, X_n$  are independent if and only if

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n F_{X_i}(x_i), \quad \text{for all } x_1, \dots, x_n \in \mathbb{R}. \quad (7.40)$$

- (ii) Assume that  $X_1, \dots, X_n$  are discrete random variables. Then  $X_1, \dots, X_n$  are independent if and only if

$$\mathbf{P}[X_1 = x_1, \dots, X_n = x_n] = \prod_{i=1}^n \mathbf{P}[X_i = x_i], \quad \text{for all } x_i \in \Omega_{X_i}, 1 \leq i \leq n. \quad (7.41)$$

- (iii) Assume that  $X_1, \dots, X_n$  are continuous random variables. Then  $X_1, \dots, X_n$  are independent if and only if

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i), \quad \text{for all } x_1, \dots, x_n \in \mathbb{R}. \quad (7.42)$$

*Proof.* For (i), we assume that  $X_1, \dots, X_n$  are independent, then

$$\mathbf{P}[X_1 \leq x_1, \dots, X_n \leq x_n] = \prod_{i=1}^n \mathbf{P}[X_i \leq x_i]. \quad (7.43)$$

The other direction requires tools from measure theory, and we omit it.

For (ii), we consider  $A_1 \in \mathcal{P}(\Omega_{X_1}), \dots, A_n \in \mathcal{P}(\Omega_{X_n})$ . Then

$$\begin{aligned} \mathbf{P}[X_1 \in A_1, \dots, X_n \in A_n] &= \sum_{x_1 \in A_1, \dots, x_n \in A_n} \mathbf{P}[X_1 = x_1, \dots, X_n = x_n] \\ &= \sum_{x_1 \in A_1, \dots, x_n \in A_n} \prod_{i=1}^n \mathbf{P}[X_i = x_i] = \prod_{i=1}^n \mathbf{P}[X_i \in A_i]. \end{aligned} \quad (7.44)$$

For (iii), assume that  $X_1, \dots, X_n$  are independent, then by part (i) we have (7.40). The claim then follows by taking the partial derivatives with respect to  $x_1, \dots, x_n$ . For the other direction, let us assume that (7.42) holds. Take  $x_1, \dots, x_n \in \mathbb{R}$ , then by Proposition 7.9,

$$\begin{aligned} F_{X_1, \dots, X_n}(x_1, \dots, x_n) &= \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f_{X_1, \dots, X_n}(t_1, \dots, t_n) dt_n \dots dt_1 \\ &= \prod_{i=1}^n \int_{-\infty}^{x_i} f_{X_i}(t_i) dt_i = \prod_{i=1}^n F_{X_i}(x_i). \end{aligned} \quad (7.45)$$

It follows from part (i) that  $X_1, \dots, X_n$  are independent.  $\square$

*Example 7.14.* Consider again Example 7.10. The random variables  $X$  and  $Y$  in part (i) are not independent. On the other hand, the random variables  $X$  and  $Y$  in part (ii), since  $f_{X,Y}(x, y) = f_X(x)f_Y(y)$  for all  $x, y \in \mathbb{R}$ .

Intuitively, if the random variables  $X$  and  $Y$  are independent, then  $Z = f(X)$  and  $W = g(Y)$  should also be independent for any functions  $f, g : \mathbb{R} \rightarrow \mathbb{R}$ . This is indeed the case, and we have the *propagation of independence*:

**Proposition 7.15.** *Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  independent real random variables. Furthermore, assume that  $h_i : (\mathbb{R}, \mathcal{B}(\mathbb{R})) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  are measurable functions. Then the functions*

$$h_1(X_1), \dots, h_n(X_n) \text{ are independent.} \quad (7.46)$$

*Proof.* For all  $B_i \in \mathcal{B}(\mathbb{R})$ , we have

$$\begin{aligned} \mathbf{P}[h_1(X_1) \in B_1, \dots, h_n(X_n) \in B_n] &= \mathbf{P}[X_1 \in h_1^{-1}(B_1), \dots, X_n \in h_n^{-1}(B_n)] \\ &= \prod_{i=1}^n \mathbf{P}[X_i \in h_i^{-1}(B_i)] = \prod_{i=1}^n \mathbf{P}[h_i(X_i) \in B_i]. \end{aligned} \quad (7.47)$$

$\square$

*Example 7.16.* If  $X$  and  $Y$  are independent, then  $Z = X^2 - \sin(X)$  and  $W = \cos(Y^2) \mathbb{1}_{\{Y \geq 0\}}$  are independent as well.

The proof of Proposition 7.15 works in a similar way for vector-valued random variables. With this, we can show that any functions of two disjoint subsets of  $\{X_1, \dots, X_n\}$  are again independent, if  $X_1, \dots, X_n$  are independent.

**Example 7.17.** If  $X_1, X_2, X_3, X_4, X_5$  are independent, then  $Z = \sin(X_1^2 + X_3^2)$  and  $Y = X_5 \exp(X_2)$  are independent.

**Theorem 7.18.** Let  $X$  and  $Y$  be independent real random variables with  $\mathbf{E}[X^2] < \infty$  and  $\mathbf{E}[Y^2] < \infty$ . Then

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]. \quad (7.48)$$

*Proof.* We first assume that  $X$  and  $Y$  are both discrete. Then, by applying (7.33) with the (measurable) function  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $g(x, y) = xy$ :

$$\begin{aligned} \mathbf{E}[XY] &= \sum_{x \in \Omega_X, y \in \Omega_Y} xy \underbrace{p_{X,Y}(x, y)}_{=p_X(x)p_Y(y)} \\ &= \sum_{x \in \Omega_X} xp_X(x) \sum_{y \in \Omega_Y} yp_Y(y) \\ &= \mathbf{E}[X]\mathbf{E}[Y]. \end{aligned} \quad (7.49)$$

Now let  $X$  and  $Y$  be jointly continuous. Then

$$\begin{aligned} \mathbf{E}[XY] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \underbrace{xy f_{X,Y}(x, y)}_{f_X(x)f_Y(y)} dx dy \\ &= \int_{-\infty}^{\infty} x f_X(x) \left( \int_{-\infty}^{\infty} y f_Y(y) dy \right) dx \\ &= \mathbf{E}[X]\mathbf{E}[Y]. \end{aligned} \quad (7.50)$$

We omit the proof of all other cases.  $\square$

**Example 7.19.** Let  $X$  and  $Y$  be the side-lengths of a random rectangle. We assume that  $X, Y \sim \mathcal{U}([0, 1])$  and  $X$  is independent of  $Y$ . The expected value of the area of the rectangle is

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y] = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}. \quad (7.51)$$

Clearly, independence is crucial for (7.48) to hold. Indeed, if we look at  $\tilde{Y} = X$ , then

$$\mathbf{E}[X\tilde{Y}] = \mathbf{E}[X^2] = \int_0^1 x^2 dx = \frac{1}{3}. \quad (7.52)$$

The previous example is a special case of a general concept, that will be very important in the rest of the course.

**Definition 7.20.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. A sequence  $(X_n)_{n \in \mathbb{N}}$  of random variables  $X_1, X_2, \dots : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  is called *independent and identically distributed*, or *i.i.d.*, if  $(X_i)_{i \in \mathbb{N}}$  are independent and for every  $i, j \in \mathbb{N}$ ,  $\mathbf{P}_{X_i} = \mathbf{P}_{X_j}$ . Being i.i.d. is defined similarly for finitely many random variables  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ .

Can we always find an i.i.d. sequence of random variables with a given distribution?

**Theorem 7.21.** Let  $\mathbf{Q}$  be a probability measure on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . Then there exists a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$  and a sequence  $(X_n)_{n \in \mathbb{N}}$  of random variables  $X_1, X_2, \dots : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  that are i.i.d. with  $\mathbf{P}_{X_n} = \mathbf{Q}$  for all  $n \in \mathbb{N}$ .

---

End of Lecture 16

## 8. Operations with random variables

(Reference: [1, Sections 6.3, 6.6])

In this chapter we will study some important ways to combine independent random variables.

### 8.1. Extremes

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X_1, \dots, X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be i.i.d. real random variables. We are interested in the distribution of the maximum

$$X^{(n)}(\omega) = \max\{X_1, \dots, X_n\}(\omega) = \max\{X_1(\omega), \dots, X_n(\omega)\} \quad (8.1)$$

and the minimum

$$X_{(n)}(\omega) = \min\{X_1, \dots, X_n\}(\omega) = \min\{X_1(\omega), \dots, X_n(\omega)\}. \quad (8.2)$$

Note that  $X^{(n)}$  and  $X_{(n)}$  are indeed random variables, by Remark 5.2, (i).

**Proposition 8.1.** *The distribution functions of  $X^{(n)}$  and  $X_{(n)}$  are given by*

$$\begin{aligned} F_{X^{(n)}}(x) &= (F_{X_1}(x))^n, \\ F_{X_{(n)}}(x) &= 1 - (1 - F_{X_1}(x))^n, \quad x \in \mathbb{R}, \end{aligned} \quad (8.3)$$

respectively.

*Proof.* We have for  $x \in \mathbb{R}$ ,

$$\begin{aligned} F_{X^{(n)}}(x) &= \mathbf{P}[\max\{X_1, \dots, X_n\} \leq x] = \mathbf{P}\left[\bigcap_{i=1}^n \{X_i \leq x\}\right] \\ &= \prod_{i=1}^n \mathbf{P}[X_i \leq x] \quad (\text{by independence}) \\ &= (F_{X_1}(x))^n. \end{aligned} \quad (8.4)$$

Similarly, we see that for  $x \in \mathbb{R}$ ,

$$\begin{aligned} 1 - F_{X_{(n)}}(x) &= \mathbf{P}[\min\{X_1, \dots, X_n\} > x] = \mathbf{P}\left[\bigcap_{i=1}^n \{X_i > x\}\right] \\ &= \prod_{i=1}^n \mathbf{P}[X_i > x] \quad (\text{by independence}) \\ &= (1 - F_{X_1}(x))^n. \end{aligned} \quad (8.5)$$

□



As an example, consider  $X_1, \dots, X_n \sim \mathcal{E}(\lambda)$  i.i.d. random variables. These may be interpreted as waiting times for independent, exponentially distributed events. What is the law of the maximum and minimum? We have

$$\begin{aligned} F_{X(n)}(x) &= (F_{X_1}(x))^n = (1 - e^{-\lambda x})^n \mathbb{1}_{[0, \infty)}(x), \\ F_{X(n)}(x) &= 1 - (1 - F_{X_1}(x))^n = 1 - (1 - (1 - e^{-\lambda x}) \mathbb{1}_{[0, \infty)}(x))^n = (1 - e^{-\lambda n x}) \mathbb{1}_{[0, \infty)}(x). \end{aligned} \quad (8.6)$$

Here we used the expression for the cumulative distribution function of an exponentially distributed random variable (5.16). In particular, we see that  $X(n) \sim \mathcal{E}(\lambda n)$ .

## 8.2. Sums of independent random variables

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Assume that  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  are independent real random variables with known distributions  $\mathbf{P}_X$  and  $\mathbf{P}_Y$ . What is the distribution of  $X + Y$ ? We will assume that  $X$  and  $Y$  are both discrete or both continuous. Let us start with the former (easier) case.

**Proposition 8.2.** *Assume that  $X, Y$  are independent discrete real random variables with probability mass functions  $(p_X(k))_{k \in \Omega_X}$  and  $(p_Y(\ell))_{\ell \in \Omega_Y}$ . Then the sum  $Z = X + Y$  is again discrete and its probability mass function is given by*

$$p_Z(k) = \sum_{\ell \in \Omega_Y} p_X(k - \ell) p_Y(\ell), \quad (8.7)$$

for  $k \in \Omega_X + \Omega_Y = \{a + \ell; a \in \Omega_X, \ell \in \Omega_Y\}$ .

*Proof.* Clearly,  $Z$  is discrete, since it takes only values in  $\Omega_X + \Omega_Y$ , which is countable. Now take  $k \in \Omega_X + \Omega_Y$ . Then (since  $\bigcup_{\ell \in \Omega_Y} \{Y = \ell\} = \Omega$ ):

$$\begin{aligned} \mathbf{P}[Z = k] &= \mathbf{P} \left[ \bigcup_{\ell \in \Omega_Y} \{X = k - \ell, Y = \ell\} \right] \\ &= \sum_{\ell \in \Omega_Y} \underbrace{\mathbf{P}[X = k - \ell]}_{=p_X(k-\ell)} \cdot \underbrace{\mathbf{P}[Y = \ell]}_{=p_Y(\ell)}, \end{aligned} \quad (8.8)$$

where we used the independence assumption in the second step.  $\square$

*Example 8.3.* (i) Let  $X \sim \text{Bin}(n, p)$  and  $Y \sim \text{Bin}(m, p)$  for  $n, m \in \mathbb{N}$  and  $p \in (0, 1)$  be independent. Then  $X + Y \sim \text{Bin}(n + m, p)$ . Indeed, the random variable  $X + Y$  can

attain values in  $\{0, \dots, n + m\}$ , so let  $0 \leq k \leq n + m$ . Then

$$\begin{aligned}
 p_Z(k) &= \sum_{\ell=0}^m p_X(\ell) p_Y(k - \ell) \\
 &= \sum_{\ell=0}^k \binom{n}{k - \ell} p^{k - \ell} (1 - p)^{n - (k - \ell)} \binom{m}{\ell} p^{\ell} (1 - p)^{m - \ell} \\
 &= \sum_{\ell=0}^k \binom{n}{k - \ell} \binom{m}{\ell} p^k (1 - p)^{n + m - k} \\
 &= \binom{n + m}{k} p^k (1 - p)^{n + m - k}.
 \end{aligned} \tag{8.9}$$

In the last line we used the *Vandermonde identity*

$$\binom{n + m}{k} = \sum_{\ell=0}^k \binom{n}{k - \ell} \binom{m}{\ell}, \quad n, m, k \in \mathbb{N}_0. \tag{8.10}$$

There are various ways to prove this identity. A particularly simple proof goes as follows: Consider the polynomial functions  $x \mapsto (1 + x)^r$ , for  $r \in \mathbb{N}$ . We see that

$$\begin{aligned}
 (1 + x)^{n + m} &= (1 + x)^n \cdot (1 + x)^m \\
 \Rightarrow \sum_{k=0}^{n + m} \binom{n + m}{k} x^k &= \left( \sum_{i=0}^n \binom{n}{i} x^i \right) \left( \sum_{j=0}^m \binom{m}{j} x^j \right) \\
 &= \sum_{k=0}^{n + m} \left( \sum_{\ell=0}^k \binom{n}{k - \ell} \binom{m}{\ell} \right) x^k.
 \end{aligned} \tag{8.11}$$

By extracting the coefficient  $k$  on both sides, we obtain (8.10).

- (ii) Let  $X \sim \text{Pois}(\lambda)$  and  $Y \sim \text{Pois}(\mu)$ ,  $\lambda, \mu > 0$  be independent. Then  $X + Y \sim \text{Pois}(\lambda + \mu)$ . We leave this as an Exercise.

In particular, part (i) shows that for i.i.d.  $X_1, \dots, X_n \sim \text{Ber}(p)$ , the random variable  $S_n = \sum_{i=1}^n X_i$  fulfills  $S_n \sim \text{Bin}(n, p)$ . In particular, we see

$$\mathbf{E}[S_n] = n\mathbf{E}[X_1] = np, \tag{8.12}$$

where we used (6.32), and Theorem 6.7.

We now move to the sum of continuous, independent random variables.

**Proposition 8.4.** Assume that  $X, Y$  are independent continuous real random variables with probability density functions  $f_X$  and  $f_Y$ . Then the sum  $Z = X + Y$  is again continuous and its probability density function is given by

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(z - y) f_Y(y) dy. \tag{8.13}$$

*Proof.* Let  $z \in \mathbb{R}$ . We have

$$\begin{aligned}
 F_Z(z) &= \mathbf{P}[X + Y \leq z] \\
 &= \mathbf{P}[(X, Y) \in \{(x, y) \in \mathbb{R}^2; x + y \leq z\}] \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) \mathbb{1}_{\{(x,y) \in \mathbb{R}^2; x+y \leq z\}} dx dy \\
 &= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{z-y} f_{X,Y}(x, y) dx \right) dy \\
 &= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{z-y} f_X(x) dx \right) f_Y(y) dy \\
 &= \int_{-\infty}^{\infty} F_X(z - y) f_Y(y) dy.
 \end{aligned} \tag{8.14}$$

We differentiate with respect to  $z$  on both sides to obtain

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(z - y) f_Y(y) dy. \tag{8.15}$$

□

*Remark 8.5.* If  $X$  and  $Y$  are independent real random variables defined on  $(\Omega, \mathcal{F}, \mathbf{P})$  with laws  $\mathbf{P}_1 = \mathbf{P}_X$  and  $\mathbf{P}_2 = \mathbf{P}_Y$ , we say that the law  $\mathbf{Q} = \mathbf{P}_{X+Y}$  of the sum of  $X$  and  $Y$  is the *convolution* of the laws of  $X$  and  $Y$ , and we write

$$\mathbf{Q} = \mathbf{P}_1 * \mathbf{P}_2 = \mathbf{P}_X * \mathbf{P}_Y. \tag{8.16}$$

For instance, one can show that if  $X, Y \sim \mathcal{E}(\lambda)$  are independent, then  $X + Y \sim \Gamma(2, \lambda)$  (the *Gamma*-distribution with parameters  $\alpha = 2$  and  $\lambda$ , characterized by the density  $f(x) = \frac{1}{\Gamma(\alpha)} \lambda^\alpha x^{\alpha-1} e^{-\lambda x} \mathbb{1}_{[0, \infty)}(x)$  for  $\lambda, \alpha > 0$ , with  $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$ ). Thus one has

$$\mathcal{E}(\lambda) * \mathcal{E}(\lambda) = \Gamma(2, \lambda). \tag{8.17}$$

For other distributions, one has (see Exercises and Example 8.3)

$$\begin{aligned}
 \text{Bin}(n, p) * \text{Bin}(m, p) &= \text{Bin}(n + m, p), & n, m \in \mathbb{N}, p \in (0, 1), \\
 \text{Pois}(\lambda) * \text{Pois}(\mu) &= \text{Pois}(\lambda + \mu), & \lambda, \mu > 0, \\
 \mathcal{N}(\mu_1, \sigma_1^2) * \mathcal{N}(\mu_2, \sigma_2^2) &= \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2), & \mu_1, \mu_2 \in \mathbb{R}, \sigma_1, \sigma_2 > 0, \\
 \Gamma(\alpha_1, \lambda) * \Gamma(\alpha_2, \lambda) &= \Gamma(\alpha_1 + \alpha_2, \lambda), & \alpha_1, \alpha_2, \lambda > 0.
 \end{aligned} \tag{8.18}$$

---

*End of Lecture 17*

## 9. More on expectation

(Reference: [1, Section 8.5])

In this short section, we state two useful inequalities that are relevant for applications.

### 9.1. Jensen's inequality

In this subsection, we will show a useful inequality for the expectation. We start with a reminder on convexity.

**Definition 9.1.** Consider  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ . We say that  $\varphi$  is *convex* if for every  $x, y \in \mathbb{R}$  and  $\lambda \in [0, 1]$ , one has

$$\varphi(\lambda x + (1 - \lambda)y) \leq \lambda\varphi(x) + (1 - \lambda)\varphi(y). \quad (9.1)$$

We say that  $\varphi$  is *concave* if  $-\varphi$  is convex.

In particular, if  $\varphi$  is twice differentiable, it is convex if and only if  $\varphi''(x) \geq 0$  for every  $x \in \mathbb{R}$ .

