

# Probability Theory

Kalle Kytölä



## Contents

Foreword	iv
Glossary of notations	vi
Chapter O. Introduction	ix
O.1. What are the basic objects of probability theory?	ix
O.2. Informal examples of the basic objects in random phenomena	x
O.3. Probability theory vs. measure theory	xiii
Chapter I. Structure of event spaces	1
I.1. Set operations on events	1
I.2. Definition of sigma algebra	2
I.3. Generating sigma algebras	3
Chapter II. Measures and probability measures	7
II.1. Measurable spaces	7
II.2. Definition of measures and probability measures	8
II.3. Properties of measures and probability measures	13
II.4. Identification and construction of measures	16
Chapter III. Random variables	21
III.1. Measurable functions and random variables	22
III.2. Indicator random variables	23
III.3. Constructing random variables	24
III.4. Simple functions	30
Chapter IV. Information generated by random variables	35
IV.1. Definition of $\sigma$ -algebra generated by random variables	36
IV.2. Doob's representation theorem	37
Chapter V. Independence	39
V.1. Definition of independence	39
V.2. Verifying independence	42
V.3. Borel – Cantelli lemmas	43
Chapter VI. Events of the infinite horizon	47
VI.1. Definition of the tail $\sigma$ -algebra	47
VI.2. Kolmogorov's 0-1 law	50
Chapter VII. Integration theory	55
VII.1. Integral for non-negative simple functions	57
VII.2. Integral for non-negative measurable functions	59
VII.3. Integral for integrable functions	63
VII.4. Convergence theorems for integrals	65
VII.5. Integrals over subsets and restriction of measures	69
VII.6. Riemann integral vs. Lebesgue integral	70

Chapter VIII. Expected values	71
VIII.1. Expected values in terms of laws	72
VIII.2. Applications of convergence theorems for expected values	75
VIII.3. Space of $p$ -integrable random variables	78
Chapter IX. Product spaces and Fubini's theorem	81
IX.1. Product sigma algebra	82
IX.2. Product measure	84
Chapter X. Probability on product spaces	93
X.1. Joint laws	93
X.2. Variances and covariances of square integrable random variables	97
X.3. Independence and products	99
X.4. A formula for the expected value	103
Chapter XI. Probabilistic notions of convergence	107
XI.1. Notions of convergence in stochastics	107
XI.2. Weak and strong laws of large numbers	110
XI.3. Proof of the weak law of large numbers	111
XI.4. Proof of the strong law of large numbers	113
XI.5. Kolmogorov's strong law of large numbers	115
Chapter XII. Central limit theorem and convergence in distribution	117
XII.1. Characteristic functions	118
XII.2. Convergence in distribution	125
XII.3. Central limit theorem	125
Appendix A. Set theory preliminaries	127
A.1. Intersections and unions of sets	127
A.2. Set differences and complements	127
A.3. Images and preimages of sets under functions	128
A.4. Cartesian products	128
A.5. Power set	129
A.6. Sequences of sets	129
A.7. Countable and uncountable sets	130
Appendix B. Topological preliminaries	137
B.1. Topological properties of the real line	137
B.2. Metric space topology	140
Appendix C. Dynkin's identification and monotone class theorem	145
C.1. Monotone class theorem	145
C.2. Auxiliary results	145
C.3. Proof of Dynkin's identification theorem	147
C.4. Proof of Monotone class theorem	148
Appendix D. Monotone convergence theorem	149
D.1. Monotone convergence theorem for simple functions	149
D.2. Monotone convergence theorem for general non-negative functions	150
Appendix E. Orthogonal projections and conditional expected values	155
E.1. Geometry of the space of square integrable random variables	155
E.2. Conditional expected values	160

Appendix F. Characteristic functions	165
F.1. Lévy's inversion theorem	165
F.2. Equivalent conditions for convergence in distribution	169
Appendix. Index	175
Appendix. References	179

## Foreword

These lecture notes are primarily intended for the regular M.Sc. level course **MS-E1600** *Probability Theory* at Aalto University.

The principal aim of the course is to familiarize the students with the mathematical foundations of randomness. The reasons why one should study such a theoretical formalism vary according to the ambitions of the individual. The development of a logically solid theory of random phenomena should perhaps be seen as worthwhile in its own right. We will stick strictly to the fundamentals, so this course offers quite ideal practice on precise mathematical reasoning, including formulating proofs. A more pragmatic motivation might be that these theoretical foundations are necessary for following many subsequent courses in probability and statistics, and for understanding more advanced topics. In any case, whether one plans to work in statistics, machine learning, or pure mathematics, relevant research literature often requires familiarity of this theory as a language, e.g., being able to distinguish between convergence in probability, convergence in distribution, convergence almost surely, or other notions of convergence of random variables. The present course attempts to provide just enough of the core mathematical theory to develop an appreciation of such differences.

The course in its current format is very concise: 12 lectures and six sets of exercises during a six weeks period. One of the regrettable consequences is that there is almost no time to enter any of the interesting applications of the theory that is being developed. There are other courses devoted to more specific topics in probability which build on the theoretical foundations of the present course and come closer to actual applications.

In order to be prepared to internalize the theory during this concise course, the student must have a little bit of mathematical maturity to begin with. Besides some calculus of infinite series, differentiation, and integration, it is crucial to have a working knowledge of set theory, especially the notion of countability, and a little bit of familiarity with continuous functions in the context of metric spaces or topological spaces, say. Appendices A and B serve as quick reminders of such prerequisites, and before engaging in this text beyond the introduction, one should make sure to grasp their content.

The material in the chapters which correspond to the 12 lectures has been kept to minimum. A number of basic results that do not fit in these are left to Appendices C – F. The material in the main chapters can be considered as the minimum of what the students are expected to internalize during the short course, while the material in these appendices is something that one can expect to encounter soon after this basic course, and it can then be quickly picked up with a little bit of further effort.

There is already a vast number of textbooks in probability theory, and the contents of mathematics courses at advanced B.Sc. or early M.Sc. level have become quite standard. For the students of the present course we recommend in particular [JP04], because it is a very concise account of probability theory quickly covering very much the same topics as the present course. A slightly more challenging alternative is [Wil91], which is a remarkably well-written, mathematically elegant account of probability that manages incorporate fascinating and important probabilistic insights into a brief text. Both the theory and a significant number of interesting

and relevant examples and applications are covered in [Dur10]. The present lecture notes borrow shamelessly from all of the above sources. And the purpose of these notes is not to replace the best textbooks on the subject, but rather to provide the students an account that follows the structure and scope of the concise six weeks probability theory course as closely as possible.

The structure of these notes is largely based on an earlier version of the course taught by Lasse Leskelä, and on parts of the textbook [Wil91]. I have received very valuable comments, especially by Joonas Karjalainen, Alex Karrila, and Niko Lietzén, as well as many students, which have led to improvements to the notes. I am, of course, responsible for all remaining mistakes. Still, you could help me — and perhaps more importantly the students who will use this material — by sending comments about mistakes, misprints, needs for clarification, etc., to me by email ([kalle.kytola@aalto.fi](mailto:kalle.kytola@aalto.fi)).

## Glossary of notations

For convenience, we provide here a list of some of the used mathematical notation and abbreviations, together with brief explanations or references to the appropriate definitions during the course.

### Numbers

$\mathbb{Z}$	the set of integers	$\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$
$\mathbb{Z}_{\geq 0}$	the set of non-negative integers	$\mathbb{Z}_{\geq 0} = \{0, 1, 2, 3, \dots\}$
$\mathbb{N}$	the set of natural numbers	$\mathbb{N} = \{1, 2, 3, 4, \dots\}$
$\mathbb{Q}$	the set of rational numbers	$\mathbb{Q} = \{\frac{n}{m} \mid n, m \in \mathbb{Z}, m \neq 0\}$
$\mathbb{R}$	the set of real numbers	
$\mathbb{C}$	the set of complex numbers	$\mathbb{C} = \{x + iy \mid x, y \in \mathbb{R}\}$
$i$	imaginary unit	$i = \sqrt{-1} \in \mathbb{C}$
$(a, b)$	open interval	$(a, b) = \{x \in \mathbb{R} \mid a < x < b\}$
$[a, b]$	closed interval	$[a, b] = \{x \in \mathbb{R} \mid a \leq x \leq b\}$
$(a, b]$		$(a, b] = \{x \in \mathbb{R} \mid a < x \leq b\}$
$[a, b)$		$[a, b) = \{x \in \mathbb{R} \mid a \leq x < b\}$

### Logical notation

$\Rightarrow$	logical implication (“only if”)	$P \Rightarrow Q$ means: if $P$ is true then also $Q$ is true
$\Leftarrow$	reverse logical implication (“if”)	$P \Leftarrow Q$ means: if $Q$ is true then also $P$ is true
$\Leftrightarrow$	logical equivalence	$P \Leftrightarrow Q$ means: $P$ is true if and only if $Q$ is true
$\forall$	“for all” (logical quantifier)	
$\exists$	“there exists” (logical quantifier)	
s.t.	such that	
iff	if and only if	



**Set operations**

$s \in S$	$s$ is an element of set $S$	
$\#S$	number of elements in the set $S$	
$\subset$	subset relation	$A \subset B$ means: if $x \in A$ then also $x \in B$
$\emptyset$	the empty set	
$\mathcal{P}$	power set / collection of all subsets	$\mathcal{P}(S) = \{A \subset S\}$
$\cup, \bigcup$	union	$A \cup B = \{x \mid x \in A \text{ or } x \in B\}$ $\bigcup_j A_j = \{x \mid x \in A_j \text{ for some } j\}$
$\cap, \bigcap$	intersection	$A \cap B = \{x \mid x \in A \text{ and } x \in B\}$ $\bigcap_j A_j = \{x \mid x \in A_j \text{ for all } j\}$
$\times$	Cartesian product	$A \times B = \{(a, b) \mid a \in A, b \in B\}$
$\setminus$	set difference	$A \setminus B = \{x \mid x \in A, x \notin B\}$
$(\dots)^c$	complement of $\dots$	
$\limsup$		$\limsup_n A_n := \bigcap_{m \in \mathbb{N}} \bigcup_{n \geq m} A_n$
$\liminf$		$\liminf_n A_n := \bigcup_{m \in \mathbb{N}} \bigcap_{n \geq m} A_n$

**Probability and measure theory**

$\Omega$	sample space / the set of outcomes	
$\omega$	an outcome	$\omega \in \Omega$
$\mathcal{F}$	sigma-algebra / the collection of events	see Lecture I
$\mathbf{P}$	probability measure	see Lecture II
$\mathbf{E}$	expected value	see Lecture VII
$\mathbb{I}_A$	indicator of event $A$	
$\perp$	independent	see Lecture V
a.s.	almost surely / with probability one	Remark II.6

*Specific sigma algebras*

$\mathcal{S}$	generic sigma-algebra on a set $S$	Definition I.1
$\mathcal{F}$	generic sigma-algebra on the sample space $\Omega$	
$\mathcal{B}(\mathfrak{X})$	Borel sigma-algebra on a topological space $\mathfrak{X}$	Definition I.10
$\mathcal{B} = \mathcal{B}(\mathbb{R})$	Borel sigma-algebra on the real line $\mathbb{R}$	Section I.3.1
$\mathcal{T}_\infty$	tail sigma-algebra	Equation (VI.4)
$\sigma(\dots)$	the sigma-algebra generated by $\dots$	Definitions I.6 and IV.1

*Specific measures*

$\mu$	generic measure on measurable space $(S, \mathcal{S})$	Definition II.4
$\mu_\#$	counting measure (on some set)	Example II.10
$\mathbf{P}$	generic probability measure on sample space $\Omega$	Definition II.5
$\Lambda$	Lebesgue measure on the real line $\mathbb{R}$	Example II.12

*Spaces of random variables*

$m\mathcal{F}$	random variables measurable w.r.t. $\mathcal{F}$
$m\mathcal{F}^+$	non-negative random variables measurable w.r.t. $\mathcal{F}$
$s\mathcal{F}$	simple random variables measurable w.r.t. $\mathcal{F}$
$s\mathcal{F}^+$	non-negative simple random variables measurable w.r.t. $\mathcal{F}$
$b\mathcal{F}$	bounded random variables measurable w.r.t. $\mathcal{F}$
$b\mathcal{F}^+$	non-negative bounded random variables measurable w.r.t. $\mathcal{F}$
$\mathcal{L}^p$	$p$ -integrable random variables

*Notions of convergence*

$\xrightarrow{\text{a.s.}}$	convergence almost surely	Definition XI.1
$\xrightarrow{\mathbf{P}}$	convergence in probability	Definition XI.2
$\xrightarrow{\mathcal{L}^1}$	convergence in $\mathcal{L}^1(\mathbf{P})$	Definition XI.8
$\xrightarrow{\text{law}}$	convergence in distribution (law)	Definition XII.10

## Introduction

### O.1. What are the basic objects of probability theory?

Probability theory forms the mathematically precise and powerful foundations for the study of randomness. Its most basic objects — defined and studied in the rest of this course — are:

$\Omega$  — **Outcomes (of a random experiment)**

An *outcome*  $\omega$  of a random experiment represents a single *realization* of the randomness involved. The *sample space*  $\Omega$  is the set consisting of all possible outcomes.

$\mathcal{F}$  — **Events**<sup>1</sup>

An *event*  $E$  is a subset  $E \subset \Omega$  of the set of possible outcomes. The event  $E$  is said to occur if the randomly realized outcome  $\omega \in \Omega$  belongs to this subset, i.e., if  $\omega \in E$ . Generally we can not allow all subsets of  $\Omega$  as events, but instead we have to select a suitable collection  $\mathcal{F}$  of subsets on which it is possible to have consistent rules of probability.

$P$  — **Probability (measure)**<sup>2</sup>

To each event  $E$  we assign the *probability*  $P[E]$  of the event, which is a real number between 0 and 1.

In addition to the three basic objects  $(\Omega, \mathcal{F}, P)$  above, the following two fundamental notions will also be indispensable:

**Random variable**<sup>3</sup>

Random variables are the quantities of interest in our probabilistic model. A *random variable* is a suitable function  $X: \Omega \rightarrow S$ , associating to each possible outcome  $\omega \in \Omega$  a value  $X(\omega) \in S$ . You may think of the *Goddess of Chance* choosing the outcome  $\omega$  at random, and the chosen outcome subsequently determining the value  $X(\omega)$  of any quantity of interest.

**Expected value**<sup>4</sup>

The *expected value*  $E[X]$  of a real-valued quantity of interest  $X$ , i.e., a random variable  $X: \Omega \rightarrow S \subset \mathbb{R}$ , represents an average of the possible values of  $X$  over all randomness, weighted according to probabilities  $P$ . The expected value is an integral with respect to the probability measure  $P$  in the sense of Lebesgue, and we will correspondingly use the

---

<sup>1</sup>We will address the precise axioms required of the collection  $\mathcal{F}$  of events in Lecture I.

<sup>2</sup>We will address the precise axioms required of the probability measure  $P$  in Lecture II.

<sup>3</sup>Random variables will be defined precisely in Lecture III.

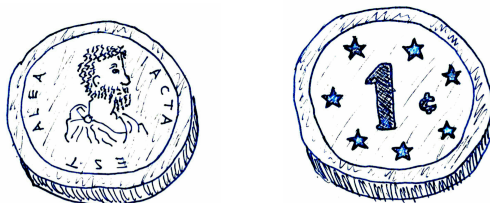
<sup>4</sup>Expected values will be defined precisely in Lecture VII.

following notations interchangeably

$$\mathbb{E}[X] = \int_{\Omega} X(\omega) \, d\mathbb{P}(\omega).$$

The purpose of this course is to make precise mathematical sense of the above notions. Before rushing into the theory, however, we continue with a brief informal introduction. For the informal examples below, it suffices to have an intuitive idea of the above notions.

## O.2. Informal examples of the basic objects in random phenomena



### Example O.1 (One coin toss).

The possible outcomes of a single coin toss are “heads” and “tails”, abbreviated H and T, respectively. The sample space of a single coin toss experiment would thus be

$$\Omega = \{H, T\}.$$

As events, we can in this case allow all subsets of  $\Omega$ , so the collection of events is

$$\mathcal{F} = \left\{ \emptyset, \{H\}, \{T\}, \{H, T\} \right\},$$

with interpretations of the events:

$\{H\}$	the event that the coin toss results in “heads”
$\{T\}$	the event that the coin toss results in “tails”
$\emptyset$	the event that the coin toss results in neither “heads” nor “tails”
$\{H, T\}$	the event that the coin toss results in either “heads” or “tails”

The last two events may appear perplexingly trivial, but we want to allow them as events, because logical reasoning with other events may result in impossibilities or certainties. In fact,  $\emptyset \subset \Omega$  always corresponds to the impossible event, which is not realized by any possible outcome  $\omega \in \Omega$ , whereas  $\Omega \subset \Omega$  always corresponds to the sure event which is realized by any outcome  $\omega \in \Omega$  of the randomness.

A fair coin toss is considered equally likely to result in “heads” or “tails”, and a single fair coin toss is thus unsurprisingly governed by the probability measure  $\mathbb{P}$  which assigns the following probabilities to the above events:

$$\mathbb{P}[\{H\}] = \frac{1}{2}, \quad \mathbb{P}[\{T\}] = \frac{1}{2}, \quad \mathbb{P}[\emptyset] = 0, \quad \mathbb{P}[\{H, T\}] = 1.$$

This example is not overly exciting, but the distinct roles of the three basic objects  $\Omega$ ,  $\mathcal{F}$ , and  $\mathbb{P}$  should be recognized here!

### Example O.2 (Repeated coin tossing).

In an experiment where coin tosses are repeated ad infinitum, the possible outcomes are all possible sequences of “heads” and “tails”, i.e., functions from  $\mathbb{N}$  to  $\{H, T\}$ . The sample space of such a repeated coin tossing experiment would be the space of all such functions

$$\Omega = \{H, T\}^{\mathbb{N}} = \left\{ \omega: \mathbb{N} \rightarrow \{H, T\} \right\},$$



which is uncountably infinite (it can be identified with the set of infinite binary sequences, Example A.16). This uncountable cardinality in a rather innocent probabilistic model can be taken as the first warning that some care is needed in a proper mathematical treatment of probability.

Let  $X_n$  denote the relative frequency of heads in the first  $n$  coin tosses. The relative frequency is a function of the outcome  $\omega$  (as random variables generally are!), given by the ratio

$$X_n(\omega) = \frac{\#\{j \mid j \leq n \text{ and } \omega(j) = \text{H}\}}{n}.$$

We may ask whether the frequency  $X_n$  tends to  $\frac{1}{2}$  in the long run, as  $n \rightarrow \infty$ . This is certainly not true for all  $\omega \in \Omega = \{\text{H}, \text{T}\}^{\mathbb{N}}$ . For example for the sequence  $\omega' = (\text{H}, \text{H}, \text{H}, \text{H}, \dots)$  of all heads, we have  $X_n(\omega') = 1$  for all  $n$  and therefore  $\lim_{n \rightarrow \infty} X_n(\omega') = 1 \neq \frac{1}{2}$ . Even worse, for the sequence

$$\omega'' = (\underbrace{\text{H}, \text{T}, \text{T}, \text{T}}_{3 \text{ times}}, \underbrace{\text{H}, \text{H}, \text{H}, \text{H}, \text{H}, \text{H}, \text{H}, \text{H}, \text{H}, \text{H}, \text{H}}_{3^2 \text{ times}}, \underbrace{\text{T}, \dots, \text{T}}_{3^3 \text{ times}}, \dots)$$

the limit  $\lim_{n \rightarrow \infty} X_n(\omega'')$  does not even exist. In view of these “counterexamples”, it should be clear that any statement about  $X_n$  tending to  $\frac{1}{2}$  in the long run must be probabilistic in nature (somehow referring to  $\mathbb{P}$ ), rather than pointwise on the entire sample space  $\Omega$ .

There are various possible probabilistic notions of  $X_n$  tending to  $\frac{1}{2}$  as  $n \rightarrow \infty$ . Compare for example the following two statements:

(a) For any  $\varepsilon > 0$  we have

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ \left| X_n - \frac{1}{2} \right| > \varepsilon \right] = 0.$$

(b) We have

$$\mathbb{P} \left[ \lim_{n \rightarrow \infty} X_n = \frac{1}{2} \right] = 1.$$

As we develop the mathematical foundations of probability, we should ask, for example:

- Can we make precise mathematical sense of statements such as (a) and (b) above?
- Does one of these two statements imply the other?
- Is either of the statements actually true?  
(say for the usual probability  $\mathbb{P}$  governing fair, repeated coin tossing)

The next example further illustrates the types of questions one may want to answer with the mathematical theory of probability. It concerns a branching population model known as Galton-Watson process.<sup>5</sup> In the footnotes we indicate which parts of the present course are relevant for the questions that arise, but the reader interested in the detailed solutions should look them up in books on stochastic processes.

<sup>5</sup>The idea of using this branching process to illustrate the applicability of probability theory is borrowed from the excellent textbook [Wil91].

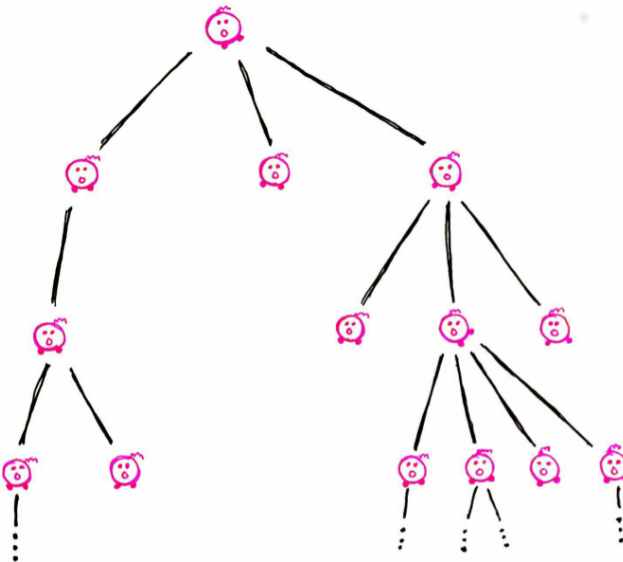
**Example O.3** (Branching process).

Consider a population producing offspring randomly as follows (this population model is known as the Galton-Watson process).

The population is started from a single ancestor, who has a random number of descendants according to some given distribution. These descendants of the ancestor form the first generation in the population, and we denote the random number of individuals in the first generation by  $Z_1$ .

Each of the  $Z_1$  individuals in the first generation then independently of each other and the ancestor has a random number of its own descendants according to the same distribution. These descendants of the descendants of the ancestor form the second generation, and we denote the number of individuals in the second generation by  $Z_2$ .

The process continues branching in the same way, with all individuals in each generation having descendants randomly, which together form the next generation. The numbers of descendants of each individual are assumed to follow the same distribution and to be independent of each other. The random number of individuals in the  $n$ :th generation is denoted by  $Z_n$ .



The first question that this model poses in the foundations of probability theory is:

Is there a well defined mathematical model of this branching process?

In particular, what are the sample space  $\Omega$ , the collection of events  $\mathcal{F}$ , and the probability measure  $P$  describing the process?<sup>6</sup> At least each  $Z_n$  should be a random variable, i.e., a function of the realization  $\omega \in \Omega$  of all involved randomness.

Admitting that the model can be defined, we can start asking more interesting questions about the model itself. One such question concerns the survival of the population, or conversely the extinction of the population:

Does the lineage of the ancestor ever terminate?

The lineage terminates, i.e., the population becomes extinct, if in some generation  $n \in \mathbb{N}$  there are no individuals,  $Z_n = 0$ . The event  $E_n$  that the generation  $n$  contains no individuals consists of those outcomes  $\omega$  of the randomness for which  $Z_n(\omega) = 0$ , i.e.,

$$E_n = \left\{ \omega \in \Omega \mid Z_n(\omega) = 0 \right\}.$$

<sup>6</sup>It turns out that countable product spaces of rather simple discrete probability spaces are sufficient to precisely construct the model, cf. Lecture IX.

We may observe that if the generation  $n$  contains no individuals, then the next generation  $n + 1$  can not contain any either, so

$$Z_n(\omega) = 0 \implies Z_{n+1}(\omega) = 0.$$

Equivalently, for the corresponding events we have the inclusion

$$E_n \subset E_{n+1}.$$

These events thus form an increasing sequence (of subsets of  $\Omega$ )

$$E_1 \subset E_2 \subset E_3 \subset \dots$$

The probabilities of the events should increase correspondingly,

$$\mathbb{P}[E_1] \leq \mathbb{P}[E_2] \leq \mathbb{P}[E_3] \leq \dots$$

and as a bounded increasing sequence of real numbers, these probabilities have a limit

$$\lim_{n \rightarrow \infty} \mathbb{P}[E_n].$$

One natural question is: can we use the events  $E_n$  to construct the event  $E$  that extinction occurs eventually?<sup>7</sup> And is its probability  $\mathbb{P}[E]$  equal to the above limit  $\lim_{n \rightarrow \infty} \mathbb{P}[E_n]$ ?<sup>8</sup> A more ambitious version of the question is: can we concretely calculate the probability  $\mathbb{P}[E]$  of eventual extinction?<sup>9</sup>

The expected size  $\mathbb{E}[Z_n]$  of generation  $n$  turns out to be  $d^n$ , where  $d$  is the expected number of descendants of one individual. Imagine then that we define the “renormalized size” of generation  $n$  as  $R_n = d^{-n} Z_n$ , so as to have expected value one,  $\mathbb{E}[R_n] = 1$ . By techniques that go just a little bit beyond the present course (martingale convergence theorem) one can show that there exists a limit of the random variables  $R_n$  as  $n \rightarrow \infty$ . The limit

$$\lim_{n \rightarrow \infty} R_n$$

is itself a random variable, which can be interpreted to describe the asymptotic long term size of the population in units that make the expected sizes equal to one. Given this, it is natural to wonder whether the expected value of this asymptotic quantity can be calculated by interchanging the limit and the expected value<sup>10</sup>

$$\mathbb{E} \left[ \lim_{n \rightarrow \infty} R_n \right] \stackrel{?}{=} \lim_{n \rightarrow \infty} \underbrace{\mathbb{E} [R_n]}_{=1} = 1.$$

Probability theory provides the means to make sense of and answer the many questions that arise in this branching process example — as well as in other interesting models.

### O.3. Probability theory vs. measure theory

The present course may appear to involve not so much of random phenomena themselves, but more of dry and formal measure theory instead. Our justification for this is that virtually all advanced probability and statistics builds on the measure theoretical foundations covered in the present course. Frequently used measure theoretical

<sup>7</sup>It turns out that  $\sigma$ -algebras, studied in Lecture I, permit just flexible enough logical operations to allow such a construction.