*Example 9.2.* The functions  $\varphi_j : \mathbb{R} \rightarrow \mathbb{R}$  with  $1 \leq j \leq 3$  with  $\varphi_1(x) = x^2$ ,  $\varphi_2(x) = e^{rx}$  for  $r > 0$ ,  $\varphi_3(x) = |x|$  are convex.

We now state *Jensen's inequality*.

**Lemma 9.3.** Let  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  be convex and  $X$  a real random variable with  $\mathbf{E}[|X|] < \infty$  and  $\mathbf{E}[|\varphi(X)|] < \infty$ . Then

$$\mathbf{E}[\varphi(X)] \geq \varphi(\mathbf{E}[X]). \quad (9.2)$$

*Proof.* We only present the proof in the case that  $\varphi$  is smooth. In this case, one can perform a Taylor expansion of  $\varphi$  around the value  $\mathbf{E}[X]$ :

$$\varphi(x) = \varphi(\mathbf{E}[X]) + \varphi'(\mathbf{E}[X])(x - \mathbf{E}[X]) + \frac{\varphi''(\xi)(x - \mathbf{E}[X])^2}{2}, \quad (9.3)$$

for some  $\xi$  between  $x$  and  $\mathbf{E}[X]$ . By convexity,  $\varphi''(\xi) \geq 0$ , and so

$$\varphi(X) \geq \varphi(\mathbf{E}[X]) + \varphi'(\mathbf{E}[X])(X - \mathbf{E}[X]). \quad (9.4)$$

Upon taking expectation, we have

$$\mathbf{E}[\varphi(X)] \geq \varphi(\mathbf{E}[X]) + \varphi'(\mathbf{E}[X])\mathbf{E}[X - \mathbf{E}[X]] = \varphi(\mathbf{E}[X]). \quad (9.5)$$

□

A special case is the function  $\varphi(x) = x^2$ , where one has the (known) inequality

$$\mathbf{E}[X^2] \geq (\mathbf{E}[X])^2. \quad (9.6)$$

## 9.2. Hölder's inequality

Suppose that  $X, Y : (\Omega, \mathcal{F}, \mathbf{P}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  are real random variables. By considering the example that  $X = Y$  with law on  $\mathbb{N}$  given by

$$\mathbf{P}[X = n] = \frac{1}{c_3} \cdot \frac{1}{n^3}, \quad c_3 = \sum_{n=1}^{\infty} \frac{1}{n^3} \simeq 1.202057..., \quad (9.7)$$

we see that even though  $\mathbf{E}[X] = \mathbf{E}[Y]$  exist, the expectation of  $\mathbf{E}[X \cdot Y] = \mathbf{E}[X^2]$  is infinite. Are there situations where we can conclude that  $X \cdot Y$  has a finite expectation, based on information on the moments of  $X$  and  $Y$ ? It turns out that this is the case, and we have the *Hölder inequality*:

**Theorem 9.4.** *Let  $p, q > 1$  with  $\frac{1}{p} + \frac{1}{q} = 1$ . Let  $X, Y$  be two real random variables, defined on the same probability space and assume that the expectations of  $|X|^p$  and  $|Y|^q$  exist. Then  $\mathbf{E}[|XY|] < \infty$  and*

$$\mathbf{E}[|X \cdot Y|] \leq (\mathbf{E}[|X|^p])^{\frac{1}{p}} \cdot (\mathbf{E}[|Y|^q])^{\frac{1}{q}}. \quad (9.8)$$

*Proof.* We first show that for all  $a, b \geq 0$ , one has

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q} \quad (\text{Young's inequality}). \quad (9.9)$$

Indeed, fix  $b > 0$  (the inequality is trivial if  $a = 0$  or  $b = 0$ ) and set  $f(x) = \frac{x^p}{p} + \frac{b^q}{q} - xb$ , then  $f$  is twice differentiable, and

$$f'(x) = x^{p-1} - b, \quad f''(x) = (p-1)x^{p-2}. \quad (9.10)$$

In particular,  $f$  attains its (unique) minimum at  $x_0 = b^{\frac{1}{p-1}}$ . Since  $q = \frac{p}{p-1}$ ,  $x_0^p = b^q$ , so

$$f(x_0) = \left(\frac{1}{p} + \frac{1}{q}\right) b^q - b^{\frac{1}{p-1}} b = 0, \quad (9.11)$$

showing (9.9). Suppose now that the expectations of  $|X|^p$  and  $|Y|^q$  are not zero (otherwise the claim becomes trivial). We can then apply (9.9) to

$$a = \frac{|X(\omega)|}{(\mathbf{E}[|X|^p])^{\frac{1}{p}}}, \quad b = \frac{|Y(\omega)|}{(\mathbf{E}[|Y|^q])^{\frac{1}{q}}}. \quad (9.12)$$

This yields

$$\frac{|X(\omega)Y(\omega)|}{(\mathbf{E}[|X|^p])^{\frac{1}{p}} \cdot (\mathbf{E}[|Y|^q])^{\frac{1}{q}}} \leq \frac{|X(\omega)|^p}{p\mathbf{E}[|X|^p]} + \frac{|Y(\omega)|^q}{q\mathbf{E}[|Y|^q]}. \quad (9.13)$$

In particular, the expectation of  $|XY|$  exists, and the inequality (9.8) follows from taking the expectation and using that  $\frac{1}{p} + \frac{1}{q} = 1$ .  $\square$

---

*End of Lecture 18*

## 10. Covariance and correlation

(Reference: [1, Section 7.4], or [2, Section 4.3])

As we have already seen, the expectation is additive. We want to find a similar expression for the variance of the sum of two random variables. This motivates the following definition.

**Definition 10.1.** Let  $X$  and  $Y$  be two real random variables defined on some probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ , fulfilling  $\mathbf{E}[X^2] < \infty$  and  $\mathbf{E}[Y^2] < \infty$ .

(i) The *covariance* of  $X$  and  $Y$  is defined as

$$\text{Cov}[X, Y] = \mathbf{E}[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])]. \quad (10.1)$$

(ii) The *correlation coefficient* of  $X$  and  $Y$  is defined as

$$\rho(X, Y) = \begin{cases} \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}}, & \text{Var}[X] \neq 0, \text{Var}[Y] \neq 0, \\ 0, & \text{else.} \end{cases} \quad (10.2)$$

We collect some properties of the covariance and correlation coefficient.

**Proposition 10.2.** Let  $X$  and  $Y$  be two real random variables defined on the same probability space, fulfilling  $\mathbf{E}[X^2] < \infty$  and  $\mathbf{E}[Y^2] < \infty$ .

(i) The covariance can be written as

$$\text{Cov}[X, Y] = \mathbf{E}[XY] - \mathbf{E}[X]\mathbf{E}[Y] = \text{Cov}[Y, X]. \quad (10.3)$$

(ii)

$$\text{Var}[X] = \text{Cov}[X, X]. \quad (10.4)$$

(iii) If  $X$  and  $Y$  are independent, then  $\text{Cov}[X, Y] = \rho(X, Y) = 0$ .

(iv) If  $X = \pm Y$ , then

$$\text{Cov}[X, Y] = \pm \text{Var}[X] \text{ and } \rho(X, Y) = \pm 1 \text{ (if } \text{Var}[X] \neq 0 \text{ and } \text{Var}[Y] \neq 0). \quad (10.5)$$

*Proof.* For (i), we use

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbf{E}[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])] \\ &= \mathbf{E}[XY] - \mathbf{E}[X\mathbf{E}[Y]] - \mathbf{E}[Y\mathbf{E}[X]] + \mathbf{E}[X]\mathbf{E}[Y] \\ &= \mathbf{E}[XY] - \mathbf{E}[X]\mathbf{E}[Y] = \text{Cov}[Y, X]. \end{aligned} \quad (10.6)$$

The claim (ii) follows directly from the definition of the variance (6.26).

For (iii), we see that

$$\text{Cov}[X, Y] = \mathbf{E}[XY] - \mathbf{E}[X]\mathbf{E}[Y] \stackrel{(7.48)}{=} \mathbf{E}[X]\mathbf{E}[Y] - \mathbf{E}[X]\mathbf{E}[Y] = 0. \quad (10.7)$$

(iv) The first part is obvious. Now

$$\rho(X, \pm X) = \frac{\pm \text{Var}[X]}{\left(\sqrt{\text{Var}[X]}\right)^2} \stackrel{\text{Var}[X] > 0}{=} \pm 1. \quad (10.8)$$

□

We now prove the bilinearity of the covariance.

**Proposition 10.3.** *Let  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_m$  random variables on the same probability space, with finite second moment. For  $a, b, c_1, \dots, c_n, d_1, \dots, d_m \in \mathbb{R}$  one has*

$$\text{Cov}\left[a + \sum_{i=1}^n c_i X_i, b + \sum_{j=1}^m d_j Y_j\right] = \sum_{i=1}^n \sum_{j=1}^m c_i d_j \text{Cov}[X_i, Y_j]. \quad (10.9)$$

*Proof.* We use (10.3) to see that for every random variables  $Z_1, Z_2, Z_3$  with finite second moment:

$$\begin{aligned} \text{Cov}[Z_1 + Z_2, Z_3] &= \mathbf{E}[(Z_1 + Z_2)Z_3] - \mathbf{E}[Z_1 + Z_2]\mathbf{E}[Z_3] \\ &= \mathbf{E}[Z_1 Z_3] + \mathbf{E}[Z_2 Z_3] - \mathbf{E}[Z_1]\mathbf{E}[Z_3] - \mathbf{E}[Z_2]\mathbf{E}[Z_3] \\ &= \text{Cov}[Z_1, Z_3] + \text{Cov}[Z_2, Z_3]. \end{aligned} \quad (10.10)$$

Moreover, for  $\lambda \in \mathbb{R}$ :

$$\begin{aligned} \text{Cov}[\lambda Z_1, Z_2] &= \mathbf{E}[\lambda Z_1 Z_2] - \mathbf{E}[\lambda Z_1]\mathbf{E}[Z_2] \\ &= \lambda \text{Cov}[Z_1, Z_2]. \end{aligned} \quad (10.11)$$

Finally, we have

$$\text{Cov}[\lambda, Z_1] = \mathbf{E}[\lambda Z_1] - \mathbf{E}[\lambda]\mathbf{E}[Z_1] = 0. \quad (10.12)$$

Equation (10.9) follows from the previous three displays and the symmetry of Cov. □

A special case of this is the *Bienaymé formula* for the variance of the sum of random variables.

**Corollary 10.4.** *Let  $X_1, \dots, X_n$  be real random variables defined on the same probability space, with finite second moment.*

(i)

$$\text{Var}\left[\sum_{i=1}^n X_i\right] = \sum_{i,j=1}^n \text{Cov}[X_i, X_j] = \sum_{i=1}^n \text{Var}[X_i] + \sum_{\substack{i=1 \\ i \neq j}}^n \text{Cov}[X_i, X_j]. \quad (10.13)$$

(ii) If  $X_1, \dots, X_n$  are uncorrelated (i.e.  $\text{Cov}[X_i, X_j] = 0$  for every  $1 \leq i \neq j \leq n$ ), then

$$\text{Var} \left[ \sum_{i=1}^n X_i \right] = \sum_{i=1}^n \text{Var}[X_i]. \quad (10.14)$$

Let us combine part (ii) with Example 8.3, (i): Since for i.i.d.  $X_1, \dots, X_n \sim \text{Ber}(p)$ , the random variable  $S_n = \sum_{i=1}^n X_i$  fulfills  $S_n \sim \text{Bin}(n, p)$ , we have

$$\text{Var}[S_n] = n\text{Var}[X_1] = np(1-p), \quad (10.15)$$

where we used (6.32).

**Theorem 10.5.** Let  $X, Y$  be real random variables defined on some probability space with  $\mathbf{E}[X^2], \mathbf{E}[Y^2] < \infty$ .

(i)  $|\rho(X, Y)| \leq 1$ .

(ii)  $\rho(X, Y) = \pm 1$  if and only if there are  $a, b \in \mathbb{R}$ ,  $a \neq 0$  with

$$\mathbf{P}[Y = aX + b] = 1, \quad (10.16)$$

and we have

$$\begin{aligned} a &> 0, & \text{if } \rho(X, Y) &= 1, \\ a &< 0, & \text{if } \rho(X, Y) &= -1. \end{aligned} \quad (10.17)$$

*Proof.* We first show (i). The variance of any random variables is nonnegative, so

$$\begin{aligned} 0 &\leq \text{Var} \left( \frac{X}{\sqrt{\text{Var}(X)}} + \frac{Y}{\sqrt{\text{Var}(Y)}} \right) \\ &= \text{Var} \left( \frac{X}{\sqrt{\text{Var}(X)}} \right) + \text{Var} \left( \frac{Y}{\sqrt{\text{Var}(Y)}} \right) + 2\text{Cov} \left( \frac{X}{\sqrt{\text{Var}(X)}}, \frac{Y}{\sqrt{\text{Var}(Y)}} \right) \\ &= \frac{\text{Var}(X)}{\text{Var}(X)} + \frac{\text{Var}(Y)}{\text{Var}(Y)} + \frac{2\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = 2(1 + \rho(X, Y)), \end{aligned} \quad (10.18)$$

where we used Corollary 10.4, (i). We see that  $\rho(X, Y) \geq -1$ . By replacing  $X$  by  $-X$ , we also have  $\rho(X, Y) \leq 1$ .

For (ii), we first assume that  $\rho(X, Y) = -1$ . Then, by (10.18),

$$\mathbf{P} \left[ \frac{X}{\sqrt{\text{Var}(X)}} + \frac{Y}{\sqrt{\text{Var}(Y)}} = c \right] = 1, \quad (10.19)$$

where we used Corollary 6.13. In other words, we have

$$\mathbf{P}[Y = aX + b] = 1, \quad a = -\frac{\sqrt{\text{Var}(Y)}}{\sqrt{\text{Var}(X)}} < 0. \quad (10.20)$$



If  $\rho(X, Y) = 1$ , we can do the same calculation but with  $\text{Var} \left( \frac{X}{\sqrt{\text{Var}(X)}} - \frac{Y}{\sqrt{\text{Var}(Y)}} \right) = \dots = 0$ .  $\square$

**Definition 10.6.** Let  $X_1, \dots, X_n$  be random variables with finite second moment, defined on the same probability space. The matrix

$$\Sigma = \Sigma(X) = (\text{Cov}(X_i, X_j))_{i,j=1}^n \quad (10.21)$$

is called the *covariance matrix* of the random vector  $X = (X_1, \dots, X_n)$ . For  $u \in \mathbb{R}^n$ , one has

$$\text{Var}(u \cdot X) = u \cdot \Sigma u. \quad (10.22)$$

Here for  $u, v \in \mathbb{R}^n$ ,  $u \cdot v = \sum_{i=1}^n u_i v_i$  denotes the standard scalar product in  $\mathbb{R}^n$ .

Note that the matrix  $\Sigma$  is symmetric and non-negative definite.

---

*End of Lecture 19*

# 11. Conditional distributions and conditional expectation

(Reference: [1, Sections 6.4–6.5, 7.5–7.6],

In this chapter, we define conditional distributions and conditional expectations of random variables. Recall that we already discussed the *joint* distribution of random variables  $X$  and  $Y$ . A natural question is then:

*What is the distribution of  $X$ , if the value of  $Y$  is known?*

For instance, we may think of the following:

- $X$  is the value of the first outcome when rolling a fair die three times,  $Y$  is the sum of outcomes, what is  $\mathbf{P}[X = x|Y = y]$ ?
- $X$  is the length of a randomly chosen fish of species  $A$ ,  $Y$  is its age. What is  $\mathbf{P}[X \in B|Y = y]$ , for  $B \in \mathcal{B}(\mathbb{R})$ ?

## 11.1. Discrete conditional distributions

Suppose that  $X, Y$  are two discrete real random variables with joint probability mass function  $p_{X,Y}(x, y)$  for  $x \in \Omega_X, y \in \Omega_Y$ . Clearly, one has

$$\mathbf{P}[X = x|Y = y] = \frac{\mathbf{P}[X = x, Y = y]}{\mathbf{P}[Y = y]} = \frac{p_{X,Y}(x, y)}{p_Y(y)}, \quad (11.1)$$

if  $p_Y(y) > 0$ . This leads to the following fact.

**Proposition 11.1.** *Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  two discrete real random variables with probability mass functions  $(p_X(x))_{x \in \Omega_X}$  and  $(p_Y(y))_{y \in \Omega_Y}$  and joint probability mass function  $p_{X,Y}(x, y)_{(x,y) \in \Omega_X \times \Omega_Y}$ . We define*

$$p_{X|Y=y}(x) = \begin{cases} \frac{p_{X,Y}(x, y)}{p_Y(y)}, & \text{if } p_Y(y) > 0, \\ 0, & \text{else.} \end{cases} \quad (11.2)$$

*Then  $p_{X|Y=y}(x) \geq 0$ ,  $\sum_{x \in \Omega_X} p_{X|Y=y}(x) = 1$  for any  $y \in \Omega_Y$  with  $p_Y(y) > 0$ , and*

$$p_X(x) = \sum_{y; p_Y(y) > 0} p_{X|Y=y}(x) p_Y(y). \quad (11.3)$$

*Proof.* This all follows from the fact that  $\mathbf{P}[\cdot | \{Y = y\}]$  is a probability measure as long as  $\mathbf{P}[Y = y] = p_Y(y) > 0$  (see Proposition 2.4), and the law of  $X$  under  $\mathbf{P}[\cdot | \{Y = y\}]$  is discrete with probability mass function  $\mathbf{P}[X = x | Y = y] = \frac{p_{X,Y}(x,y)}{p_Y(y)}$ . Equation (11.3) follows from the law of total probability.  $\square$

We say that  $\mathbf{P}[\cdot | Y = y]_X$  is the *conditional distribution / law of  $X$  given  $Y = y$* , and its probability mass function  $p_{X|Y=y}$  is the *conditional probability mass function of (the law of)  $X$  given  $Y = y$* . We also record that if  $X$  and  $Y$  are independent, then

$$p_{X|Y=y}(x) = p_X(x) \quad \text{for all } x \in \Omega_X, p_Y(y) > 0. \quad (11.4)$$

*Example 11.2.* (i) Recall the example of tossing a coin three times from (7.2) in Chapter 7, i.e.

$$\begin{aligned} \Omega &= \{0, 1\}^3 = \{(\omega_1, \omega_2, \omega_3); \omega_i \in \{0, 1\}\}, \\ X(\omega) &= \omega_1, \quad \Omega_X = \{0, 1\}, \\ Y(\omega) &= \sum_{i=1}^3 \omega_i, \quad \Omega_Y = \{0, 1, 2, 3\}. \end{aligned} \quad (11.5)$$

We found for the joint probability mass function:

$x_i \setminus y_j$	0	1	2	3	$p_X(x_i)$
0	1/8	2/8	1/8	0	1/2
1	0	1/8	2/8	1/8	1/2
$p_Y(y_j)$	1/8	3/8	3/8	1/8	1

With this, we have for instance:

$$\begin{aligned} p_{X|Y=1}(0) &= \frac{p_{X,Y}(0,1)}{p_Y(1)} = \frac{2/8}{2/8 + 1/8} = \frac{2}{3}, \\ p_{X|Y=1}(1) &= \frac{p_{X,Y}(1,1)}{p_Y(1)} = \frac{1/8}{2/8 + 1/8} = \frac{1}{3}. \end{aligned}$$

(ii) Suppose that  $X$  and  $Y$  are independent Poisson random variables with parameters  $\lambda > 0$  and  $\mu > 0$ , respectively. We want to calculate the conditional distribution of  $X$  given  $X + Y = n$ .

$$\begin{aligned} \mathbf{P}[X = k | X + Y = n] &= \frac{\mathbf{P}[X = k] \mathbf{P}[Y = n - k]}{\mathbf{P}[X + Y = n]} \\ &= \frac{e^{-\lambda} \lambda^k}{k!} \cdot \frac{e^{-\mu} \mu^{n-k}}{(n-k)!} \cdot \frac{n!}{e^{-(\lambda+\mu)} (\lambda + \mu)^n} \\ &= \binom{n}{k} \cdot \left( \frac{\lambda}{\lambda + \mu} \right)^k \cdot \left( \frac{\mu}{\lambda + \mu} \right)^{n-k}. \end{aligned}$$

Here we used that  $X + Y \sim \text{Pois}(\lambda + \mu)$ , see Example 8.3, (ii). In other words,  $X \sim \text{Bin}(n, \frac{\lambda}{\lambda + \mu})$  under  $\mathbf{P}[\cdot | X + Y = n]$ . This property is known as *splitting of Poisson random variables*.

## 11.2. Continuous conditional distributions

We wish to extend the discussion of the previous section to the case where  $X, Y$  are jointly continuous real random variables with joint density  $f_{X,Y}$ . Unfortunately,  $\mathbf{P}[Y = y] = 0$ , and we cannot really “condition” on the event  $\{Y = y\}$ . However, the following heuristics still gives a valuable definition of a conditional density: Suppose  $a < b$  and  $f_Y(y) > 0$ , then

$$\begin{aligned} \mathbf{P}[X \in [a, b] | Y = y] &\approx \mathbf{P}[X \in [a, b] | Y \in [y, y + \Delta y)], \quad \Delta y \text{ small} \\ &= \frac{\mathbf{P}[X \in [a, b], Y \in [y, y + \Delta y)]}{\mathbf{P}[Y \in [y, y + \Delta y)]} \\ &= \frac{\int_a^b \int_y^{y+\Delta y} f_{X,Y}(x, z) dz dx}{\int_{y+\Delta y} f_Y(z) dz} \\ &\approx \frac{\int_a^b f_{X,Y}(x, y) \Delta y dx}{f_Y(y) \Delta y} = \frac{\int_a^b f_{X,Y}(x, y) dx}{f_Y(y)}. \end{aligned}$$

This suggests that we may define the conditional density of  $X$  given  $Y = y$  as a quotient of the joint density  $f_{X,Y}$  and the marginal density  $f_Y$ .

**Proposition 11.3.** *Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  two jointly continuous real random variables joint probability density function  $f_{X,Y}$ . We define*

$$f_{X|Y=y}(x) = \begin{cases} \frac{f_{X,Y}(x,y)}{f_Y(y)}, & \text{if } f_Y(y) > 0, \\ 0, & \text{else.} \end{cases} \quad (11.6)$$

Then  $f_{X|Y=y}(x) \geq 0$ ,  $\int_{-\infty}^{\infty} f_{X|Y=y}(x) dx = 1$  for any  $y \in \mathbb{R}$  with  $f_Y(y) > 0$ , and

$$f_X(x) = \int_{-\infty}^{\infty} f_{X|Y=y}(x) f_Y(y) dy. \quad (11.7)$$

The quantity  $f_{X|Y=y}$  is called the *conditional probability density function of (the law of)  $X$  given  $Y = y$* .

---

End of Lecture 20

## 11.3. Conditional expectation

We have seen in the previous section how to define the conditional probability mass function  $p_{X|Y=y}$  or the conditional probability density function  $f_{X|Y=y}$ . We can also calculate expectations with these quantities.

**Definition 11.4.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  two real random variables.

- (i) Suppose both  $X$  and  $Y$  are discrete with values in  $\Omega_X$  and  $\Omega_Y$  and joint probability mass function  $p_{X,Y}$ . For any  $y \in \Omega_Y$  with  $p_Y(y) > 0$ , we define the *conditional expectation of  $X$  given  $Y = y$*  as

$$\mathbf{E}[X|Y = y] = \sum_{x \in \Omega_X} x \cdot p_{X|Y=y}(x), \quad (11.8)$$

if  $\sum_{x \in \Omega_X} |x| \cdot p_{X|Y=y}(x) < \infty$ .

- (ii) Suppose that  $X, Y$  are jointly continuous with joint probability density function  $f_{X,Y}$ . For  $y \in \mathbb{R}$  with  $f_Y(y) > 0$ , we define the *conditional expectation of  $X$  given  $Y = y$*  as

$$\mathbf{E}[X|Y = y] = \int_{-\infty}^{\infty} x f_{X|Y=y}(x), \quad (11.9)$$

if  $\int_{-\infty}^{\infty} |x| \cdot f_{X|Y=y}(x) < \infty$ .

**Remark 11.5.** In the discrete case, the definition of the conditional expectation of  $X$  given  $Y = y$  is exactly the expectation of  $X$  under the conditional probability measure  $\mathbf{P}[\cdot | Y = y]$ , since the the distribution of  $X$  under this probability measure  $(\mathbf{P}[\cdot | Y = y])_X$  has probability mass function  $p_{X|Y=y}(x)$ .

**Example 11.6.** Suppose that the joint probability density function of  $X$  and  $Y$  is given by

$$f_{X,Y}(x, y) = \frac{e^{-\frac{x}{y}} e^{-y}}{y} \mathbb{1}_{(0, \infty)^2}(x, y). \quad (11.10)$$

We want to calculate  $\mathbf{E}[X|Y = y]$  and start by computing the conditional density. For  $x, y > 0$ :

$$\begin{aligned} f_{X|Y=y}(x) &= \frac{f_{X,Y}(x, y)}{f_Y(y)} \\ &= \frac{(1/y) e^{-\frac{x}{y}} e^{-y}}{\int_0^{\infty} (1/y) e^{-\frac{x}{y}} e^{-y} dx} = \frac{1}{y} e^{-\frac{x}{y}}. \end{aligned} \quad (11.11)$$

So the conditional distribution of  $X$ , given  $Y = y$  is the exponential distribution  $\mathcal{E}(\frac{1}{y})$ , and thus

$$\mathbf{E}[X|Y = y] = y. \quad (11.12)$$

Note that in the discrete case, the set of  $\omega \in \Omega$  for which  $p_y(Y(\omega)) = 0$  has probability zero. We then set

$$\mathbf{E}[X|Y](\omega) = \sum_{y \in \Omega_Y, p_Y(y) > 0} \mathbf{E}[X|Y = y] \mathbb{1}_{\{y=Y(\omega)\}}. \quad (11.13)$$

With this definition,  $\mathbf{E}[X|Y] : \Omega \rightarrow \mathbb{R}$  is also a random variable. A similar argument can be performed in the continuous case (requiring some measure theory).

**Theorem 11.7.** Let  $X, Y, Z : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be real random variables. Assume that  $\mathbf{E}[|X|] < \infty$  and  $\mathbf{E}[|Y|] < \infty$ . Then

- (i)  $\mathbf{E}[\mathbf{E}[X|Y]] = \mathbf{E}[X]$ ,
- (ii)  $X, Y$  stochastically independent, then  $\mathbf{E}[X|Y] = \mathbf{E}[X]$ ,<sup>1</sup>
- (iii)  $\mathbf{E}[Xh(Y)|Y] = h(Y)\mathbf{E}[X|Y]$  for any (measurable) function  $h$ , in particular

<sup>1</sup>Apart from a set  $N$  with  $\mathbf{P}[N] = 0$ .

- $\mathbf{E}[\text{const.}|Y] = \text{const.},$
- $\mathbf{E}[h(Y)|Y] = h(Y),$

(iv)  $\mathbf{E}[\alpha X + \beta Y|Z] = \alpha \mathbf{E}[X|Z] + \beta \mathbf{E}[Y|Z]$  (*linearity*) for  $\alpha, \beta \in \mathbb{R}$ .

*Proof.* We only consider the case where  $X$  and  $Y$  are both discrete.

(i) We calculate

$$\begin{aligned} \mathbf{E}[\mathbf{E}[X|Y]] &= \sum_{y \in \Omega_Y, p_Y(y) > 0} \mathbf{E}[X|Y = y] p_Y(y) \\ &= \sum_{y \in \Omega_Y, p_Y(y) > 0} \sum_{x \in \Omega_X} x \underbrace{p_{X|Y=y}(x) p_Y(y)}_{=p_{X,Y}(x,y)} = \mathbf{E}[X]. \end{aligned} \quad (11.14)$$

(ii) Note that we have

$$\mathbf{E}[X|Y = y] = \sum_{x \in \Omega_X} x \underbrace{p_{X|Y=y}(x)}_{=p_X(x)} = \mathbf{E}[X]. \quad (11.15)$$

(iii) Suppose that  $\omega \in \Omega$  and  $y = Y(\omega)$ . Then we have (for  $p_Y(y) > 0$ ):

$$\begin{aligned} \mathbf{E}[X \cdot h(Y)|Y](\omega) &= \mathbf{E}[X \cdot h(Y)|Y = y] = h(y) \mathbf{E}[X|Y = y] \\ &= h(Y(\omega)) \mathbf{E}[X|Y](\omega). \end{aligned} \quad (11.16)$$

(iv) Follows from a calculation. □

---

*End of Lecture 21*

As an application, we discuss *random sums*.

*Example 11.8.* Suppose that in an insurance company, a random (integer) number  $N$  of claims is made during a year. We assume that the size of the claims form an i.i.d. sequence  $(X_n)_{n \in \mathbb{N}}$  of (non-negative) real random variables which we assume to be independent from  $N$ . The total amount of claims to the company is

$$S = \sum_{j=1}^N X_j. \quad (11.17)$$

Note that this is really a random variable of the form

$$\omega \mapsto S(\omega) = \sum_{j=1}^{N(\omega)} X_j(\omega) = \sum_{j=1}^{\infty} X_j(\omega) \mathbb{1}_{\{N(\omega) \geq j\}} \quad (11.18)$$

We are interested in  $\mathbf{E}[S]$  and  $\text{Var}[S]$ .

$$\begin{aligned}\mathbf{E}[S|N = n] &= \mathbf{E}\left[\sum_{j=1}^N X_j \middle| N = n\right] = \sum_{j=1}^n \mathbf{E}[X_j] = n\mathbf{E}[X_1], \\ \Rightarrow \quad \mathbf{E}[S|N] &= N\mathbf{E}[X_1] \\ \Rightarrow \quad \mathbf{E}[S] &= \mathbf{E}[\mathbf{E}[S|N]] = \mathbf{E}[N]\mathbf{E}[X_1].\end{aligned}\tag{11.19}$$

Moreover, we have

$$\begin{aligned}\mathbf{E}[S^2|N = n] &= \mathbf{E}\left[\left(\sum_{j=1}^N X_j\right)^2 \middle| N = n\right] = \mathbf{E}\left[\left(\sum_{j=1}^n X_j\right)^2\right] \\ &= \text{Var}\left[\sum_{j=1}^n X_j\right] + \left(\mathbf{E}\left[\sum_{j=1}^n X_j\right]\right)^2 \\ &= n\text{Var}[X_1] + n^2(\mathbf{E}[X_1])^2.\end{aligned}\tag{11.20}$$

From here, it follows that

$$\begin{aligned}\mathbf{E}[S^2|N] &= N\text{Var}[X_1] + N^2(\mathbf{E}[X_1])^2 \\ \Rightarrow \quad \mathbf{E}[S^2] &= \mathbf{E}[N]\text{Var}[X_1] + \mathbf{E}[N^2](\mathbf{E}[X_1])^2 \\ \Rightarrow \quad \text{Var}[S] &= \mathbf{E}[S^2] - (\mathbf{E}[S])^2 = \mathbf{E}[N]\text{Var}[X_1] + \text{Var}[N](\mathbf{E}[X_1])^2.\end{aligned}\tag{11.21}$$

Compare this to the variance of a *fixed* sum  $\tilde{S} = \sum_{j=1}^n X_j$ , which is  $\text{Var}[\tilde{S}] = n\text{Var}[X_1]$ .

The conditional expectation  $\mathbf{E}[Y|X]$  gives the “average” over many independent realizations of  $Y$ , when the value  $X$  is given. We can make this more rigorous, and discuss as a final part of this section the *best prediction* / *best linear prediction* of random variables.

**Lemma 11.9.** *Suppose  $X, Y$  are real random variables, defined on the same probability space and  $\mathbf{E}[X^2] < \infty$ ,  $\mathbf{E}[Y^2] < \infty$ .*

- (i) *The conditional expectation  $\mathbf{E}[Y|X] =: h_*(X)$  minimizes the expression  $\mathbf{E}[(Y - h(X))^2]$  among all  $h : \mathbb{R} \rightarrow \mathbb{R}$  (measurable).*
- (ii) *The values  $a, b \in \mathbb{R}$  that minimize  $\mathbf{E}[(Y - (a + bX))^2]$  are given by*

$$b = \frac{\text{Cov}[X, Y]}{\text{Var}[X]}, \quad a = \mathbf{E}[Y] - \frac{\text{Cov}[X, Y]}{\text{Var}[X]}\mathbf{E}[X].\tag{11.22}$$

*Proof.* (i) is given as an exercise. For (ii), consider

$$\begin{aligned}\mathbf{E}[(Y - a - bX)^2] &= \text{Var}[Y - a - bX] + (\mathbf{E}[Y] - a - b\mathbf{E}[X])^2 \\ &= \text{Var}[Y - bX] + (\mathbf{E}[Y] - a - b\mathbf{E}[X])^2.\end{aligned}$$

Both terms are non-negative, so we see immediately (by minimizing the second term) that

$$a = \mathbf{E}[Y] - b\mathbf{E}[X].\tag{11.23}$$

On the other hand:

$$\begin{aligned} \text{Var}[Y - bX] &= \text{Var}[Y] + b^2\text{Var}[X] - 2b\text{Cov}[X, Y] \\ \frac{\partial}{\partial b}\text{Var}[Y - bX] &= 2b\text{Var}[X] - 2\text{Cov}[X, Y] = 0 \quad \Leftrightarrow \quad b = \frac{\text{Cov}[X, Y]}{\text{Var}[X]}. \end{aligned}$$

This concludes the proof.  $\square$

We also remark that the *mean square error* for these optimal values of  $a$  and  $b$  is given by

$$\mathbf{E}[(Y - a - bX)^2] = \text{Var}[Y] - \frac{\text{Cov}[X, Y]^2}{\text{Var}[X]} = \text{Var}[Y](1 - \rho(X, Y)^2). \quad (11.24)$$

Note that  $\rho(X, Y)$  close to  $\pm 1$  means that we typically only make a small error when approximating  $Y$  as a linear function of  $X$ .

Lemma 11.9 thus tells us: if we want to *predict* the value of  $Y$ , *knowing* the value of  $X$ :

- The best way in general to do it is to choose  $\mathbf{E}[Y|X]$ ,
- The best way to do it with a linear function is  $a + bX$  with  $a, b$  given in (11.22).

*Example 11.10.* Let  $Y = X^2 + X + Z$  with  $X, Z \sim \mathcal{N}(0, 1)$  i.i.d., suppose we want to find a good prediction of  $Y$  knowing  $X$ .

- The best prediction is given by the conditional expectation:

$$\mathbf{E}[Y|X] = X^2 + X + \mathbf{E}[Z|X] = X^2 + X,$$

and the mean square error is

$$\mathbf{E}[(Y - X^2 - X)^2] = \mathbf{E}[Z^2] = 1.$$

- The best *linear* prediction of  $Y$  is given by  $a + bX$  with

$$\begin{aligned} b &= \frac{\text{Cov}[X, Y]}{\text{Var}[X]} = \text{Cov}[X, X^2 + X] = \text{Cov}[X, X^2] + \text{Var}[X] = 1, \\ a &= \mathbf{E}[Y] - \mathbf{E}[X] = \mathbf{E}[X^2] = 1, \end{aligned}$$

so  $1 + X$  is the best linear prediction. Its mean square error is

$$\begin{aligned} \mathbf{E}[(Y - (1 + X))^2] &= \mathbf{E}[(X^2 + Z - 1)^2] \\ &= \mathbf{E}[X^4] + 2\mathbf{E}[X^2Z] + \mathbf{E}[Z^2] - 2\mathbf{E}[X^2] - 2\mathbf{E}[Z] + 1 = 3. \end{aligned}$$



## 12. Generating functions

(Reference: [1, Sections 7.7], or [2, Section 4.4])

In this chapter, we will define three functions which characterize probability distributions  $\mathbf{P}_X$  of a real random variable  $X$  and are also helpful in the calculation of moments.

**Definition 12.1.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  a real random variable. We define the *moment generating function*  $\psi_X$  of  $X$  by

$$\psi_X(t) = \mathbf{E}[e^{tX}] = \begin{cases} \sum_{k \in \Omega_X} e^{tk} \mathbf{P}[X = k], & \text{if } X \text{ is discrete,} \\ \int_{-\infty}^{\infty} e^{tx} f_X(x) dx, & \text{if } X \text{ is continuous and its law has density } f_X, \end{cases} \quad (12.1)$$

for  $t \in \mathbb{R}$ , if this expectation exists.

Note that  $\psi_X(t)$  always exists for  $t = 0$ , but may not exist for general  $t \neq 0$ . Here are some properties of moment generating functions:

**Lemma 12.2.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X, Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  independent real random variables, and  $\lambda, \mu \in \mathbb{R}$ .

(i) Assume that  $\psi_X(t)$ ,  $\psi_Y(t)$  and  $\psi_{X+Y}(t)$  exist for some  $t \in \mathbb{R}$ , then

$$\psi_{X+Y}(t) = \psi_X(t)\psi_Y(t). \quad (12.2)$$

(ii) Assume that  $\psi_X(\lambda t)$  exists, then also  $\psi_{\lambda X + \mu}(t)$  exists and

$$\psi_{\lambda X + \mu}(t) = e^{\mu t} \psi_X(\lambda t). \quad (12.3)$$

(iii) Suppose that  $\psi_X(t)$  exists for all  $t \in (-\varepsilon, \varepsilon)$ ,  $\varepsilon > 0$ . Then

$$\psi_X^{(n)}(0) = \left. \frac{d^n}{dt^n} \psi_X(t) \right|_{t=0} = \mathbf{E}[X^n]. \quad (12.4)$$

*Proof.* (i) We have

$$\psi_{X+Y}(t) = \mathbf{E}[e^{t(X+Y)}] = \mathbf{E}[e^{tX} e^{tY}] \stackrel{(7.49)}{=} \mathbf{E}[e^{tX}] \mathbf{E}[e^{tY}] = \psi_X(t) \psi_Y(t), \quad (12.5)$$

where we used that since  $X$  and  $Y$  are independent, so are  $e^{tX}$  and  $e^{tY}$  (see Proposition 7.15).