<sup>8</sup>This indeed turns out to be true, by monotonicity properties of probability measures established in Lecture II.

<sup>9</sup>This can be done using generating functions — a close cousin of the characteristic functions studied in Lecture XII.

<sup>10</sup>Whether this can be done turns out to be subtle. In Lectures VII, VIII, and IX we will learn under which conditions one can interchange the order of limits and expected values (and integrals and sums, etc.). Such interchanges of order of operations are tremendously useful in many calculations in practice.

tools in stochastics include Dominated Convergence Theorem, Monotone Convergence Theorem, Fubini's theorem,  $L^p$ -spaces, etc. The formal foundations also serve as a common and well-defined language across different branches of stochastics. For example, various different notions of convergence of random variables used in mathematical statistics are what we will be ready to introduce in the last two of our lectures. The role of the present course is to develop the mathematical foundations of probability mainly for future use!

Since a large number of basic definitions and results in measure theory and probability theory are literally identical, anyone who has already studied measure theory will recognize many familiar notions. The overlap may raise the question: are measure theory and probability theory really separate topics in their own right, and is it necessary to study them separately? It would, in fact, be possible to combine measure theory and probability theory in one extended and coherent course, but the scope of that could easily become daunting. As the topics are currently taught in separate courses, it does not matter much whether one first studies measure theory and later learns about its applicability to probability, or if one proceeds in the opposite order. Even concerning abstract measure theoretic notions and results, probability theory in fact offers very interesting and useful interpretations. In this course, such interpretations include, e.g.,

- product measures interpreted as probabilistic independence
- push-forward measures interpreted as laws of random variables
- (sub-)sigma-algebras interpreted as (partial) information.

At its best, probabilistic thinking leads to entirely new techniques in mathematics, such as

- coupling arguments
- existence proofs relying on random choice.

And finally, probability theory includes inherently stochastic results which do not belong to the domain of measure theory. Some such results covered in the present course are:

- zero-one laws (Borel-Cantelli lemmas, Kolmogorov's 0-1 law)
- laws of large numbers
- central limit theorems.

And despite the similarities, measure theory and probability theory are ultimately concerned with different questions. To highlight just one difference in emphasis, note that the identification of the law of a random variable occupies a much more central place in probability theory than the corresponding question does in analysis. In developments beyond this first theoretical course, it becomes even more apparent that probability theory is not just a subset of measure theory<sup>11</sup> — consider, e.g., martingales, ergodic theory, large deviations, stochastic calculus, optimal stopping, etc. It is also in such further studies, which build on the present foundational course, that the advantages of the theory will become clearer.

---

<sup>11</sup>Vice versa, of course, there are topics covered in courses of measure theory such as [Kin16], which are not covered in this course of probability theory, and there are aspects of the theory into which one gains valuable insights from analysis. Therefore, especially for serious mathematicians, it is highly recommended to study both topics!



We hope that these reassurances and the occasional genuinely probabilistic interpretations and results included in the lectures provide a sufficient motivation to seriously study also the formal (measure theoretical) aspects of probability!



## Lecture I

### Structure of event spaces

#### I.1. Set operations on events

Recall that an event is a subset  $E \subset \Omega$  of the set of all possible outcomes, and the event is said to occur if the outcome  $\omega \in \Omega$  that is realized belongs to this subset,  $\omega \in E$ . We then have the following interpretations of set operations on events:

		interpretation
the whole sample space	$\Omega$	sure event (contains all possible outcomes)
the empty set	$\emptyset$	impossible event (contains no outcomes)
intersection	$E_1 \cap E_2$	“events $E_1$ and $E_2$ both occur”
union	$E_1 \cup E_2$	“event $E_1$ or event $E_2$ occurs”
complement	$E^c = \Omega \setminus E$	“event $E$ does not occur”
subset	$E_1 \subset E_2$	“occurrence of $E_1$ implies $E_2$ ”

In other words, set theoretic operations enable logical reasoning with events. The logical operations **and**, **or**, and **not** are implemented by intersections  $\cap$ , unions  $\cup$ , and complements  $(\dots)^c$ , respectively.<sup>1</sup> We therefore at least want that the collection  $\mathcal{F}$  of events is stable under such operations, i.e.,

- $\Omega \in \mathcal{F}$  and  $\emptyset \in \mathcal{F}$
- if  $E_1, E_2 \in \mathcal{F}$  then  $E_1 \cap E_2 \in \mathcal{F}$  and  $E_1 \cup E_2 \in \mathcal{F}$
- if  $E \in \mathcal{F}$  then  $E^c = \Omega \setminus E \in \mathcal{F}$ .

In fact, for a meaningful mathematical theory, we need to be able to form also countably infinite intersections and unions of events. This is the reason for the following definition.

---

<sup>1</sup>Intersection  $\cap$  is indeed the equivalent of the logical quantifier “for all”,  $\forall$ , and union  $\cup$  is the equivalent of the logical quantifier “for some”, i.e., “there exists”,  $\exists$ .

## I.2. Definition of sigma algebra

**Definition I.1** (Sigma algebra).

A collection  $\mathcal{F} \subset \mathcal{P}(\Omega)$  of subsets of a set  $\Omega$  is a  $\sigma$ -algebra on  $\Omega$  if

$$\begin{aligned} (\Sigma-1) : & \quad \Omega \in \mathcal{F} \\ (\Sigma-c) : & \quad \text{if } E \in \mathcal{F} \text{ then } E^c = \Omega \setminus E \in \mathcal{F} \\ (\Sigma-\cup) : & \quad \text{if } E_1, E_2, \dots \in \mathcal{F} \text{ then } \bigcup_{n \in \mathbb{N}} E_n \in \mathcal{F}. \end{aligned}$$

**Remark I.2** (A sigma algebra is stable under countable set operations).

Note that properties  $(\Sigma-1)$  and  $(\Sigma-c)$  imply that  $\emptyset \in \mathcal{F}$ , since  $\emptyset = \Omega^c$ .

Likewise, properties  $(\Sigma-\cup)$  and  $(\Sigma-c)$  imply that if  $E_1, E_2, \dots \in \mathcal{F}$  then also  $\bigcap_{n=1}^{\infty} E_n \in \mathcal{F}$ , since  $\bigcap_{n=1}^{\infty} E_n = \left(\bigcup_{n=1}^{\infty} E_n^c\right)^c$  by De Morgan's laws, Proposition A.1.

Also, since  $\emptyset \in \mathcal{F}$ , we can extend any finite sequence  $E_1, E_2, \dots, E_k \in \mathcal{F}$  of members of the collection to an infinite sequence by setting  $E_n = \emptyset$  for all  $n > k$ , and we thus deduce from  $(\Sigma-\cup)$  that the finite union  $E_1 \cup \dots \cup E_n \in \mathcal{F}$  also belongs to the collection. Similarly, finite intersections  $E_1 \cap \dots \cap E_n$  of members of the collection remain in the collection.

In view of the definition and remark above,  $\sigma$ -algebras are stable under countable set operations. Since we will always assume the collection of events to be a  $\sigma$ -algebra, we are thus allowed to perform rather flexible logical constructions with events.

**Example I.3** (Examples and counterexamples of sigma algebras).

- (i)  $\mathcal{F} = \{\emptyset, \Omega\}$  is a  $\sigma$ -algebra on  $\Omega$ , albeit not a very interesting one: it only contains the impossible event  $\emptyset$  and the sure event  $\Omega$ .
- (ii)  $\mathcal{F} = \mathcal{P}(\Omega)$ , the collection of all subsets of  $\Omega$ , is a  $\sigma$ -algebra on  $\Omega$ . However, when  $\Omega$  is uncountably infinite, consistent rules of probability can typically only be given on a smaller collection of events.
- (iii)  $\mathcal{F} = \mathcal{T}(\mathbb{R})$ , the collection of all open subsets of  $\mathbb{R}$ , is not a  $\sigma$ -algebra on  $\mathbb{R}$ ! You should recall that arbitrary unions of open sets are open, and finite intersections of open sets are also open. However, for example the countable intersection  $\bigcap_{n=1}^{\infty} \left(-\frac{1}{n}, \frac{1}{n}\right) = \{0\}$  of open intervals consists of a single point and is not open. The collection in fact satisfies  $(\Sigma-1)$  and  $(\Sigma-\cup)$ , but fails to satisfy  $(\Sigma-c)$ .

**Exercise I.1** (Sigma algebras on small finite sets).

Let  $a, b, c$  be three distinct points.

- (a) Write down all  $\sigma$ -algebras on  $\Omega = \{a, b\}$ .
- (b) Write down all  $\sigma$ -algebras on  $\Omega = \{a, b, c\}$ .
- (c) Give an explicit counterexample which shows that the union of two  $\sigma$ -algebras is not necessarily a  $\sigma$ -algebra.

In probability theory, we require the collection of events  $\mathcal{F}$  to be a  $\sigma$ -algebra on the sample space  $\Omega$ . The next examples illustrate what sorts of countable set operations we might encounter in practice. These examples also give a fair idea of the expressive power of such operations, when used iteratively.

**Example I.4** (The event of branching process extinction).

Let us revisit Example O.3 about the branching process. Let  $Z_n$  denote the random size of the population in generation  $n \in \mathbb{N}$ , which, as any random variable, depends on the outcome  $\omega$  of the underlying randomness. Consider first an event defined by the condition

that the generation  $n$  contains no individuals

$$E_n = \left\{ \omega \in \Omega \mid Z_n(\omega) = 0 \right\}.$$

Extinction happens in some generation in the future if there exists some  $n \in \mathbb{N}$  such that  $Z_n = 0$ . The corresponding event is

$$\begin{aligned} E &= \left\{ \omega \in \Omega \mid \exists n \in \mathbb{N} : Z_n(\omega) = 0 \right\} \\ &= \bigcup_{n \in \mathbb{N}} \left\{ \omega \in \Omega \mid Z_n(\omega) = 0 \right\} = \bigcup_{n \in \mathbb{N}} E_n. \end{aligned}$$

We see that the event  $E$  of eventual extinction is the union over the countably many generations  $n \in \mathbb{N}$  of the events  $E_n$  that generation  $n$  is already extinct.

Countable set operations came to our rescue!

**Example I.5** (Long term frequency of heads in coin tossing).

Let us revisit Example O.2 about repeated coin tossing. Let  $X_n$  be the relative frequency of heads in the first  $n$  coin tosses. Observe that the following are logically equivalent ways of expressing the property that the frequency tends to  $\frac{1}{2}$  in the long run:

$$\begin{aligned} \lim_{n \rightarrow \infty} X_n = \frac{1}{2} &\iff \forall \varepsilon > 0 \exists k \in \mathbb{N} \text{ such that } \forall n \geq k \text{ we have } \left| X_n - \frac{1}{2} \right| < \varepsilon \\ &\iff \forall m \in \mathbb{N} \exists k \in \mathbb{N} \text{ such that } \forall n \geq k \text{ we have } \left| X_n - \frac{1}{2} \right| < \frac{1}{m}. \end{aligned}$$

The last expression requires only quantifiers over countable collections, and is therefore good for our purposes. Consider first an event by the condition  $\left| X_n - \frac{1}{2} \right| < \frac{1}{m}$ ,

$$E_n^{(m)} := \left\{ \omega \in \{H, T\}^{\mathbb{N}} \mid \frac{1}{2} - \frac{1}{m} < X_n(\omega) < \frac{1}{2} + \frac{1}{m} \right\},$$

we can now express the event  $E$  that the frequency tends to  $\frac{1}{2}$  in the long run by the following countable set operations:

$$E = \bigcap_{m \in \mathbb{N}} \bigcup_{k \in \mathbb{N}} \bigcap_{n \geq k} E_n^{(m)}.$$

So at least provided that each  $E_n^{(m)}$  belongs to the collection  $\mathcal{F}$  of admissible events (seems reasonable) the properties of  $\sigma$ -algebras allow us to construct the more complicated but much more interesting event  $E \in \mathcal{F}$ , which contains precise information about the long term behavior of the frequencies of heads.

Iteratively constructed countable set operations came to our rescue!

### I.3. Generating sigma algebras

**Definition I.6** (Sigma algebra generated by a collection of subsets).

Let  $\mathcal{C} \subset \mathcal{P}(\Omega)$  be a collection of subsets of  $\Omega$ . Then we define  $\sigma(\mathcal{C})$  as the smallest  $\sigma$ -algebra on  $\Omega$  which contains the collection  $\mathcal{C}$ . We call  $\sigma(\mathcal{C})$  the  $\sigma$ -algebra generated by the collection  $\mathcal{C}$ .

**Remark I.7.** The language of the above definition is intended to be as accessible as possible, but let us make sure that the precise meanings are clear as well:

- We say that a  $\sigma$ -algebra  $\mathcal{F}$  contains the collection  $\mathcal{C}$ , if each member of  $\mathcal{C}$  is a member in  $\mathcal{F}$ , i.e., if we have the inclusion  $\mathcal{C} \subset \mathcal{F}$  of the collections of sets.
- For two  $\sigma$ -algebras  $\mathcal{F}_1$  and  $\mathcal{F}_2$ , we say that  $\mathcal{F}_1$  is smaller than  $\mathcal{F}_2$  if  $\mathcal{F}_1 \subset \mathcal{F}_2$ .

With these clarifications, the meaning should be unambiguous, but one still has to verify that  $\sigma(\mathcal{C})$  becomes well defined. If we want to define  $\sigma(\mathcal{C})$  as the smallest  $\sigma$ -algebra containing  $\mathcal{C}$ ,

we first need to know that such a  $\sigma$ -algebra exists and that it is unique!<sup>2</sup> These concerns will be settled in Corollary I.9 below.

One standard use of generated  $\sigma$ -algebras is the following. If  $\mathcal{C}$  is a collection of some basic events that we want to be able to discuss, in our definition of a probabilistic model we could set  $\mathcal{F} = \sigma(\mathcal{C})$ , which is exactly the smallest possible collection of events that contains the basic events and behaves well under countable operations. Isn't this convenient!

In Lecture IV we discuss the interpretation of  $\sigma$ -algebras as describing information. We will realize that the notion of generated  $\sigma$ -algebras corresponds to what information can be deduced from some initially given pieces of information (the knowledge about events in the generating collection  $\mathcal{C}$ ).

Finally, generated  $\sigma$ -algebras can also be used as a technical tool. It is often very difficult to describe explicitly all members of even very common and reasonable  $\sigma$ -algebras. Working with suitably chosen generating collections can bring about significant simplifications.

Having thus motivated the notion of generated  $\sigma$ -algebras, let us finally address their well-definedness. The key observation is that intersections of  $\sigma$ -algebras are themselves  $\sigma$ -algebras.<sup>3</sup>

**Lemma I.8.** *Suppose that  $(\mathcal{F}_\alpha)_{\alpha \in I}$  is a non-empty collection (indexed by  $I \neq \emptyset$ ) of  $\sigma$ -algebras  $\mathcal{F}_\alpha$  on  $\Omega$ . Then also the intersection  $\mathcal{F} = \bigcap_{\alpha \in I} \mathcal{F}_\alpha$  is a  $\sigma$ -algebra on  $\Omega$ .*

*Proof.* By requiring the collection to be non-empty, we ensured that the intersection is well-defined.

We need to verify that the intersection  $\bigcap_{\alpha \in I} \mathcal{F}_\alpha$  satisfies the three properties in Definition I.1. Note that for a subset  $E \subset \Omega$ , we have  $E \in \mathcal{F} = \bigcap_{\alpha \in I} \mathcal{F}_\alpha$  if and only if  $E \in \mathcal{F}_\alpha$  for all  $\alpha \in I$ .

We have  $\Omega \in \mathcal{F}_\alpha$  for all  $\alpha \in I$ , and therefore  $\Omega \in \mathcal{F}$ . Thus condition  $(\Sigma-1)$  holds for  $\mathcal{F}$ .

Suppose that  $E \in \mathcal{F}$ . Then for all  $\alpha \in I$  we have  $E \in \mathcal{F}_\alpha$ . By property  $(\Sigma-c)$  for the  $\sigma$ -algebra  $\mathcal{F}_\alpha$  we get that  $E^c \in \mathcal{F}_\alpha$ . Since this holds for all  $\alpha$ , we conclude  $E^c \in \mathcal{F}$ . Thus condition  $(\Sigma-c)$  holds for  $\mathcal{F}$ .

Suppose that  $E_1, E_2, \dots \in \mathcal{F}$ . Then for all  $\alpha \in I$  we have  $E_1, E_2, \dots \in \mathcal{F}_\alpha$ . By property  $(\Sigma-\cup)$  for the  $\sigma$ -algebra  $\mathcal{F}_\alpha$  we get that  $\bigcup_{n=1}^{\infty} E_n \in \mathcal{F}_\alpha$ . Since this holds for all  $\alpha$ , we conclude  $\bigcup_{n=1}^{\infty} E_n \in \mathcal{F}$ . Thus condition  $(\Sigma-\cup)$  holds for  $\mathcal{F}$ .  $\square$

We are now ready to conclude that Definition I.1 indeed made sense.

**Corollary I.9** (Well-definedness of the generated sigma algebra).

*Let  $\mathcal{C} \subset \mathcal{P}(\Omega)$  be a collection of subsets of  $\Omega$ . Then the smallest  $\sigma$ -algebra on  $\Omega$  which contains the collection  $\mathcal{C}$  exists and is unique.*

<sup>2</sup>Such issues must be taken seriously. To illustrate the existence issue, imagine trying to define  $s > 0$  as the smallest real number which is strictly positive: no such thing exists, and if we disregard that fact, we will soon run into logical contradictions. To illustrate the uniqueness issue, suppose that  $a, b, c$  are three distinct elements, and imagine trying to define  $S \subset \{a, b, c\}$  as the smallest subset which contains an odd number of elements: any of the three singleton subsets  $\{a\}, \{b\}, \{c\} \subset \{a, b, c\}$  are equally small, so which one should  $S$  be?

<sup>3</sup>By contrast, in Exercise I.1(c) you showed that the union of  $\sigma$ -algebras may fail to be a  $\sigma$ -algebra.

*Proof.* The uniqueness part is usual abstract nonsense. Suppose we had two different smallest  $\sigma$ -algebras which contain the collection  $\mathcal{C}$ , say  $\mathcal{F}_1$  and  $\mathcal{F}_2$ . Then we would have  $\mathcal{F}_1 \subset \mathcal{F}_2$  because  $\mathcal{F}_1$  is smallest and  $\mathcal{F}_2 \subset \mathcal{F}_1$  because  $\mathcal{F}_2$  is smallest, so we get that  $\mathcal{F}_1 = \mathcal{F}_2$ .

It remains to show that a smallest  $\sigma$ -algebra which contains the collection  $\mathcal{C}$  exists. Consider the collection  $\mathcal{S}$  of all  $\sigma$ -algebras  $\mathcal{F}$  on  $\Omega$  such that  $\mathcal{C} \subset \mathcal{F}$ . This collection  $\mathcal{S}$  is non-empty, because as in Example I.3(ii), the power set of  $\Omega$  is such a  $\sigma$ -algebra, i.e.,  $\mathcal{P}(\Omega) \in \mathcal{S}$ . Let us define  $\sigma(\mathcal{C})$  as the intersection

$$\sigma(\mathcal{C}) := \bigcap_{\mathcal{F} \in \mathcal{S}} \mathcal{F}$$

of all these  $\sigma$ -algebras. By Lemma I.8,  $\sigma(\mathcal{C})$  is itself a  $\sigma$ -algebra. Since for each  $\mathcal{F} \in \mathcal{S}$  we have  $\mathcal{C} \subset \mathcal{F}$ , the intersection also has this property,  $\mathcal{C} \subset \sigma(\mathcal{C})$ . If  $\mathcal{F}$  is any  $\sigma$ -algebra which contains  $\mathcal{C}$ , then clearly  $\sigma(\mathcal{C})$  is smaller, since  $\mathcal{F} \in \mathcal{S}$  appears in the intersection and thus  $\sigma(\mathcal{C}) \subset \mathcal{F}$ . Thus the intersection  $\sigma(\mathcal{C})$  is smallest.  $\square$

### I.3.1. Borel sigma algebra

**Definition I.10** (Borel sigma algebra).

For a topological space  $\mathfrak{X}$ , the *Borel  $\sigma$ -algebra* on  $\mathfrak{X}$  is the  $\sigma$ -algebra  $\mathcal{B}(\mathfrak{X})$  generated by the collection  $\mathcal{T}(\mathfrak{X})$  of open sets in  $\mathfrak{X}$ .

Arguably the most important  $\sigma$ -algebra in all of probability theory is the Borel  $\sigma$ -algebra on the real line  $\mathbb{R}$ , because it is needed whenever we consider real valued random variables. We denote simply  $\mathcal{B} = \mathcal{B}(\mathbb{R})$ . By definition,  $\mathcal{B}$  is the smallest  $\sigma$ -algebra on  $\mathbb{R}$  which contains all open sets  $V \subset \mathbb{R}$ . The following proposition establishes that  $\mathcal{B}$  can alternatively be generated by various convenient collections of subsets of the real line.

**Proposition I.11** (Generating the Borel sigma algebra on the real line).

*The Borel  $\sigma$ -algebra  $\mathcal{B}$  on the real line  $\mathbb{R}$  is generated by any of the following collections of subsets of  $\mathbb{R}$ :*

$$\begin{aligned} \text{(i)} : \quad \mathcal{C} &= \left\{ (-\infty, x] \mid x \in \mathbb{R} \right\}, & \text{(iii)} : \quad \mathcal{C} &= \left\{ (x, y) \mid x, y \in \mathbb{R}, x < y \right\}, \\ \text{(ii)} : \quad \mathcal{C} &= \left\{ [x, y] \mid x, y \in \mathbb{R}, x \leq y \right\}, & \text{(iv)} : \quad \mathcal{C} &= \left\{ (x, y] \mid x, y \in \mathbb{R}, x < y \right\}. \end{aligned}$$

**Remark I.12.** The reader can certainly imagine further variations of generating collections of intervals, and is invited to think about the modifications needed in the proof below.

*Proof of Proposition I.11.* We will only explicitly check that the collection (i) generates  $\mathcal{B}$ , the other cases are similar.

So for the case (i), let  $\mathcal{C}$  be the collection of all intervals of the form  $(-\infty, x]$ , with  $x \in \mathbb{R}$ . In order to show that  $\sigma(\mathcal{C}) = \mathcal{B}$ , we will separately check the two converse inclusions  $\sigma(\mathcal{C}) \subset \mathcal{B}$  and  $\sigma(\mathcal{C}) \supset \mathcal{B}$ .

*inclusion  $\sigma(\mathcal{C}) \subset \mathcal{B}$ :* To show that  $\sigma(\mathcal{C}) \subset \mathcal{B}$ , it is sufficient to show that  $\mathcal{B}$  contains all intervals of the form  $(-\infty, x]$ , because  $\sigma(\mathcal{C})$  is by definition the smallest such  $\sigma$ -algebra.

Note that the set  $(x, +\infty)$  is open, and thus is contained in the Borel  $\sigma$ -algebra by definition. The complement of it is  $((x, +\infty))^c = \mathbb{R} \setminus (x, +\infty) = (-\infty, x]$ . Since  $\mathcal{B}$  is a  $\sigma$ -algebra on  $\mathbb{R}$  which contains  $(x, +\infty)$ , by property  $(\Sigma\text{-c})$  it contains also the complement  $(-\infty, x]$ . The inclusion  $\sigma(\mathcal{C}) \subset \mathcal{B}$  follows.

*inclusion*  $\sigma(\mathcal{C}) \supset \mathcal{B}$ : To show that  $\sigma(\mathcal{C}) \supset \mathcal{B}$ , it is sufficient to show that  $\sigma(\mathcal{C})$  contains all open sets  $V \subset \mathbb{R}$ , because  $\mathcal{B}$  is by definition the smallest such  $\sigma$ -algebra. Let us show step by step that  $\sigma(\mathcal{C})$  contains all sets of the forms

- (a) semi-open intervals  $(x, y]$ , for  $x, y \in \mathbb{R}$
- (b) open intervals  $(x, z)$ , for  $x, z \in \mathbb{R}$
- (c) open sets  $V \subset \mathbb{R}$ .

For case (a), note that

$$(x, y] = (-\infty, y] \setminus (-\infty, x] = (-\infty, y] \cap (-\infty, x]^c,$$

so  $(x, y]$  is obtained from members of the collection  $\mathcal{C}$  by countable (in fact finite) intersections and complements. Therefore we have  $(x, y] \in \sigma(\mathcal{C})$ .

For case (b), note that

$$(x, z) = \bigcup_{n=1}^{\infty} \left(x, z - \frac{1}{n}\right],$$

so the open interval  $(x, z)$  is obtained from intervals of type (a) by a countable union. Since intervals of type (a) are already known to belong to  $\sigma(\mathcal{C})$ , we get that  $(x, z) \in \sigma(\mathcal{C})$ .

Finally, for case (c), note that any open set  $V \subset \mathbb{R}$  is a countable union of open intervals — see Proposition B.5. Since open intervals are already known to belong to  $\sigma(\mathcal{C})$  by case (b), we also get  $V \in \sigma(\mathcal{C})$ . This concludes the proof.  $\square$

The Borel  $\sigma$ -algebra on  $\mathbb{R}$  will be needed in particular for real valued random variables. Likewise, for vector valued random variables, we will need the Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R}^d)$  on the vector spaces  $\mathbb{R}^d$ . Recall that by definition  $\mathcal{B}(\mathbb{R}^d)$  is the smallest  $\sigma$ -algebra on  $\mathbb{R}^d$  which contains all open sets  $V \subset \mathbb{R}^d$  of the  $d$ -dimensional Euclidean space. The next exercise gives a concrete and useful generating collection for this  $\sigma$ -algebra, in the case  $d = 2$ .

**Exercise I.2** (Borel  $\sigma$ -algebra on the two-dimensional Euclidean space).

Denote by

$$\mathcal{C} = \left\{ (-\infty, x] \times (-\infty, y] \mid x \in \mathbb{R}, y \in \mathbb{R} \right\}.$$

the collection of closed south-west quadrants in  $\mathbb{R}^2$ . Prove that the collection  $\mathcal{C}$  of closed south-west quadrants generates  $\mathcal{B}(\mathbb{R}^2)$ , that is, show that  $\mathcal{B}(\mathbb{R}^2) = \sigma(\mathcal{C})$ .

**Hint:** Compare with the proof of Proposition I.11. You may use the fact that every open set in  $\mathbb{R}^2$  can be written as a countable union  $\bigcup_{n=1}^{\infty} R_n$  of open rectangles of the form  $R_n = (a_n, b_n) \times (a'_n, b'_n)$ .



## Lecture II

### Measures and probability measures

Recall that the basic objects of probability theory are:

$\Omega$  — the set of all possible outcomes (sample space)

$\mathcal{F}$  — the collection of all events

$\mathbf{P}$  — the probability (measure).

The sample space  $\Omega$  can be any (non-empty) set, which we in our probabilistic modelling deem representative of the possible outcomes of the randomness involved.

In the previous lecture we explained why the collection  $\mathcal{F}$  of events should be stable under countable set operations, i.e., why it must be a  $\sigma$ -algebra on  $\Omega$ .

In this lecture we examine the last remaining basic object,  $\mathbf{P}$ , the probability itself. We give the axiomatic properties that  $\mathbf{P}$  is required to satisfy, and we begin studying the consequences. By the axioms,  $\mathbf{P}$  is a special case of a mathematical object called a measure, so there is a large amount of overlap between probability theory and measure theory. We in fact choose to first develop measure theory in the general setup up to some point, because even in stochastics we make use of also other measures besides just probability measures.

#### II.1. Measurable spaces

In the previous lecture, we emphasized the importance of being able to perform countable set operations. This merits a definition in its own right.

**Definition II.1** (Measurable space).

If  $S$  is a set and  $\mathcal{S}$  is a  $\sigma$ -algebra on  $S$ , then we call the pair  $(S, \mathcal{S})$  a *measurable space*. A subset  $A \subset S$  is called *measurable* if  $A \in \mathcal{S}$ .

Think of measurable spaces as spaces which are ready to accommodate measures. They come equipped with a good collection  $\mathcal{S}$  of subsets, which behaves well under set operations as discussed in Lecture I, and a measure will assign to each of these good subsets a numerical value appropriately quantifying the size of the subset.

For convenience, let us once more unravel the definition and summarize what a measurable space is:

- $S$  is a set
- $\mathcal{S} \subset \mathcal{P}(S)$  is a collection of subsets which satisfies
  - ( $\Sigma$ -1):  $S \in \mathcal{S}$
  - ( $\Sigma$ -c): if  $A \in \mathcal{S}$  then  $A^c = S \setminus A \in \mathcal{S}$ .
  - ( $\Sigma$ - $\cup$ ): if  $A_1, A_2, \dots \in \mathcal{S}$  then  $\bigcup_{n=1}^{\infty} A_n \in \mathcal{S}$ .

Let us then give some examples of measurable spaces.

The following simple example is primarily relevant when  $S$  is finite or countably infinite.

**Example II.2** (Measurable spaces where all subsets are measurable).

If  $S$  is any set and  $\mathcal{P}(S)$  is the collection of all subsets of  $S$ , then by Example I.3(ii),  $\mathcal{P}(S)$  is a  $\sigma$ -algebra on  $S$ . Thus the pair  $(S, \mathcal{P}(S))$  is a measurable space.

The following example is extremely important: in integration theory of real valued functions (or real valued random variables) one needs the set of real numbers  $\mathbb{R}$  to carry the structure of a measurable space.

**Example II.3** (Real line as a measurable space).

Consider  $S = \mathbb{R}$  and let  $\mathcal{S} = \mathcal{B}$  be the Borel  $\sigma$ -algebra on  $\mathbb{R}$  as in Section I.3.1. Then the pair  $(\mathbb{R}, \mathcal{B})$  is a measurable space.

The measurable space of Example II.3 will in particular accommodate the usual measure of length on the real line  $\mathbb{R}$  (cf. Example II.12 below).

Finally, the class of examples most relevant for probability theory: any pair  $(\Omega, \mathcal{F})$  of a sample space  $\Omega$  together with the collection of events  $\mathcal{F}$  on it has to be a measurable space — ready to accommodate a probability measure in the next step!

## II.2. Definition of measures and probability measures

The final basic object of probability theory is the probability measure  $\mathbb{P}$ . In fact, even in probability theory we actually very often use also measures which are not necessarily probability measures. For example, counting measures are used when handling infinite sums, and the (intuitively) familiar measures on  $\mathbb{R}$  and  $\mathbb{R}^d$  are used as references when talking about densities of real valued or vector valued random variables, respectively.

Let us therefore first define measures in general.

**Definition II.4** (Measure).

Let  $(S, \mathcal{S})$  be a measurable space. A *measure*  $\mu$  on  $(S, \mathcal{S})$  is a function

$$\mu: \mathcal{S} \rightarrow [0, +\infty]$$

such that

$$\mu[\emptyset] = 0 \tag{M-0}$$

and if  $A_1, A_2, \dots \in \mathcal{S}$  are disjoint, then

$$\mu\left[\bigcup_{n=1}^{\infty} A_n\right] = \sum_{n=1}^{\infty} \mu[A_n]. \tag{M-U}$$

A probability measure has just one further requirement added: that the total probability must be equal to one.

**Definition II.5** (Probability measure).

Let  $(\Omega, \mathcal{F})$  be a measurable space. A *probability measure*  $P$  on  $(\Omega, \mathcal{F})$  is a measure on  $(\Omega, \mathcal{F})$  such that  $P[\Omega] = 1$ .

**Remark II.6** (Sure vs. almost sure).

The event  $\Omega$  is the *sure event*: it contains all possible outcomes. The additional requirement in the above definition merely says that the probability of the sure event is one:  $P[\Omega] = 1$ .

It is worth noting that there may be also other events  $E \in \mathcal{F}$ ,  $E \subsetneq \Omega$ , which have probability one,  $P[E] = 1$ . We say that such an event  $E$  is *almost sure*, or alternatively we say that the event  $E$  occurs *almost surely*. The notion of *sure* only depends on the sample space  $\Omega$  itself, but the notion of *almost sure* depends on the probability measure  $P$  as well, so occasionally it is appropriate to use the more specific terminology *P-almost sure* and *P-almost surely*.

**Remark II.7** (Abuse of terminology).

Often the  $\sigma$ -algebra of the underlying measurable space is clear from the context. Then, rather than saying that  $\mu$  is a measure on  $(S, \mathcal{S})$ , we simply say that  $\mu$  is a measure on  $S$ . Likewise, rather than saying that  $P$  is a probability measure on  $(\Omega, \mathcal{F})$ , we simply say that  $P$  is a probability measure on  $\Omega$ .

If  $\mu$  is a measure on  $(S, \mathcal{S})$ , then we call the triple  $(S, \mathcal{S}, \mu)$  a *measure space*. Likewise, if  $P$  is a probability measure on  $(\Omega, \mathcal{F})$ , then we call the triple  $(\Omega, \mathcal{F}, P)$  a *probability space*. In particular, any probability space is a measure space, and all results for measure spaces can be used for probability spaces. For this reason, especially in Section II.3 below, we content ourselves to stating basic properties only for measure spaces in general. There are some results which are valid for probability measures, but not for general measures. Many such results would actually hold under a milder assumption, described below.

**Definition II.8** (Total mass).

Let  $\mu$  be a measure on  $(S, \mathcal{S})$ . The value  $\mu[S] \in [0, +\infty]$  that the measure assigns to the whole space  $S$  is called the *total mass* of  $\mu$ .

**Definition II.9** (Finite measure).

We say that the measure  $\mu$  on  $(S, \mathcal{S})$  is *finite* if its total mass is finite,  $\mu[S] < +\infty$ . We then also say that the corresponding measure space  $(S, \mathcal{S}, \mu)$  is finite.

Probability measures, in particular, are finite measures (since  $P[\Omega] = 1 < +\infty$ ).

Let us now give a few examples of measures and probability measures.

**Example II.10** (Counting measure).

Let  $S$  be any set. Equip  $S$  with the  $\sigma$ -algebra  $\mathcal{P}(S)$  consisting of all subsets of  $S$ . Then the *counting measure* on  $S$  is the measure  $\mu_{\#}$ , which associates to any subset  $A \subset S$  the number of elements  $\#A$  in the subset,

$$\mu_{\#}[A] = \#A.$$

In particular, for all infinite subsets  $A \subset S$  we have  $\mu_{\#}[A] = +\infty$ . Property (M- $\emptyset$ ) holds for  $\mu_{\#}$ , since the empty set has no elements. Property (M- $\cup$ ) holds since the number of elements in a disjoint union of sets is obtained by adding up the numbers of elements in each set.

The counting measure  $\mu_{\#}$  is a finite measure if and only if the underlying set  $S$  is a finite set.

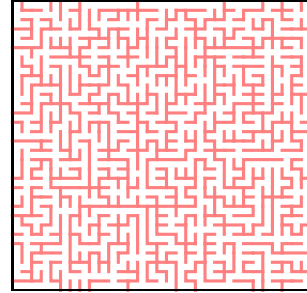
**Example II.11** (Discrete uniform probability measure).

Let  $\Omega$  be any finite non-empty set. Equip it with the  $\sigma$ -algebra  $\mathcal{P}(\Omega)$  consisting of all subsets of  $\Omega$ . Then the (discrete) uniform probability measure on  $\Omega$  is the measure  $\mathbb{P}_{\text{unif}}$  given by

$$\mathbb{P}_{\text{unif}}[E] = \frac{\#E}{\#\Omega} \quad \text{for all } E \subset \Omega.$$

In other words, the (discrete) uniform probability measure is just the counting measure normalized to have total mass one,  $\mathbb{P}_{\text{unif}} = \frac{1}{\#\Omega} \mu_{\#}$ .

The figure on the right shows a uniform random sample from the finite set of all  $30 \times 30$  labyrinths, illustrating that despite its simplicity, the discrete uniform probability measure can give rise to intricate behavior.

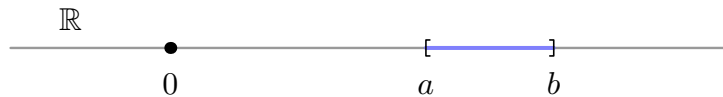


**Example II.12** (Lebesgue measure on the real line).

The natural notion of “length” on the real line  $\mathbb{R}$  corresponds to the *Lebesgue measure*  $\Lambda$  on  $(\mathbb{R}, \mathcal{B})$ . For instance a closed interval  $[a, b] \subset \mathbb{R}$ , with  $a \leq b$ , has measure

$$\Lambda[[a, b]] = b - a$$

equal to the length of the interval, and this property in fact is sufficient to characterize the measure  $\Lambda$ . The length of the entire real axis, on the other hand, is infinite:  $\Lambda[\mathbb{R}] = +\infty$ .



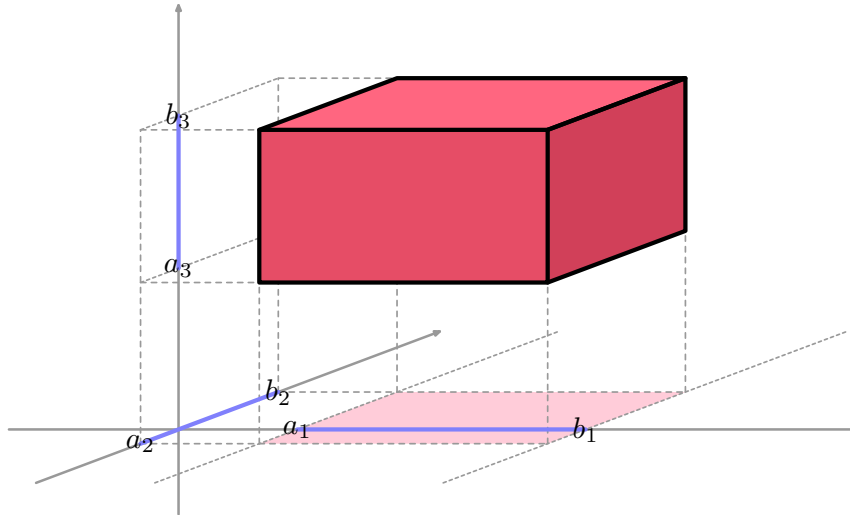
**Example II.13** (Higher dimensional Lebesgue measure).

Example II.12 on the one-dimensional space  $\mathbb{R}$  has a  $d$ -dimensional generalization — a measure on the Euclidean space  $\mathbb{R}^d$ . The cases  $d = 1$ ,  $d = 2$ , and  $d = 3$  have the interpretation of “length on the line  $\mathbb{R}$ ”, “area in the plane  $\mathbb{R}^2$ ”, and “volume in the space  $\mathbb{R}^3$ ”, respectively.

The Euclidean space  $\mathbb{R}^d$  is equipped with the Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R}^d)$  (see Definition I.10). The  $d$ -dimensional Lebesgue measure  $\Lambda^d$  is a measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ , which is characterized by the property that any rectangular box  $[a_1, b_1] \times \cdots \times [a_d, b_d] \subset \mathbb{R}^d$  has measure

$$\Lambda^d[[a_1, b_1] \times \cdots \times [a_d, b_d]] = \prod_{j=1}^d (b_j - a_j)$$

given by the product of the side lengths of the box.



Any course in traditional measure theory covers the particular example of Lebesgue measures  $\Lambda$  and  $\Lambda^d$  in great detail, so we choose not to elaborate on them too extensively here.

**Exercise II.1** (Truncation of measures and conditioning of probability measures).

- (a) Let  $\mu$  be a measure on  $(S, \mathcal{S})$  and let  $B \in \mathcal{S}$ . Show that also  $A \mapsto \mu[A \cap B]$  defines a measure on  $(S, \mathcal{S})$ .
- (b) Let  $P$  be a probability measure on  $(\Omega, \mathcal{F})$ , and let  $B \in \mathcal{F}$  be an event such that  $P[B] > 0$ . Show that the conditional probability  $A \mapsto P[A | B] := \frac{P[A \cap B]}{P[B]}$  is a probability measure on  $(\Omega, \mathcal{F})$ .

**Example II.14** (Uniform probability measure on an interval).

Consider the truncation of the Lebesgue measure to the unit interval  $[0, 1] \subset \mathbb{R}$ . Let  $\Lambda$  be the Lebesgue measure on  $\mathbb{R}$  as in Example II.12. Define a new measure  $P$  on  $\mathbb{R}$  by truncation as in part (a) of Exercise II.1:

$$P[A] = \Lambda[A \cap [0, 1]] \quad \text{for all } A \in \mathcal{B}.$$

Then we have  $P[\mathbb{R}] = \Lambda[[0, 1]] = 1$ , so  $P$  is a probability measure on  $\mathbb{R}$ . For subsets  $A \subset [0, 1]$  of the unit interval,  $P$  coincides with the Lebesgue measure,  $P[A] = \Lambda[A]$ . For subsets outside the unit interval, on the other hand, we have  $A \cap [0, 1] = \emptyset$  and thus these sets carry no probability mass:  $P[A] = \Lambda[\emptyset] = 0$ . We call  $P$  the *uniform probability measure* on the unit interval.

More generally, the uniform probability measure on any interval  $[a, b] \subset \mathbb{R}$  of positive length is defined by the formula  $A \mapsto \frac{1}{b-a} \Lambda[A \cap [a, b]]$ , where we first truncate the Lebesgue measure to  $[a, b]$  and then normalize the total mass by the length  $b - a$  of the interval.

Before starting to examine general results about measures, we still look into one further class of examples of probability measures which is relatively easy, yet important in practice.

### II.2.1. Probability distributions on countable spaces

Many probabilistic models concern distributions on the natural numbers  $\mathbb{N}$ , the integers  $\mathbb{Z}$ , or other countable sets. It turns out that on such countable spaces,

we can characterize probability measures in an intuitive way using probability mass functions.

For the rest of this section we therefore assume that  $\Omega$  is a non-empty countable set. Then there exists an enumeration<sup>1</sup>  $\Omega = \{\omega_1, \omega_2, \omega_3, \dots\}$  with distinct elements  $\omega_1, \omega_2, \omega_3, \dots \in \Omega$ . Summation over  $\Omega$  can be defined using the enumeration,

$$\sum_{\omega \in \Omega} a(\omega) := \sum_j a(\omega_j),$$

and if the terms of the sum are non-negative,  $a(\omega) \geq 0$ , then the result of the sum is independent of the chosen enumeration.

**Definition II.15** (Probability mass function).

A *probability mass function* (p.m.f.) on  $\Omega$  is a function

$$p: \Omega \rightarrow [0, 1]$$

such that

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

To a probability mass function  $p$  it is natural to associate a measure defined by

$$\mathbf{P}[E] = \sum_{\omega \in E} p(\omega) \quad \text{for all } E \subset \Omega, \quad (\text{II.1})$$

and conversely to a probability measure  $\mathbf{P}$  on  $(\Omega, \mathcal{P}(\Omega))$  it is natural to associate masses of singleton events  $\{\omega\} \subset \Omega$

$$p(\omega) = \mathbf{P}[\{\omega\}] \quad \text{for all } \omega \in \Omega. \quad (\text{II.2})$$

The following exercise shows that on countable spaces, probability mass functions are in one-to-one correspondence with probability measures via the above formulas.

**Exercise II.2** (Probability distributions on countable spaces).

Let  $\Omega$  be a finite or a countably infinite set, and denote by  $\mathcal{P}(\Omega)$  the collection of all subsets of  $\Omega$ .

- (a) Show that if  $p$  is a probability mass function on  $\Omega$ , then the set function  $\mathbf{P}$  defined by (II.1) is a probability measure on  $(\Omega, \mathcal{P}(\Omega))$ .
- (b) Show that if  $\mathbf{P}$  is a probability measure on  $(\Omega, \mathcal{P}(\Omega))$ , then the function  $p$  defined by (II.2) is a probability mass function on  $\Omega$ .

**Example II.16** (Poisson distribution).

Let  $\lambda > 0$ . Recalling the power series  $\sum_{k=0}^{\infty} \frac{1}{k!} \lambda^k = e^\lambda$  of the exponential function, it is easy to see that the function  $p$  given by

$$p(k) = e^{-\lambda} \frac{\lambda^k}{k!} \quad \text{for } k \in \mathbb{Z}_{\geq 0} = \{0, 1, 2, \dots\} \quad (\text{II.3})$$

is a probability mass function on  $\mathbb{Z}_{\geq 0}$ .

The Poisson distribution with parameter  $\lambda$  is the probability measure on  $\mathbb{Z}_{\geq 0}$  with the above probability mass function.

---

<sup>1</sup>If  $\Omega$  is finite, the enumeration terminates,  $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$ . The more interesting case is if  $\Omega$  is countably infinite.

**Example II.17** (Geometric distribution).

Let  $q \in (0, 1)$ . Using the geometric series  $\sum_{k=0}^{\infty} r^k = \frac{1}{1-r}$  with  $r = 1 - q$ , it is easy to see that the function  $p$  given by

$$p(k) = (1 - q)^{k-1}q \quad \text{for } k \in \mathbb{N} = \{1, 2, 3, \dots\} \quad (\text{II.4})$$

is a probability mass function on  $\mathbb{N}$ .

The geometric distribution with parameter  $q$  is the probability measure on  $\mathbb{N}$  with the above probability mass function.

**Example II.18** (Binomial distribution).

Let  $n \in \mathbb{N}$  and  $q \in (0, 1)$ . Using the binomial formula  $\sum_{k=0}^n \binom{n}{k} a^k b^{n-k} = (a + b)^n$  with  $a = q$  and  $b = 1 - q$  it is easy to see that the function  $p$  given by

$$p(k) = \binom{n}{k} q^k (1 - q)^{n-k} \quad \text{for } k \in \{0, 1, \dots, n - 1, n\} \quad (\text{II.5})$$

is a probability mass function on the finite set  $\{0, 1, \dots, n - 1, n\}$ .

The binomial distribution with parameters  $n$  and  $q$  is the probability measure on the finite set  $\{0, 1, \dots, n - 1, n\}$  with the above probability mass function.

### II.3. Properties of measures and probability measures

Let us now discuss some of the first properties of measures and probability measures.

#### Subadditivity of measures and the union bound

For repeated later use, we start by proving the following additivity properties for measures of disjoint sets, subadditivity properties of measures of (not necessarily disjoint) sets, as well as related monotonicity and monotone convergence properties of measures.

**Lemma II.19** (First properties of measures).

Let  $\mu$  be a measure on a measurable space  $(S, \mathcal{S})$ . Then we have the following:

- (a) Finite additivity: If  $A_1, \dots, A_n \in \mathcal{S}$  are disjoint measurable sets, then we have:

$$\mu[A_1 \cup \dots \cup A_n] = \mu[A_1] + \dots + \mu[A_n]. \quad (\text{II.6})$$

- (b) Monotonicity: If  $A, B \in \mathcal{S}$  and  $A \subset B$ , then we have:

$$\mu[A] \leq \mu[B]. \quad (\text{II.7})$$

- (c) Finite subadditivity: If  $A_1, \dots, A_n \in \mathcal{S}$  are any measurable sets, then we have:

$$\mu[A_1 \cup \dots \cup A_n] \leq \mu[A_1] + \dots + \mu[A_n]. \quad (\text{II.8})$$

- (d) Monotone convergence of measures: Let  $A_1 \subset A_2 \subset \dots$  be an increasing sequence of measurable sets,  $A_n \in \mathcal{S}$  for all  $n \in \mathbb{N}$ . Then the measures of the increasing limit  $A_n \uparrow A = \bigcup_{j=1}^{\infty} A_j$  of sets constitute the increasing limit

$$\mu[A_n] \uparrow \mu[A]. \quad (\text{II.9})$$

(e) Countable subadditivity: If  $A_1, A_2, \dots \in \mathcal{S}$  is a sequence of measurable sets (not necessarily disjoint), then we have:

$$\mu\left[\bigcup_{j=1}^{\infty} A_j\right] \leq \sum_{j=1}^{\infty} \mu[A_j]. \quad (\text{II.10})$$

*Proof of (a):* Given the disjoint measurable sets  $A_1, \dots, A_n \in \mathcal{S}$ , let us extend this finite sequence by empty sets: define  $A_{n+1} = A_{n+2} = \dots = \emptyset$ . Since  $\emptyset \in \mathcal{S}$  by properties of  $\sigma$ -algebras, we thus obtain an infinite sequence  $A_1, A_2, \dots \in \mathcal{S}$  of measurable sets. This sequence of sets is disjoint (the newly added empty sets do not have common elements with the already disjoint  $A_1, \dots, A_n$ ). Thus from axiom (M- $\cup$ ) it follows that

$$\mu\left[\bigcup_{j=1}^{\infty} A_j\right] = \sum_{j=1}^{\infty} \mu[A_j].$$

But on the left hand side, the union is simply  $\bigcup_{j=1}^{\infty} A_j = A_1 \cup \dots \cup A_n$ , because the empty sets do not contribute to the union. On the right hand side, the sum is  $\sum_{j=1}^{\infty} \mu[A_j] = \mu[A_1] + \dots + \mu[A_n]$ , because the measures of the empty sets  $\mu[\emptyset] = 0$  do not contribute to the sum. Assertion (a) follows.

*Proof of (b):* Assume that  $A, B \in \mathcal{S}$ . Note that then also  $B \setminus A = B \cap A^c \in \mathcal{S}$  by properties of  $\sigma$ -algebras. If  $A \subset B$ , then  $B = A \cup (B \setminus A)$  is a disjoint union. For these two disjoint sets, we can use part (a) to get

$$\mu[B] = \mu[A \cup (B \setminus A)] \stackrel{(a)}{=} \mu[A] + \mu[B \setminus A].$$

Since  $\mu[B \setminus A] \geq 0$  by properties of measures, the assertion  $\mu[B] \geq \mu[A]$  follows.

*Proof of (c):* We will prove the inequality

$$\mu[A_1 \cup \dots \cup A_n] \leq \mu[A_1] + \dots + \mu[A_n].$$

for all  $A_1, \dots, A_n \in \mathcal{S}$  by induction on the number  $n$  of sets in the union.

The case  $n = 1$  is clear — the two sides of the inequality are in fact equal. Now assume the inequality for unions of  $n$  sets, and consider  $A_1, \dots, A_{n+1} \in \mathcal{S}$ . Define  $A = A_1 \cup \dots \cup A_n$  and  $B = A_{n+1} \setminus A$ . Then we have  $A_1 \cup \dots \cup A_{n+1} = A \cup B$ , where the sets  $A$  and  $B$  are disjoint. Thus by part (a) we get

$$\mu[A \cup B] = \mu[A] + \mu[B].$$

The first term on the right hand side is

$$\mu[A] = \mu[A_1 \cup \dots \cup A_n] \leq \mu[A_1] + \dots + \mu[A_n]$$

by induction assumption. The second term on the right hand side is  $\mu[B] \leq \mu[A_{n+1}]$  by monotonicity proven in part (b), since  $B \subset A_{n+1}$ . Notice that the right hand side  $\mu[A \cup B] = \mu[A_1 \cup \dots \cup A_{n+1}]$  is the measure we are interested in. Therefore, by combining the observations, we conclude

$$\mu[A_1 \cup \dots \cup A_{n+1}] \leq \left(\mu[A_1] + \dots + \mu[A_n]\right) + \mu[A_{n+1}],$$

which finishes the proof of assertion (c) by induction.

*Proof of (d):* Suppose that  $A_1 \subset A_2 \subset \dots$  is an increasing sequence of measurable sets, and denote its limit by  $A = \bigcup_{j=1}^{\infty} A_j$ . Then  $A$  is also measurable by properties of  $\sigma$ -algebras. Now write first  $B_1 = A_1$ , and then  $B_2 = A_2 \setminus A_1, \dots, B_n = A_n \setminus A_{n-1}, \dots$ . These sets  $B_1, B_2, \dots$  are disjoint and  $A_n = B_1 \cup \dots \cup B_n$  for all  $n \in \mathbb{N}$ . From part (a) we get

$$\mu[A_n] = \mu[B_1 \cup \dots \cup B_n] = \sum_{j=1}^n \mu[B_j].$$



The right hand sides are the partial sums of an infinite sum with non-negative terms, so they form a sequence increasing to that infinite sum, and we conclude

$$\mu[A_n] \uparrow \sum_{j=1}^{\infty} \mu[B_j] \quad \text{as } n \rightarrow \infty.$$

On the other hand, by disjointness of  $B_1, B_2, \dots$  and axiom (M- $\cup$ ), this infinite sum equals

$$\sum_{j=1}^{\infty} \mu[B_j] = \mu\left[\bigcup_{j=1}^{\infty} B_j\right] = \mu\left[\bigcup_{j=1}^{\infty} A_j\right] = \mu[A].$$

This proves the assertion (d),  $\mu[A_n] \uparrow \mu[A]$  as  $n \rightarrow \infty$ .