(ii) We see that

$$\psi_{\lambda X + \mu}(t) = \mathbf{E}[e^{\lambda t X} e^{\mu t}] = e^{\mu t} \psi_X(\lambda t). \quad (12.6)$$

(iii) Let  $X$  be discrete, and we assume that we can interchange differentiation and summation<sup>1</sup>:

$$\begin{aligned}\psi_X^{(n)}(0) &= \sum_{k \in \Omega_X} \frac{d^n}{dt^n} e^{tk} \mathbf{P}[X = k] \Big|_{t=0} \\ &= \sum_{k \in \Omega_X} k^n e^{tk} \mathbf{P}[X = k] \Big|_{t=0} = \sum_{k \in \Omega_X} k^n \mathbf{P}[X = k] = \mathbf{E}[X^n].\end{aligned}\tag{12.7}$$

If  $X$  is continuous, one has (interchanging integration and differentiation)

$$\begin{aligned}\psi_X^{(n)}(0) &= \int_{-\infty}^{\infty} \frac{d^n}{dt^n} e^{tx} f_X(x) dx \Big|_{t=0} \\ &= \int_{-\infty}^{\infty} x^n e^{tx} f_X(x) dx \Big|_{t=0} = \int_{-\infty}^{\infty} x^n f_X(x) dx = \mathbf{E}[X^n].\end{aligned}\tag{12.8}$$

□

The third property justifies the name moment generating function. Note that the assumption in part (i) the above Lemma can typically be justified by Hölder's inequality: Indeed, if  $\psi_X(2t)$  and  $\psi_Y(2t)$  exist for  $t > 0$ , then so does  $\psi_{X+Y}(t)$  (see Theorem 9.4 with  $p = q = 2$ ), similarly for (ii). Also note that the moment generating function exists for all  $t \in \mathbb{R}$  if  $X$  is bounded.

*Example 12.3.* Consider  $S \sim \text{Bin}(n, p)$  for  $n \in \mathbb{N}$  and  $p \in [0, 1]$ , then  $\psi_S(t)$  exists for every  $t \in \mathbb{R}$  and

$$\psi_S(t) = \sum_{k=0}^n e^{tk} \binom{n}{k} p^k (1-p)^{n-k} = (pe^t + 1 - p)^n.\tag{12.9}$$

By differentiation, we have

$$\psi'_S(t) = n(pe^t + 1 - p)^{n-1} pe^t \stackrel{(12.4)}{\Rightarrow} \mathbf{E}[S] = np.\tag{12.10}$$

With this we reproduce (8.12).

We now quote without proof another fundamental property of the moment generating function, which is that it characterizes the law of  $X$ .

**Theorem 12.4.** *Suppose that  $X$  and  $Y$  are two real random variables with moment generating functions  $\psi_X$  and  $\psi_Y$ , which both exist in  $(-\varepsilon, \varepsilon)$  for some  $\varepsilon > 0$ . Then*

$$\psi_X(t) = \psi_Y(t) \text{ for all } t \in (-\varepsilon, \varepsilon) \quad \Leftrightarrow \quad \mathbf{P}_X = \mathbf{P}_Y.\tag{12.11}$$

---

*End of Lecture 22*

This gives another characterization of equality in law, besides the characterization with cumulative distribution functions (see (5.30)), or the equality of the probability mass functions / probability density functions. To give an example, we consider the case of Binomial distributions.

---

<sup>1</sup>This is immediate if  $\Omega_X$  is finite, and can be justified generally.

**Example 12.5.** Let  $X \sim \text{Bin}(n, p)$  and  $Y \sim \text{Bin}(m, p)$  with  $n, m \in \mathbb{N}$  and  $p \in [0, 1]$ . The moment generating functions of  $X$  and  $Y$  exist for all  $t \in \mathbb{R}$ . Then we see, using (12.9) and Lemma 12.2, (i),

$$\psi_{X+Y}(t) = (pe^t + (1-p))^n \cdot (pe^t + (1-p))^m = (pe^t + (1-p))^{n+m}, \quad (12.12)$$

and this is the moment generating function for a  $\text{Bin}(n+m, p)$ , so  $X+Y \sim \text{Bin}(n+m, p)$  (which we already saw in Example 8.3, (i)).

For completeness, we also give the definitions and properties for two other relevant functions.

**Definition 12.6.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  a real random variable. We define the *characteristic function*  $\varphi_X : \mathbb{R} \rightarrow \mathbb{C}$  of  $X$  by

$$\varphi_X(t) = \mathbf{E}[e^{itX}] = \begin{cases} \sum_{k \in \Omega_X} e^{itk} \mathbf{P}[X = k], & \text{if } X \text{ is discrete,} \\ \int_{-\infty}^{\infty} e^{itx} f_X(x) dx, & \text{if } X \text{ is continuous and its law has density } f_X, \end{cases} \quad (12.13)$$

for  $t \in \mathbb{R}$ . Here,  $i$  is the imaginary unit (with  $i^2 = -1$ ).

The characteristic function has a similar uniqueness property as we saw for the moment generating function (Theorem 12.4), and it has the advantage that it always exists for real random variables. We can say

$$\varphi_X(t) = \varphi_Y(t) \text{ for all } t \in \mathbb{R} \quad \Leftrightarrow \quad \mathbf{P}_X = \mathbf{P}_Y. \quad (12.14)$$

**Definition 12.7.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  a discrete real random variable with  $\Omega_X \subseteq \mathbb{N}_0$  (i.e.  $X$  only attains values in  $\mathbb{N}_0$ ). We define the *probability generating function*  $G_X$  of  $X$  by

$$G_X(t) = \mathbf{E}[t^X] = \sum_{k=0}^{\infty} t^k \mathbf{P}[X = k], \quad (12.15)$$

for all  $t \geq 0$  for which it exists.

Clearly,  $G_X(t)$  exists for  $t \in [0, 1]$  and is differentiable at least on  $[0, 1)$ . We have similar properties as for the moment generating functions, which we state here without proof (see [2, Theorem 4.33] for a proof).

**Lemma 12.8.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  a discrete real random variable with  $\Omega_X \subseteq \mathbb{N}_0$ . Then

(i) For every  $k \in \mathbb{N}$ , one has

$$\mathbf{P}[X = k] = \frac{G_X^{(k)}(0)}{k!}. \quad (12.16)$$

In particular,  $\mathbf{P}_X$  is uniquely determined by  $G_X$ .<sup>2</sup>

(ii)  $\mathbf{E}[X]$  exists if and only if  $G'(1) = \lim_{t \uparrow 1} G'(t)$  exists and then

$$\mathbf{E}[X] = G'(1). \quad (12.17)$$

<sup>2</sup>So we have, just as in (12.11) and (12.14), that for two  $\mathbb{N}_0$ -valued random variables  $X$  and  $Y$ , that  $G_X(t) = G_Y(t)$  for every  $t \in [0, 1] \Leftrightarrow \mathbf{P}_X = \mathbf{P}_Y$ .

## 13. Convergence in probability, almost sure convergence and the law of large numbers

(Reference: [1, Sections 8.2, 8.4], or [2, Section 5.1])

### 13.1. Convergence in probability and the weak law of large numbers

**Definition 13.1.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $(Z_n)_{n \in \mathbb{N}}$  a sequence of random variables  $Z_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . We say that  $(Z_n)_{n \in \mathbb{N}}$  *converges in probability* to some random variable  $Z : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , if

$$\lim_{n \rightarrow \infty} \mathbf{P}[|Z_n - Z| > \varepsilon] = 0, \quad \text{for all } \varepsilon > 0. \quad (13.1)$$

We write this as

$$Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} Z. \quad (13.2)$$

We immediately state the (weak) law of large numbers.

**Theorem 13.2.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $(X_n)_{n \in \mathbb{N}}$  a sequence of i.i.d. random variables  $X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  with  $\mathbf{E}[X_1^2] < \infty$  and  $\mu = \mathbf{E}[X_1]$ . Then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \mu. \quad (13.3)$$

*Proof.* Let  $\sigma^2 = \text{Var}[X_1]$ . By Chebyshev's inequality (6.42), we have for any given  $\varepsilon > 0$ :

$$\begin{aligned} \mathbf{P}[|\bar{X}_n - \mu| > \varepsilon] &\leq \frac{\text{Var}[\bar{X}_n]}{\varepsilon^2} \\ &= \frac{1}{\varepsilon^2} \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \stackrel{(6.29), (10.14)}{=} \frac{1}{n^2 \varepsilon^2} \sum_{i=1}^n \underbrace{\text{Var}[X_i]}_{=\sigma^2} \\ &= \frac{\sigma^2}{n \varepsilon^2} \xrightarrow[n \rightarrow \infty]{} 0. \end{aligned} \quad (13.4)$$

□

**Remark 13.3.** The assumption that  $\mathbf{E}[X_1^2] < \infty$  can be relaxed to  $\mathbf{E}[|X_1|] < \infty$  (see Theorem 5.7 in [2]). However, if  $\mathbf{E}[|X_1|] = \infty$ , the weak law of large numbers does not necessarily hold anymore.

---

End of Lecture 23

## 13.2. Almost sure convergence and the strong law of large numbers

We have seen that for  $X_1, \dots, X_n$  i.i.d. with  $\mathbf{E}[X_1^2] < \infty$ , the probability to see a deviation from the average by at least  $\varepsilon > 0$  converges to zero. In fact, something stronger is true: We will see that not only does this probability go to zero, but in fact with probability 1,  $\bar{X}_n$  converges to  $\mathbf{E}[X_1]$ . We formalize this kind of convergence.

**Definition 13.4.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $(Z_n)_{n \in \mathbb{N}}$  a sequence of random variables  $Z_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . We say that  $(Z_n)_{n \in \mathbb{N}}$  *converges  $\mathbf{P}$ -almost surely* (abbreviated as  *$\mathbf{P}$ -a.s.*) to some random variable  $Z : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , if

$$\mathbf{P} \left[ \left\{ \omega \in \Omega ; \lim_{n \rightarrow \infty} Z_n(\omega) = Z(\omega) \right\} \right] = 1 \quad (13.5)$$

We write this as

$$Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}\text{-a.s.}} Z. \quad (13.6)$$

One can show that the set  $\{\omega \in \Omega ; \lim_{n \rightarrow \infty} Z_n(\omega) = Z(\omega)\}$  is in  $\mathcal{F}$ . Almost sure convergence is a stronger notion than convergence in probability. Indeed, one has:

**Proposition 13.5.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $(Z_n)_{n \in \mathbb{N}}$  a sequence of random variables  $Z_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ .

$$Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}\text{-a.s.}} Z \quad \Rightarrow \quad Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} Z. \quad (13.7)$$

*Proof.* Note that

$$\begin{aligned} \mathbf{P}[|Z_n - Z| > \varepsilon] &\leq \mathbf{P} \left[ \bigcup_{k=n}^{\infty} \{|Z_k - Z| > \varepsilon\} \right] \\ &\xrightarrow[n \rightarrow \infty]{} \mathbf{P} \left[ \bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} \{|Z_k - Z| > \varepsilon\} \right] \leq \mathbf{P} \left[ \left\{ \omega \in \Omega ; \lim_{n \rightarrow \infty} Z_n(\omega) \neq Z(\omega) \right\} \right]. \end{aligned} \quad (13.8)$$

We have used Proposition 1.15, (vii) in the second step. Now if  $Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}\text{-a.s.}} Z$ , then the right-hand side of the above inequality is zero, and thus  $\lim_{n \rightarrow \infty} \mathbf{P}[|Z_n - Z| > \varepsilon] = 0$  follows.  $\square$

We now state the (strong) law of large numbers

**Theorem 13.6.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $(X_n)_{n \in \mathbb{N}}$  a sequence of i.i.d. random variables  $X_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  with  $\mathbf{E}[X_1^4] < \infty$  and  $\mu = \mathbf{E}[X_1]$ . Then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbf{P}\text{-a.s.}} \mu. \quad (13.9)$$

*Proof.* We follow the proof in [1, Section 8.4]. Suppose that  $\mu = \mathbf{E}[X_1] = 0$  and define  $K = \mathbf{E}[X_1^4] (< \infty)$ . Then (for  $n \geq 4$ ):

$$\begin{aligned} \mathbf{E}[\bar{X}_n^4] &= \frac{1}{n^4} \mathbf{E} \left[ \sum_{i=1}^n X_i^4 + \sum_{\{i,j\} \subseteq \{1,\dots,n\}, |\{i,j\}|=2} \left( \binom{4}{1} X_i^3 X_j + \binom{4}{2} X_i^2 X_j^2 \right) \right. \\ &\quad \left. + \sum_{\{i,j,k\} \subseteq \{1,\dots,n\}, |\{i,j,k\}|=3} \binom{4}{3} X_i^2 X_j X_k + \sum_{\{i,j,k,\ell\} \subseteq \{1,\dots,n\}, |\{i,j,k,\ell\}|=4} X_i X_j X_k X_\ell \right]. \end{aligned} \quad (13.10)$$

Now note that due to independence, we have

$$\begin{aligned} \mathbf{E}[X_i^3 X_j] &= \mathbf{E}[X_i^3] \mathbf{E}[X_j] = 0, \\ \mathbf{E}[X_i^2 X_j X_k] &= \mathbf{E}[X_i^2] \mathbf{E}[X_j] \mathbf{E}[X_k] = 0, \\ \mathbf{E}[X_i X_j X_k X_\ell] &= \mathbf{E}[X_i] \mathbf{E}[X_j] \mathbf{E}[X_k] \mathbf{E}[X_\ell] = 0. \end{aligned}$$

Inserting this into (13.10) shows that only the  $n$  terms  $\mathbf{E}[X_i^4]$  and the  $6\binom{n}{2}$  terms of the form  $\mathbf{E}[X_i^2] \mathbf{E}[X_j^2]$  contribute to the sum, and so:

$$\mathbf{E}[\bar{X}_n^4] \leq nK + 3n(n-1)\mathbf{E}[X_1^2]^2 \leq \frac{K}{n^3} + \frac{3K}{n^2}. \quad (13.11)$$

In the second step, we used Jensen's inequality (9.2):  $\mathbf{E}[X_2^2]^2 \leq \mathbf{E}[X_1^4] = K$ . Now consider the random variable

$$\Xi(\omega) = \sum_{n=1}^{\infty} \bar{X}_n(\omega)^4 \quad (13.12)$$

(note that all terms in the sum are non-negative, so we get a  $[0, \infty) \cup \{\infty\}$ -valued random variable). Suppose now that  $\mathbf{P}[\Xi = \infty] > 0$ . Since  $\Xi$  is non-negative, we must have  $\mathbf{E}[\Xi] = \infty$ , but on the other hand:

$$\mathbf{E} \left[ \sum_{n=1}^{\infty} \bar{X}_n^4 \right] = \lim_{N \rightarrow \infty} \sum_{n=1}^N \mathbf{E}[\bar{X}_n^4] \leq \lim_{N \rightarrow \infty} \sum_{n=1}^N \left( \frac{K}{n^3} + \frac{3K}{n^2} \right) < \infty. \quad (13.13)$$

Since this is a contradiction, we must have that  $\mathbf{P}[\Xi = \infty] = 0$ . But on the event  $\{\Xi < \infty\}$ , we necessarily have that  $\bar{X}_n^4 \rightarrow 0$ , and therefore  $\bar{X}_n \rightarrow 0$ . Finally, if  $\mathbf{E}[X_1] = \mu \neq 0$ , consider  $Y_j = X_j - \mu$ , then we see from the proof that  $\bar{Y}_n \rightarrow 0$  **P**-a.s., and therefore the claim follows.

*Remark 13.7.* (i) The proof has some caveats: One first has to define the notion of a  $[0, \infty) \cup \{\infty\}$ -valued random variable. The bigger issue is that we use in (13.13) the *monotone convergence theorem* to exchange limit and expectation, see [2, Theorem 4.7].

(ii) Another proof of the strong law of large numbers can be found in [2, Section 5.1.3], under weaker assumptions. In fact the weakest assumptions under which the strong law of large numbers is still valid are quite surprising: Etemadi (1981) proved that  $(X_n)_{n \in \mathbb{N}}$  *pairwise* independent, identically distributed and  $\mathbf{E}[|X_1|] < \infty$  already suffices to guarantee the validity of (13.9). Note that this is a much (!) weaker assumption than the i.i.d. assumption. Clearly, this version of the strong law of large numbers implies the weak law of large numbers, by Proposition 13.5.

To close this section, we give the following result known as the *continuous mapping theorem*.

**Proposition 13.8.** Consider a family  $(Z_n)_{n \in \mathbb{N}}$  of random variables  $Z_n : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  and a random variable  $Z : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . Furthermore, let  $g : (\mathbb{R}, \mathcal{B}(\mathbb{R})) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  be a measurable function such that

$$\mathbf{P}[Z \in D_g] = 0, \quad D_g = \{x \in \mathbb{R}; g \text{ is not continuous in } x\}. \quad (13.14)$$

(i) If  $Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} Z$ , then also  $g(Z_n) \xrightarrow[n \rightarrow \infty]{\mathbf{P}} g(Z)$ ,

(ii) If  $Z_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}\text{-a.s.}} Z$ , then also  $g(Z_n) \xrightarrow[n \rightarrow \infty]{\mathbf{P}\text{-a.s.}} g(Z)$ .

For instance, consider  $(X_n)_{n \in \mathbb{N}}$  be an i.i.d. sequence with  $X_1 \sim \mathcal{E}(\lambda)$ . By the strong law of large numbers,

$$\bar{X}_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \mathbf{E}[X_1] = \frac{1}{\lambda}. \quad (13.15)$$

Therefore, by the continuous mapping theorem, Proposition 13.8, (ii), we have

$$\frac{1}{\bar{X}_n} \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \lambda. \quad (13.16)$$

□

### 13.3. Application: Monte-Carlo integration

The law of large numbers essentially states that the average  $\bar{X}_n(\omega)$  of i.i.d. random variables  $X_1, \dots, X_n$  is close to  $\mu$  with high probability of large  $n$ . An important application of the law of large numbers are *Monte-Carlo methods* to calculate integrals or sums.

*Example 13.9.* Consider a piecewise continuous function  $h : [0, 1] \rightarrow \mathbb{R}$  we want to calculate the integral

$$I = \int_0^1 h(x) dx, \quad (13.17)$$

which cannot be calculated elementary. We also assume that  $\int_0^1 h^4(x) dx < \infty$ <sup>1</sup>. The *Monte-Carlo integration* gives us a way to find an estimate of  $I$ .

Let  $X_1, X_2, \dots \sim \mathcal{U}([0, 1])$  be i.i.d. random variables. Then also  $h(X_1), h(X_2), \dots$  are i.i.d. by Proposition 7.15 and the fact (not proved here) that the function  $h$  is measurable, since it is piecewise continuous. Note that

$$\mathbf{E}[h(X_1)] \stackrel{(6.19)}{=} \int_0^1 h(x) dx = I. \quad (13.18)$$

<sup>1</sup>For instance since we know some bound. Note in particular that  $I$  is therefore well-defined and finite.

By the strong law of large numbers, Theorem 13.6, we have

$$\frac{1}{n} \sum_{i=1}^n h(X_i) \xrightarrow[n \rightarrow \infty]{\text{P-a.s.}} I. \quad (13.19)$$

This means that:

$$\mathbf{P} \left[ \left\{ \omega \in \Omega; \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n h(X_i) = I \right\} \right] = 1. \quad (13.20)$$

In other words, by simulating i.i.d. uniform random variable, we have a method to estimate the integral  $I$ .

*Remark 13.10.* (i) In the same way, we can calculate approximately other integrals or series by considering i.i.d. random variables  $X_1, X_2, \dots$  with a different distribution. For instance, if  $X_1, X_2, \dots \sim \mathcal{N}(0, 1)$  i.i.d., we can approximate

$$\int_{-\infty}^{\infty} e^{-\frac{x^2}{2}} h(x) dx \approx \frac{\sqrt{2\pi}}{n} \sum_{i=1}^n h(X_i), \quad (13.21)$$

if  $h : \mathbb{R} \rightarrow \mathbb{R}$  is piecewise continuous and  $\int_{-\infty}^{\infty} h^4(x) e^{-\frac{x^2}{2}} dx < \infty$ .