*Proof of (e):* Let  $A_1, A_2, \dots \in \mathcal{S}$  be a sequence of measurable sets. Form their finite unions  $C_n = A_1 \cup \dots \cup A_n$ , for all  $n \in \mathbb{N}$ . Then, as a union of measurable sets, each  $C_n$  is also measurable. This sequence is clearly increasing  $C_1 \subset C_2 \subset \dots$ , and its limit is the countably infinite union  $C = \bigcup_{j=1}^{\infty} A_j$ . We can therefore apply part (d) to get

$$\mu[C_n] \uparrow \mu[C] \quad \text{as } n \rightarrow \infty. \tag{II.11}$$

On the other hand, by part (c) we have for all  $n \in \mathbb{N}$

$$\mu[C_n] = \mu[A_1 \cup \dots \cup A_n] \leq \sum_{j=1}^n \mu[A_j] \leq \sum_{j=1}^{\infty} \mu[A_j].$$

If we denote the value of the infinite sum by  $c := \sum_{j=1}^{\infty} \mu[A_j]$ , then this bound  $\mu[C_n] \leq c$  for all  $n$  implies that for the limit (II.11) we also have  $\mu[C] \leq c$ . Recalling what  $C$  and  $c$  are, we have now obtained

$$\mu\left[\bigcup_{j=1}^{\infty} A_j\right] = \mu[C] \leq c = \sum_{j=1}^{\infty} \mu[A_j].$$

This finishes the proof. □

Especially part (e) of the above lemma is, despite its simplicity, so useful in probability theory that it has been given an affectionate name: “the union bound”. Because of its importance, we record this fact once more in the probabilistic context.

**Theorem II.20** (The union bound).

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and let  $E_1, E_2, \dots \in \mathcal{F}$  be a sequence of events. Then we have

$$\mathbb{P}\left[\bigcup_{j=1}^{\infty} E_j\right] \leq \sum_{j=1}^{\infty} \mathbb{P}[E_j]. \tag{II.12}$$

In other words, the probability that at least one event in a sequence occurs can not exceed the sum of the probabilities of the events in the sequence.

Probability measures enjoy some properties that may not be valid for measures of infinite total mass. The following exercise gives a few of them.

**Exercise II.3** (Properties specific to probability measures).

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space.

- (a) Show that for any event  $E \in \mathcal{F}$  we have

$$\mathbb{P}[E^c] = 1 - \mathbb{P}[E].$$

- (b) Show that for any two events  $E_1, E_2 \in \mathcal{F}$  we have

$$\mathbb{P}[E_1 \cup E_2] = \mathbb{P}[E_1] + \mathbb{P}[E_2] - \mathbb{P}[E_1 \cap E_2].$$

### Monotone convergence of probability measures

Part (d) of Lemma II.19 is a monotone convergence statement of measures for increasing sequences of sets. For general measures we do not have the corresponding monotone convergence for decreasing sequences of sets, as the following counterexample shows.

**Example II.21.** (Decreasing monotone convergence of measures can fail in general)

Consider the set  $\mathbb{N} = \{1, 2, 3, \dots\}$  of natural numbers with the counting measure  $\mu_{\#}$ , as defined in Example II.10:

$$\mu_{\#}[A] = \#A \quad \text{for all } A \subset \mathbb{N}.$$

Consider the subsets  $A_n = \{n, n+1, n+2, \dots\} \subset \mathbb{N}$ . Each of these is an infinite set, so their counting measures are infinite,  $\mu_{\#}[A_n] = +\infty$ . These sets form a decreasing sequence,

$$A_1 \supset A_2 \supset A_3 \supset \dots,$$

and the limit is the intersection

$$A = \bigcap_{n \in \mathbb{N}} A_n.$$

But no natural number  $m \in \mathbb{N}$  belongs to all of  $A_n$ ,  $n \in \mathbb{N}$ , (indeed,  $m \notin A_n$  as soon as  $n > m$ ). Therefore the intersection is empty,  $A = \emptyset$ . The number of elements in this empty set is zero,  $\#A = 0$ . In particular, the sequence of counting measures  $\mu_{\#}[A_n] = +\infty$  does *not* tend to the counting measure  $\mu_{\#}[A] = 0$  of the decreasing limit set  $A$ .

**Exercise II.4.** Construct a similar counterexample with the Lebesgue measure  $\Lambda$  on  $\mathbb{R}$ .

For probability measures, however, monotone convergence of measures holds for both increasing and decreasing sequences of events.

**Theorem II.22** (Monotone convergence of probability measures).

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space.

- (a) If  $E_1 \subset E_2 \subset \dots$  is an increasing sequence of events with limit  $E = \bigcup_{n \in \mathbb{N}} E_n$ , then we have  $\mathbb{P}[E_n] \uparrow \mathbb{P}[E]$  as  $n \rightarrow \infty$ .
- (b) If  $E_1 \supset E_2 \supset \dots$  is a decreasing sequence of events with limit  $E = \bigcap_{n \in \mathbb{N}} E_n$ , then we have  $\mathbb{P}[E_n] \downarrow \mathbb{P}[E]$  as  $n \rightarrow \infty$ .

*Proof.* Part (a) follows from part (d) of Lemma II.19, since any probability measure is a measure. Part (b) is left as an exercise.  $\square$

**Exercise II.5.** Prove part (b) of Theorem II.22 above.

## II.4. Identification and construction of measures

We now turn to the following questions:

- Does a measure with some desired properties exist?  
(How can we construct measures?)
- How can we check whether two measures are the same?  
(What does one need to know to uniquely identify a measure?)

As we noted in Section I.3, it can be complicated to work with  $\sigma$ -algebras. For the purposes of the above two questions, in particular, we often prefer to work with simpler collections. For the identification part, the appropriate simpler collections are called  $\pi$ -systems (Definition II.23 below)

Courses on traditional measure theory focus quite a lot on the construction of measures. We refer the interested reader to such measure theory courses for details — the dedicated reader will find for example Carathéodory's extension theorem, with which it is possible to construct, e.g., the Lebesgue measure  $\Lambda$  on  $\mathbb{R}$  (Example II.12) and its multi-dimensional analogue  $\Lambda^d$  on  $\mathbb{R}^d$  (Example II.13).

Instead, identification of probability measures is of practical relevance in stochastics and statistics, so we focus more on this latter question. The proof of the main identification result (Dynkin's identification theorem, Theorem II.26), is not given immediately. It can be found in Appendix C, and may be most natural to study together with the topics of Lecture IV.

### Identification of probability measures

Collections of the following type are good for the purposes of identification of measures.

#### Definition II.23 (Pi-system).

A collection  $\mathcal{J}$  of subsets of  $S$  is called a  $\pi$ -system if the following holds:

$$(\text{II-}\cap) : \quad \text{if } A, B \in \mathcal{J}, \text{ then also } A \cap B \in \mathcal{J}.$$

**Remark II.24.** Any  $\sigma$ -algebra  $\mathcal{S}$  is also a  $\pi$ -system (since the intersection of two sets in a  $\sigma$ -algebra belongs to the  $\sigma$ -algebra). However, not every  $\pi$ -system is a  $\sigma$ -algebra, as the following example shows.

#### Example II.25 (A pi-system of semi-infinite intervals).

Consider the collection

$$\mathcal{J}(\mathbb{R}) := \left\{ (-\infty, x] \mid x \in \mathbb{R} \right\} \tag{II.13}$$

of semi-infinite intervals  $(-\infty, x] \subset \mathbb{R}$ . Then  $\mathcal{J}(\mathbb{R})$  is a  $\pi$ -system on  $\mathbb{R}$ : indeed it is clearly a non-empty collection of subsets of  $\mathbb{R}$ , and given any two intervals from the collection,  $(-\infty, x_1]$  and  $(-\infty, x_2]$ , their intersection is the interval

$$(-\infty, x_1] \cap (-\infty, x_2] = (-\infty, z], \quad \text{where } z = \min\{x_1, x_2\}$$

which itself belongs to the collection  $\mathcal{J}(\mathbb{R})$ . Thus property (II- $\cap$ ) holds for  $\mathcal{J}(\mathbb{R})$ .

In Proposition I.11 we saw that  $\mathcal{J}(\mathbb{R})$  generates the Borel  $\sigma$ -algebra  $\mathcal{B}$  on  $\mathbb{R}$ , i.e.,  $\sigma(\mathcal{J}(\mathbb{R})) = \mathcal{B}$ . It is one of the simplest such  $\pi$ -systems, and for this reason we will use  $\mathcal{J}(\mathbb{R})$  over and over again, especially when dealing with real valued random variables.

The main result which is used for identification of measures is the following.

#### Theorem II.26 (Dynkin's identification theorem).

Let  $\mathbb{P}_1$  and  $\mathbb{P}_2$  be two probability measures on a measurable space  $(\Omega, \mathcal{F})$ . Assume that  $\mathcal{J}$  is a  $\pi$ -system on  $\Omega$  such that the  $\sigma$ -algebra  $\sigma(\mathcal{J})$  generated by it coincides with the  $\sigma$ -algebra  $\mathcal{F}$  of measurable sets in the measurable space, i.e.,  $\sigma(\mathcal{J}) = \mathcal{F}$ . Then the following are equivalent:

- (i)  $P_1[E] = P_2[E]$  for all  $E \in \mathcal{J}$
- (ii) the two probability measures are equal,  $P_1 = P_2$ .

**Remark II.27.** Condition (ii) above is clearly stronger than condition (i). Namely, the equality of probability measures  $P_1 = P_2$  means that  $P_1[E] = P_2[E]$  for all  $E \in \mathcal{F}$ , and there are in general more sets  $E$  in this  $\sigma$ -algebra  $\mathcal{F}$  than in the  $\pi$ -system  $\mathcal{J}$ . The nontrivial part of the proof is therefore that condition (i) implies (ii). This will be proven in Appendix C.3.

### Cumulative distribution function

In the following we consider a probability measure on the real axis  $\mathbb{R}$ . Typically such a probability measure could appear as the law of a real-valued random variable, as we will discuss later on in the course. In order to have a suitable notation for such common situations, let us denote the probability measure in this case by  $\nu$  instead of  $P$ .

**Definition II.28** (Cumulative distribution function).

If  $\nu$  is a probability measure on  $(\mathbb{R}, \mathcal{B})$ , then the *cumulative distribution function* (c.d.f.) of  $\nu$  is the function  $F: \mathbb{R} \rightarrow [0, 1]$  defined by

$$F(x) := \nu\left[(-\infty, x]\right].$$

A simple but important application of Dynkin's identification theorem is the following. This case is applicable, e.g., to the identification of the laws of real-valued random variables.

**Corollary II.29** (Cumulative distribution function identify distributions).

Let  $\nu_1$  and  $\nu_2$  be two probability measures on  $(\mathbb{R}, \mathcal{B})$ , and  $F_1$  and  $F_2$  their cumulative distribution functions, respectively. Then the following are equivalent:

- (i) The cumulative distribution functions are equal,  $F_1 = F_2$ .
- (ii) The probability measures are equal,  $\nu_1 = \nu_2$ .

*Proof:* Equivalence is proved by establishing both implications (i)  $\Rightarrow$  (ii) and (ii)  $\Rightarrow$  (i).

*proof of (ii)  $\Rightarrow$  (i):* Assuming the probability measures are equal,  $\nu_1 = \nu_2$ , we get for any  $x \in \mathbb{R}$

$$F_1(x) := \nu_1[(-\infty, x]] = \nu_2[(-\infty, x]] =: F_2(x).$$

*proof of (i)  $\Rightarrow$  (ii):* Assume the equality  $F_1 = F_2$  of cumulative distribution functions, i.e., that  $F_1(x) = F_2(x)$  for all  $x \in \mathbb{R}$ . Consider the  $\pi$ -system  $\mathcal{J}(\mathbb{R})$  of Example II.25. A set  $A \in \mathcal{J}(\mathbb{R})$  of this  $\pi$ -system is by definition of the form  $A = (-\infty, x]$  for some  $x \in \mathbb{R}$ . For such a set, we get

$$\nu_1[(-\infty, x]] =: F_1(x) = F_2(x) =: \nu_2[(-\infty, x]],$$

so we have that  $\nu_1$  and  $\nu_2$  coincide on  $\mathcal{J}(\mathbb{R})$ . The  $\sigma$ -algebra  $\sigma(\mathcal{J}(\mathbb{R}))$  generated by the  $\pi$ -system  $\mathcal{J}(\mathbb{R})$  coincides with the Borel  $\sigma$ -algebra  $\mathcal{B}$  by Proposition I.11(i). Theorem II.26 then guarantees that  $\nu_1$  and  $\nu_2$  coincide on the entire Borel  $\sigma$ -algebra  $\mathcal{B}$ .  $\square$

Since cumulative distribution functions characterize probability measures on  $(\mathbb{R}, \mathcal{B})$  by Corollary II.29 above, it is natural to next ask which functions  $F$  can qualify as cumulative distribution functions. There is indeed a rather explicit characterization.

Deriving the necessary conditions below is an instructive application of the basic properties of measures which we established in Lemma II.19.

**Proposition II.30** (Properties of cumulative distribution functions).

*If  $F: \mathbb{R} \rightarrow [0, 1]$  is the cumulative distribution function of a probability measure  $\nu$  on  $(\mathbb{R}, \mathcal{B})$ , then it satisfies the following properties:*

- (a)  *$F$  is increasing: if  $x \leq y$  then  $F(x) \leq F(y)$*
- (b)  *$F$  is right-continuous: if  $x_n \downarrow x \in \mathbb{R}$  as  $n \rightarrow \infty$ , then  $F(x_n) \downarrow F(x)$*
- (c)  *$\lim_{x \rightarrow +\infty} F(x) = 1$  and  $\lim_{x \rightarrow -\infty} F(x) = 0$ .*

**Exercise II.6.** Prove Proposition II.30.

**Hint:** Use appropriate parts of Lemma II.19.



## Lecture III

### Random variables

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space, i.e.,

- $\Omega$  — the set of all possible outcomes
- $\mathcal{F}$  — the collection of events (see Lecture I)
- $\mathbf{P}$  — the probability measure (see Lecture II).

The key idea of a random variable is the following two step procedure by which randomness is thought to have effect:

- 1.) “Chance determines the random outcome  $\omega \in \Omega$ .”
- 2.) “The outcome  $\omega$  determines various quantities of interest.”  
(random variables)

Therefore, a random variable will be a function  $X$  defined on  $\Omega$ , which to an outcome  $\omega \in \Omega$  associates the value  $X(\omega)$  of some quantity of interest. The function

$$X: \Omega \rightarrow S'$$

takes values in a suitable set  $S'$  of possible values of the quantity of interest (some examples are given below, in Example III.4). Crucially, this function has to be sufficiently well-behaved so that we can talk about probabilities with which the quantity assumes certain values. So whenever  $A' \subset S'$  is a reasonable enough subset of the possible values, the set of outcomes  $\omega$  for which  $X(\omega)$  belongs to  $A'$  should constitute an event, i.e.,

$$\{\omega \in \Omega \mid X(\omega) \in A'\} \in \mathcal{F}. \quad (\text{III.1})$$

This requirement of well-behavedness of the function  $X: \Omega \rightarrow S'$  is what the notion of *measurable function* captures (cf. Definition III.1 below).

The set in (III.1) above is just the preimage of  $A'$  under the function  $X: \Omega \rightarrow S'$ : indeed by definition we have

$$X^{-1}(A') = \{\omega \in \Omega \mid X(\omega) \in A'\} \subset \Omega.$$

For simplicity, we will often abbreviate this just as

$$\{X \in A'\} \subset \Omega.$$

This last slight abuse of notation is not only shorter, but it also has the advantage that the probabilistic interpretation

“the value of our (random) quantity of interest  $X$  lies in  $A'$ ”

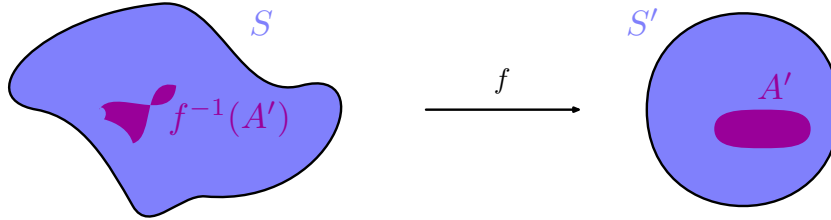
of the event becomes apparent at a glance.

### III.1. Measurable functions and random variables

Let  $(S, \mathcal{S})$  and  $(S', \mathcal{S}')$  be two measurable spaces, i.e.,  $S$  and  $S'$  are two sets and  $\mathcal{S}$  and  $\mathcal{S}'$  are  $\sigma$ -algebras on these two respectively.

**Definition III.1** (Measurable function).

A function  $f: S \rightarrow S'$  is called  $\mathcal{S}/\mathcal{S}'$ -measurable<sup>1</sup> if for all  $A' \in \mathcal{S}'$  we have  $f^{-1}(A') \in \mathcal{S}$ .



Let now  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space, and  $(S', \mathcal{S}')$  a measurable space.

**Definition III.2** (Random variable).

A *random variable* with values in  $S'$  is a  $\mathcal{F}/\mathcal{S}'$ -measurable function  $X: \Omega \rightarrow S'$ .

**Remark III.3.** This definition precisely requires that  $\{X \in A'\} \subset \Omega$  is an event whenever the subset  $A' \subset S'$  is  $\mathcal{S}'$ -measurable.

**Example III.4** (Examples of types of random variables).

Depending on our quantity of interest, the set  $S'$  of allowed values of the random variable can be for example one of the following:

**random numbers:** The case  $S' = \mathbb{R}$ ,  $\mathcal{S}' = \mathcal{B}$  (Borel  $\sigma$ -algebra on the real line) corresponds to real-valued random variables, i.e., *random numbers*.

**random vectors:** e.g.,  $S' = \mathbb{R}^d$ ,  $\mathcal{S}' = \mathcal{B}(\mathbb{R}^d)$  (Borel sigma-algebra on  $\mathbb{R}^d$ )

**random matrices:** e.g.,  $S' = \mathbb{R}^{m \times n}$ ,  $\mathcal{S}' = \mathcal{B}(\mathbb{R}^{m \times n})$  (Borel sigma-algebra on  $\mathbb{R}^{m \times n}$ )

**random graphs:**  $S'$  some set of graphs,  $\mathcal{S}'$  a suitably chosen  $\sigma$ -algebra (often simply the power set  $\mathcal{P}(S')$ )

**etc.:** ...

Usually  $S$  and  $S'$  are at least topological spaces, so we can and will equip them with their Borel  $\sigma$ -algebras  $\mathcal{S} = \mathcal{B}(S)$  and  $\mathcal{S}' = \mathcal{B}(S')$ , generated by their open subsets (see Definition I.10). A  $\mathcal{B}(S)/\mathcal{B}(S')$ -measurable function  $f: S \rightarrow S'$  will be called a *Borel-measurable function* or simply a *Borel function*. Let us give one interesting example of a random variable with values in a topological space which is not as simple as the finite dimensional Euclidean spaces in the previous example.

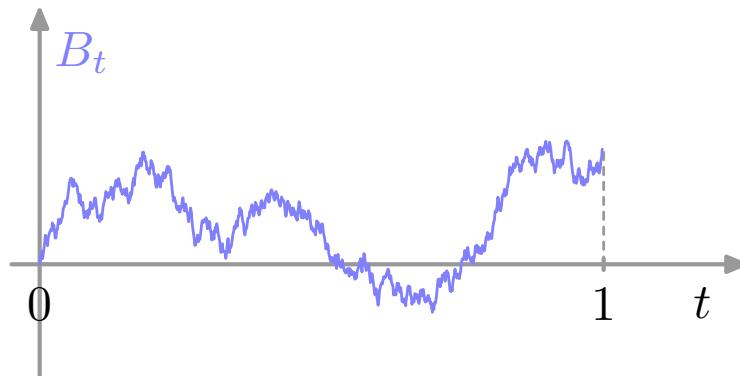
**Example III.5** (Brownian motion as a random variable).

Brownian motion is one of the most important stochastic processes: it is a continuous real valued Markov process in continuous time, which is used to model many things from thermal motion of microscopic particles to stock prices in finance. Mathematically, the Brownian

<sup>1</sup>When the two  $\sigma$ -algebras are clear from the context, we usually just say that the function is *measurable*.



motion on the unit time interval is a certain random variable taking values in the space  $S' = \mathcal{C}([0, 1])$  of continuous functions  $h: [0, 1] \rightarrow \mathbb{R}$ , with the topology induced by the uniform norm  $\|h\|_\infty = \sup_{t \in [0, 1]} |h(t)|$  and the corresponding Borel  $\sigma$ -algebra  $\mathcal{B}(\mathcal{C}([0, 1]))$ .



### The law of a random variable

Suppose that  $X: \Omega \rightarrow S'$  is a random variable. Then there is a probability measure on  $S'$  which describes how the values of the random variable are distributed.

**Definition III.6** (The law of a random variable).

The *law* (or the *distribution*) of the random variable  $X: \Omega \rightarrow S'$  is the probability measure  $P_X$  on  $(S', \mathcal{S}')$  defined by

$$P_X[A'] = \mathbf{P}[X^{-1}(A')], \quad \text{for } A' \in \mathcal{S}'. \quad (\text{III.2})$$

**Exercise III.1.** Verify that  $P_X$  given by (III.2) is indeed a probability measure on  $(S', \mathcal{S}')$ .

With a slight abuse of notation, we usually write the defining equation (III.2) of the law of  $X$  in the more descriptive form

$$P_X[A'] = \mathbf{P}[X \in A'].$$

If we were to insist on carefully following the notation that was introduced in Lecture II, then instead of “ $X \in A'$ ” we should in principle write “ $\{\omega \in \Omega \mid X(\omega) \in A'\}$ ” (this subset of the sample space is the event whose probability concerns us). But it is clear that doing so would become really cumbersome in actual practice, so shorthand notations of the above kind are commonplace in probabilistic literature.

## III.2. Indicator random variables

Constant functions provide trivial examples of random variables (verify the measurability directly from the definition!). They in fact correctly model situations when a quantity of interest contains no randomness — such random variables are usually called “deterministic”.

Constant functions have only one possible value (the constant in question). Arguably the next simplest example of a random variable would be one which assumes one of two possible values (depending on the random outcome). It is convenient to take 0

and 1 as those two values, in which case one speaks of an *indicator random variable*. For  $E \subset \Omega$  a subset, we define the *indicator function*  $\mathbb{I}_E: \Omega \rightarrow \mathbb{R}$  of  $E$  by

$$\mathbb{I}_E(\omega) = \begin{cases} 1 & \text{if } \omega \in E \\ 0 & \text{if } \omega \notin E. \end{cases} \quad (\text{III.3})$$

**Exercise III.2** (Measurability of indicators).

Prove that  $\mathbb{I}_E$  is  $\mathcal{F}$ -measurable if and only if  $E \in \mathcal{F}$ , i.e., if  $E \subset \Omega$  is an event.

When  $\mathbb{I}_E$  is  $\mathcal{F}$ -measurable, we call it the *indicator random variable* of the event  $E$ . It “indicates” the occurrence of the event  $E$  in the sense that it takes the value 1 if the event  $E$  occurs and value 0 otherwise.

**Exercise III.3** (Indicators of intersections and unions).

Let  $A, B \subset \Omega$ .

- (a) Show that  $\mathbb{I}_{A \cap B} = \mathbb{I}_A \mathbb{I}_B$ .
- (b) When is it true that  $\mathbb{I}_{A \cup B} = \mathbb{I}_A + \mathbb{I}_B$ ?

### III.3. Constructing random variables

So far we have given the definition of random variables, but we have hardly addressed the issue of constructing them — we have only left it as an exercise to characterize the measurability of random variables with values 0 and 1. At this stage one might therefore still worry that perhaps the requirements of a measurable function are too stringent for any interesting examples to exist. . . Fortunately, this is not the case: almost all functions that you ever encounter turn out to be measurable. The rest of this lecture is devoted to understanding why.

#### A very easy case

The example here concerns one case in which we do not have to worry about the existence of random variables at all — when all subsets of the sample space  $\Omega$  are events, then every function on  $\Omega$  is measurable.

**Example III.7** (A case when all functions are measurable).

Suppose that  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space in which all subsets of  $\Omega$  are events  $\mathcal{F} = \mathcal{P}(\Omega)$ . Then any function  $X: \Omega \rightarrow \mathbb{R}$  is a random variable. Indeed, for any  $B \in \mathcal{B}$ , the preimage  $X^{-1}(B) \subset \Omega$  is a subset, and therefore  $X^{-1}(B) \in \mathcal{P}(\Omega)$ .

As we have mentioned before, we can usually only take  $\mathcal{F} = \mathcal{P}(\Omega)$  when  $\Omega$  is countable, so Example III.7 only reassures us of the existence of plenty of random variables on countable sample spaces. It remains to convince ourselves that interesting random variables exist in other common situations.

### Composition of measurable functions

Most of the time in practice, we construct more complicated random variables from simpler ones. The measurability of the simpler building blocks should ideally directly imply the measurability of the more complicated construction.

Composition of functions is one such construction.

**Proposition III.8** (Composition of measurable functions is measurable).

*Suppose that  $(S, \mathcal{S})$ ,  $(S', \mathcal{S}')$ , and  $(S'', \mathcal{S}'')$  are three measurable spaces and that we are given two measurable functions:*

$$\begin{aligned} f: S &\rightarrow S' && \text{is } \mathcal{S}/\mathcal{S}'\text{-measurable} \\ g: S' &\rightarrow S'' && \text{is } \mathcal{S}'/\mathcal{S}''\text{-measurable.} \end{aligned}$$

*Then the composite function  $s \mapsto g(f(s))$  is also measurable:*

$$g \circ f: S \rightarrow S'' \quad \text{is } \mathcal{S}/\mathcal{S}''\text{-measurable.}$$

*Proof.* We check the measurability of  $g \circ f$  directly using the definition. Let  $A'' \in \mathcal{S}''$  be a measurable subset of  $S''$ . Then first, by the  $\mathcal{S}'/\mathcal{S}''$ -measurability of  $g$ , the preimage under  $g$  of  $A''$  is measurable: we have  $g^{-1}(A'') \in \mathcal{S}'$ . Using the  $\mathcal{S}/\mathcal{S}'$ -measurability of  $f$ , it then follows that also the preimage under  $f$  of the measurable set  $g^{-1}(A'')$  is measurable: we have  $f^{-1}(g^{-1}(A'')) \in \mathcal{S}$ . It is an easy fact about preimages of composite functions that  $f^{-1}(g^{-1}(A'')) = (g \circ f)^{-1}(A'')$ . Therefore we have checked the measurability of the preimage  $(g \circ f)^{-1}(A'')$  of any measurable subset  $A'' \subset S''$  under the composite function  $g \circ f: S \rightarrow S''$ , which by definition says that the composite function  $g \circ f$  is measurable.  $\square$

We postpone examples until after we have some more tools at our disposal.