- (ii) The speed of convergence of Monte-Carlo is quite poor (we can find the speed by the central limit theorem later). Therefore, Monte-Carlo methods are typically only used if  $h$  is very irregular, or to perform high-dimensional integration. For regular, one-dimensional functions  $h$  numerical methods are typically much better than Monte-Carlo methods.

---

*End of Lecture 24*



## 14. The central limit theorem

(Reference: [1, Section 8.3], or [2, Section 5.2, 5.3])

Consider rolling a die  $n$  times. We are interested in the sum of all numbers that came during the  $n$  rolls. By now we know that such a process can be described in different ways, and one of them is to consider  $X_1, X_2, \dots, X_n$  i.i.d. random variables with  $X_1 \sim \mathcal{U}(\{1, \dots, 6\})$ . From the previous section we learned that

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \mathbf{E}[X_1]. \quad (14.1)$$

This means that the average of the numbers shown should be close to  $\mathbf{E}[X_1] = 3.5$  for large  $n$ , consequently the sum of the numbers should be close to  $3.5 \cdot n$ . What is not clear however is *how* close it really is. For instance, we could ask

$$\text{How likely is it that the sum of outcomes when rolling a die 100 times exceeds 400?} \quad (14.2)$$

In other words, we are interested in *fluctuations* around  $\mathbf{E}[X_1]$ . This will be the main statement of the central limit theorem. For this, we need another notion of convergence.

### 14.1. Convergence in distribution

**Definition 14.1.** Let  $(Z_n)_{n \in \mathbb{N}}$  be a sequence of real random variables with cumulative distribution functions  $F_{Z_n}$  and  $Z$  a real random variable with cumulative distribution function  $F_Z$ . We say that  $Z_n$  *converges in law / in distribution* if

$$F_{Z_n}(z) = \mathbf{P}[Z_n \leq z] \xrightarrow[n \rightarrow \infty]{} F_Z(z), \quad (14.3)$$

for all points of continuity  $z \in \mathbb{R}$  of  $F_Z$ . We write

$$Z_n \xrightarrow[n \rightarrow \infty]{d} Z. \quad (14.4)$$

**Remark 14.2.** (i) Since  $F_Z$  is monotone and bounded, one can show that there are at most countably many points of discontinuity.

(ii) Note that the above notion of convergence does not depend on the random variables  $Z_n$  and  $Z$ , but only on their laws  $\mathbf{P}_{Z_n}$  and  $\mathbf{P}_Z$ . If  $Z_n \xrightarrow[n \rightarrow \infty]{d} Z$  also say that  $\mathbf{P}_{Z_n}$  converges *weakly* to  $\mathbf{P}_Z$ , also written as

$$\mathbf{P}_{Z_n} \xrightarrow[n \rightarrow \infty]{w} \mathbf{P}_Z. \quad (14.5)$$

- (iii) To understand why we exclude discontinuity points of  $F_Z$ , consider the deterministic functions  $Z_n = \frac{1}{n}$  and  $Z = 0$ . Then

$$F_{Z_n}(z) = \begin{cases} 1, & z \geq \frac{1}{n}, \\ 0, & z < \frac{1}{n}, \end{cases} \quad F_Z(z) = \begin{cases} 1, & z \geq 0, \\ 0, & z < 0. \end{cases} \quad (14.6)$$

We see that  $F_{Z_n}(0) = 0$  does not converge to  $F_Z(0) = 1$ , but of course we would like to say that  $Z_n \xrightarrow[n \rightarrow \infty]{d} Z$ , so we exclude the convergence at this point.

## 14.2. The central limit theorem

Let us start with a motivation: As explained at the beginning of this chapter, we want to be able to calculate probabilities of fluctuations for sums of i.i.d. random variables. Naturally we should ask, how large these fluctuations typically are.

*Example 14.3.* Let  $X_1, \dots, X_n$  be i.i.d. real random variables with  $X_1 \sim \mathcal{N}(\mu, \sigma^2)$ . Then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}(\mu, \frac{\sigma^2}{n}) \quad \Rightarrow \quad \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \sim \mathcal{N}(0, 1). \quad (14.7)$$

In other words, we have that

$$F_{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}(x) = \Phi(x), \quad x \in \mathbb{R}, \quad (14.8)$$

where  $\Phi$  is the cumulative distribution function of a  $\mathcal{N}(0, 1)$ -distributed random variable.

From this, we see that the scaling  $\sqrt{n}$  is the “correct” scaling to measure the fluctuations of  $\bar{X}_n$ , at least if  $X_1, \dots, X_n$  are i.i.d. with  $X_1 \sim \mathcal{N}(\mu, \sigma^2)$ . That this is true in general is the statement of the celebrated *central limit theorem*:

**Theorem 14.4.** *Let  $X_1, X_2, \dots$  be a sequence of i.i.d. real random variables with  $\mu = E[X_1]$  and  $\sigma^2 = \text{Var}[X_1] \in (0, \infty)$ . Then,*

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, 1). \quad (14.9)$$

The notation in (14.9) means that

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow[n \rightarrow \infty]{d} Z, \text{ where } Z \sim \mathcal{N}(0, 1). \quad (14.10)$$

There are different ways to establish the central limit theorem. In these notes, two are presented, one using moment generating functions and the other more direct one, is presented in the Appendix A.2.

For the first method, we need the following *continuity theorem for moment generating functions*, which we state without proof:

**Lemma 14.5.** Let  $(Z_n)_{n \in \mathbb{N}}$  be a sequence of real random variables with  $Z$  a real random variable. Suppose the moment generating functions  $\psi_{Z_n}(t)$  and  $\psi_Z(t)$  exist for all  $n \geq 1$ ,  $|t| < \varepsilon$  for some  $\varepsilon > 0$ . Then

$$Z_n \xrightarrow[n \rightarrow \infty]{d} Z \quad \Leftrightarrow \quad \psi_{Z_n}(t) \xrightarrow[n \rightarrow \infty]{} \psi_Z(t) \quad \text{for all } t \in (-\varepsilon, \varepsilon). \quad (14.11)$$

Since we will need it in the rest, we calculate the moment generating function of  $Z \sim \mathcal{N}(0, 1)$ :

$$\begin{aligned} \psi_Z(t) &= \mathbf{E}[e^{tZ}] = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{tz} e^{-\frac{z^2}{2}} dz \\ &= \frac{1}{\sqrt{2\pi}} e^{\frac{t^2}{2}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(z-t)^2} dz = e^{\frac{t^2}{2}}, \end{aligned} \quad (14.12)$$

for every  $t \in \mathbb{R}$ .

*Proof of Theorem 14.4.* To show that claim (14.9), we show that the moment generating function of  $\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}$  converges to that of  $Z \sim \mathcal{N}(0, 1)$ . We also assume that all moment generating functions exist, at least for  $t$  in a neighborhood of 0.

Now we have:

$$\begin{aligned} \psi_{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}(t) &= \psi_{\frac{1}{\sqrt{n}\sigma} (\sum_{i=1}^n (X_i - \mu))}(t) = \prod_{i=1}^n \psi_{X_i - \mu} \left( \frac{t}{\sqrt{n}\sigma} \right) \\ &= \prod_{i=1}^n \left( e^{-\mu \frac{t}{\sqrt{n}\sigma}} \psi_{X_i} \left( \frac{t}{\sqrt{n}\sigma} \right) \right) = e^{-\mu \sqrt{n} \frac{t}{\sigma}} \left( \psi_{X_1} \left( \frac{t}{\sqrt{n}\sigma} \right) \right)^n. \end{aligned} \quad (14.13)$$

We consider  $L(u) = \log \psi_{X_1}(u)$ . Note that

$$\begin{aligned} L(0) &= 0, \\ L'(0) &= \frac{\psi'_{X_1}(0)}{\psi_{X_1}(0)} = \mathbf{E}[X_1] = \mu, \\ L''(0) &= \frac{\psi_{X_1}(0)\psi''_{X_1}(0) - (\psi'_{X_1}(0))^2}{(\psi_{X_1}(0))^2} = \text{Var}[X_1] = \sigma^2. \end{aligned} \quad (14.14)$$

Now we take the logarithm in (14.13):

$$\log \psi_{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}(t) = -\mu \sqrt{n} \frac{t}{\sigma} + n \log \psi_{X_1} \left( \frac{t}{\sqrt{n}\sigma} \right) = -\mu \sqrt{n} \frac{t}{\sigma} + nL \left( \frac{t}{\sqrt{n}\sigma} \right). \quad (14.15)$$

One can then perform a Taylor expansion of the term on the right-hand side:

$$\begin{aligned} \log \psi_{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}(t) &= -\mu \sqrt{n} \frac{t}{\sigma} + n \left( L(0) + \frac{t}{\sqrt{n}\sigma} L'(0) + \frac{t^2}{n\sigma^2} L''(0) + O \left( \frac{t^3}{n^{3/2}\sigma^3} \right) \right) \\ &\stackrel{(14.14)}{=} \frac{t^2}{2} + O \left( \frac{t^3}{\sqrt{n}\sigma^3} \right). \end{aligned} \quad (14.16)$$

Sending  $n \rightarrow \infty$  gives us

$$\log \psi_{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}(t) \rightarrow \frac{t^2}{2} \quad \Rightarrow \quad \psi_{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}(t) \rightarrow e^{\frac{t^2}{2}} \stackrel{(14.12)}{=} \psi_Z(t), \quad (14.17)$$

by the continuity of exp. The claim then follows by applying Lemma 14.5.  $\square$

Note that the proof above relied on Lemma 14.5 (which we did not prove) and does not work if the moment generating function of  $X_1$  fails to exist in a neighborhood of 0. For completeness, a more direct (but technically more complicated) proof is sketched in the Appendix A.2.

Let us turn to some applications of the central limit theorem.

*Example 14.6.* (i) Recall the question (14.2) posed at the beginning of this chapter, namely suppose we roll a die  $n = 100$  times and we are interested in the probability that the sum of the outcomes exceeds 400. To model this, let  $X_1, X_2, \dots$  be i.i.d. random variables with  $X_1 \sim \mathcal{U}(\{1, \dots, 6\})$ . An easy calculation shows that

$$\mu = \mathbf{E}[X_1] = 3.5, \quad \sigma^2 = \text{Var}[X_1] = \frac{35}{12}. \quad (14.18)$$

Thus, the assumptions of the central limit theorem (Theorem 14.4) are fulfilled and thus for every  $x \in \mathbb{R}$ ,

$$\lim_{n \rightarrow \infty} \mathbf{P} \left[ \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \leq x \right] = \lim_{n \rightarrow \infty} \mathbf{P} \left[ \frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n}\sigma} \leq x \right] = \Phi(x). \quad (14.19)$$

For  $n = 100$ , we therefore have the approximate probability

$$\begin{aligned} \mathbf{P} \left[ \sum_{i=1}^{100} X_i > 400 \right] &= \mathbf{P} \left[ \frac{\sum_{i=1}^{100} X_i - 350}{\sqrt{100} \cdot \sqrt{\frac{35}{12}}} > \frac{400 - 350}{\sqrt{100} \cdot \sqrt{\frac{35}{12}}} \right] \\ &\approx \mathbf{P} \left[ Z > \frac{400 - 350}{\sqrt{100} \cdot \sqrt{\frac{35}{12}}} \right], \quad Z \sim \mathcal{N}(0, 1) \\ &= 1 - \Phi \left( \frac{400 - 350}{\sqrt{100} \cdot \sqrt{\frac{35}{12}}} \right) = 1 - \Phi(2.927) \\ &\approx 1 - 0.9983 = 0.27\%. \end{aligned} \quad (14.20)$$

(ii) As another application, we present a general formula to approximate the binomial distribution  $\text{Bin}(n, p)$  by a normal distribution, if  $n \cdot p$  is not too small and  $n$  is large. We rely on the fact that for  $X_1, X_2, \dots$  i.i.d. random variables with  $X_1 \sim \text{Ber}(p)$ ,  $p \in (0, 1)$ , we have

$$S_n = \sum_{i=1}^n X_i \sim \text{Bin}(n, p). \quad (14.21)$$

Also recall that we have

$$\mathbf{E}[X_1] = p, \quad \text{Var}[X_1] = p(1 - p). \quad (14.22)$$

Suppose now we want to approximate the probability mass function of  $\text{Bin}(n, p)$  for  $0 \leq k \leq n$ , which is given by  $p_{S_n}(k) = \binom{n}{k} p^k (1 - p)^{n-k}$ . Note that

$$\begin{aligned} \mathbf{P}[S_n = k] &= \mathbf{P}[k - 0.5 < S_n \leq k + 0.5] \\ &= \mathbf{P}\left[\frac{k - 0.5 - np}{\sqrt{np(1 - p)}} < \underbrace{\frac{S_n - np}{\sqrt{np(1 - p)}}}_{=\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1 - p)}}} \leq \frac{k + 0.5 - np}{\sqrt{np(1 - p)}}\right] \\ &\approx \mathbf{P}\left[\frac{k - 0.5 - np}{\sqrt{np(1 - p)}} < Z \leq \frac{k + 0.5 - np}{\sqrt{np(1 - p)}}\right], \quad Z \sim \mathcal{N}(0, 1) \\ &= \Phi\left(\frac{k + 0.5 - np}{\sqrt{np(1 - p)}}\right) - \Phi\left(\frac{k - 0.5 - np}{\sqrt{np(1 - p)}}\right). \end{aligned} \quad (14.23)$$

In the approximation we have chosen the interval of length 1, which makes the approximation better for finite  $n$ .

---

Some other useful properties of convergence in distribution are given in the Appendix [A.3](#).  
*End of Lecture 25; End of relevant material for the Final Exam.*

# 15. The Poisson Process

(Reference: [1, Section 9.1], or [2, Section 3.5])

In this chapter, we will introduce a stochastic process in time that models the arrival of random events of a discrete nature, such as arrival times of radioactive particles at a Geiger counter, the times at which electronic components fail, or the arrival of customers in a store.

**Definition 15.1.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space,  $I \neq \emptyset$  an index set and  $E$  a set equipped with a  $\sigma$ -algebra  $\mathcal{E}$ . A *stochastic process* with *time parameter set*  $I$  and *state space*  $E$  is a collection  $(X_t)_{t \in I}$  of random variables  $X_t : (\Omega, \mathcal{F}) \rightarrow (E, \mathcal{E})$ . For  $\omega \in \Omega$ , we say that  $t \mapsto X_t(\omega)$  is a *sample path* of the process.

For us,  $I = \mathbb{N}$  or  $I = [0, \infty)$  are natural choices. The easiest examples of stochastic processes would include  $(X_n)_{n \in \mathbb{N}}$ , where  $X_n = Z : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  for all  $n \in \mathbb{N}$  (a constant process, that only depends on the initial randomness  $Z(\omega)$ ), or  $(Y_n)_{n \in \mathbb{N}}$ , where  $Y_1, Y_2, \dots : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  are i.i.d. random variables. In these examples,  $I = \mathbb{N}$  and  $(E, \mathcal{E}) = (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . We now introduce the Poisson process (with rate  $\lambda > 0$ ), which is a stochastic process with  $I = [0, \infty)$  and  $(E, \mathcal{E}) = (\mathbb{N}_0, \mathcal{P}(\mathbb{N}_0))$ .

**Definition 15.2.** A stochastic process  $(N_t)_{t \geq 0}$  of random variables  $N_t : (\Omega, \mathcal{F}) \rightarrow (\mathbb{N}_0, \mathcal{P}(\mathbb{N}_0))$  is called a *Poisson process of rate  $\lambda > 0$*  if it fulfills the following properties:

- (i)  $N_0 = 0$ ,
- (ii) the paths  $t \mapsto N_t(\omega)$  are right-continuous for every  $\omega \in \Omega$ ,
- (iii) for  $0 = t_0 < t_1 < \dots < t_n$ , the increments  $(N_{t_k} - N_{t_{k-1}})_{k=1, \dots, n}$  are independent,
- (iv) for  $0 \leq s < t$ ,  $N_t - N_s \sim \text{Pois}(\lambda(t - s))$ .

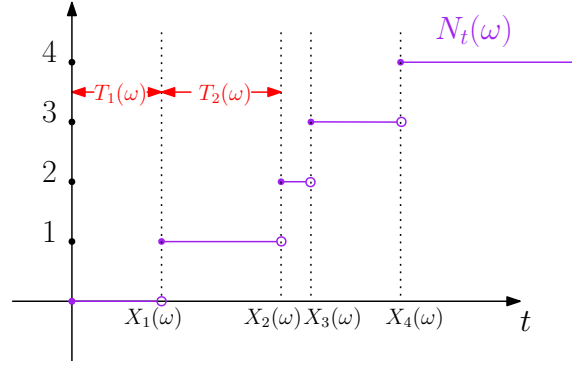
*Example 15.3.* Let  $N_t$  be the number of emitted particles of a radioactive source during the time interval  $[0, t]$ . Then (iii) means that the numbers of particles emitted during disjoint time intervals are independent.

Since  $N_t = N_s + (N_t - N_s)$ , where  $N_t - N_s \geq 0$ , we see that  $t \mapsto N_t(\omega)$  is nondecreasing. Moreover, by definition, the process only attains values in  $\mathbb{N}_0$ . Therefore, sample paths will be piecewise constant, with “jumps” occurring at random times (see figure 15). How long do we have to wait for the first jump? Let us define

$$T = \inf\{t \geq 0; N_t > 0\}, \quad (15.1)$$

then we have that

$$\mathbf{P}[T > t] = \mathbf{P}[N_t = 0] = e^{-\lambda t}, \quad (15.2)$$

Figure 15.1.: Plot for one trajectory for  $t \mapsto N_t(\omega)$ .

since  $N_t = N_t - N_0 \sim \text{Pois}(\lambda t)$ . More generally, we have that

$$\mathbf{P}[N_{s+t} - N_s = 0] = e^{-\lambda t}, \quad s \geq 0, t > 0. \quad (15.3)$$

This suggests that the “random jumps” occur after independent, exponentially distributed times, the *interarrival times*. We use this idea to construct a Poisson process, and briefly explain afterwards, that this property actually characterizes Poisson processes.

**Theorem 15.4.** *Let  $T_1, T_2, \dots \sim \mathcal{E}(\lambda)$  be i.i.d. random variables with  $\lambda > 0$ . Define*

$$X_n = \sum_{k=1}^n T_k, \quad N_t = |\{n \in \mathbb{N}_0; X_n \leq t\}|. \quad (15.4)$$

*The family  $(N_t)_{t \geq 0}$  is a Poisson process with intensity  $\lambda$ .*

*Proof.* The proof follows [2, Theorem 3.34]. We must show that for any  $n \in \mathbb{N}$  and any sequence  $0 = t_0 < t_1 < \dots < t_n$ , we have that  $(N_{t_i} - N_{t_{i-1}})_{i=1}^n$  are independent and  $N_{t_i} - N_{t_{i-1}} \sim \text{Pois}(\lambda(t_i - t_{i-1}))$ . For simplicity, we show it only sketch the proof in the case of  $n = 2$ . Thus we want to show that for  $0 < s < t$  and  $\ell, k \in \mathbb{N}_0$ :

$$\mathbf{P}[N_s = k, N_t - N_s = \ell] = \left( e^{-\lambda s} \frac{(\lambda s)^k}{k!} \right) \left( e^{-\lambda(t-s)} \frac{(\lambda(t-s))^\ell}{\ell!} \right). \quad (15.5)$$

This implies that  $N_s$  and  $N_t - N_s$  are independent by Proposition 7.13, (ii). Moreover, by summing over  $k \in \mathbb{N}_0$ , we have  $N_t - N_s \sim \text{Pois}(\lambda(t-s))$ . Now the joint law of  $T_1, \dots, T_{k+\ell+1}$ ,  $\mathbf{P}_{T_1, \dots, T_{k+\ell+1}}$ , has the density

$$f_{T_1, \dots, T_{k+\ell+1}}(x_1, \dots, x_{k+\ell+1}) = \lambda^{k+\ell+1} e^{-\lambda(x_1 + \dots + x_{k+\ell+1})} \prod_{j=1}^{k+\ell+1} \mathbb{1}_{[0, \infty)}(x_j). \quad (15.6)$$

Let us only consider the cases where  $k, \ell \geq 1$  (the other cases, i.e. where either  $k = 0$  or  $\ell = 0$  are similar). Now we have

$$\begin{aligned} \mathbf{P}[N_s = k, N_t - N_s = \ell] &= \mathbf{P}[X_k \leq s < X_{k+1} \leq X_{k+\ell} \leq t < X_{k+\ell+1}] \\ &= \int_0^\infty \dots \int_0^\infty dx_1 \dots dx_{k+\ell+1} \lambda^{k+\ell+1} e^{-\lambda(x_1 + \dots + x_{k+\ell+1})} \\ &\quad \cdot \mathbb{1}_{s \in [x_1 + \dots + x_k, x_1 + \dots + x_{k+1})} \cdot \mathbb{1}_{t \in [x_1 + \dots + x_{k+\ell}, x_1 + \dots + x_{k+\ell+1})}. \end{aligned} \quad (15.7)$$

We now integrate now first over  $x_{k+\ell+1}$ : This yields for fixed  $x_1, \dots, x_{k+\ell}$

$$\int_0^\infty dx_{k+\ell+1} \lambda e^{-\lambda(x_1 + \dots + x_{k+\ell+1})} \mathbb{1}_{\{x_1 + \dots + x_{k+\ell} + x_{k+\ell+1} > t\}} = \int_t^\infty dz \lambda e^{-\lambda z} = e^{-\lambda t}. \quad (15.8)$$

We then fix  $x_1, \dots, x_k$  and integrate over  $x_{k+1}, \dots, x_{k+\ell}$ :

$$\begin{aligned} &\int_0^\infty \dots \int_0^\infty dx_{k+1} \dots dx_{k+\ell} \mathbb{1}_{\{s < x_1 + \dots + x_{k+1} \leq x_1 + \dots + x_{k+\ell} \leq t\}} \\ &= \int_0^\infty \dots \int_0^\infty dy_1 \dots dy_\ell \mathbb{1}_{\{y_1 + \dots + y_\ell \leq t - s\}} = \frac{(t - s)^\ell}{\ell!}. \end{aligned} \quad (15.9)$$

We have used the substitution  $y_1 = x_1 + \dots + x_{k+1} - s$ ,  $y_2 = x_{k+2}$ , ...,  $y_\ell = x_{k+\ell}$ . Finally, we need to integrate over the remaining variables  $x_1, \dots, x_k$ :

$$\int_0^\infty \dots \int_0^\infty dx_1 \dots dx_k \mathbb{1}_{\{x_1 + \dots + x_k \leq s\}} = \frac{s^k}{k!}. \quad (15.10)$$

This means that we have

$$\mathbf{P}[N_s = k, N_t - N_s = \ell] = e^{-\lambda t} \lambda^{k+\ell} \frac{s^k}{k!} \frac{(t - s)^\ell}{\ell!}. \quad (15.11)$$

This is exactly  $\square$

Theorem 15.4 gives us the construction of a Poisson process: If we can simulate  $T_1, T_2, \dots \sim \mathcal{E}(\lambda)$  i.i.d. random variables, then we can define  $(N_t)_{t \geq 0}$  as in (15.4) and obtain a Poisson process with rate  $\lambda > 0$ . In fact, the converse is also true, which we state here without proof.

**Proposition 15.5.** *Let  $(N_t)_{t \geq 0}$  be a Poisson process of rate  $\lambda > 0$ . We define the jump times  $(X_i)_{i \geq 1}$  by*

$$X_i(\omega) = \inf\{t \geq 0; N_t(\omega) = i\}, \quad i \in \mathbb{N}. \quad (15.12)$$

*For a set  $\tilde{\Omega} \in \mathcal{F}$  with  $\mathbf{P}[\tilde{\Omega}] = 1$ , we have that  $X_k(\omega) < \infty$  for  $\omega \in \tilde{\Omega}$  for every  $k \in \mathbb{N}$ , and the interarrival times*

$$T_i = X_i - X_{i-1}, \quad i \geq 1, \quad (15.13)$$

*are i.i.d. random variables with  $T_1 \sim \mathcal{E}(\lambda)$ .*

Let us consider a worked example.



*Example 15.6.* Queuing at a post office can be modelled by a Poisson process with rate  $\lambda = \frac{1}{2}$  (in units  $\frac{\text{customers}}{\text{min}}$ ).

- (i) The expected number of arrivals during the first 10 minutes of an hour is given by

$$\mathbf{E}[N_{10}] = \lambda \cdot 10 = \frac{1}{2} \cdot 10 = 5, \quad (15.14)$$

since  $N_{10} \sim \text{Pois}(\lambda \cdot 10)$ .

- (ii) The probability to have 4 or less arrivals during the first 10 minutes of an hour is

$$\mathbf{P}[N_{10} \leq 4] = \sum_{k=0}^4 \frac{(\lambda \cdot 10)^k}{k!} e^{-\lambda \cdot 10} = e^{-5} \sum_{k=0}^4 \frac{5^k}{k!} \approx 0.4405. \quad (15.15)$$

- (iii) What is the probability that there are 3 customers in the first 10 minutes and 5 customers in the next 20 minutes? We have:

$$\begin{aligned} \mathbf{P}[N_{10} = 3, N_{30} - N_{10} = 5] &= \mathbf{P}[N_{10} = 3] \cdot \mathbf{P}[N_{30} - N_{10} = 5] \quad (\text{independent increments}) \\ &= \left( e^{-5} \frac{5^3}{3!} \right) \cdot \left( e^{-10} \frac{10^5}{5!} \right) \quad (\text{since } N_t - N_s \sim \text{Pois}(\lambda(t-s))) \\ &\approx 0.0053. \end{aligned} \quad (15.16)$$

---

*End of Lecture 26*

## 16. Markov chains: An overview

(Reference: [1, Section 9.2], or [2, Chapter 6])

In this chapter, we will define *Markov chains* and give a very brief overview over some classical properties. To keep things as simple as possible, we will only consider Markov chains

- in discrete time,
- with a finite state space.<sup>1</sup>

Throughout the entire chapter, we will set

$$E = \{1, 2, \dots, N\}, \quad (16.1)$$

with  $N \in \mathbb{N}$ . This will be the *state space*. A (discrete-time) Markov chain will be a stochastic process  $(X_n)_{n \in \mathbb{N}_0}$  with time parameter set  $I = \mathbb{N}_0$  and state space  $E$  in the sense of Definition 15.1. Its characterizing feature is that at time  $n$ , the next position  $X_{n+1} \in E$  of the process only depends on the current position  $X_n$ .

**Definition 16.1.** A  $\mathbb{R}^{n \times n}$ -matrix  $\Pi = (\Pi(x, y))_{(x, y) \in E^2}$  is called a *matrix of transition probabilities* if for every  $x \in E$ , the function

$$\Pi(x, \cdot) : E \rightarrow \mathbb{R}, y \mapsto \Pi(x, y) \quad (16.2)$$

is a probability mass function.

We will interpret the value  $\Pi(x, y)$  as

$$\text{the probability that the chain is at } y \text{ at time } n + 1, \text{ if it was at } x \text{ at time } n. \quad (16.3)$$

*Example 16.2.* Let  $E = \{1, 2, 3\}$  and consider the three matrices

$$\Pi_1 = \begin{pmatrix} 0.1 & 0.2 & 0.7 \\ 0.3 & 0.3 & 0.4 \\ 0 & 0.5 & 0.5 \end{pmatrix}, \quad \Pi_2 = \begin{pmatrix} 0.1 & 0.1 & 0.7 \\ 0.3 & 0.9 & 0.4 \\ 0 & 0 & 0.5 \end{pmatrix}, \quad \Pi_3 = \begin{pmatrix} 0.2 & -0.1 & 0.9 \\ -0.3 & 0.9 & 0.4 \\ 0.6 & 0 & 0.5 \end{pmatrix}. \quad (16.4)$$

Out of the three matrices, only  $\Pi_1$  is a valid matrix of transition probabilities. Indeed, in  $\Pi_2$ , the rows do not sum to one, and  $\Pi_3$  has negative entries.

To visualize this, one often uses *transition graphs*, which are formally directed weighted graphs on the set  $E$ , and the weight of the directed edge  $\overrightarrow{(x, y)}$  is exactly  $\Pi(x, y)$ . Edges with weight 0 are often omitted completely. An example for such a transition graph is given below.

We now come to the formal definition of a Markov chain.

---

<sup>1</sup>Markov chains with an at most countable state space can be treated very similarly, and we refer to [2, Chapter 6] for more on this

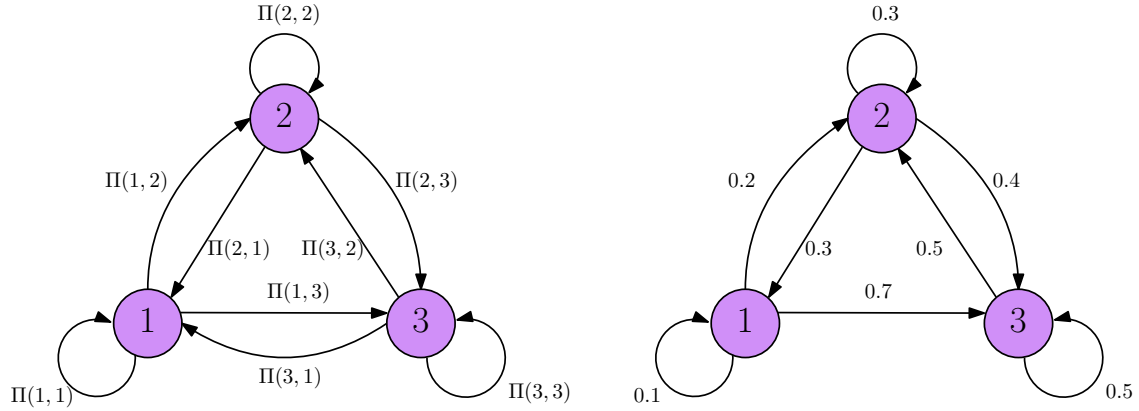


Figure 16.1.: Transition graph for a general Markov chain on  $E = \{1, 2, 3\}$  on the left, and for the choice  $\Pi_1$  from Example 16.2 on the right.

**Definition 16.3.** Let  $X_0, X_1, \dots$  be a sequence of random variables defined on a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ , all taking values in  $E$ . The stochastic process  $(X_n)_{n \in \mathbb{N}_0}$  is called a *Markov chain* on  $E$  with transition probabilities  $\Pi$  if

$$\mathbf{P}[X_{n+1} = x_{n+1} | X_0 = x_0, \dots, X_n = x_n] = \Pi(x_n, x_{n+1}), \quad (16.5)$$

for every  $n \geq 0$  and  $x_0, \dots, x_{n+1} \in E$  such that  $\mathbf{P}[X_0 = x_0, \dots, X_n = x_n] > 0$ . The law  $\mu = \mathbf{P}_{X_0}$  of  $X_0$  is called the *initial distribution* of the Markov chain.

*Remark 16.4.* (i) From the preceding remarks, it is not immediate that Markov chains actually exist. One can show however, that given any initial distribution and any matrix of transition probabilities  $\Pi$ , one can construct a corresponding Markov chain. This construction can be found, e.g., in [2, Remark 6.2 (d)].

(ii) Note that the right-hand side of (16.5) does not explicitly depend on  $n$ , and also not on  $x_0, \dots, x_{n-1}$ . This formalizes the intuition given earlier: The next step taken by a Markov chain only depends on the *current position*  $\{X_n = x_n\}$ .

A Markov chain is in fact characterized completely by its initial distribution and its transition probabilities. To understand this, we make the following observation: Suppose that for some fixed  $x \in E$ :

$$\mathbf{P}[X_0 = y] = \mathbb{1}_{\{x=y\}}. \quad (16.6)$$

This means that at time 0, the chain starts in the point  $x$ . We have:

$$\begin{aligned}
\mathbf{P}[X_1 = x_1, \dots, X_n = x_n] &= \mathbf{P}[X_0 = x, X_1 = x_1, \dots, X_n = x_n] \\
&= \underbrace{\mathbf{P}[X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_1 = x_1, X_0 = x]}_{=\Pi(x_{n-1}, x_n)} \\
&\quad \cdot \mathbf{P}[X_{n-1} = x_{n-1}, \dots, X_1 = x_1, X_0 = x] \\
&= \Pi(x_{n-1}, x_n) \underbrace{\mathbf{P}[X_{n-1} = x_{n-1} | X_{n-2} = x_{n-2}, \dots, X_1 = x_1, X_0 = x]}_{=\Pi(x_{n-2}, x_{n-1})} \\
&\quad \cdot \mathbf{P}[X_{n-2} = x_{n-2}, \dots, X_1 = x_1, X_0 = x] \\
&= \dots = \Pi(x_{n-1}, x_n) \cdot \dots \cdot \Pi(x_2, x_1) \mathbf{P}[X_1 = x_1 | X_0 = x] \cdot \mathbf{P}[X_0 = x] \\
&= \prod_{j=1}^n \Pi(x_{j-1}, x_j).
\end{aligned} \tag{16.7}$$

Let us first exemplify this: Suppose in Example 16.2 (with  $\Pi_1$ ) we also know that  $\mathbf{P}[X_0 = 1]$ . What is the probability that the chain first jumps to 2, then to 3 and then to 3 again? Clearly:

$$\mathbf{P}[X_1 = 2, X_2 = 3, X_3 = 3] = \Pi(1, 2)\Pi(2, 3)\Pi(3, 3) = 0.2 \cdot 0.4 \cdot 0.5 = 0.04.$$

A similar calculation can be performed in the case where  $X_0$  has a non-trivial distribution  $\mu$  under  $\mathbf{P}$ . We summarize this in the following Proposition:

**Proposition 16.5.** *Suppose that  $(X_n)_{n \in \mathbb{N}_0}$  is a Markov chain on  $E$  with transition probabilities  $\Pi$  and initial distribution  $\mu$ . Then we have*

(i) *We have*

$$\mathbf{P}[X_0 = x_0, \dots, X_n = x_n] = \mu(x_0)\Pi(x_0, x_1) \cdot \dots \cdot \Pi(x_{n-1}, x_n). \tag{16.8}$$

(ii) *The Chapman-Kolmogorov equation holds: Define the  $n$ -step transition probabilities as*

$$\Pi^{(n)}(x, y) = \mathbf{P}[X_{n+m} = y | X_n = x]. \tag{16.9}$$

*Then*

$$\Pi^{(n)}(x, y) = \sum_{z \in E} \Pi^{(r)}(x, z) \Pi^{(n-r)}(z, y), \quad \text{for } 0 < r < n. \tag{16.10}$$

*In particular (16.10) means that in the  $n$ -step transition probability is the  $n$ -fold matrix product of  $\Pi$ , i.e.  $\Pi^{(n)} = \Pi^n = \Pi \cdot \dots \cdot \Pi$  ( $n$  times).*

(iii) *The probability for the chain to be at  $y$  at time  $n$  is given by*

$$\mathbf{P}[X_n = y] = \sum_{z \in E} \mu(z) (\Pi^n)(z, y) \tag{16.11}$$

Is it possible that the law of  $X_{n+1}$  coincides with the law of  $X_n$  for all  $n$ , i.e. for the Markov chain to be *stationary*? The answer is provided by (16.11): It states that in order to obtain the law of  $X_1$ , we can write the probability mass function of the law of  $X_0$  as a row vector and multiply it with  $\Pi$  from the left. For instance, suppose that  $E = \{1, 2, 3\}$  and  $\mu = (1 \ 0 \ 0)$ , i.e. at time 0 the chain starts in 1. Then considering again Example 16.2 (with  $\Pi_1$ ), clearly

$$(\mathbf{P}[X_1 = 1] \ \mathbf{P}[X_1 = 2] \ \mathbf{P}[X_1 = 3]) = (1 \ 0 \ 0) \cdot \begin{pmatrix} 0.1 & 0.2 & 0.7 \\ 0.3 & 0.3 & 0.4 \\ 0 & 0.5 & 0.5 \end{pmatrix} = (0.1 \ 0.2 \ 0.7).$$

This leads us to the following idea.

**Definition 16.6.** Let  $\Pi$  be a matrix of transition probabilities of a Markov chain on  $E = \{1, \dots, n\}$ . We say that a probability mass function  $(\pi(x))_{x \in E}$  is a *stationary* distribution for  $\Pi$  if

$$\pi = \pi \cdot \Pi \tag{16.12}$$

(as matrices), or in components:

$$\pi(x) = \sum_{z \in E} \pi(z) \Pi(z, x), \quad \text{for all } x \in E. \tag{16.13}$$

In other words: A stationary distribution is an *eigenvector* to the matrix  $\Pi^\top$  corresponding to eigenvalue 1, with positive entries and  $\ell^1$ -norm equal to 1. Finding a stationary distribution is therefore a problem in linear algebra! One may now ask, whether stationary distributions always exist and whether they are unique. The answer is given by the following theorem, which we state without proof.

**Theorem 16.7.** Let  $\Pi$  be a matrix of transition probabilities for a Markov chain on  $E = \{1, \dots, n\}$ . Suppose that there exists  $n \in \mathbb{N}$  with the property that<sup>2</sup>

$$\Pi^n(x, y) > 0 \quad \text{for all } x, y \in E. \tag{16.14}$$

Then there exists a unique stationary distribution for the Markov chain.

For many further topics, including

- classification of states (absorbing, recurrent, transient, ...),
- return times and their distribution,
- convergence of Markov chains,

we refer to the literature (see [2, Section 6], and [3] for a much more complete treatment).

---

<sup>2</sup>The following condition means that the chain is *irreducible*: In other words, every state  $y$  may be reached by starting from any other state  $x$  in some finite time.

## A. Appendix

### A.1. Multiple integrals

We give some details and intuition about multiple integrals. Such multiple integrals appear in Chapter 7. The focus of this appendix is *not* to study the most general multiple integrals, but rather to give some examples and calculation rules, and to act as a kind of *survival kit* if multiple integrals are unfamiliar.

*Remark A.1.* (i) Let  $f : [a, b] \rightarrow [0, \infty)$ ,  $a < b$ , a real function. Under some regularity conditions on  $f$  (for instance, if  $f$  is continuous) we know that for  $a \leq a_1 < b_1 \leq b$  the integral

$$\int_{[a_1, b_1]} f(x_1) dx_1 = \int_{a_1}^{b_1} f(x_1) dx_1 \quad (\text{A.1})$$

gives us the area that is enclosed by the graph of the function  $f$ , the two lines  $x_1 = a_1$  and  $x_2 = b_1$  and the  $x_1$ -axis (which can be expressed as the line  $y = 0$ ). Moreover, we remark that in (A.1), we could also have expressed the integral  $\int_{a_1}^{b_1} f(x_1) dx_1$  as  $\int_{[a_1, b_1]} f(x_1) dx_1$ ,  $\int_{(a_1, b_1)} f(x_1) dx_1$  or  $\int_{(a_1, b_1)} f(x_1) dx_1$ , since the boundary points do not contribute to the area.