### Practical verification of measurability

To verify measurability, it is in fact enough to check the defining condition for just some collection of subsets that generates the  $\sigma$ -algebra on the target space.

**Lemma III.9** (A sufficient condition for measurability).

*Let  $\mathcal{C}' \subset \mathcal{P}(S')$  be a collection of subsets of  $S'$  that generates the  $\sigma$ -algebra  $\mathcal{S}'$ , i.e.,  $\sigma(\mathcal{C}') = \mathcal{S}'$ . Then a function  $f: S \rightarrow S'$  is  $\mathcal{S}/\mathcal{S}'$ -measurable if and only if  $f^{-1}(C') \in \mathcal{S}$  for all  $C' \in \mathcal{C}'$ .*

*Proof.* The condition is clearly necessary for measurability of  $f$ : if  $C' \in \mathcal{C}' \subset \sigma(\mathcal{C}') = \mathcal{S}'$  then the definition of measurability requires that  $f^{-1}(C') \in \mathcal{S}$ . It therefore remains only to prove that the condition is also sufficient.

Assume now that  $f^{-1}(C') \in \mathcal{S}$  for all  $C' \in \mathcal{C}'$ . We must prove that then the function  $f$  is  $\mathcal{S}/\mathcal{S}'$ -measurable. Define

$$\mathcal{G}' = \left\{ G' \in \mathcal{S}' \mid f^{-1}(G') \in \mathcal{S} \right\}$$

as the collection of the “good” subsets  $G'$  of  $S'$ , whose preimages are measurable. By assumption we have  $\mathcal{C}' \subset \mathcal{G}'$ . Now since preimages satisfy the properties

$$\begin{aligned} f^{-1}(S') &= S \\ f^{-1}((G')^c) &= (f^{-1}(G'))^c \\ f^{-1}\left(\bigcup_{n=1}^{\infty} G'_n\right) &= \bigcup_{n=1}^{\infty} f^{-1}(G'_n) \end{aligned}$$

(see Exercise A.2) and  $\mathcal{S}$  is a  $\sigma$ -algebra (on  $S$ ), we see that the collection  $\mathcal{G}'$  of subsets of  $S'$  is a  $\sigma$ -algebra (on  $S'$ ).

The fact that the  $\sigma$ -algebra  $\mathcal{G}'$  contains the collection  $\mathcal{C}'$  implies that it contains also the  $\sigma$ -algebra generated by that collection, i.e.,  $\sigma(\mathcal{C}') \subset \mathcal{G}'$ . But we have assumed that  $\sigma(\mathcal{C}') = \mathcal{S}'$ , so we conclude that  $\mathcal{S}' \subset \mathcal{G}'$ . By definition of  $\mathcal{G}'$  this means that every  $\mathcal{S}'$ -measurable subset  $A'$  has the property that  $f^{-1}(A') \in \mathcal{S}$ . This shows the measurability of the function  $f$ .  $\square$

As an application, we see that all continuous functions are good for our purposes. This gives us quite a lot of measurable functions already.

**Corollary III.10** (Continuous functions are Borel-measurable).

*Let  $\mathfrak{X}$  and  $\mathfrak{X}'$  be two topological spaces (e.g., metric spaces). Then any continuous function  $f: \mathfrak{X} \rightarrow \mathfrak{X}'$  is Borel-measurable.*

*Proof.* Let  $f: \mathfrak{X} \rightarrow \mathfrak{X}'$  be a continuous function and let  $\mathcal{B}(\mathfrak{X})$  and  $\mathcal{B}(\mathfrak{X}')$  be the Borel  $\sigma$ -algebras on  $\mathfrak{X}$  and  $\mathfrak{X}'$ , respectively. Recall that Borel-measurability of the function  $f$  means that  $f$  is  $\mathcal{B}(\mathfrak{X})/\mathcal{B}(\mathfrak{X}')$ -measurable.

Recall also that the Borel  $\sigma$ -algebra  $\mathcal{B}(\mathfrak{X}')$  is generated by the collection  $\mathcal{T}(\mathfrak{X}')$  of all open subsets  $V' \subset \mathfrak{X}'$ . If  $f: \mathfrak{X} \rightarrow \mathfrak{X}'$  is continuous, then for any open subset  $V' \subset \mathfrak{X}'$  the preimage  $f^{-1}(V') \subset \mathfrak{X}$  is also open:  $f^{-1}(V') \in \mathcal{T}(\mathfrak{X})$ . As an open set, the preimage is in particular Borel-measurable:  $f^{-1}(V') \in \mathcal{B}(\mathfrak{X})$ . Therefore we have checked that  $f^{-1}(V') \in \mathcal{B}(\mathfrak{X})$  holds for all  $V'$  in the collection  $\mathcal{T}(\mathfrak{X}')$  which generates  $\mathcal{B}(\mathfrak{X}')$ . By Lemma III.9, this is sufficient to show that the function  $f$  is  $\mathcal{B}(\mathfrak{X})/\mathcal{B}(\mathfrak{X}')$ -measurable.  $\square$

For the case  $S' = \mathbb{R}$  (e.g., random numbers), we can in fact use even simpler collections that generate the Borel  $\sigma$ -algebra.

**Corollary III.11** (Measurability of real valued functions).

*A function  $f: S \rightarrow \mathbb{R}$  is  $\mathcal{S}/\mathcal{B}$ -measurable if and only if for all  $c \in \mathbb{R}$  we have  $\{f \leq c\} \in \mathcal{S}$ .*

**Remark III.12.** Recall that the notation  $\{f \leq c\}$  is shorthand for

$$\{f \leq c\} := \{s \in S \mid f(s) \leq c\} \subset S,$$

which is also the preimage  $f^{-1}((-\infty, c])$  of the interval  $(-\infty, c]$ .

*Proof of Corollary III.11.* The collection  $\mathcal{C} = \{(-\infty, c] \mid c \in \mathbb{R}\}$  generates the Borel  $\sigma$ -algebra  $\mathcal{B}$  on  $\mathbb{R}$  by Proposition I.11, so the assertion follows from Lemma III.9 above.  $\square$

### Pointwise operations on measurable functions

Since real numbers have addition and multiplication, we can perform such operations on real-valued functions pointwise. If  $f_1, f_2: S \rightarrow \mathbb{R}$  are two functions, then the

pointwise sum  $f_1 + f_2: S \rightarrow \mathbb{R}$  is defined by

$$(f_1 + f_2)(s) = f_1(s) + f_2(s) \quad \forall s \in S,$$

the pointwise product  $f_1 f_2: S \rightarrow \mathbb{R}$  by

$$(f_1 f_2)(s) = f_1(s) f_2(s) \quad \forall s \in S,$$

and if  $f: S \rightarrow \mathbb{R}$  is a function and  $\lambda \in \mathbb{R}$  is a scalar, then the pointwise scalar multiple  $\lambda f: S \rightarrow \mathbb{R}$  is defined by

$$(\lambda f)(s) = \lambda f(s) \quad \forall s \in S.$$

Let us denote by  $\mathfrak{m}\mathcal{S}$  the set of all  $\mathcal{S}/\mathcal{B}$ -measurable functions  $S \rightarrow \mathbb{R}$ . The pointwise operations allow us to construct new measurable functions from old ones.

**Proposition III.13** (Pointwise sums and products preserve measurability).

*Let  $(S, \mathcal{S})$  be a measurable space. Then the set  $\mathfrak{m}\mathcal{S}$  of all measurable real valued functions on it is stable under taking pointwise sums, pointwise products, and pointwise scalar multiples, i.e., the following hold:*

- (i)  $f \in \mathfrak{m}\mathcal{S}, \lambda \in \mathbb{R} \implies \lambda f \in \mathfrak{m}\mathcal{S}$
- (ii)  $f_1, f_2 \in \mathfrak{m}\mathcal{S} \implies f_1 + f_2 \in \mathfrak{m}\mathcal{S}$
- (iii)  $f_1, f_2 \in \mathfrak{m}\mathcal{S} \implies f_1 f_2 \in \mathfrak{m}\mathcal{S}$ .

*Proof:* We will prove part (ii), and leave parts (i) and (iii) as exercises.

*proof of (ii):* Suppose that  $f_1, f_2 \in \mathfrak{m}\mathcal{S}$ . Note that for  $c \in \mathbb{R}$  and  $s \in S$ , the condition  $f_1(s) + f_2(s) > c$  holds if and only if there exists a rational number  $q \in \mathbb{Q}$  such that we have  $f_1(s) > q$  and  $f_2(s) > c - q$ . In other words, the following two subsets of  $S$  are equal

$$\{f_1 + f_2 > c\} = \bigcup_{q \in \mathbb{Q}} \left( \{f_1 > q\} \cap \{f_2 > c - q\} \right).$$

By measurability of the functions  $f_1$  and  $f_2$ , the subsets  $\{f_1 > q\}, \{f_2 > c - q\} \subset S$  are measurable. By properties of  $\sigma$ -algebras, then, the set on the right hand side above is measurable (note that the set  $\mathbb{Q}$  of rational numbers is countable). We have thus shown that  $\{f_1 + f_2 > c\} \in \mathcal{S}$ , and by taking complements we get

$$\{f_1 + f_2 \leq c\} = \{f_1 + f_2 > c\}^c \in \mathcal{S}.$$

By Corollary III.11, this is sufficient to show the measurability of the pointwise sum function  $f_1 + f_2$ .  $\square$

**Exercise III.4.** Prove the assertions (i) and (iii) in Proposition III.13.

We can also use more subtle pointwise operations. If we are given a sequence  $f_1, f_2, \dots$  of functions  $f_n: S \rightarrow \mathbb{R}$  for  $n \in \mathbb{N}$ , then we can define functions  $\sup_n f_n$  and  $\inf_n f_n$  on  $S$  by the pointwise supremum and infimum

$$\left( \sup_{n \in \mathbb{N}} f_n \right)(s) = \sup_{n \in \mathbb{N}} f_n(s) \quad \text{and} \quad \left( \inf_{n \in \mathbb{N}} f_n \right)(s) = \inf_{n \in \mathbb{N}} f_n(s).$$

These functions, however, may assume values  $+\infty$  and  $-\infty$  even if each  $f_n$  takes only finite real values. Thus the functions

$$\sup_{n \in \mathbb{N}} f_n: S \rightarrow [-\infty, +\infty] \quad \text{and} \quad \inf_{n \in \mathbb{N}} f_n: S \rightarrow [-\infty, +\infty]$$

are defined so that their allowed range of values is the *extended real line*

$$[-\infty, +\infty] = \mathbb{R} \cup \{+\infty\} \cup \{-\infty\}.$$

The extended real line is equipped with the topology of a closed interval, so that for example the tangent function

$$\tan: \left[-\frac{\pi}{2}, \frac{\pi}{2}\right] \rightarrow [-\infty, +\infty]$$

is a homeomorphism, when we interpret  $\tan(\pi/2) = +\infty$  and  $\tan(-\pi/2) = -\infty$  (see Example B.19). From this topology, we get the Borel  $\sigma$ -algebra  $\mathcal{B}([-\infty, +\infty])$  on the extended real line, in the usual way (Definition I.10). It is so convenient to occasionally allow  $+\infty$  and  $-\infty$  as possible values that we will slightly abuse the notation, and write  $f \in \mathfrak{m}\mathcal{S}$  also for functions

$$f: S \rightarrow [-\infty, +\infty]$$

that are  $\mathcal{S}/\mathcal{B}([-\infty, +\infty])$ -measurable.

Likewise, for a sequence  $f_1, f_2, \dots$  of functions  $S \rightarrow \mathbb{R}$  we can define the pointwise lim sup and lim inf,

$$\left(\limsup_n f_n\right)(s) = \limsup_n f_n(s) \quad \text{and} \quad \left(\liminf_n f_n\right)(s) = \liminf_n f_n(s).$$

**Proposition III.14** (Pointwise supremum and infimum preserve measurability).

Let  $(S, \mathcal{S})$  be a measurable space. Suppose that  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}$ . Then we have:

- (i)  $\sup_n f_n \in \mathfrak{m}\mathcal{S}$
- (ii)  $\inf_n f_n \in \mathfrak{m}\mathcal{S}$
- (iii)  $\limsup_n f_n \in \mathfrak{m}\mathcal{S}$
- (iv)  $\liminf_n f_n \in \mathfrak{m}\mathcal{S}$ .

*Proof:* We will prove only parts (i) and (iii) — parts (ii) and (iv) are entirely similar.<sup>2</sup>

*proof of part (i).* Note that we have  $\sup_n f_n(s) \leq c$  if and only if  $f_n(s) \leq c$  for all  $n \in \mathbb{N}$ . Therefore we can write

$$\left\{ \sup_n f_n \leq c \right\} = \bigcap_{n \in \mathbb{N}} \{f_n \leq c\}.$$

For each  $n \in \mathbb{N}$  we have  $\{f_n \leq c\} = f_n^{-1}((-\infty, c]) \in \mathcal{S}$ , since  $f_n \in \mathfrak{m}\mathcal{S}$ . Therefore we have  $\{\sup_n f_n \leq c\} \in \mathcal{S}$  as a countable intersection of sets in the  $\sigma$ -algebra  $\mathcal{S}$ . By Corollary III.11 we conclude that  $\sup_n f_n$  is a  $\mathcal{S}$ -measurable function.

*proof of part (iii).* Recall that we have

$$\limsup_n f_n(s) = \inf_{n \in \mathbb{N}} \left( \sup_{k \geq n} f_k(s) \right),$$

because the sequence of functions  $g_1, g_2, \dots$  defined by

$$g_n(s) = \sup_{k \geq n} f_k(s)$$

is decreasing (for larger  $n$  the supremum contains fewer terms)

$$g_1(s) \geq g_2(s) \geq g_3(s) \geq \dots$$

The function  $g_n$  is a pointwise supremum of the measurable functions  $f_k \in \mathfrak{m}\mathcal{S}$ ,  $k \geq n$ , and thus itself  $\mathcal{S}$ -measurable by part (i). Consequently,  $\limsup_n f_n$  is the pointwise infimum of the measurable functions  $g_n \in \mathfrak{m}\mathcal{S}$ ,  $n \in \mathbb{N}$ , and thus itself  $\mathcal{S}$ -measurable by part (ii).  $\square$

<sup>2</sup>In the case when values  $\pm\infty$  do not appear, parts (ii) and (iv) in fact directly follow from parts (i) and (iii) and Proposition III.13(i) by observing that  $\inf_n f_n = -\sup_n(-f_n)$  and  $\liminf_n f_n = -\limsup_n(-f_n)$ .

With the above we can conclude the extremely useful fact that pointwise limits of measurable functions are measurable.

**Corollary III.15** (Pointwise limits of measurable functions are measurable).

Let  $(S, \mathcal{S})$  be a measurable space. Suppose that  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}$ , and suppose that for all  $s \in S$  the limit

$$\lim_{n \rightarrow \infty} f_n(s)$$

exists. Then the pointwise limit function is measurable:  $\lim_{n \rightarrow \infty} f_n \in \mathfrak{m}\mathcal{S}$ .

*Proof.* When the assumed limit exists, it coincides with both the lim sup and lim inf (cf. Proposition B.4),

$$\liminf_n f_n(s) = \lim_{n \rightarrow \infty} f_n(s) = \limsup_n f_n(s).$$

Therefore the assertion follows from Proposition III.14, part (iii) or (iv).  $\square$

Let us now revisit earlier examples about repeated coin tossing. The following illustrates how the operations we studied above allow us to construct rather nontrivial random variables starting from very basic ones.

**Example III.16** (The question of existence of limit frequency in coin tossing).

Let  $\Omega = \{\text{H}, \text{T}\}^{\mathbb{N}}$  be the sample space for repeated coin tossing as in Examples O.2 and I.5. Let  $\mathcal{F}$  be the  $\sigma$ -algebra on  $\Omega$  generated by the events

$$E_j := \{\omega \in \Omega \mid \omega(j) = \text{H}\} = \text{“the } j\text{:th coin toss is heads”}$$

for  $j \in \mathbb{N}$ , i.e.,  $\mathcal{F} = \sigma(\{E_j \mid j \in \mathbb{N}\})$ . Then the indicator random variable of  $E_j$ ,

$$\mathbb{I}_{E_j}(\omega) = \begin{cases} 1 & \text{if } \omega(j) = \text{H} \\ 0 & \text{if } \omega(j) = \text{T}, \end{cases}$$

is  $\mathcal{F}$ -measurable by Exercise III.2. The relative frequency of heads in the first  $n$  coin tosses,

$$X_n(\omega) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}_{E_j}(\omega),$$

is then also  $\mathcal{F}$ -measurable by Proposition III.13 (the relative frequency  $X_n$  is constructed from the indicators  $\mathbb{I}_{E_j}$  by finite sums and a scalar multiple). Therefore also the upper and lower limits

$$L^+(\omega) = \limsup_n X_n(\omega) \quad \text{and} \quad L^-(\omega) = \liminf_n X_n(\omega)$$

are random variables, by Proposition III.14. Knowing that these are random variables, do we learn something interesting about events?

For example, for any  $r \in [0, 1]$ , we should hope to be able to form the event

$$\text{“relative frequencies of heads tend to } r\text{”} = \left\{ \lim_{n \rightarrow \infty} X_n = r \right\}$$

which in more careful notation is the subset  $\{\omega \in \Omega \mid \lim_{n \rightarrow \infty} X_n(\omega) = r\} \subset \Omega$ . Note that we have

$$\left\{ \lim_{n \rightarrow \infty} X_n = r \right\} = \{L^+ = r\} \cap \{L^- = r\},$$

and since  $L^+$  and  $L^-$  are  $\mathcal{F}$ -measurable random variables, the two preimages

$$\{L^+ = r\} = (L^+)^{-1}(\{r\}) \quad \text{and} \quad \{L^- = r\} = (L^-)^{-1}(\{r\})$$

are  $\mathcal{F}$ -measurable events. Thus also their intersection is indeed an event,

$$\left\{ \lim_{n \rightarrow \infty} X_n = r \right\} \in \mathcal{F},$$

and so it will at least be meaningful to talk about the probability (once we introduce a probability measure  $\mathbb{P}$ ) of the events that the relative frequency tends to a particular limiting value  $r \in [0, 1]$ .

In Example I.5 we in fact already reached essentially the same conclusion by direct manipulations with events, without using random variables. The present approach, however, has some advantages. Consider for instance the slight variation:

$$\text{“relative frequencies have a limit”} = \left\{ \exists \lim_{n \rightarrow \infty} X_n \right\}.$$

To show that this is an event, we may now just notice that the limit exists if and only if the corresponding upper and lower limits coincide,  $L^+ = L^-$ , i.e. if their difference vanishes

$$\left\{ \exists \lim_{n \rightarrow \infty} X_n \right\} = \{L^+ - L^- = 0\}.$$

This is the preimage of  $\{0\} \subset \mathbb{R}$  under the pointwise difference  $L^+ - L^-$  of random variables, and as such it is measurable,

$$\left\{ \exists \lim_{n \rightarrow \infty} X_n \right\} \in \mathcal{F}.$$

It would be more cumbersome to try to conclude the same directly by manipulating events using countable set operations in the spirit of Example I.5.

### III.4. Simple functions

In Section III.2 we discussed the measurability of functions assuming two possible values, the indicator functions of subsets  $A \subset S$  defined by

$$\mathbb{I}_A: S \rightarrow \mathbb{R}, \quad \mathbb{I}_A(s) = \begin{cases} 1 & \text{if } s \in A \\ 0 & \text{if } s \notin A. \end{cases}$$

In Exercise III.2 the measurability of the indicator functions was characterized: we have  $\mathbb{I}_A \in \mathfrak{m}\mathcal{S}$  if and only if  $A \in \mathcal{S}$ .

From two possible values we next proceed modestly to finitely many possible values.

**Definition III.17** (Simple function).

A real valued measurable function assuming only finitely many different values is called a *simple function* (or a *simple random variable* in the probabilistic context).

Any finite linear combination of indicator functions

$$f(s) = \sum_{k=1}^m a_k \mathbb{I}_{A_k} \tag{III.4}$$

of indicators of measurable sets  $A_1, \dots, A_m \in \mathcal{S}$  is a simple function: it is measurable by Proposition III.13, and there are only finitely many real numbers that can be expressed as a sum of some of the coefficients  $a_1, \dots, a_m \in \mathbb{R}$ . In fact, if a simple function  $f: S \rightarrow \mathbb{R}$  can only assume  $m$  different real values  $a_1, \dots, a_m$ , then we can write it as a linear combination (III.4) with the disjoint measurable sets  $A_k := f^{-1}(\{a_k\})$ . Occasionally, choosing this minimal linear combination is very convenient, but at other times we might not want to insist on disjointness.



In Corollary III.15 we saw that any pointwise limit of measurable functions is measurable. An important version of the converse statement is also true: any measurable function can be obtained as a pointwise limit of simple functions!

When  $(S, \mathcal{S})$  is a measurable space, denote by  $m\mathcal{S}^+$  the set of all functions

$$f: S \rightarrow [0, +\infty]$$

which are  $\mathcal{S}/\mathcal{B}([0, +\infty])$  measurable.

**Lemma III.18** (Approximation of non-negative measurable functions).

Let  $f \in m\mathcal{S}^+$ . Then there exists a sequence  $f_1, f_2, \dots : S \rightarrow [0, +\infty)$  of non-negative simple functions such that  $f_n \uparrow f$  pointwise as  $n \rightarrow \infty$ .

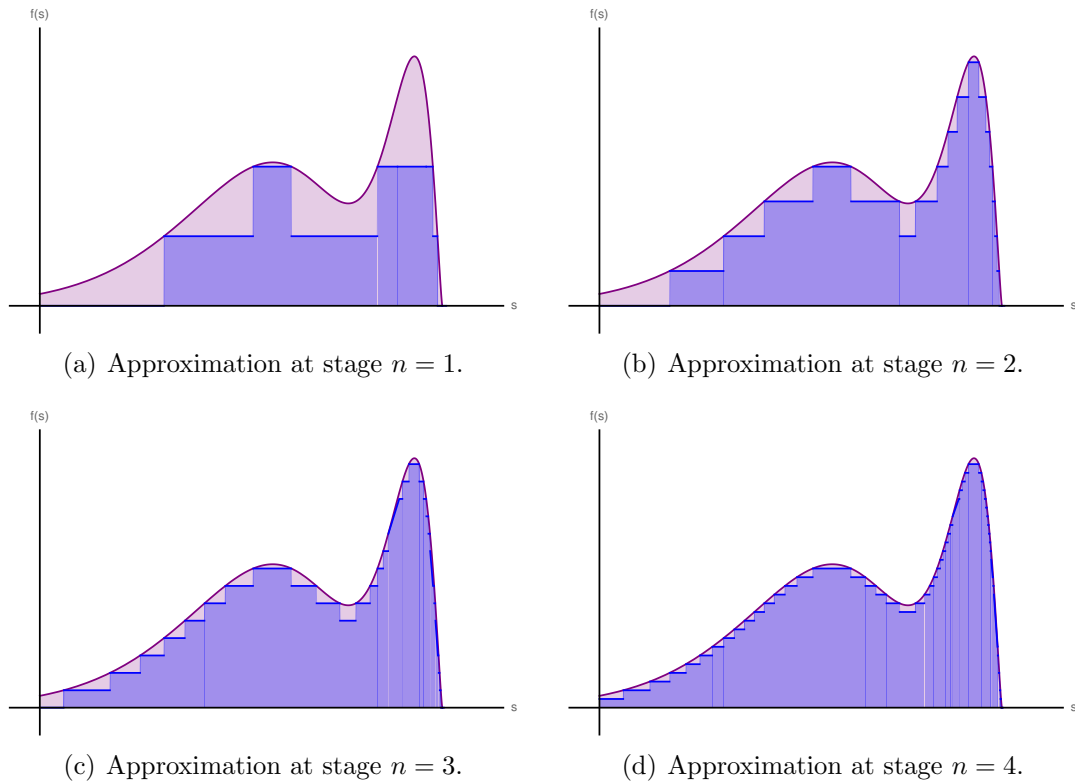


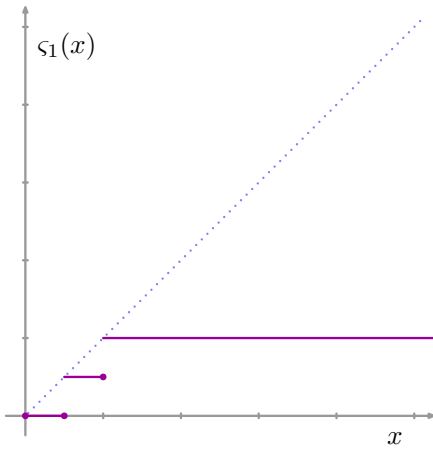
FIGURE III.1. Pointwise increasing approximation of non-negative measurable functions by simple functions.

An illustration of the approximation of a non-negative measurable function is given in Figure III.1. To prove the above approximation lemma, we will construct the approximating sequence explicitly. The idea is that at the  $n$ :th stage of approximation, we truncate the values that exceed level  $n$  to exactly  $n$ , and we replace values below level  $n$  by the nearby values on a grid of mesh  $2^{-n}$ . As  $n \rightarrow \infty$ , the truncation is done ever further away, and the grid becomes ever finer. The truncation and discretization at the  $n$ :th stage are achieved with the *staircase functions*

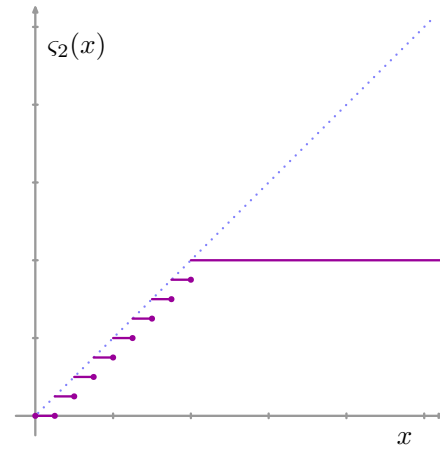
$$\varsigma_n: [0, +\infty] \rightarrow [0, n]$$

illustrated in Figure III.2 and defined piecewise by the formula

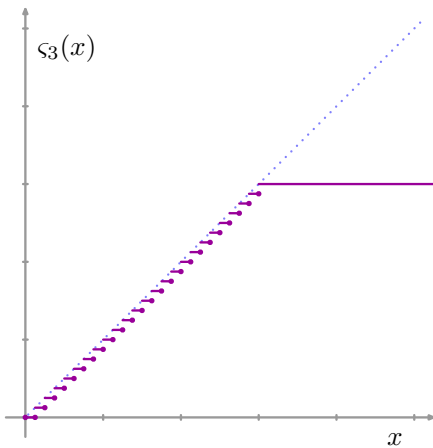
$$\varsigma_n(x) = \begin{cases} 0 & \text{if } 0 \leq x \leq \frac{1}{2^n} \\ \frac{1}{2^n} & \text{if } \frac{1}{2^n} < x \leq \frac{2}{2^n} \\ \frac{2}{2^n} & \text{if } \frac{2}{2^n} < x \leq \frac{3}{2^n} \\ \vdots & \\ \frac{n2^n-1}{2^n} & \text{if } \frac{n2^n-1}{2^n} < x \leq n \\ n & \text{if } n < x. \end{cases} \quad (\text{III.5})$$



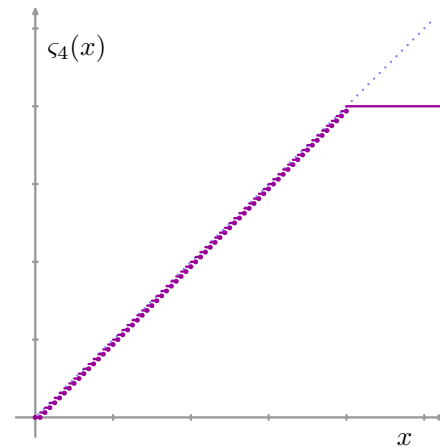
(a) The staircase function  $\varsigma_n$  for  $n = 1$ .



(b) The staircase function  $\varsigma_n$  for  $n = 2$ .



(c) The staircase function  $\varsigma_n$  for  $n = 3$ .



(d) The staircase function  $\varsigma_n$  for  $n = 4$ .

FIGURE III.2. The staircase functions  $\varsigma_n: [0, +\infty] \rightarrow [0, n]$ .

We first check that these staircase functions provide a good approximation of the identity function of  $[0, +\infty]$ .

**Lemma III.19** (Properties of the staircase functions).

- (a) *The staircase functions  $\varsigma_n$  are simple, Borel-measurable, and left-continuous.*
- (b) *For every  $x \in [0, +\infty]$ , we have  $\varsigma_n(x) \uparrow x$  as  $n \rightarrow \infty$ .*

*Proof.* From the definition (III.5) it is clear that the set of possible values of  $\varsigma_n$  is the finite set  $\{j2^{-n} \mid j \in \{0, 1, \dots, n2^{-n}\}\}$ . Each of the preimages  $\varsigma_n^{-1}(\{j2^{-n}\})$  is an interval, therefore a Borel set, so  $\varsigma_n$  is Borel-measurable and simple. Left continuity

$$\lim_{x' \uparrow x} \varsigma_n(x') = \varsigma_n(x) \quad \forall x \in [0, +\infty] \quad (\text{III.6})$$

is also clear from the definition (III.5). This proves part (a).

For part (b), we consider the cases  $x = +\infty$  and  $x \in [0, +\infty)$  separately.

Consider first  $x = +\infty$ . For any  $n \in \mathbb{N}$ , we have  $\varsigma_n(+\infty) = n$  by definition (III.5). These values form an increasing sequence tending to infinity, i.e., we have  $\varsigma_n(+\infty) \uparrow +\infty$  as claimed.

Consider then  $x \in [0, +\infty)$ . For any  $n > x$ , we have  $|\varsigma_n(x) - x| \leq 2^{-n}$  by definition (III.5). This shows that  $\varsigma_n(x) \rightarrow x$  as  $n \rightarrow \infty$ . It is also easy to see that the sequence of values is increasing,  $\varsigma_1(x) \leq \varsigma_2(x) \leq \dots$ . Thus we have  $\varsigma_n(x) \uparrow x$  as claimed. This finishes the proof of part (b).  $\square$

With this, we are ready to prove the approximation lemma.

*Proof of Lemma III.18* Let  $f \in \mathfrak{m}\mathcal{S}^+$  be a non-negative measurable function. For each  $n \in \mathbb{N}$ , define the function  $f_n = \varsigma_n \circ f$ , i.e.,

$$f_n(s) = \varsigma_n(f(s)) \quad \text{for } s \in S.$$

This  $f_n$  is measurable as the composition of the measurable function  $f: S \rightarrow [0, +\infty]$  and the Borel function  $\varsigma_n: [0, +\infty] \rightarrow [0, n]$ . The possible values of  $f_n$  are contained in the finite set of possible values of  $\varsigma_n$ . We conclude that  $f_n$  is a non-negative simple function.

At any  $s \in S$  we have  $f(s) \in [0, +\infty]$ , and therefore

$$f_n(s) = \varsigma_n(f(s)) \uparrow f(s) \quad \text{as } n \rightarrow \infty$$

by part (b) of the previous lemma. We have thus constructed the desired sequence  $f_1, f_2, \dots$  of simple functions, which approximates  $f$  pointwise in a monotone increasing way.  $\square$

**Remark III.20** (Approximating approximations).

For the proof of the approximation lemma itself, it was not really important whether we made the staircase functions left-continuous or not. Where left-continuity is actually convenient is the following situation, which appears in particular in the proof of the Monotone convergence theorem (Theorem VII.8) in Appendix D.

If  $g_1, g_2, \dots \in \mathfrak{m}\mathcal{S}^+$  is any sequence of non-negative functions such that  $g_n \uparrow g$  as  $n \rightarrow \infty$ , then we may construct the simple approximations  $g_n^{(r)} := \varsigma_r \circ g_n$ ,  $r \in \mathbb{N}$ , of each  $g_n$ ,

$$g_n^{(r)} \uparrow g_n \quad \text{as } r \rightarrow \infty$$

as well as the simple approximations  $g^{(r)} := \varsigma_r \circ g$ ,  $r \in \mathbb{N}$ , of the limit function  $g$ ,

$$g^{(r)} \uparrow g \quad \text{as } r \rightarrow \infty.$$

In this setup, the approximations of the limit function  $g$  are the limits of the approximations of each  $g_n$ : for every  $s \in S$  we have by assumption  $g_n(s) \uparrow g(s)$  as  $n \rightarrow \infty$ , and therefore by left-continuity (III.6) of  $\varsigma_n$  we get, for any  $r \in \mathbb{N}$ ,

$$g_n^{(r)}(s) = \varsigma_r(g_n(s)) \uparrow \varsigma_r(g(s)) = g^{(r)}(s) \quad \text{as } n \rightarrow \infty.$$



## Lecture IV

### Information generated by random variables

Probability theory offers an important interpretation of  $\sigma$ -algebras: they describe information.

Let us first mention a few contexts in which this notion of information is used in stochastics. In this course we will first use of the notion in the next Lecture V in relation to independence: the elementary notion of independence of events is generalized to the notion of independence of information. Another common and fruitful use of information (see Appendix E) is conditional expected value, which represents the best estimate of a random number given some (partial) information about it. Finally, stochastic processes are random time-dependent phenomena, and it is often relevant to model how information accumulates as we observe the phenomenon over a period of time — the mathematical notion suitable for this is refining collections of  $\sigma$ -algebras indexed by time known as filtrations. There would be yet other contexts, but it is in fact useful to interpret all  $\sigma$ -algebras as describing information, and relate the notion of measurability of functions to this interpretation. So let us start with an informal description of the ideas and then proceed to precise definition and properties.

To describe the idea informally, suppose that  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space and we have an indexed collection  $(Y_\gamma)_{\gamma \in \Gamma}$  of random variables  $Y_\gamma: \Omega \rightarrow \mathbb{R}$  on it.<sup>1</sup> To understand what information is contained in these random variables, recall the two step procedure by which randomness is thought to arise:

- 1.) “Chance determines the random outcome  $\omega \in \Omega$ .”
- 2.) “The outcome  $\omega$  determines the values  $Y_\gamma(\omega)$  of quantities of interest  $Y_\gamma$ .”

The motivating question about information is then:

“If you do not know the outcome  $\omega$  of all randomness, but someone tells you the values of the quantities of interest  $Y_\gamma$  for all  $\gamma \in \Gamma$ , then for which events  $E \in \mathcal{F}$  are you able to decide whether  $E$  occurs or not?”

Formulated in this way it makes sense that the information contained in the collection  $(Y_\gamma)_{\gamma \in \Gamma}$  of random variables is some collection of events — namely those events whose occurrence can be decided based on the random variables. Evidently, any event of type  $\{Y_\gamma \in A'\}$  concerning the value of any one of the random variables  $Y_\gamma$  can be decided, and thus belongs to the collection. However, in deciding about

---

<sup>1</sup>There is no fundamental reason to require that the random variables are real-valued, but we assume this for the sake of concreteness — and in order to avoid the very awkward notation that would arise if each random variable  $Y_\gamma: \Omega \rightarrow S_\gamma$  in the collection would have a different set  $S_\gamma$  of allowed values (necessarily then also equipped with its own  $\sigma$ -algebra  $\mathcal{S}_\gamma$ ).

events you are furthermore allowed to use logical reasoning (e.g., the logical operations described in Section I.1), so the collection of events that can be decided should itself be a  $\sigma$ -algebra. This motivates the following definition.

#### IV.1. Definition of $\sigma$ -algebra generated by random variables

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and  $(Y_\gamma)_{\gamma \in \Gamma}$  an indexed collection of random variables  $Y_\gamma: \Omega \rightarrow \mathbb{R}$  on it.

**Definition IV.1** (Sigma-algebra generated by random variables).

The  $\sigma$ -algebra generated by the collection  $(Y_\gamma)_{\gamma \in \Gamma}$  of random variables is the smallest  $\sigma$ -algebra  $\mathcal{Y}$  on  $\Omega$  such that for each  $\gamma \in \Gamma$  the random variable  $Y_\gamma$  is  $\mathcal{Y}$ -measurable. We denote the  $\sigma$ -algebra generated by the collection by  $\mathcal{Y} = \sigma((Y_\gamma)_{\gamma \in \Gamma})$ .

**Remark IV.2.** Similarly to Section I.3, the smallest  $\sigma$ -algebra with the above property exists and is unique: it is the intersection of all  $\sigma$ -algebras satisfying the property.

**Remark IV.3.** Since each random variable  $Y_\gamma$  is by definition at least  $\mathcal{F}$ -measurable, we obviously have  $\sigma((Y_\gamma)_{\gamma \in \Gamma}) \subset \mathcal{F}$ . According to the information interpretation,  $\mathcal{F}$  represents “full information” (all events on our probability space), so no amount of random variables could contain more information than that.

**Remark IV.4.** Although we use the notation  $\sigma(\dots)$  both for the  $\sigma$ -algebra generated by a collection of subsets (Definition I.6) and for the  $\sigma$ -algebra generated by a collection of random variables (Definition IV.1), we trust that there is no risk of confusion: it should always be clear from the context which of these two closely related notions is meant.

An already interesting special case is a collection which contains only one random variable: the  $\sigma$ -algebra generated by  $Y: \Omega \rightarrow \mathbb{R}$  is the smallest  $\sigma$ -algebra with respect to which  $Y$  is measurable. We denote it simply by  $\sigma(Y)$ .

**Exercise IV.1** (The  $\sigma$ -algebra generated by a random number).

Let  $Y$  be a real-valued random variable defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ .

- Show that the  $\sigma$ -algebra  $\sigma(Y)$  generated by the random variable  $Y$  coincides with the  $\sigma$ -algebra  $\sigma(Y^{-1}(\mathcal{B}))$  generated by the collection of events  $Y^{-1}(\mathcal{B}) = \{Y^{-1}(B) \mid B \in \mathcal{B}\}$ .
- Show that we in fact have the equality  $\sigma(Y) = Y^{-1}(\mathcal{B})$ .

**Exercise IV.2** (A  $\pi$ -system to generate the  $\sigma$ -algebra generated by a random number).

Let  $Y$  be a real-valued random variable defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . Let  $\mathcal{J}(\mathbb{R}) = \{(-\infty, x] \mid x \in \mathbb{R}\}$  be the  $\pi$ -system on  $\mathbb{R}$  as in Example II.25 and define

$$\mathcal{J} = Y^{-1}(\mathcal{J}(\mathbb{R})) = \left\{ Y^{-1}((-\infty, x]) \mid x \in \mathbb{R} \right\}.$$

Show that  $\mathcal{J}$  is a  $\pi$ -system on  $\Omega$  which generates the  $\sigma$ -algebra  $\sigma(Y)$  generated by the random variable  $Y$ , i.e.,  $\sigma(\mathcal{J}) = \sigma(Y)$ .

## IV.2. Doob's representation theorem

The definition of information contained in random variables may seem abstract. We next state and prove a theorem, which offers a good interpretation for the information contained in one real-valued random variable.

**Theorem IV.5** (Doob's representation theorem).

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and  $Y: \Omega \rightarrow \mathbb{R}$  and  $Z: \Omega \rightarrow \mathbb{R}$  two real valued random variables on it. Then  $Z$  is  $\sigma(Y)$ -measurable<sup>2</sup> if and only if there exists a Borel function  $f: \mathbb{R} \rightarrow \mathbb{R}$  such that

$$Z = f(Y).$$

**Remark IV.6.** The precise meaning of

$$Z = f(Y)$$

is that for all possible outcomes  $\omega \in \Omega$  we have

$$Z(\omega) = f(Y(\omega)).$$

In other words, keeping in mind that random variables are ultimately just functions on the sample space  $\Omega$ , the function  $Z: \Omega \rightarrow \mathbb{R}$  is the composition

$$Z = f \circ Y$$

of functions  $Y: \Omega \rightarrow \mathbb{R}$  and  $f: \mathbb{R} \rightarrow \mathbb{R}$

$$\Omega \xrightarrow{\begin{array}{c} Y \\ \searrow Z=f \circ Y \end{array}} \mathbb{R} \xrightarrow{f} \mathbb{R}.$$

Theorem IV.5 gives the following interpretation: to say that  $Z$  is measurable with respect to the information contained in  $Y$  means that the value of  $Z$  could be obtained from the value of  $Y$  by applying some deterministic function  $f$ .

In the proof we use the Monotone class theorem, which can be found in Appendix C.

*Proof of Theorem IV.5.* The “if” direction of the statement is easy. Namely, if we have  $Z = f \circ Y$  for some  $\mathcal{B}/\mathcal{B}$ -measurable function  $f: \mathbb{R} \rightarrow \mathbb{R}$ , then since  $Y: \Omega \rightarrow \mathbb{R}$  is by definition  $\sigma(Y)/\mathcal{B}$ -measurable, it follows from Proposition III.8 that the composition  $Z = f \circ Y: \Omega \rightarrow \mathbb{R}$  is  $\sigma(Y)/\mathcal{B}$ -measurable. It therefore remains to prove the “only if” direction. We first prove the “only if” direction assuming that  $Z$  is bounded, and afterwards extend to full generality.

*bounded  $Z$ :* Let us define  $\mathcal{H}$  as the collection of all those bounded functions  $Z: \Omega \rightarrow \mathbb{R}$  which can be written as  $Z = f \circ Y$  for some bounded Borel function  $f$ . Our goal is to show that the collection  $\mathcal{H}$  contains all  $\sigma(Y)$ -measurable bounded functions.

Denote  $\mathcal{J} = \sigma(Y)$ . We start by checking that at least the indicator  $\mathbb{I}_E$  of any event  $E \in \mathcal{J}$  belongs to  $\mathcal{H}$ . Recall the fact  $\mathcal{J} = \sigma(Y) = \{Y^{-1}(B) \mid B \in \mathcal{B}\}$  from Exercise IV.1: any  $E \in \mathcal{J}$  is of the form  $E = Y^{-1}(B)$  for some  $B \in \mathcal{B}$ . Therefore the indicator of  $E$  is the function

$$\mathbb{I}_E(\omega) = \mathbb{I}_{Y^{-1}(B)}(\omega) = \begin{cases} 1 & \text{if } \omega \in Y^{-1}(B) \\ 0 & \text{otherwise.} \end{cases} \quad (\text{IV.1})$$

In contrast, the indicator function of  $B \subset \mathbb{R}$  is a function defined on  $\mathbb{R}$ : in fact,  $\mathbb{I}_B: \mathbb{R} \rightarrow \mathbb{R}$  is a bounded Borel function. We notice that

$$\mathbb{I}_B(Y(\omega)) = \begin{cases} 1 & \text{if } Y(\omega) \in B \\ 0 & \text{otherwise.} \end{cases} \quad (\text{IV.2})$$

<sup>2</sup>More precisely,  $\sigma(Y)/\mathcal{B}$ -measurable.

Comparing (IV.1) and (IV.2) we see that

$$\mathbb{I}_E(\omega) = \mathbb{I}_B(Y(\omega)).$$

This conclusion  $\mathbb{I}_E = \mathbb{I}_B \circ Y$  shows that the indicator of any  $E \in \mathcal{G}$  belongs to the collection we are considering,  $\mathbb{I}_E \in \mathcal{H}$ .

We now show that  $\mathcal{H}$  is a *monotone class* as defined in Appendix C (see Definition C.1), i.e., it satisfies the three properties (MC-1), (MC- $\mathbb{R}$ ), and (MC- $\uparrow$ ).

The first of these properties is obvious: if we take  $f$  to be the constant function  $f(x) = 1$  for all  $x \in \mathbb{R}$  then  $Z = f \circ Y$  is the constant random variable  $Z(\omega) = f(Y(\omega)) = 1$  for all  $\omega \in \Omega$ . Thus indeed the constant function 1 on  $\Omega$  belongs to  $\mathcal{H}$ . This is property (MC-1) for  $\mathcal{H}$ .

The second property is easy, too. If we have  $Z_1, Z_2 \in \mathcal{H}$ , then we can write  $Z_1 = f_1 \circ Y$  and  $Z_2 = f_2 \circ Y$  for some bounded Borel-measurable functions  $f_1, f_2: \mathbb{R} \rightarrow \mathbb{R}$ . Then for  $c_1, c_2 \in \mathbb{R}$  we have  $c_1 Z_1 + c_2 Z_2 = f \circ Y$ , where  $f = c_1 f_1 + c_2 f_2$  is a pointwise linear combination function  $\mathbb{R} \rightarrow \mathbb{R}$ , which is also bounded and Borel-measurable by Proposition III.13. This is property (MC- $\mathbb{R}$ ) for  $\mathcal{H}$ .

The last property is checked as follows. Suppose that  $Z_n \uparrow Z$  as  $n \rightarrow \infty$ , and that  $Z_n \in \mathcal{H}$  for all  $n \in \mathbb{N}$ , and that  $0 \leq Z_n(\omega) \leq K$  for all  $\omega \in \Omega$  and some constant  $K < \infty$ . Then we have  $Z_n = f_n \circ Y$  for some bounded Borel-measurable functions  $f_n: \mathbb{R} \rightarrow \mathbb{R}$ . We may assume that  $0 \leq f_n \leq K$  pointwise<sup>3</sup>. Now define  $f = \limsup_n f_n$ . Then  $f: \mathbb{R} \rightarrow \mathbb{R}$  is Borel measurable by Proposition III.14, and it is also bounded:  $0 \leq f \leq K$ . Moreover, from our assumptions it now follows that

$$Z = \lim_{n \rightarrow \infty} Z_n = \lim_{n \rightarrow \infty} f_n \circ Y = \limsup_n f_n \circ Y = f \circ Y.$$

We thus obtain that  $Z \in \mathcal{H}$ . This is property (MC- $\uparrow$ ) for  $\mathcal{H}$ .

We have shown that  $\mathcal{H}$  is a monotone class which contains the indicator functions of all sets  $E$  in the collection  $\mathcal{G} = \sigma(Y)$ . The collection  $\mathcal{G} = \sigma(Y)$  is a  $\sigma$ -algebra and thus in particular a  $\pi$ -system (Remark II.24). The Monotone class theorem (Theorem C.2) therefore guarantees that  $\mathcal{H}$  contains all bounded  $\sigma(Y)$ -measurable functions. In other words, every bounded  $\sigma(Y)$ -measurable  $Z: \Omega \rightarrow \mathbb{R}$  is of the form  $Z = f \circ Y$  for some bounded Borel function  $f: \mathbb{R} \rightarrow \mathbb{R}$ .

*unbounded  $Z$ :* Consider now the general case, where  $Z: \Omega \rightarrow \mathbb{R}$  can be unbounded. In that case we can apply the function

$$\arctan: \mathbb{R} \rightarrow \left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$$

and get the bounded random variable

$$\tilde{Z} = \arctan \circ Z.$$

By the composition property, if  $Z$  is  $\sigma(Y)$ -measurable, then  $\tilde{Z}$  is also ( $\arctan$  is continuous and therefore Borel). From the first step of the proof, we thus know that  $\tilde{Z} = \tilde{f} \circ Y$  for some bounded Borel function  $\tilde{f}: \mathbb{R} \rightarrow \mathbb{R}$ . Then we have

$$Z = \tan \circ \tilde{Z} = \tan \circ \tilde{f} \circ Y = f \circ Y,$$

where  $f = \tan \circ \tilde{f}$ . This  $f$  is Borel measurable and the proof is complete.  $\square$

<sup>3</sup>If  $f_n$  do not already satisfy this property, we may truncate them. Using the functions

$$\tilde{f}_n(x) := \begin{cases} 0 & \text{when } f_n(x) < 0 \\ f_n(x) & \text{when } 0 \leq f_n(x) \leq K \\ K & \text{when } K < f_n(x) \end{cases}$$

instead, we still have  $Z_n = \tilde{f}_n \circ Y$ .



## Lecture V

### Independence

In the previous lecture (Lecture IV) we saw that probability theory offers an important interpretation of  $\sigma$ -algebras: they describe information. One of the first uses of the notion of information pertains to probabilistic independence: we will generalize the elementary notion of independence of events to the notion of independence of information.

The intuitive interpretation of probabilistic independence should be familiar from basic courses in probability and statistics, so we recall the idea only very briefly. Suppose that  $A$  and  $B$  are two events, and assume also that  $\mathbb{P}[B] > 0$  so that the conditional probability  $\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]}$  of the event  $A$  given the occurrence of  $B$  can be defined. If the probabilities of the two events satisfy

$$\mathbb{P}[A \cap B] = \mathbb{P}[A] \mathbb{P}[B],$$

then we get

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} = \frac{\mathbb{P}[A] \mathbb{P}[B]}{\mathbb{P}[B]} = \mathbb{P}[A],$$

i.e., the conditional probability of  $A$  given  $B$  is just the probability of  $A$ . We interpret this as saying that the knowledge of the occurrence of  $B$  does not reveal anything that could be used to improve our estimate about the occurrence of  $A$ , and we therefore consider the event  $A$  independent of the event  $B$ .<sup>1</sup>

We start this lecture by introducing the abstract and general notion of probabilistic independence, and we then show its relation to the more familiar elementary notion which was also used in the interpretation above. After the definitions and basic properties, we also discuss the first profound techniques related to independence: the Borel-Cantelli lemmas.

#### V.1. Definition of independence

Throughout, let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space.

Let  $(\mathcal{G}_j)_{j \in J}$  be a collection of  $\sigma$ -algebras,  $\mathcal{G}_j \subset \mathcal{F}$  for all  $j \in J$ . For each  $j$ , we think of  $\mathcal{G}_j$  representing some information (for example available to a person  $j$ ). Whether these informations, for different  $j$ , are independent of each other with respect to the underlying probability  $\mathbb{P}$  is captured by the following definition.

**Definition V.1** (Independence of sigma-algebras).

The collection  $(\mathcal{G}_j)_{j \in J}$  of  $\sigma$ -algebras is *independent* if for any distinct  $j_1, \dots, j_n \in$

---

<sup>1</sup> Since the condition  $\mathbb{P}[A \cap B] = \mathbb{P}[A] \mathbb{P}[B]$  is symmetric under interchange of  $A$  and  $B$ , we also consider the event  $B$  independent of  $A$ . It is, in fact, better to use symmetric terminology and say that events  $A$  and  $B$  are independent.

$J$  and any events  $A_{j_1} \in \mathcal{G}_{j_1}, \dots, A_{j_n} \in \mathcal{G}_{j_n}$  we have

$$\mathbb{P}[A_{j_1} \cap A_{j_2} \cap \dots \cap A_{j_n}] = \mathbb{P}[A_{j_1}] \mathbb{P}[A_{j_2}] \cdots \mathbb{P}[A_{j_n}]. \quad (\text{V.1})$$

This is the abstract and general notion of probabilistic independence, on which other notions of independence are founded. It is important to note that while the collection of  $\sigma$ -algebras here may be infinite, the formula (V.1) is only ever required for intersections of finitely many events.

To get from the abstract and general notion to a more concrete case, let now  $(X_j)_{j \in J}$  be a collection of random variables. We need to define what it means for these random variables to be independent. For this purpose, we consider the  $\sigma$ -algebra  $\sigma(X_j)$  generated by  $X_j$ , for each  $j \in J$  separately (see Definition IV.1).

**Definition V.2** (Independence of random variables).

The collection  $(X_j)_{j \in J}$  of random variables is *independent* if the collection  $(\sigma(X_j))_{j \in J}$  of  $\sigma$ -algebras generated by them is independent in the sense of Definition V.1.

Finally, let  $(E_j)_{j \in J}$  be a collection of events,  $E_j \in \mathcal{F}$ . To define what it means for these events to be independent, we consider the indicator random variables  $\mathbb{I}_{E_j}: \Omega \rightarrow \mathbb{R}$  of the events  $E_j$ , for  $j \in J$  (see Equation (III.3)).

**Definition V.3** (Independence of events).

The collection  $(E_j)_{j \in J}$  of events is *independent* if the collection  $(\mathbb{I}_{E_j})_{j \in J}$  of the corresponding indicator random variables is independent in the sense of Definition V.2.

**Remark V.4** (The elementary notion of independence of events).

The  $\sigma$ -algebra generated by the indicator random variable  $\mathbb{I}_E: \Omega \rightarrow \mathbb{R}$  of an event  $E \subset \Omega$  is

$$\sigma(\mathbb{I}_E) = \{\emptyset, E, E^c, \Omega\}.$$

Thus events  $(E_j)_{j \in J}$  are independent if and only if the  $\sigma$ -algebras  $(\{\emptyset, E_j, E_j^c, \Omega\})_{j \in J}$  are independent. In view of Definition V.1, this amounts to the equalities

$$\mathbb{P}[E_{j_1}^* \cap E_{j_2}^* \cap \dots \cap E_{j_n}^*] = \mathbb{P}[E_{j_1}^*] \mathbb{P}[E_{j_2}^*] \cdots \mathbb{P}[E_{j_n}^*], \quad (\text{V.2})$$

where each  $E_{j_k}^*$  stands for either  $E_{j_k}$  or its complement  $E_{j_k}^c$  (note that including impossible events  $\emptyset$  in the intersection is unnecessary, since both sides of the equation would then vanish automatically, and including sure events  $\Omega$  just amounts to having fewer terms in both the intersection and the product on the two sides of the equation).