(ii) Let  $f : [a, b] \times [a', b'] \rightarrow [0, \infty)$ , for some  $a < b$ ,  $a' < b'$ , be a real function depending on two variables. One can visualize the expression  $z = f(x_1, x_2)$  as a kind of “mountainous landscape”, where the value  $z$  describes the height with respect to zero over the point  $(x_1, x_2) \in [a, b] \times [a', b']$ . Consider the expression

$$\int_{[a_1, b_1] \times [a_2, b_2]} f(x_1, x_2) d^2x = \int_{a_1}^{b_1} \left( \int_{a_2}^{b_2} f(x_1, x_2) dx_2 \right) dx_1. \quad (\text{A.2})$$

If  $f$  is sufficiently regular (again, if  $f$  is continuous<sup>1</sup>, this is fulfilled), the expression in (A.2) can be considered as the volume enclosed by the graph of the function  $f$ , the four planes  $x_1 = a_1$ ,  $x_1 = b_1$ ,  $x_2 = a_2$ ,  $x_2 = b_2$  and the  $x_1$ - $x_2$ -plane (which can be expressed as  $z = 0$ ). A graphical representation is given in Figure A.1 below. Similarly as in the one-dimensional case, we could also write the expression in (A.2) as  $\int_{[a_1, b_1] \times [a_2, b_2]} f(x_1, x_2) d^2x$ ,  $\int_{(a_1, b_1) \times (a_2, b_2)} f(x_1, x_2) d^2x$ , etc., since planes like  $x_1 = a_1$ ,  $x_1 = b_1$ ,  $x_2 = a_2$  and  $x_2 = b_2$  give no contribution to the volume.

(iii) The expression in (A.2) is evaluated from the inner integral to the outer integral. This means that we first perform the integration with respect to  $x_2$ , treating  $x_1$  as a constant.

<sup>1</sup>Continuity means that for every two sequences  $(x_1^{(n)})_{n \geq 1} \subseteq [a, b]$  and  $(x_2^{(n)})_{n \geq 1} \subseteq [a', b']$  that fulfill  $x_1^{(n)} \rightarrow x_1$ ,  $x_2^{(n)} \rightarrow x_2$  as  $n \rightarrow \infty$ , we necessarily have  $\lim_{n \rightarrow \infty} f(x_1^{(n)}, x_2^{(n)}) = f(x_1, x_2)$ . Intuitively this means that the graph of  $f$  does not have any “fault lines”.

The resulting function  $x_1 \mapsto \int_{a_2}^{b_2} f(x_1, x_2) dx_2$  is then integrated with respect to  $x_1$ , over the interval  $[a_1, b_1]$ . We demonstrate this in Example A.2 below.

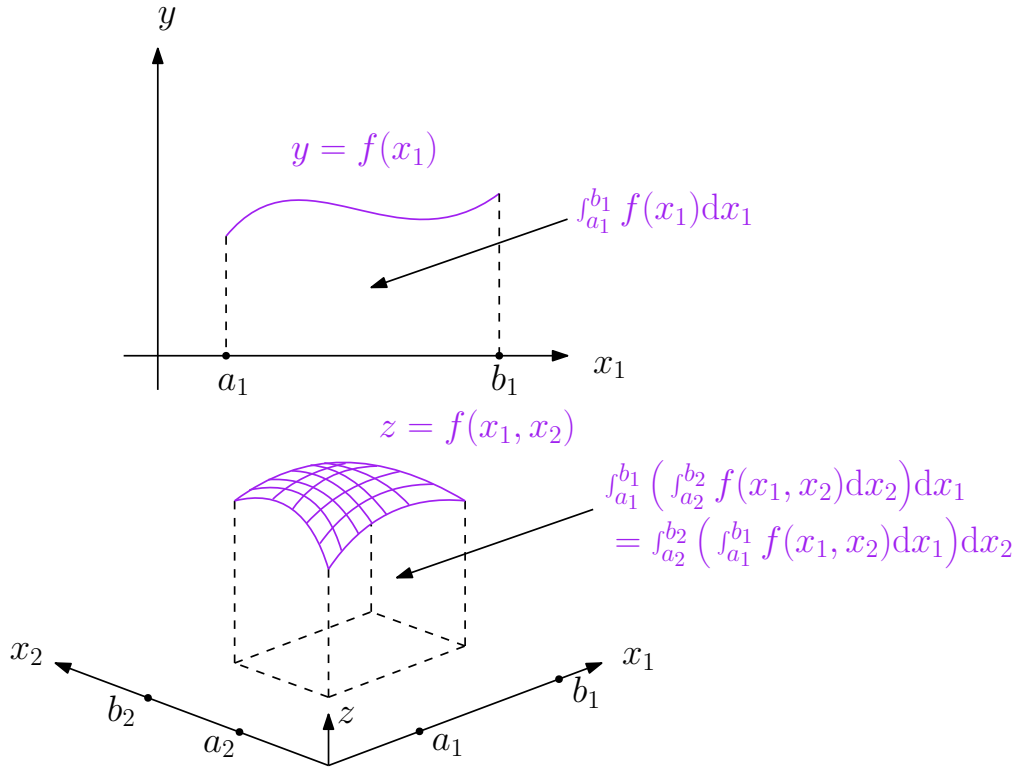


Figure A.1.: Upper panel: The area under a curve for a function  $f$  depending on one variable. Lower panel: The volume under the graph for a function  $f$  depending on two variables.

- (iv) It can be shown that under certain conditions on  $f$  (which are always fulfilled in this course, again continuity on  $[a, b] \times [a', b']$  suffices), the expression in (A.2) can also be rewritten as follows:

$$\int_{a_1}^{b_1} \left( \int_{a_2}^{b_2} f(x_1, x_2) dx_2 \right) dx_1 = \int_{a_2}^{b_2} \left( \int_{a_1}^{b_1} f(x_1, x_2) dx_1 \right) dx_2. \quad (\text{A.3})$$

In other words, the order in which we perform the integration does not matter. This fact is known as *Fubini's theorem*. Some care is required if we do not integrate  $f$  over a rectangle, but over more complicated domains (see the discussion below).

(v) Now let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . For  $a_i < b_i$ , we can look at the expression<sup>2</sup>

$$\begin{aligned} \int_{[a_1, b_1] \times [a_2, b_2] \times \dots \times [a_n, b_n]} f(x_1, \dots, x_n) d^n x \\ = \int_{a_1}^{b_1} \left( \int_{a_2}^{b_2} \dots \left( \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_n \right) \dots dx_2 \right) dx_1. \end{aligned} \quad (\text{A.4})$$

Again, this expression is evaluated from the inside to the outside. Note that is nothing else but (A.1) in the case of  $n = 1$  or (A.2) in the case of  $n = 2$ . In the latter cases, we can also make sense of negative values of the function  $f$ , by assigning the area below the  $x_1$ -axis a negative value and the volume below the  $x_1$ - $x_2$ -plane a negative value. As in (A.3), we can perform the integrals in any order we like. Mathematically, one has for every permutation  $\sigma : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$  (a bijective map from  $\{1, 2, \dots, n\}$  to itself, which is nothing else but a reordering of the numbers from 1 to  $n$ ):

$$\begin{aligned} \int_{a_1}^{b_1} \left( \int_{a_2}^{b_2} \dots \left( \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_n \right) \dots dx_2 \right) dx_1 \\ = \int_{a_{\sigma(1)}}^{b_{\sigma(1)}} \left( \int_{a_{\sigma(2)}}^{b_{\sigma(2)}} \dots \left( \int_{a_{\sigma(n)}}^{b_{\sigma(n)}} f(x_1, \dots, x_n) dx_{\sigma(n)} \right) \dots dx_{\sigma(2)} \right) dx_{\sigma(1)}. \end{aligned} \quad (\text{A.5})$$

Finally, we can define the improper integral

$$\begin{aligned} \int_{\mathbb{R}^n} f(x_1, \dots, x_n) d^n x &= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} \dots \left( \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_n \right) \dots dx_2 \right) dx_1 \\ &= \lim_{a_1 \rightarrow -\infty} \lim_{b_1 \rightarrow \infty} \dots \lim_{a_n \rightarrow -\infty} \lim_{b_n \rightarrow \infty} \int_{a_1}^{b_1} \left( \int_{a_2}^{b_2} \dots \left( \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_n \right) \dots dx_2 \right) dx_1 \end{aligned} \quad (\text{A.6})$$

Improper integrals in which only one boundary is infinity (which appear for instance when calculating joint cumulative distribution functions, see Proposition 7.9) are defined similarly. For instance:

$$\begin{aligned} \int_{(-\infty, b_1] \times \dots \times (-\infty, b_n]} f(x_1, \dots, x_n) d^n x \\ = \int_{-\infty}^{b_1} \left( \int_{-\infty}^{b_2} \dots \left( \int_{-\infty}^{b_n} f(x_1, \dots, x_n) dx_n \right) \dots dx_2 \right) dx_1 \\ = \lim_{a_1 \rightarrow -\infty} \dots \lim_{a_n \rightarrow -\infty} \int_{a_1}^{b_1} \left( \int_{a_2}^{b_2} \dots \left( \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_n \right) \dots dx_2 \right) dx_1 \end{aligned} \quad (\text{A.7})$$

---

<sup>2</sup>Or equivalently  $\int_{[a_1, b_1) \times [a_2, b_2) \times \dots \times [a_n, b_n)} f(x_1, \dots, x_n) d^n x$ ,  $\int_{[a_1, b_1) \times [a_2, b_2) \times \dots \times [a_n, b_n)} f(x_1, \dots, x_n) d^n x$  or any other number of expressions, where we only modify whether boundary values  $a_j$  or  $b_j$  are included.



*Example A.2.* (i) Consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(x_1, x_2) = x_1 x_2^2$ . We want to integrate  $f$  over the rectangle  $[0, 2] \times [0, 1] \subseteq \mathbb{R}^2$ . This integral is

$$\begin{aligned} \int_{[0,2] \times [0,1]} f(x_1, x_2) d^2x &= \int_0^2 \left( \int_0^1 f(x_1, x_2) dx_2 \right) dx_1 \\ &= \int_0^2 \left( \int_0^1 x_1 x_2^2 dx_2 \right) dx_1 \\ &= \int_0^2 x_1 \left[ \frac{1}{3} x_2^3 \right]_0^1 dx_1 \\ &= \frac{1}{3} \int_0^2 x_1 dx_1 = \frac{1}{3} \left[ \frac{1}{2} x_1^2 \right]_0^2 = \frac{2}{3}. \end{aligned} \tag{A.8}$$

We have marked in the relevant inner integral all terms that play a role ( $x_1$  is treated as a constant when integrating over  $x_2$ ).

(ii) Consider the function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $f(x_1, x_2, x_3) = \cos(x_1) x_2 \exp(x_2 x_3)$ . We integrate  $f$  over the cuboid  $[0, 1] \times [2, 3] \times [1, 4]$ .

$$\begin{aligned} \int_{[0,1] \times [2,3] \times [1,4]} f(x_1, x_2, x_3) d^3x &= \int_0^1 \left( \int_2^3 \left( \int_1^4 f(x_1, x_2, x_3) dx_3 \right) dx_2 \right) dx_1 \\ &= \int_0^1 \left( \int_2^3 \left( \int_1^4 \cos(x_1) x_2 \exp(x_2 x_3) dx_3 \right) dx_2 \right) dx_1 \\ &= \int_0^1 \left( \int_2^3 \cos(x_1) x_2 \left[ \frac{1}{x_2} \exp(x_2 x_3) \right]_{x_3=1}^4 dx_2 \right) dx_1 \\ &= \int_0^1 \cos(x_1) \left( \int_2^3 (\exp(4x_2) - \exp(x_2)) dx_2 \right) dx_1 \\ &= \int_0^1 \cos(x_1) \left[ \frac{1}{4} \exp(4x_2) - \exp(x_2) \right]_{x_2=2}^3 dx_1 \\ &= \left( \frac{1}{4} e^8 (e^4 - 1) - (e - 1) e^2 \right) [\sin(x_1)]_{x_1=0}^1 \\ &= \left( \frac{1}{4} e^8 (e^4 - 1) - (e - 1) e^2 \right) \sin(1). \end{aligned} \tag{A.9}$$

Note that in the innermost integral in the second line, although  $x_2$  was treated as constant, one has to be a bit careful since we cannot divide by 0. This is however excluded, since  $x_2 \in [2, 3]$ .

So far, we only discussed integrals over generalized rectangles  $[a_1, b_1] \times \dots \times [a_n, b_n]$ . In one dimension, this is what we really need. However, in dimension  $n \geq 2$ , we may want to integrate over more general sets  $A \subseteq \mathbb{R}^n$ , i.e. calculate integrals of the form

$$\int_A f(x_1, \dots, x_n) d^n x. \tag{A.10}$$

In general, this is a hard problem (first we would have to clarify which kinds of sets can be used to define integrals: In essence, one can use here any  $A \in \mathcal{B}(\mathbb{R}^n)$ ). We will only need sets of specific types, for instance those that we encountered in Example 7.10.

*Remark A.3.* (i) Suppose we are given the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , assumed to be sufficiently regular, which we want to integrate over the following area  $A$ , given by

$$A = \{(x_1, x_2) ; x_1 \in [a_1, b_1], x_2 \in [u(x_1), v(x_1)]\}. \quad (\text{A.11})$$

Here  $u : [a_1, b_1] \rightarrow \mathbb{R}$  and  $v : [a_1, b_1] \rightarrow \mathbb{R}$  are piecewise continuous functions that fulfill  $u(x_1) \leq v(x_1)$  for every  $x_1 \in [a_1, b_1]$ . We can define the integral

$$\int_A f(x_1, x_2) d^2x = \int_{a_1}^{b_1} \left( \int_{u(x_1)}^{v(x_1)} f(x_1, x_2) dx_2 \right) dx_1. \quad (\text{A.12})$$

In the case where  $f \geq 0$ , we can interpret this integral again as a volume, namely of the set bounded by the graph of the function  $f|_A : A \rightarrow [0, \infty)$ , the planes  $x_1 = a_1$ ,  $x_1 = b_1$ , the  $x_1$ - $x_2$ -plane and the “curved planes” given by the relations  $x_2(x_1) = u(x_1)$  and  $x_2(x_1) = v(x_1)$ . Figure A.3 below gives a graphical representation of the situation.

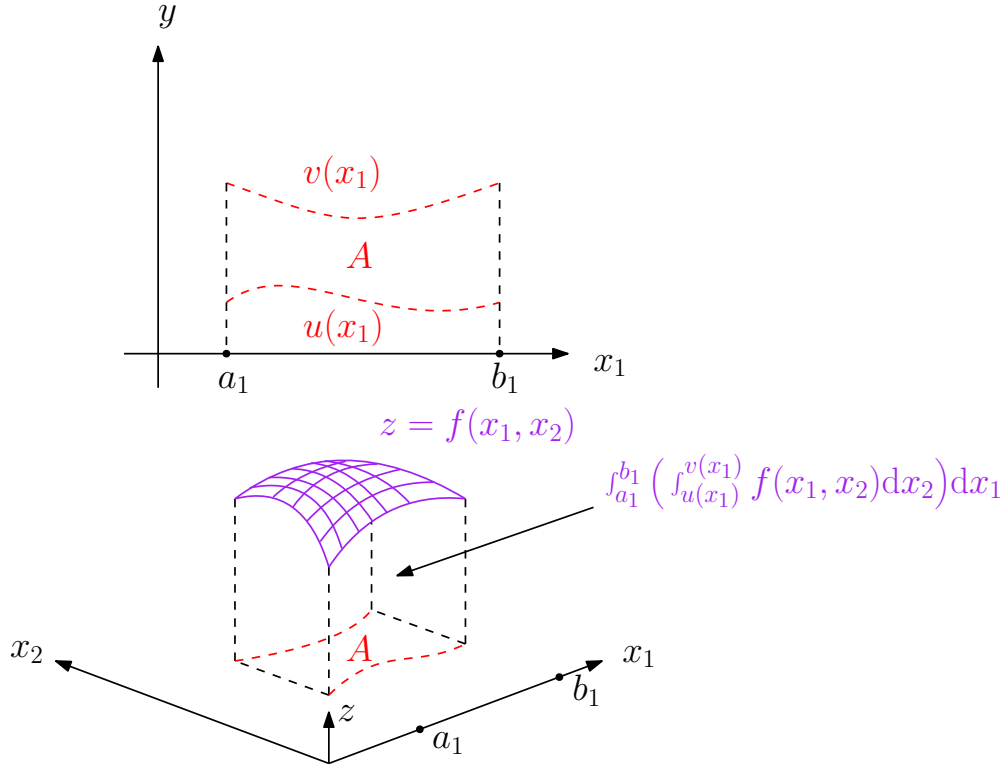


Figure A.2.: Upper panel: The area  $A$  over which we want to integrate over. Lower panel: The volume under the graph for a function  $f$  depending on two variables, integrated over  $A$ . The upper panel is the picture seen in the  $x_1$ - $x_2$ -plane of the lower panel.

(ii) More generally, we can consider integrals of (sufficiently regular) functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  over sets of the form

$$A = \{(x_1, \dots, x_n) ; x_1 \in [a_1, b_1], x_2 \in [u_2(x_1), v_2(x_2)], \dots, x_n \in [u_n(x_1, \dots, x_{n-1}), v_n(x_1, \dots, x_{n-1})]\}, \quad (\text{A.13})$$

where the functions  $u_2, v_2, u_3, v_3, \dots, u_n, v_n$  have to be sufficiently regular themselves and fulfill  $u_j \leq v_j$  for every  $j = 2, \dots, n$ . We set

$$\begin{aligned} & \int_A f(x_1, x_2, \dots, x_n) d^n x \\ &= \int_{a_1}^{b_1} \left( \int_{u_2(x_1)}^{v_2(x_1)} \left( \dots \left( \int_{u_n(x_1, \dots, x_{n-1})}^{v_n(x_1, \dots, x_{n-1})} f(x_1, x_2, \dots, x_n) dx_n \right) \dots \right) dx_2 \right) dx_1. \end{aligned} \quad (\text{A.14})$$

(iii) Let us remark that in dimension  $n = 2$  the expression  $\int_A 1 d^2 x = \int_{a_1}^{b_1} \left( \int_{u_2(x_1)}^{v_2(x_1)} dx_2 \right) dx_1$  (corresponding to setting  $f = 1$  in (A.12)) is exactly the area of the set  $A$  in (A.11). Formally this corresponds to the volume of a prism with base area  $A$  and height 1. In a similar manner, in dimension  $n = 3$ , the expression  $\int_V 1 d^3 x = \int_{a_1}^{b_1} \left( \int_{u_2(x_1)}^{v_2(x_1)} \left( \int_{u_3(x_1, x_2)}^{v_3(x_1, x_2)} f(x_1, x_2, x_3) dx_3 \right) dx_2 \right) dx_1$  gives us the volume of the set

$$V = \{(x_1, x_2, x_3) ; x_1 \in [a_1, b_1], x_2 \in [u_2(x_1), v_2(x_1)], x_3 \in [u_3(x_1, x_2), v_3(x_1, x_2)]\}. \quad (\text{A.15})$$

*Example A.4.* (i) Let us first calculate the integral (we write  $x, y, z$  instead of  $x_1, x_2, x_3$ ):

$$\begin{aligned} \int_0^1 \left( \int_0^{2x} \left( \int_0^{x+y} 1 dz \right) dy \right) dx &= \int_0^1 \left( \int_0^{2x} (x+y) dy \right) dx \\ &= \int_0^1 \left[ xy + \frac{1}{2} y^2 \right]_{y=0}^{2x} dx \\ &= 4 \int_0^1 x^2 dx \\ &= 4 \left[ \frac{1}{3} x^3 \right]_0^1 = \frac{4}{3}. \end{aligned} \quad (\text{A.16})$$

This calculation gives us the volume of the set

$$V = \{(x, y, z) ; 0 \leq x \leq 1, 0 \leq y \leq 2x, 0 \leq z \leq x+y\}. \quad (\text{A.17})$$

(ii) We give some more details in the calculations appearing in (7.29) and (7.32). Here, we given the information that  $X, Y$  are random variables with a joint density  $f_{X,Y}$ . This means that for a given set  $B \in \mathcal{B}(\mathbb{R}^2)$ , one has

$$\mathbf{P}[(X, Y) \in B] = \mathbf{P}_{(X,Y)}[B] = \int_B f_{X,Y}(x_1, x_2) d^2 x, \quad (\text{A.18})$$

see (7.17). We wanted to find the value  $\mathbf{P}[X^2 \leq Y]$ . Define the set

$$A = \{(x, y) ; x \in [0, 1], y \in [x^2, 1]\}, \quad (\text{A.19})$$

which is of the form (A.11) for  $u(x) = x^2$  and  $v(x) = 1$ . Moreover, we know that  $\mathbf{P}[(X, Y) \in [0, 1]^2] = 1$ . This means that

$$\begin{aligned} \mathbf{P}[X^2 \leq Y] &= \mathbf{P}[\{X^2 \leq Y\} \cap \{(X, Y) \in [0, 1]^2\}] \\ &= \mathbf{P}_{(X, Y)}[\{(x, y) ; x^2 \leq y \text{ and } 0 \leq x, y \leq 1\}] \\ &= \mathbf{P}_{(X, Y)}[A] \\ &= \int_0^1 \left( \int_{x^2}^1 f_{X, Y}(x, y) dy \right) dx. \end{aligned} \quad (\text{A.20})$$

The rest of the calculation proceeds as in (7.29) or (7.32).