For the sake of concreteness, consider now just two events,  $E_1$  and  $E_2$ . If  $E_1$  and  $E_2$  are independent, then as a special case of (V.2) we at least have the familiar defining equality

$$\mathbb{P}[E_1 \cap E_2] = \mathbb{P}[E_1] \mathbb{P}[E_2]. \quad (\text{V.3})$$

But also conversely, if (V.3) holds then one can derive the equations (V.2) involving possible complements — the reader is invited to directly derive at least the equality  $\mathbb{P}[E_1 \cap E_2^c] = \mathbb{P}[E_1] \mathbb{P}[E_2^c]$  from (V.3). Of course, case by case checking the equations (V.2) would be impractical as well as inelegant, so we instead develop systematic tools in Section V.2 below.

**Exercise V.1.** Show that the equality  $\mathbb{P}[E_1 \cap E_2^c] = \mathbb{P}[E_1] \mathbb{P}[E_2^c]$  follows from (V.3).

**Exercise V.2** (Independence is preserved under measurable functions).

Let  $X$  and  $Y$  be two  $\mathbb{R}$ -valued random variables defined on a common probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ .

- (a) Show that if  $X$  and  $Y$  are independent, then  $f(X)$  and  $g(Y)$  are independent for any Borel-measurable functions  $f, g: \mathbb{R} \rightarrow \mathbb{R}$ .
- (b) Show that if  $f(X)$  and  $g(Y)$  are independent for all Borel-measurable functions  $f, g: \mathbb{R} \rightarrow \mathbb{R}$ , then  $X$  and  $Y$  are independent.

**Exercise V.3.** Assume that  $X$  and  $Y$  are two  $\mathbb{R}$ -valued random variables such that we have  $\mathbb{P}[X + Y = 42] = 1$ . Is it possible that  $X$  and  $Y$  are independent?

**Exercise V.4** (Calculations with two independent geometrically distributed numbers).

Let  $X, Y: \Omega \rightarrow \mathbb{N}$  be two independent random variables with

$$\mathbb{P}[X = j] = \mathbb{P}[Y = j] = \frac{1}{2^j} \quad \text{for all } j \in \mathbb{N} = \{1, 2, \dots\}.$$

- (a) Show that  $\mathbb{P}[Y > n] = \frac{1}{2^n}$  for any  $n \in \mathbb{N}$ .

Calculate the following probabilities

- (b):  $\mathbb{P}[X = Y]$
- (c):  $\mathbb{P}[\min(X, Y) \leq k]$ , where  $k \in \mathbb{N}$
- (d):  $\mathbb{P}[Y > X]$
- (e):  $\mathbb{P}[X > kY]$ , where  $k \in \mathbb{N}$
- (f):  $\mathbb{P}[X \text{ divides } Y]$

**Hint:** The correct final results are among the following:

$$\frac{1}{3}, \quad \frac{1}{2^{k+1} - 1}, \quad 1 - \frac{1}{4^k}, \quad \sum_{k=1}^{\infty} \frac{1}{2^{k+1} - 1} = \sum_{k=1}^{\infty} \frac{1}{4^k - 2^k}.$$

### Notation

We abbreviate independence by the symbol  $\perp$ . We thus denote:

- $(\mathcal{G}_j)_{j \in J} \perp$  if the collection  $(\mathcal{G}_j)_{j \in J}$  of  $\sigma$ -algebras is independent
- $(X_j)_{j \in J} \perp$  if the collection  $(X_j)_{j \in J}$  of random variables is independent
- $(E_j)_{j \in J} \perp$  if the collection  $(E_j)_{j \in J}$  of events is independent.

In the case of enumerated (countable) collections, we use the notation:

- $\mathcal{G}_1, \mathcal{G}_2, \dots \perp$  if the collection  $(\mathcal{G}_j)_{j \in \mathbb{N}}$  of  $\sigma$ -algebras is independent
- $X_1, X_2, \dots \perp$  if the collection  $(X_j)_{j \in \mathbb{N}}$  of random variables is independent
- $E_1, E_2, \dots \perp$  if the collection  $(E_j)_{j \in \mathbb{N}}$  of events is independent.

In the case of collections of just two members, we use notation:

- $\mathcal{G}_1 \perp \mathcal{G}_2$  if the collection  $(\mathcal{G}_j)_{j \in \{1,2\}}$  of  $\sigma$ -algebras is independent
- $X_1 \perp X_2$  if the collection  $(X_j)_{j \in \{1,2\}}$  of random variables is independent
- $E_1 \perp E_2$  if the collection  $(E_j)_{j \in \{1,2\}}$  of events is independent.

**Exercise V.5** (Pairwise independence does not imply independence).

Construct an example in which three  $\sigma$ -algebras  $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3$  are pairwise independent,

$$\mathcal{G}_1 \perp \mathcal{G}_2, \quad \mathcal{G}_1 \perp \mathcal{G}_3, \quad \mathcal{G}_2 \perp \mathcal{G}_3,$$

but the collection  $(\mathcal{G}_j)_{j \in \{1,2,3\}}$  of all three is not independent.

## V.2. Verifying independence

To check independence, one usually does not want to work directly with the general definition, but rather use an easier sufficient condition.

**Proposition V.5** (A sufficient condition for independence).

Suppose that  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are two  $\pi$ -systems on  $\Omega$  and let  $\mathcal{G}_1 = \sigma(\mathcal{I}_1)$  and  $\mathcal{G}_2 = \sigma(\mathcal{I}_2)$  be the  $\sigma$ -algebras generated by them. Then the following conditions are equivalent:

- (i)  $\mathcal{G}_1 \perp \mathcal{G}_2$
- (ii) For all  $I_1 \in \mathcal{I}_1$  and  $I_2 \in \mathcal{I}_2$ , we have  $\mathbb{P}[I_1 \cap I_2] = \mathbb{P}[I_1] \mathbb{P}[I_2]$ .

*Proof.* Clearly condition (i) implies (ii), because  $\mathcal{I}_1 \subset \sigma(\mathcal{I}_1) = \mathcal{G}_1$  and  $\mathcal{I}_2 \subset \sigma(\mathcal{I}_2) = \mathcal{G}_2$ .  $\sigma$ -algebras It remains to show that (ii) implies (i). Let us therefore assume (ii). According to Definition V.1, we must prove that then

$$\mathbb{P}[E_1 \cap E_2] = \mathbb{P}[E_1] \mathbb{P}[E_2]$$

holds for all  $E_1 \in \mathcal{G}_1$ ,  $E_2 \in \mathcal{G}_2$ . We do this in two steps: first assuming that  $E_1$  is in the  $\pi$ -system  $\mathcal{I}_1$ , and then for a general  $E_1 \in \mathcal{G}_1$ .

*step 1:* Let  $I_1 \in \mathcal{I}_1$ . In this first step we seek to prove that

$$\mathbb{P}[I_1 \cap E_2] = \mathbb{P}[I_1] \mathbb{P}[E_2] \tag{V.4}$$

holds for all  $E_2 \in \mathcal{G}_2$ . If  $\mathbb{P}[I_1] = 0$ , then equations (V.4) hold trivially, because both sides vanish. We may therefore assume that  $\mathbb{P}[I_1] > 0$ . Then we define a new probability measure  $\tilde{\mathbb{P}}_2$  on  $(\Omega, \mathcal{G}_2)$  by the formula

$$\tilde{\mathbb{P}}_2[E_2] = \frac{\mathbb{P}[I_1 \cap E_2]}{\mathbb{P}[I_1]} \quad \text{for } E_2 \in \mathcal{G}_2 \tag{V.5}$$

(this is indeed a probability measure by Exercise II.1). From assumption (ii), it follows that the two probability measures  $\mathbb{P}$  and  $\tilde{\mathbb{P}}_2$  coincide on the  $\pi$ -system  $\mathcal{I}_2$ : if  $I_2 \in \mathcal{I}_2$  then

$$\tilde{\mathbb{P}}_2[I_2] = \frac{\mathbb{P}[I_1 \cap I_2]}{\mathbb{P}[I_1]} \stackrel{\text{(ii)}}{=} \frac{\mathbb{P}[I_1] \mathbb{P}[I_2]}{\mathbb{P}[I_1]} = \mathbb{P}[I_2].$$

Dynkin's identification theorem (Theorem II.26) thus implies that they coincide on the entire  $\sigma$ -algebra  $\mathcal{G}_2$  generated by this  $\pi$ -system, i.e.  $\tilde{\mathbb{P}}_2[E_2] = \mathbb{P}[E_2]$  for all  $E_2 \in \mathcal{G}_2$ . This shows that

$$\mathbb{P}[E_2] = \tilde{\mathbb{P}}_2[E_2] = \frac{\mathbb{P}[I_1 \cap E_2]}{\mathbb{P}[I_1]},$$

which upon multiplying by  $\mathbb{P}[I_1]$  gives the desired equality (V.4).

*step 2:* Let  $E_2 \in \mathcal{G}_2$ . In this last step we seek to prove, with a method analogous to step 1, that

$$\mathbb{P}[E_1 \cap E_2] = \mathbb{P}[E_1] \mathbb{P}[E_2] \tag{V.6}$$

holds for all  $E_1 \in \mathcal{G}_1$ . This will conclude the proof. We may assume that  $\mathbb{P}[E_2] > 0$ , because otherwise both sides of (V.6) vanish. Then we again define a new probability measure  $\tilde{\mathbb{P}}_1$ , this time on  $(\Omega, \mathcal{G}_1)$ , by the formula

$$\tilde{\mathbb{P}}_1[E_1] = \frac{\mathbb{P}[E_1 \cap E_2]}{\mathbb{P}[E_2]} \quad \text{for } E_1 \in \mathcal{G}_1. \tag{V.7}$$

By step 1, the two probability measures  $\mathbb{P}$  and  $\tilde{\mathbb{P}}_1$  coincide on the  $\pi$ -system  $\mathcal{I}_1$ . From Dynkin's identification theorem (Theorem II.26) it thus follows that they coincide on the  $\sigma$ -algebra  $\mathcal{G}_1$  generated by this  $\pi$ -system. This shows that for all  $E_1 \in \mathcal{G}_1$  we have

$$\mathbb{P}[E_1] = \tilde{\mathbb{P}}_1[E_1] = \frac{\mathbb{P}[E_1 \cap E_2]}{\mathbb{P}[E_2]}.$$

Upon multiplying by  $\mathbb{P}[E_2]$  we get the desired equality (V.6), which shows that the  $\sigma$ -algebras  $\mathcal{G}_1$  and  $\mathcal{G}_2$  are independent. This establishes the implication (ii)  $\Rightarrow$  (i), and concludes the proof.  $\square$

A typical example application of this is the following (familiar and) practical criterion for the independence of two random numbers using their cumulative distribution functions.

**Corollary V.6** (Independence using cumulative distribution functions).

Suppose that  $X_1, X_2: \Omega \rightarrow \mathbb{R}$  are two random variables. Then we have  $X_1 \perp\!\!\!\perp X_2$  if and only if

$$\mathbb{P}[X_1 \leq x_1, X_2 \leq x_2] = \mathbb{P}[X_1 \leq x_1] \mathbb{P}[X_2 \leq x_2] \quad (\text{V.8})$$

for all  $x_1, x_2 \in \mathbb{R}$ .

*Proof.* Let  $\mathcal{J}(\mathbb{R}) = \{(-\infty, x] \mid x \in \mathbb{R}\}$  be the  $\pi$ -system on  $\mathbb{R}$  as in Example II.25. Let

$$\mathcal{I}_1 = X_1^{-1}(\mathcal{J}(\mathbb{R})) = \left\{ \left\{ \omega \in \Omega \mid X_1(\omega) \leq x \right\} \mid x \in \mathbb{R} \right\}.$$

Then  $\mathcal{I}_1$  is a  $\pi$ -system on  $\Omega$  which generates  $\sigma(X_1)$  (see Exercise IV.2). Similarly  $\mathcal{I}_2 = X_2^{-1}(\mathcal{J}(\mathbb{R}))$  is a  $\pi$ -system on  $\Omega$  which generates  $\sigma(X_2)$ . The assumption (V.8) exactly says that  $\mathbb{P}[A_1 \cap A_2] = \mathbb{P}[A_1] \mathbb{P}[A_2]$  for all  $A_1 \in \mathcal{I}_1$  and  $A_2 \in \mathcal{I}_2$ . The statement therefore follows from Proposition V.5  $\square$

Similar statements hold for finite collections of  $\pi$ -systems and finite collections of random numbers. To indicate how to deal with more than two  $\sigma$ -algebras, consider the following exercise.

**Exercise V.6** (Independence of three sigma algebras).

Let  $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3 \subset \mathcal{F}$  be sigma-algebras on  $\Omega$ . Assume that  $\mathcal{G}_k$  is generated by a  $\pi$ -system  $\mathcal{I}_k$  which contains  $\Omega$ .

- (a) Show that  $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3$  are independent if and only if

$$\mathbb{P}[I_1 \cap I_2 \cap I_3] = \mathbb{P}[I_1] \mathbb{P}[I_2] \mathbb{P}[I_3]$$

for all  $I_1 \in \mathcal{I}_1, I_2 \in \mathcal{I}_2, I_3 \in \mathcal{I}_3$ .

- (b) Why did we here require that  $\mathcal{I}_k$  contains  $\Omega$ ?

**Hint:** Consider, e.g., the finite set  $\Omega = \{1, 2, \dots, 8\}$  and the discrete uniform probability measure  $\mathbb{P}$  on it, and three  $\pi$ -systems  $\mathcal{I}_1, \mathcal{I}_2, \mathcal{I}_3$  each consisting of a single event, suitably chosen.

### V.3. Borel – Cantelli lemmas

Given a sequence  $E_1, E_2, \dots \in \mathcal{F}$  of events we occasionally care about whether infinitely many of them occur, or whether only finitely many of them occur. Borel – Cantelli lemmas are examples of what are called 0-1 laws in probability theory — they state that under some mild conditions that are often relatively simple to verify, the probability that we are interested in is trivial: either 0 or 1.

#### Occurrence infinitely often and occurrence eventually

We use the following definitions for a sequence  $E_1, E_2, \dots \in \mathcal{F}$  of events:

- “ $E_n$  occurs infinitely often”, abbreviated “ $E_n$  i.o.”, is the event

$$\begin{aligned} \text{“}E_n \text{ i.o.} \text{”} &= \bigcap_{m \in \mathbb{N}} \bigcup_{n \geq m} E_n = \limsup_n E_n \\ &= \{ \omega \in \Omega \mid \omega \in E_n \text{ for infinitely many indices } n \}. \end{aligned}$$

- “ $E_n$  occurs eventually”, abbreviated “ $E_n$  ev.”, is the event

$$\begin{aligned} \text{“}E_n \text{ ev.} \text{”} &= \bigcup_{m \in \mathbb{N}} \bigcap_{n \geq m} E_n = \liminf_n E_n \\ &= \{ \omega \in \Omega \mid \omega \in E_n \text{ for all except finitely many } n \}. \end{aligned}$$

These two are related by taking complements: by De Morgan’s laws, we have

$$\begin{aligned} \left( \bigcap_{m \in \mathbb{N}} \bigcup_{n \geq m} E_n \right)^c &= \bigcup_{m \in \mathbb{N}} \bigcap_{n \geq m} E_n^c, \quad \text{i.e.,} \\ \left( \text{“}E_n \text{ occurs infinitely often} \text{”} \right)^c &= \text{“}E_n^c \text{ occurs eventually} \text{”}. \end{aligned}$$

**Exercise V.7** (Indicators and upper and lower limits).

Consider a sequence  $E_1, E_2, \dots \subset \Omega$  of subsets. Show that for all  $\omega \in \Omega$  we have

$$\limsup_n \mathbb{I}_{E_n}(\omega) = \mathbb{I}_{\limsup_n E_n}(\omega) \quad \text{and} \quad \liminf_n \mathbb{I}_{E_n}(\omega) = \mathbb{I}_{\liminf_n E_n}(\omega).$$

Use this to show that

$$\liminf_n E_n \subset \limsup_n E_n.$$

## The two Borel – Cantelli lemmas

The first Borel – Cantelli lemma says that whenever the probabilities of the events  $E_n$  decay fast enough, it is (almost surely) impossible for the events to occur infinitely often.

**Lemma V.7** (Borel–Cantelli lemma: convergence part).

Suppose that  $E_1, E_2, \dots \in \mathcal{F}$  are such that  $\sum_{n=1}^{\infty} \mathbb{P}[E_n] < +\infty$ . Then we have

$$\mathbb{P}[\text{“}E_n \text{ occurs infinitely often} \text{”}] = 0.$$

*Proof.* Denote  $G_m = \bigcup_{n=m}^{\infty} E_n$ , so that the event  $E = \limsup_n E_n$  of interest is the decreasing limit  $G_m \downarrow E$  as  $m \rightarrow \infty$ . By monotone convergence for probability measures, Theorem II.22, we have

$$\mathbb{P}[E] = \lim_{m \rightarrow \infty} \mathbb{P}[G_m].$$

Now we can use the union bound, Theorem II.20, to estimate

$$0 \leq \mathbb{P}[G_m] = \mathbb{P}\left[ \bigcup_{n=m}^{\infty} E_n \right] \leq \sum_{n=m}^{\infty} \mathbb{P}[E_n].$$

Since the series  $\sum_n \mathbb{P}[E_n]$  is convergent, its tail goes to zero:  $\sum_{n=m}^{\infty} \mathbb{P}[E_n] \rightarrow 0$  as  $m \rightarrow \infty$ . This shows that  $\lim_{m \rightarrow \infty} \mathbb{P}[G_m] = 0$  and thus  $\mathbb{P}[E] = 0$ .  $\square$

The second Borel–Cantelli lemma says that if the probabilities of the events  $E_n$  do not decay fast and if the events are in addition independent, then the events must (almost surely) occur infinitely often.

**Lemma V.8** (Borel – Cantelli lemma: divergence part).

Suppose that  $E_1, E_2, \dots \in \mathcal{F}$  are independent and  $\sum_{n=1}^{\infty} \mathbb{P}[E_n] = +\infty$ . Then we have

$$\mathbb{P}\left[“E_n \text{ occurs infinitely often}”\right] = 1.$$

*Proof.* Instead of the event  $E = \bigcap_{m \in \mathbb{N}} \bigcup_{n \geq m} E_n$  of interest, consider first its complement  $E^c = \bigcup_{m \in \mathbb{N}} \bigcap_{n \geq m} E_n^c$ . Denote the members in this union by  $F_m = \bigcap_{n \geq m} E_n^c$ . The event  $F_m$  is an intersection — but in order to use independence we first need to arrange things to finite intersections. To achieve this, define  $F_{m,\ell} = \bigcap_{n=m}^{\ell} E_n^c$  and note that these finite intersections decrease to the infinite intersection  $F_{m,\ell} \downarrow F_m$  as  $\ell \rightarrow \infty$ . Now by independence we have

$$\mathbb{P}[F_{m,\ell}] = \mathbb{P}\left[\bigcap_{n=m}^{\ell} E_n^c\right] = \prod_{n=m}^{\ell} \mathbb{P}[E_n^c].$$

Denote  $p_n = \mathbb{P}[E_n]$  and use the estimate  $1 - p_n \leq e^{-p_n}$  to get

$$\mathbb{P}[F_{m,\ell}] = \prod_{n=m}^{\ell} (1 - p_n) \leq \prod_{n=m}^{\ell} e^{-p_n} = \exp\left(-\sum_{n=m}^{\ell} p_n\right).$$

By the assumption  $\sum_n p_n = +\infty$  we have  $\exp\left(-\sum_{n=m}^{\ell} p_n\right) \rightarrow 0$  as  $\ell \rightarrow \infty$ . Therefore from  $F_{m,\ell} \downarrow F_m$  and monotone convergence for probability measures, Theorem II.22, we get

$$\mathbb{P}[F_m] = \lim_{\ell \rightarrow \infty} \mathbb{P}[F_{m,\ell}] = 0.$$

Now for the complement  $E^c = \bigcup_m F_m$  use the union bound, Theorem II.20, to get

$$\mathbb{P}[E^c] = \mathbb{P}\left[\bigcup_{m=1}^{\infty} F_m\right] \leq \sum_{m=1}^{\infty} \mathbb{P}[F_m] = \sum_{m=1}^{\infty} 0 = 0.$$

This result for the complement is what we wanted to show about  $E$ :

$$\mathbb{P}[E] = 1 - \mathbb{P}[E^c] = 1 - 0 = 1.$$

□

**Exercise V.8** (Independence assumption is needed in Lemma V.8).

Find a sequence  $E_1, E_2, \dots \in \mathcal{F}$  of events (do not try independent) such that we have  $\sum_{n=1}^{\infty} \mathbb{P}[E_n] = +\infty$ , but  $\mathbb{P}[“E_n \text{ occurs infinitely often}”] \neq 1$ .

**Example V.9** (Records).

Consider an annual sports contest, in which one keeps track of the winner’s scores for different years as well as the record score of all past years.

Suppose that the winner’s score for year  $n$  is a real-valued random variable  $X_n$ , and suppose that  $X_1, X_2, \dots$  are independent and identically distributed — and moreover that the cumulative distribution function of  $X_n$  is continuous. (These assumptions are not completely unreasonable!)

Consider the event

$$E_n = \left\{X_n > \max\{X_1, \dots, X_{n-1}\}\right\}$$

that a new record is made in the contest of year  $n$ . We leave it as an exercise to the reader to prove that  $\mathbb{P}[E_n] = \frac{1}{n}$  and that  $E_1, E_2, \dots$  are independent.

Then, since the harmonic series diverges  $\sum_n \mathbb{P}[E_n] = \sum_n \frac{1}{n} = +\infty$ , the divergence part of Borel – Cantelli lemmas implies that  $\mathbb{P}[E_n \text{ infinitely often}] = 1$ . In other words: almost surely new records are made infinitely many times.

Consider then the event  $F_n = E_n \cap E_{n+1}$  that records are broken in the consecutive years  $n$  and  $n + 1$ . By what we know of  $E_1, E_2, \dots$ , we get  $\mathbb{P}[F_n] = \frac{1}{n(n+1)}$ , and thus  $\sum_n \mathbb{P}[F_n] = \sum_n \frac{1}{n^2+n} < +\infty$ . Now the convergence part of Borel – Cantelli lemmas implies that  $\mathbb{P}[F_n \text{ infinitely often}] = 0$ . In other words: almost surely there are only finitely many times that new records are made in consecutive years.

**Exercise V.9** (Decimal digits of a uniform random number).

Let  $\mathbb{P}$  be the uniform probability measure on the unit interval  $[0, 1]$  (cf. Example II.12). Consider the digits  $D_k(\omega)$ ,  $k \in \mathbb{N}$ , of the decimal representation of a number  $\omega \in [0, 1]$ , so that  $\omega = \sum_{k=1}^{\infty} D_k(\omega) 10^{-k}$ . For this exercise you can consider it known (*think about it anyway!*) that the digits  $D_1, D_2, \dots$  are simple random variables (on the sample space  $\Omega = [0, 1]$ ), and that they are independent.

Let  $Z_k$  be the run length of zeroes starting from the  $k$ :th digit:

- $Z_k = 0$  if the  $k$ :th digit is not zero,  $D_k \neq 0$ ;
- $Z_k = m$  if  $D_k = D_{k+1} = \dots = D_{k+m-1} = 0$  and  $D_{k+m} \neq 0$ .

- (a) Show that  $\mathbb{P}[Z_k = m] = \frac{9}{10^{m+1}}$  for all  $m \in \mathbb{Z}_{\geq 0}$ .
- (b) Fix  $m \in \mathbb{Z}_{\geq 0}$ . Show that  $\mathbb{P}[Z_k = m \text{ for infinitely many } k] = 1$ .
- (c) Show that  $\mathbb{P}[Z_k = k \text{ for infinitely many } k] = 0$ .



## Lecture VI

### Events of the infinite horizon

Consider a sequence of random variables. The topic of this lecture is:

What information about the sequence is not sensitive to the values of any finite number of individual members of the sequence?

Which events can be decided and which random variables are determined by such information?

Although it might at first appear surprising, there are in fact many interesting events and random variables which are not affected by any finite number of individual values. We will give a number of examples.

We will moreover return to profound consequences of independence, and ask:

Assuming moreover that the sequence of random variables is independent, what can be said about the probabilities of events that are not sensitive to any finite number of individual values?

Pertaining to the last question, we will prove Kolmogorov's 0-1 law which states that under the independence assumption, any event which is not sensitive to finitely many individual values has probability either zero or one, and any random variable which is not affected by finitely many individual values is almost surely constant. This probabilistic fact underlies some surprising phenomena, in particular related to phase transitions in physical systems.

#### VI.1. Definition of the tail $\sigma$ -algebra

Throughout, let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space.

Let  $X_1, X_2, X_3, \dots$  be a sequence of random variables. Denote first of all by

$$\mathcal{X} := \sigma(X_1, X_2, \dots) = \sigma((X_n)_{n \in \mathbb{N}}) \quad (\text{VI.1})$$

the  $\sigma$ -algebra generated by the collection  $(X_n)_{n \in \mathbb{N}}$  of random variables constituting the sequence (recall Definition IV.1). It should be interpreted as the information contained in the entire sequence.

We will also consider information contained in various subsets of the collection of random variables in the sequence. Specifically,

$$\mathcal{X}_k := \sigma(X_1, X_2, \dots, X_k) = \sigma((X_n)_{n=1, \dots, k}) \quad (\text{VI.2})$$

describes the information contained in the first  $k$  members of the sequence, and

$$\mathcal{T}_k := \sigma(X_{k+1}, X_{k+2}, \dots) = \sigma((X_n)_{n > k}) \quad (\text{VI.3})$$

describes the information contained in the sequence without its first  $k$  members. In other words,  $\mathcal{T}_k$  is the information in the sequence which is not affected by the first  $k$  values.

Our main interest lies in

$$\mathcal{T}_\infty := \bigcap_{k \in \mathbb{N}} \mathcal{T}_k. \quad (\text{VI.4})$$

As an intersection of the  $\sigma$ -algebras  $\mathcal{T}_k$ , also  $\mathcal{T}_\infty$  is a  $\sigma$ -algebra by Lemma I.8. It is called the *tail  $\sigma$ -algebra* of the sequence  $X_1, X_2, \dots$ . Since  $\mathcal{T}_\infty$  contains less information than any  $\mathcal{T}_k$ ,  $k \in \mathbb{N}$ , it describes the information which is not affected by any finite number of changes in the sequence. The corresponding events  $E \in \mathcal{T}_\infty$  are called *tail events*.