## A.2. Alternative Proof of the central limit theorem 14.4

**Lemma A.5.** *Let  $X_1, \dots, X_n$  be i.i.d. random variables with  $\mathbf{E}[X_1] = 0$ ,  $\text{Var}[X_1] = 1$  and  $g : \mathbb{R} \rightarrow \mathbb{R}$  twice differentiable such that  $g''$  is uniformly continuous and bounded. For  $Z_n = \sqrt{n}\bar{X}_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$  and  $Z \sim \mathcal{N}(0, 1)$ , one has*

$$\mathbf{E}[g(Z_n)] \xrightarrow{n \rightarrow \infty} \mathbf{E}[g(Z)]. \quad (\text{A.21})$$

*Proof.* The idea is to replace  $X_1, \dots, X_n$  step-by-step by i.i.d. random variables that have a standard normal distribution, and to use Taylor's theorem.

Consider  $Y_1, \dots, Y_n \sim \mathcal{N}(0, 1)$  i.i.d. random variables such that  $X_1, \dots, X_n, Y_1, \dots, Y_n$  are stochastically independent. This means that

$$\sqrt{n}\bar{Y}_n \sim \mathcal{N}(0, 1) \quad \Rightarrow \quad \mathbf{E}[g(\sqrt{n}\bar{Y}_n)] = \mathbf{E}[g(Z)]. \quad (\text{A.22})$$

Now we replace  $X_i$  by  $Y_i$  as follows:

$$\begin{aligned} g(Z_n) - g(\sqrt{n}\bar{Y}_n) &= g\left(\frac{X_1 + \dots + X_n}{\sqrt{n}}\right) - g\left(\frac{Y_1 + \dots + X_n}{\sqrt{n}}\right) \\ &\quad + g\left(\frac{Y_1 + \dots + X_n}{\sqrt{n}}\right) - g\left(\frac{Y_1 + Y_2 + \dots + X_n}{\sqrt{n}}\right) \\ &\quad \vdots \\ &\quad + g\left(\frac{Y_1 + \dots + Y_{n-1} + X_n}{\sqrt{n}}\right) - g\left(\frac{Y_1 + \dots + Y_n}{\sqrt{n}}\right) \\ &= V_1 + \dots + V_n, \end{aligned} \quad (\text{A.23})$$

where we set

$$V_i = g\left(U_i + \frac{X_i}{\sqrt{n}}\right) - g\left(U_i + \frac{Y_i}{\sqrt{n}}\right), \quad (\text{A.24})$$

and

$$U_i = \frac{Y_1 + \dots + Y_{i-1} + X_{i+1} + \dots + X_n}{\sqrt{n}}, \quad 1 \leq i \leq n. \quad (\text{A.25})$$

A Taylor expansion of  $g$  around the point  $U_i$  yields

$$V_i = \frac{X_i - Y_i}{\sqrt{n}} g'(U_i) + \frac{1}{2n} X_i^2 g'' \left( U_i + \frac{\theta_1 X_i}{\sqrt{n}} \right) - \frac{1}{2n} Y_i^2 g'' \left( U_i + \frac{\theta_2 Y_i}{\sqrt{n}} \right), \quad (\text{A.26})$$

where  $\theta_1, \theta_2 \in [0, 1]$  are random variables. Consider the modulus of continuity

$$\begin{aligned} \delta(h) &= \sup_{|x-y| \leq h} |g''(x) - g''(y)| \\ \Rightarrow V_i &= \frac{X_i - Y_i}{\sqrt{n}} g'(U_i) + \frac{1}{2n} (X_i^2 - Y_i^2) g''(U_i) + R_i, \end{aligned} \quad (\text{A.27})$$

and we have

$$|R_i| \leq \frac{1}{2n} X_i^2 \delta \left( \frac{|X_i|}{\sqrt{n}} \right) + \frac{1}{2n} Y_i^2 \delta \left( \frac{|Y_i|}{\sqrt{n}} \right). \quad (\text{A.28})$$

Note that the random variables  $U_i$  and  $X_i - Y_i$ , resp.  $U_i$  and  $X_i^2 - Y_i^2$  are independent. By (A.26), we find that

$$\mathbf{E}[V_i] = \frac{1}{\sqrt{n}} \underbrace{\mathbf{E}[(X_i - Y_i)g'(U_i)]}_{=\mathbf{E}[X_i - Y_i]\mathbf{E}[g'(U_i)] = 0} + \frac{1}{2n} \underbrace{\mathbf{E}[(X_i^2 - Y_i^2)g''(U_i)]}_{=\mathbf{E}[X_i^2 - Y_i^2]\mathbf{E}[g''(U_i)]} + \mathbf{E}[R_i] \quad (\text{A.29})$$

Now it holds that

$$|\mathbf{E}[R_i]| \leq \mathbf{E}[|R_i|] \leq \frac{1}{2n} \mathbf{E} \left[ X_i^2 \delta \left( \frac{|X_i|}{\sqrt{n}} \right) \right] + \frac{1}{2n} \mathbf{E} \left[ X_i^2 \delta \left( \frac{|Y_i|}{\sqrt{n}} \right) \right]. \quad (\text{A.30})$$

Now we estimate

$$\begin{aligned} \mathbf{E} \left[ X_i^2 \delta \left( \frac{|X_i|}{\sqrt{n}} \right) \right] &= \mathbf{E} \left[ X_i^2 \delta \left( \frac{|X_i|}{\sqrt{n}} \right) \left( \mathbb{1}_{\{|X_i| \leq \lambda \sqrt{n}\}} + \mathbb{1}_{\{|X_i| > \lambda \sqrt{n}\}} \right) \right] \\ &\leq \delta(\lambda) \mathbf{E}[X_i^2] + C \mathbf{E}[X_i^2 \mathbb{1}_{\{|X_i| > \lambda \sqrt{n}\}}], \end{aligned} \quad (\text{A.31})$$

where  $C > 0$  is some constant. We have used that  $\delta$  is finite (since  $g''$  is bounded). Moreover, since  $g''$  is uniformly continuous, we know that  $\lim_{\lambda \downarrow 0} \delta(\lambda) = 0$ . On the other hand, for fixed  $\lambda > 0$ , we have

$$\lim_{n \rightarrow \infty} \mathbf{E}[X_i^2 \mathbb{1}_{\{|X_i| > \lambda \sqrt{n}\}}] = 0, \quad (\text{A.32})$$

(this follows from the *dominated convergence theorem* - in our case, we can only prove it if  $X_i$  is either continuous or discrete, when we can represent  $\mathbf{E}[X_i^2 \mathbb{1}_{\{|X_i| > \lambda \sqrt{n}\}}]$  either as a series or an integral). Combining everything, we see that

$$\mathbf{E}[R_i] \leq C \frac{1}{n} \cdot (\lambda + u_n), \quad (\text{A.33})$$

where  $\lim_{n \rightarrow \infty} u_n = 0$ . Thus:

$$\lim_{n \rightarrow \infty} |g(Z_n) - g(\sqrt{n} \bar{Y}_n)| \leq \lim_{n \rightarrow \infty} n |\mathbf{E}[V_1]| = \lim_{n \rightarrow \infty} u_n = 0. \quad (\text{A.34})$$

□

*Proof of Theorem 14.4.* We first observe that since  $X_1, X_2, \dots$  are i.i.d. with  $\mathbf{E}[X_1] = \mu$  and  $\text{Var}[X_1] = \sigma^2$ , the random variables  $X'_i = \frac{X_i - \mu}{\sigma}$  fulfill

$$\mathbf{E}[X'_1] = 0, \quad \text{Var}[X'_1] = 1, \quad \sqrt{n} \frac{\overline{X}_n - \mu}{\sigma} = \sqrt{n} \overline{X}'_n. \quad (\text{A.35})$$

Therefore, it suffices to show the statement of Lemma A.5 also for the functions  $g = \mathbb{1}_{(-\infty, x]}$ ,  $x \in \mathbb{R}$ . Indeed, if (A.21) holds for such  $g$ , we have

$$F_{\frac{\overline{X}_n - \mu}{\sigma}}(x) = \mathbf{E} [\mathbb{1}_{(-\infty, x]} (\sqrt{n} \overline{X}'_n)] \xrightarrow{n \rightarrow \infty} \mathbf{E} [\mathbb{1}_{(-\infty, x]}(Z)] = F_Z(x) = \Phi(x). \quad (\text{A.36})$$

So fix  $g = \mathbb{1}_{(-\infty, x]}$  and  $\varepsilon > 0$ . We choose  $g_1, g_2$  such that  $g_1''$  and  $g_2''$  are both uniformly continuous and bounded and fulfill

$$\begin{aligned} g_1(t) &\leq g(t) \leq g_2(t) \\ \int_{-\infty}^{\infty} (g_2(t) - g_1(t)) dt &\leq \varepsilon. \end{aligned} \quad (\text{A.37})$$

It follows that with  $Z_n = \sqrt{n} \overline{X}'_n$ , one has

$$\mathbf{E}[g_1(Z_n)] \leq \mathbf{E}[g(Z_n)] \leq \mathbf{E}[g_2(Z_n)]. \quad (\text{A.38})$$

The left- and rightmost items in the previous equation converge according to Lemma A.5 towards  $\mathbf{E}[g_1(Z)]$  and  $\mathbf{E}[g_2(Z)]$  respectively, and we have

$$\mathbf{E}[g_1(Z)] \leq \mathbf{E}[g(Z)] \leq \mathbf{E}[g_2(Z)]. \quad (\text{A.39})$$

By our choice of  $g_1$  and  $g_2$ , we also have

$$\mathbf{E}[g_2(Z) - g_1(Z)] = \int_{-\infty}^{\infty} (g_2(t) - g_1(t)) \varphi(t) dt \leq \varphi(0) \varepsilon. \quad (\text{A.40})$$

It follows that

$$\limsup_{n \rightarrow \infty} |\mathbf{E}[g(Z_n)] - \mathbf{E}[g(Z)]| \leq \varphi(0) \varepsilon. \quad (\text{A.41})$$

Since the latter holds for every  $\varepsilon > 0$ , we have proved the convergence (A.36).  $\square$

### A.3. More properties of convergence in distribution

**Proposition A.6.** *Let  $(X_n)_{n \in \mathbb{N}}, (Y_n)_{n \in \mathbb{N}}$  be sequences of real random variables with*

$$Y_n \xrightarrow[n \rightarrow \infty]{d} X, \quad X_n - Y_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0. \quad (\text{A.42})$$

*Then it follows that*

$$X_n \xrightarrow[n \rightarrow \infty]{d} X. \quad (\text{A.43})$$

*Proof.* Let  $Z_n = Y_n - X_n$ . We need to show that  $F_{X_n}(x) \rightarrow F_X(x)$  for all points of continuity  $x$  of  $F_X(x)$ . Let  $x \in \mathbb{R}$  and  $\varepsilon > 0$  such that  $x, x \pm \varepsilon$  are points of continuity of  $F_X$ . Then

$$\begin{aligned} F_{X_n}(x) &= \mathbf{P}[X_n \leq x] = \mathbf{P}[Y_n \leq x + Z_n] \\ &= \mathbf{P}[Y_n \leq x + Z_n, Z_n < \varepsilon] + \mathbf{P}[Y_n \leq x + Z_n, Z_n \geq \varepsilon] \\ &\leq \mathbf{P}[Y_n \leq x + \varepsilon] + \mathbf{P}[|Z_n| \geq \varepsilon] \\ \Rightarrow \quad \limsup_{n \rightarrow \infty} F_{X_n}(x) &\leq F_X(x + \varepsilon). \end{aligned} \tag{A.44}$$

Analogously (by setting  $Z_n > -\varepsilon$  instead of  $Z_n < \varepsilon$ ):

$$\liminf_{n \rightarrow \infty} F_{X_n}(x) \geq F_X(x - \varepsilon). \tag{A.45}$$

Finally, since  $x$  is a point of continuity of  $F_X$ , we find that

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x). \tag{A.46}$$

□

The following gives an overview of relations between convergence in probability and convergence in distribution.

**Theorem A.7.** *Let  $(X_n)_{n \in \mathbb{N}}$  and  $(Y_n)_{n \in \mathbb{N}}$  be sequences of real random variables,  $X$  a real random variable and  $c \in \mathbb{R}$  a constant.*

(i) *If  $(X_n)_{n \in \mathbb{N}}$  converges in probability to  $X$ , then it converges also to  $X$  in distribution:*

$$X_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} X \quad \Rightarrow \quad X_n \xrightarrow[n \rightarrow \infty]{d} X. \tag{A.47}$$

(ii) *We have that*

$$X_n \xrightarrow[n \rightarrow \infty]{d} X, Y_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0 \quad \Rightarrow \quad X_n Y_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0. \tag{A.48}$$

(iii) *One has Slutsky's theorem:*

$$\begin{aligned} X_n \xrightarrow[n \rightarrow \infty]{d} X, Y_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} c \quad \Rightarrow \quad & X_n + Y_n \xrightarrow[n \rightarrow \infty]{d} X + c, \\ & X_n Y_n \xrightarrow[n \rightarrow \infty]{d} cX, \\ & \frac{X_n}{Y_n} \xrightarrow[n \rightarrow \infty]{d} \frac{X}{c}, \text{ if } c \neq 0. \end{aligned} \tag{A.49}$$

*Proof.* For (i), set  $Y_n \equiv X$  in Proposition A.6.

For (ii), let  $k, \varepsilon > 0$ . Then

$$\begin{aligned} \mathbf{P}[|X_n Y_n| \geq \varepsilon] &= \mathbf{P}[|X_n Y_n| \geq \varepsilon, |Y_n| < \frac{\varepsilon}{k}] + \mathbf{P}[|X_n Y_n| \geq \varepsilon, |Y_n| \geq \frac{\varepsilon}{k}] \\ &\leq \mathbf{P}[|X_n| > k] + \mathbf{P}[|Y_n| \geq \frac{\varepsilon}{k}] \\ \Rightarrow \limsup_{n \rightarrow \infty} \mathbf{P}[|X_n Y_n| \geq \varepsilon] &\leq \mathbf{P}[|X| > k], \text{ if } \pm k \text{ are points of continuity of } F_X. \end{aligned} \quad (\text{A.50})$$

Replace  $k$  by a monotone sequence of points of continuity  $k_m \in \mathbb{R}$ . Since  $\sum_{m=1}^{\infty} \mathbf{P}[|X| \in (k_{m-1}, k_m]] < \infty$ , we have  $\lim_{m \rightarrow \infty} \mathbf{P}[|X| > k_m] = 0$ , thus (A.48) follows.

For (iii), we first see that

$$(X_n + Y_n) - (X_n + c) = Y_n - c \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0 \quad \Rightarrow \quad X_n + Y_n \xrightarrow[n \rightarrow \infty]{d} X + c, \quad (\text{A.51})$$

by Proposition A.6. Moreover, we have

$$X_n \xrightarrow[n \rightarrow \infty]{d} X \quad \Rightarrow \quad cX_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} cX. \quad (\text{A.52})$$

On the other hand, by (ii):

$$X_n Y_n - cX_n = X_n(Y_n - c) \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0 \quad \Rightarrow \quad X_n Y_n \xrightarrow[n \rightarrow \infty]{d} cX, \quad (\text{A.53})$$

by Proposition A.6. The last part of (iii) follows analogously.  $\square$

We illustrate the use of Proposition A.6 and Theorem A.7.

*Example A.8.* Let  $X_1, X_2, \dots$  be i.i.d. real random variables,  $\mathbf{E}[X_1] = \mu$ ,  $\text{Var}[X_1] = \sigma^2 \in (0, \infty)$  and  $\mathbf{E}[(X_1 - \mu)^4] = \mu_4^{(c)} \in (\sigma^4, \infty)$  the centered fourth moment of  $X_1$ . By the weak law of large numbers (Theorem 13.2), we know that

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \mathbf{E}[X_1] = \mu. \quad (\text{A.54})$$

If we let  $Y_i = (X_i - \mu)^2$ , then we have  $\mathbf{E}[Y_1] = \sigma^2$  and  $\text{Var}[Y_1] = \mathbf{E}[Y_1^2] - \mathbf{E}[Y_1]^2 = \mu_4^{(c)} - \sigma^4 \in (0, \infty)$ , so by the central limit theorem 14.4 we have that

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2 \right) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, \mu_4^{(c)} - \sigma^4). \quad (\text{A.55})$$

As an example consider the empirical variance given by

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2. \quad (\text{A.56})$$

We claim that

$$\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, \mu_4^{(c)} - \sigma^4). \quad (\text{A.57})$$



To see this, we consider the difference between the expression in (A.55) and the estimator  $S_n^2$  in (A.56):

$$\begin{aligned}
 & \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 - \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \right) \\
 = & - \underbrace{\sqrt{n}(\bar{X}_n - \mu)}_{\xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, \sigma^2), (14.9)} \cdot \underbrace{(\bar{X}_n - \mu)}_{\xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0, (13.3)} \\
 & \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0, \quad \text{by Theorem A.7, (ii).}
 \end{aligned} \tag{A.58}$$

Since the difference between  $\sqrt{n} \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$  and  $\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2)$  converges in probability to zero, we can combine (A.55) and Proposition A.6 to obtain (A.57).

## **Bibliography**

- [1] S. Ross, *A First Course in Probability*, 10th ed. Pearson, 2018.
- [2] H.-O. Georgii, *Stochastics: Introduction to Probability and Statistics*, 2nd ed. De Gruyter, 2012.
- [3] J.R. Norris, *Markov chains*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.