We start with examples of all of these different  $\sigma$ -algebras. For the examples, it is important to keep in mind that information in the first  $k$  members increases with  $k$ ,

$$\mathcal{X}_1 \subset \mathcal{X}_2 \subset \mathcal{X}_3 \subset \dots \subset \mathcal{X}_k \subset \mathcal{X}_{k+1} \subset \dots \subset \mathcal{X},$$

whereas information without the first  $k$  members decreases with  $k$

$$\mathcal{T}_\infty \subset \dots \subset \mathcal{T}_{k+1} \subset \mathcal{T}_k \subset \dots \subset \mathcal{T}_3 \subset \mathcal{T}_2 \subset \mathcal{T}_1 \subset \mathcal{X},$$

and of course all of these are contained in the information  $\mathcal{X}$  about the entire sequence (which in turn is contained in complete information  $\mathcal{F}$  on our probability space).

**Example VI.1** (Examples of events and random variables related to a sequence).

Suppose for concreteness that the sequence consists of real valued random variables, i.e.,

$$X_1, X_2, \dots \in m\mathcal{F}$$

are  $\mathcal{F}$ -measurable functions  $\Omega \rightarrow \mathbb{R}$ .

As a warm-up, note that by definition of  $\mathcal{X}_n = \sigma(X_1, \dots, X_n)$ , for any  $n \leq \ell$  we have

$$X_n \in m\mathcal{X}_n \subset m\mathcal{X}_\ell \subset m\mathcal{X}.$$

Because linear combinations of measurable functions are measurable by Proposition III.13, we get thus, for example, the following measurability of (finite) averages

$$\frac{X_1 + \dots + X_\ell}{\ell} \in m\mathcal{X}_\ell \subset m\mathcal{X}. \quad (\text{VI.5})$$

Also because upper and lower limits of measurable functions are measurable (Proposition III.14), we get the following

$$\liminf_n X_n \in m\mathcal{X} \quad \text{and} \quad \limsup_n X_n \in m\mathcal{X}. \quad (\text{VI.6})$$

We claim that the random variables (VI.6) are, in fact, measurable even with respect to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$ . For this, observe first that we can write, for any  $k \in \mathbb{N}$ ,

$$\liminf_n X_n = \sup_{m \in \mathbb{N}} \left( \inf_{n \geq m} X_n \right) = \sup_{m > k} \left( \inf_{n \geq m} X_n \right),$$

because  $\inf_{n \geq m} X_n$  is increasing in  $m$  (infimum over smaller collections), so omitting finitely terms corresponding to  $m = 1, 2, \dots, k$  does not affect the supremum. Now  $\inf_{n \geq m} X_n$  is the infimum of the countable collection of random variables  $(X_n)_{n \geq m}$ , which are by construction measurable with respect to  $\mathcal{T}_k = \sigma(X_{k+1}, X_{k+2}, \dots)$  when  $m > k$ . Proposition III.14 then first of all implies  $\inf_{n \geq m} X_n \in m\mathcal{T}_k$  for  $m > k$ , and the same proposition applied again to the supremum of these gives  $\liminf_n X_n = \sup_{m > k} (\inf_{n \geq m} X_n) \in m\mathcal{T}_k$ . Since this holds for any  $k$ , and the tail  $\sigma$ -algebra is defined as the intersection  $\mathcal{T}_\infty := \bigcap_k \mathcal{T}_k$ , we get  $\liminf_n X_n \in m\mathcal{T}_\infty$ . The limsup is handled similarly, and we conclude

$$\liminf_n X_n \in m\mathcal{T}_\infty \quad \text{and} \quad \limsup_n X_n \in m\mathcal{T}_\infty. \quad (\text{VI.7})$$

Provided that the limit  $\lim_{n \rightarrow \infty} X_n$  exists, it coincides with these, and is also  $\mathcal{T}_\infty$ -measurable.

We can use the above observations to check that the event of existence of the limit

$$E = \left\{ \omega \in \Omega \mid \text{the limit } \lim_{n \rightarrow \infty} X_n(\omega) \text{ exists} \right\}$$

is a tail event,  $E \in \mathcal{T}_\infty$ .<sup>1</sup> Indeed, the limit  $\lim_{n \rightarrow \infty} X_n$  does not exist if and only if the upper and lower limits are different, which happens if and only if we can find some rational number  $q \in \mathbb{Q}$  such that  $\liminf_n X_n < q < \limsup_n X_n$ . The event  $E$  of existence of limit is therefore

$$E = \left( \bigcup_{q \in \mathbb{Q}} \left( \{ \liminf_n X_n < q \} \cap \{ \limsup_n X_n > q \} \right) \right)^c, \quad (\text{VI.8})$$

where we are only using finite intersections, countable unions, and complements starting from the sets  $\{ \liminf_n X_n < q \}$  and  $\{ \limsup_n X_n > q \}$ , which belong to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$  by (VI.7). Thus we indeed have  $E \in \mathcal{T}_\infty$ .

A slightly less obvious example concerns the limit of averages (VI.5) as the number  $\ell$  of terms tends to infinity,

$$\lim_{\ell \rightarrow \infty} \frac{X_1 + \cdots + X_\ell}{\ell}.$$

The limit may or may not exist, so let us instead consider, e.g.,  $\limsup_{\ell \rightarrow \infty} \frac{X_1 + \cdots + X_\ell}{\ell}$ . The key observation is that for fixed  $k$ , we have

$$\lim_{\ell \rightarrow \infty} \frac{X_1 + \cdots + X_k}{\ell} = 0$$

(the numerator remains constant as  $\ell \rightarrow \infty$  while the denominator tends to infinity). Using this, we get

$$\begin{aligned} \limsup_{\ell \rightarrow \infty} \frac{X_1 + \cdots + X_\ell}{\ell} &= \limsup_{\ell \rightarrow \infty} \frac{(X_1 + \cdots + X_k) + (X_{k+1} + \cdots + X_\ell)}{\ell} \\ &= 0 + \limsup_{\ell \rightarrow \infty} \frac{X_{k+1} + \cdots + X_\ell}{\ell}. \end{aligned}$$

The random variables  $(X_n)_{n > k}$  are by construction measurable with respect to  $\mathcal{T}_k = \sigma(X_{k+1}, X_{k+2}, \dots)$ , so their linear combination  $\frac{X_{k+1} + \cdots + X_\ell}{\ell}$  is also  $\mathcal{T}_k$ -measurable according to Proposition III.13. By Proposition III.14 we then deduce that  $\limsup_{\ell \rightarrow \infty} \frac{X_1 + \cdots + X_\ell}{\ell} = \limsup_{\ell \rightarrow \infty} \frac{X_{k+1} + \cdots + X_\ell}{\ell}$  is  $\mathcal{T}_k$ -measurable. Since this holds for any  $k$ , and the tail  $\sigma$ -algebra is the intersection  $\mathcal{T}_\infty := \bigcap_k \mathcal{T}_k$ , we conclude

$$\limsup_{\ell} \frac{X_1 + \cdots + X_\ell}{\ell} \in \mathfrak{m}\mathcal{T}_\infty.$$

Like before, it follows that the existence of the limit of averages is a tail event

$$\left\{ \omega \in \Omega \mid \text{the limit } \lim_{\ell \rightarrow \infty} \frac{X_1(\omega) + \cdots + X_\ell(\omega)}{\ell} \text{ exists} \right\} \in \mathcal{T}_\infty,$$

and provided the limit exists, it is measurable with respect to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$ .

In particular, the example shows that there are some interesting events in the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$ . The following exercise contains a few more examples.

**Exercise VI.1** (Tail events).

Let  $X_1, X_2, X_3, \dots$  be a sequence of random numbers (i.e.,  $\mathbb{R}$ -valued random variables)

<sup>1</sup>If the sequence  $X_1, X_2, \dots$  is bounded, then there is a slick way to check this. Form the difference  $D = \limsup_n X_n - \liminf_n X_n$ , which is  $\mathcal{T}_\infty$ -measurable as a linear combination (use Proposition III.13). The limit exists if and only if the difference is zero,  $D = 0$ . In other words, the event  $E$  is the preimage  $E = D^{-1}(\{0\}) \in \mathcal{T}_\infty$  (since  $\{0\} \in \mathcal{B}$  and  $D \in \mathfrak{m}\mathcal{T}_\infty$ ).

defined on a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ . Investigate which of the following events belong to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$ .

- (a):  $\left\{ \omega \in \Omega \mid \text{the series } \sum_{n=1}^{\infty} X_n(\omega) \text{ converges} \right\}$
- (b):  $\left\{ \omega \in \Omega \mid \sum_{n=1}^{\infty} X_n(\omega) \leq -42 \right\}$
- (c):  $\left\{ \omega \in \Omega \mid \text{the sequence } X_1, X_2, \dots \text{ is bounded} \right\}$
- (d):  $\left\{ \omega \in \Omega \mid \forall \ell \in \mathbb{N} \exists n \in \mathbb{N} \text{ such that } X_n(\omega) = X_{n+1}(\omega) = \dots = X_{n+\ell}(\omega) \right\}$   
 $= \left\{ \text{there exists arbitrarily long repetitions in the sequence } X_1, X_2, \dots \right\}$

## VI.2. Kolmogorov's 0-1 law

If we assume that the sequence is independent, we get a remarkable result about the probabilities of tail events.

**Theorem VI.2** (Kolmogorov's 0-1 law).

*Suppose that  $X_1, X_2, \dots$  is a sequence of independent random variables. Then the following hold:*

- (a) *For any tail event  $E \in \mathcal{T}_\infty$  we have either*

$$\mathbf{P}[E] = 1 \quad \text{or} \quad \mathbf{P}[E] = 0.$$

- (b) *Any  $\mathbb{R}$ -valued random variable  $T$  which is measurable with respect to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$  is almost surely constant, i.e., for some  $c \in \mathbb{R}$  we have*

$$\mathbf{P}[T = c] = 1.$$

The proof strategy is almost as surprising as the statement: we will show that the information in the tail is independent of itself,  $\mathcal{T}_\infty \perp \mathcal{T}_\infty$ .

*Proof of Theorem VI.2.* We will show that  $\mathcal{T}_\infty \perp \mathcal{T}_\infty$ . The assertions (a) and (b) can then be concluded rather easily as follows.

For (a), suppose that  $E \in \mathcal{T}_\infty$ . Then by the independence property  $\mathcal{T}_\infty \perp \mathcal{T}_\infty$  we have

$$\mathbf{P}[E \cap E] = \mathbf{P}[E] \mathbf{P}[E].$$

But in view of the obvious fact  $E \cap E = E$ , this gives  $\mathbf{P}[E] = \mathbf{P}[E]^2$ . The only two solutions of this quadratic equation are 0 and 1, so we conclude  $\mathbf{P}[E] \in \{0, 1\}$ .

For (b), suppose that  $T \in \mathfrak{m}\mathcal{T}_\infty$ . Then for any  $x \in \mathbb{R}$  we have

$$\{T \leq x\} = T^{-1}((-\infty, x]) \in \mathcal{T}_\infty,$$

so by part (a) we have

$$\mathbf{P}[T \leq x] \in \{0, 1\}.$$

This shows that the cumulative distribution function of  $T$  never takes any other values except 0 and 1. If we define

$$c := \inf \{x \in \mathbb{R} \mid \mathbf{P}[T \leq x] = 1\},$$

then by right continuity of cumulative distribution functions (Proposition II.30) we have  $\mathbf{P}[T \leq c] = 1$  and  $\mathbf{P}[T < c] = 0$ . This gives

$$\mathbf{P}[T = c] = \mathbf{P}[T \leq c] - \mathbf{P}[T < c] = 1,$$

i.e.,  $T = c$  almost surely as claimed in (b).

The proof will thus be complete once we have shown that  $\mathcal{T}_\infty \perp \mathcal{T}_\infty$ . We will do this in a number of steps.

*step 1,  $\mathcal{X}_k \perp \mathcal{T}_k$ :* We claim that the two  $\sigma$ -algebras  $\mathcal{X}_k$  and  $\mathcal{T}_k$  are independent. We will verify their independence using two  $\pi$ -systems which generate these  $\sigma$ -algebras.

Let  $\mathcal{J}$  denote the collection of events of the form

$$A = \left\{ \omega \in \Omega \mid X_1(\omega) \in B_1, X_2(\omega) \in B_2, \dots, X_k(\omega) \in B_k \right\} \subset \Omega, \quad (\text{VI.9})$$

where  $B_1, B_2, \dots, B_k \in \mathcal{B}$ .

Then  $\mathcal{J}$  is a  $\pi$ -system, which generates the first of our  $\sigma$ -algebras,  $\sigma(\mathcal{J}) = \mathcal{X}_k$  (we leave it as an exercise<sup>2</sup> to the reader to verify this).

Similarly, let  $\mathcal{J}'$  denote the collection of events of the form

$$A' = \left\{ \omega \in \Omega \mid X_{k+1}(\omega) \in B_{k+1}, X_{k+2}(\omega) \in B_{k+2}, \dots, X_{k+r}(\omega) \in B_{k+r} \right\} \subset \Omega, \quad (\text{VI.10})$$

where  $r \in \mathbb{N}$  and  $B_{k+1}, B_{k+2}, \dots, B_{k+r} \in \mathcal{B}$ .

Then  $\mathcal{J}'$  is a  $\pi$ -system, which generates the second of our  $\sigma$ -algebras,  $\sigma(\mathcal{J}') = \mathcal{T}_k$  (we again leave it as an exercise to the reader to verify this).

Consider now  $A \in \mathcal{J}$  and  $A' \in \mathcal{J}'$  as in (VI.9) and (VI.10). Note that by independence of  $X_1, X_2, \dots$  we have

$$\begin{aligned} \mathbb{P}[A] &= \mathbb{P}[X_1 \in B_1, \dots, X_k \in B_k] = \prod_{j=1}^k \mathbb{P}[X_j \in B_j] \\ \mathbb{P}[A'] &= \mathbb{P}[X_{k+1} \in B_{k+1}, \dots, X_{k+r} \in B_{k+r}] = \prod_{j=k+1}^{k+r} \mathbb{P}[X_j \in B_j]. \end{aligned}$$

The probability of the intersection  $A \cap A'$  is similarly computed using independence,

$$\begin{aligned} \mathbb{P}[A \cap A'] &= \mathbb{P}[X_1 \in B_1, \dots, X_k \in B_k, X_{k+1} \in B_{k+1}, \dots, X_{k+r} \in B_{k+r}] \\ &= \prod_{j=1}^{k+r} \mathbb{P}[X_j \in B_j]. \end{aligned}$$

These expressions show that  $\mathbb{P}[A \cap A'] = \mathbb{P}[A] \mathbb{P}[A']$ . Since the  $\pi$ -systems  $\mathcal{J}$  and  $\mathcal{J}'$  generate the  $\sigma$ -algebras  $\mathcal{X}_k$  and  $\mathcal{T}_k$ , it follows from Proposition V.5 that  $\mathcal{X}_k$  and  $\mathcal{T}_k$  are independent.

*step 2,  $\mathcal{X}_k \perp \mathcal{T}_\infty$ :* In step 1 we showed  $\mathcal{X}_k \perp \mathcal{T}_k$ . But since we have  $\mathcal{T}_\infty \subset \mathcal{T}_k$ , we a fortiori have also  $\mathcal{X}_k \perp \mathcal{T}_\infty$  (there are now fewer sets for which the condition of Definition V.1 needs to be verified!).

*step 3,  $\mathcal{X} \perp \mathcal{T}_\infty$ :* We claim that the two  $\sigma$ -algebras  $\mathcal{X}$  and  $\mathcal{T}_\infty$  are independent. We will again verify their independence using  $\pi$ -systems.

Consider the collection  $\mathcal{U} = \bigcup_{k \in \mathbb{N}} \mathcal{X}_k$ . We claim that  $\mathcal{U}$  is a  $\pi$ -system which generates  $\mathcal{X}$  (note that  $\mathcal{U}$  is the union of  $\sigma$ -algebras, but  $\mathcal{U}$  is generally not a  $\sigma$ -algebra itself).

To check that  $\mathcal{U}$  is a  $\pi$ -system, suppose that  $A_1, A_2 \in \mathcal{U}$ . Then since  $\mathcal{U}$  is defined as the union of  $\mathcal{X}_k$ ,  $k \in \mathbb{N}$ , we have  $A_1 \in \mathcal{X}_{k_1}$  and  $A_2 \in \mathcal{X}_{k_2}$  for some  $k_1, k_2 \in \mathbb{N}$ . Recall that the sequence of  $\sigma$ -algebras in the union is increasing,  $\mathcal{X}_1 \subset \mathcal{X}_2 \subset \dots$ . Therefore by setting  $k_0 = \max\{k_1, k_2\}$ , we have  $A_1, A_2 \in \mathcal{X}_{k_0}$ . Since  $\mathcal{X}_{k_0}$  is itself a  $\sigma$ -algebra, we thus have also  $A_1 \cap A_2 \in \mathcal{X}_{k_0}$ . Finally, since we have  $\mathcal{X}_{k_0} \subset \bigcup_{k \in \mathbb{N}} \mathcal{X}_k = \mathcal{U}$ , we get  $A_1 \cap A_2 \in \mathcal{U}$ , which shows that  $\mathcal{U}$  is a  $\pi$ -system.

<sup>2</sup>As a hint for this exercise, it is worth noting that step 3 contains the details of a similar argument in a slightly more complicated situation.

Let us then check that  $\mathcal{U}$  generates  $\mathcal{X}$ , i.e.,  $\sigma(\mathcal{U}) = \mathcal{X}$ . We do this by showing inclusions in both directions. Note first that for any  $n \in \mathbb{N}$  we have  $X_n \in \mathfrak{m}\mathcal{X}_n$  directly from the definition, and since  $\mathcal{X}_n \subset \bigcup_{k \in \mathbb{N}} \mathcal{X}_k = \mathcal{U} \subset \sigma(\mathcal{U})$ , we also have  $X_n \in \mathfrak{m}\sigma(\mathcal{U})$ . Thus  $\sigma(\mathcal{U})$  is a  $\sigma$ -algebra with respect to which each  $X_n$  is measurable, and by definition  $\mathcal{X} = \sigma((X_k)_{k \in \mathbb{N}})$  is the smallest such  $\sigma$ -algebra, so we get  $\mathcal{X} \subset \sigma(\mathcal{U})$ . On the other hand, for each  $k \in \mathbb{N}$  we have  $\mathcal{X}_k \subset \mathcal{X}$  so also  $\mathcal{U} \subset \mathcal{X}$ , and since  $\sigma(\mathcal{U})$  is the smallest  $\sigma$ -algebra with this property, we get  $\sigma(\mathcal{U}) \subset \mathcal{X}$ . These two inclusions give the equality  $\sigma(\mathcal{U}) = \mathcal{X}$ .

Now we know that  $\mathcal{U}$  is a  $\pi$ -system which generates  $\mathcal{X}$ . We use  $\mathcal{T}_\infty$  itself as a  $\pi$ -system which generates  $\mathcal{T}_\infty$ , i.e.,  $\sigma(\mathcal{T}_\infty) = \mathcal{T}_\infty$  (recall that any  $\sigma$ -algebra is also a  $\pi$ -system). Let  $A \in \mathcal{U}$  and  $E \in \mathcal{T}_\infty$ . We have  $A \in \mathcal{X}_k$  for some  $k \in \mathbb{N}$ , because of the definition of  $\mathcal{U}$  as a union. In step 2 we showed that  $\mathcal{X}_k \perp \mathcal{T}_\infty$ , so we get  $\mathbb{P}[A \cap E] = \mathbb{P}[A] \mathbb{P}[E]$ . It therefore follows from Proposition V.5 that  $\mathcal{X}$  and  $\mathcal{T}_\infty$  are independent.

*final step,  $\mathcal{T}_\infty \perp \mathcal{T}_\infty$ :* In step 3 we showed  $\mathcal{X} \perp \mathcal{T}_\infty$ . But since we have  $\mathcal{T}_\infty \subset \mathcal{X}$ , we a fortiori have also  $\mathcal{T}_\infty \perp \mathcal{T}_\infty$  (there is less to verify!). This concludes the proof.  $\square$

Part (b) of Kolmogorov's 0-1 law is stated in Theorem above for real-valued random variables measurable with respect to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$ . This choice was made for the sake of concreteness — you may verify that the same conclusion holds much more generally.

**Exercise VI.2.** Let  $\mathfrak{X}$  be a complete separable metric space. Assume that a random variable  $T: \Omega \rightarrow \mathfrak{X}$  is  $\mathcal{T}_\infty/\mathcal{B}(\mathfrak{X})$ -measurable, where  $\mathcal{T}_\infty$  is the tail  $\sigma$ -algebra of a sequence of independent random variables  $X_1, X_2, \dots$ . Prove that  $T$  is almost surely constant, i.e., there exists a point  $x \in \mathfrak{X}$  such that  $\mathbb{P}[T = x] = 1$ .

## A few interesting examples

### *Various random series*

**Example VI.3** (Random series with independent terms).

Suppose that  $X_1, X_2, \dots$  are independent  $\mathbb{R}$ -valued random variables. Consider the random series formed from the sequence,

$$\sum_{n=1}^{\infty} X_n.$$

Then the event that this series converges,

$$E = \left\{ \omega \in \Omega \mid \sum_{n=1}^{\infty} X_n(\omega) \text{ converges} \right\},$$

belongs to the tail  $\sigma$ -algebra,  $E \in \mathcal{T}_\infty$ . Since the terms  $X_n$  were assumed independent, Kolmogorov's 0-1 law states that  $\mathbb{P}[E] \in \{0, 1\}$ , i.e., either the series  $\sum_{n=1}^{\infty} X_n$  converges almost surely or it diverges almost surely. The theorem does not tell which of these two extremes occurs.

The following example is a simple but interesting special case of random series in Example VI.3. It in particular shows that both of the two extremes in the example are indeed possible.

**Example VI.4** (Series with randomly assigned signs).

Let  $a_1, a_2, \dots \geq 0$  be a sequence of non-negative real numbers and let  $S_1, S_2, \dots$  be a sequence of independent and identically distributed  $\{\pm 1\}$ -valued random variables with

$$\mathbb{P}[S_n = +1] = \frac{1}{2} \quad \text{and} \quad \mathbb{P}[S_n = -1] = \frac{1}{2}.$$

Consider the following series of  $a_n$ 's with random signs ( $S_n = \pm 1$ ),

$$\sum_{n=1}^{\infty} a_n S_n. \quad (\text{VI.11})$$

This is a series of the type considered in Example VI.3, with terms  $X_n = a_n S_n$ . Therefore the series (VI.11) either converges almost surely or diverges almost surely. It depends on the given sequence  $(a_n)_{n \in \mathbb{N}}$ , which of these happens.

Note first that the series (VI.11) can only converge if its terms at least tend to zero. The absolute values of the terms are  $|a_n S_n| = a_n$ , so the terms tend to zero if and only if the given sequence  $(a_n)_{n \in \mathbb{N}}$  satisfies  $\lim_{n \rightarrow \infty} a_n = 0$ . In particular, if we take  $a_n = 1$  for all  $n$  (or anything else which does not tend to zero), then series (VI.11) diverges almost surely.

On the other hand, if the series  $\sum_{n=1}^{\infty} a_n$  converges, then the series (VI.11) is always absolutely convergent, whatever the random signs. In particular, if we take  $a_n = n^{-2}$  for  $n \in \mathbb{N}$  (or anything else which constitutes a convergent series), then series (VI.11) converges almost surely.

The interesting borderline cases for random series of type (VI.11) are when  $a_n \rightarrow 0$  as  $n \rightarrow \infty$ , but we have  $\sum_{n=1}^{\infty} a_n = \infty$ . An example of this is  $a_n = n^{-1}$  for  $n \in \mathbb{N}$ . Then Kolmogorov's 0-1 law still says that we either have almost sure convergence or divergence, but which one is it now?

Another variant of random series is power series with random coefficients. These also have important applications, although we will only use them for the purpose of examples.

**Example VI.5** (Random power series).

Let  $Y_0, Y_1, Y_2, \dots$  be a sequence of independent random variables. The power series

$$F(z) = \sum_{n=0}^{\infty} Y_n z^n \quad (\text{VI.12})$$

with these coefficients defines a random function  $F$  of a real (or complex) variable  $z$ . For any fixed value of  $z \in \mathbb{R}$  (or  $z \in \mathbb{C}$ ), the series (VI.12) is essentially a special case of Example VI.3, with terms  $X_n = Y_n z^n$  (the only difference is that starting the indexing from  $n = 0$  is more natural for power series).

We know from the theory of power series that the series (VI.12) has some radius of convergence  $R$  such that

$$(\text{VI.12}) \text{ is convergent for } |z| < R \quad \text{and} \quad (\text{VI.12}) \text{ is divergent for } |z| > R.$$

Since the coefficients of the series are random, it seems that the radius of convergence  $R$  should be random, too. In fact, the Cauchy–Hadamard formula in the theory of power series tells us that  $R$  can be written in terms of the coefficients as

$$R = \frac{1}{\limsup_n \sqrt[n]{|Y_n|}}. \quad (\text{VI.13})$$

Note that the function  $y \mapsto \sqrt[n]{|y|}$  is continuous  $\mathbb{R} \rightarrow \mathbb{R}$ , and in particular Borel-measurable (by Corollary III.10), so  $\sqrt[n]{|Y_n|}$  is a (measurable) random variable (by Proposition III.8). Then also the denominator,  $\limsup_n \sqrt[n]{|Y_n|}$ , is a random variable (by Proposition III.14). Finally, taking the reciprocal  $s \mapsto \frac{1}{s}$  is continuous  $[0, \infty] \rightarrow [0, \infty]$ , so we see that the radius of convergence  $R$  of the random power series (VI.12) is indeed a random variable!

However, the radius of convergence  $R$  in (VI.13) is measurable with respect to the tail  $\sigma$ -algebra  $\mathcal{T}_{\infty}$  (the limsup is insensitive to finite number of changes in the coefficients). Therefore, having assumed independence of the coefficients, we get from Kolmogorov's 0-1 law that  $R$  is almost surely equal to some constant  $c$ . In conclusion, the radius of convergence of this random power series is in fact essentially deterministic (non-random)!

*Random walks*

**Example VI.6** (Escape probability of asymmetric simple random walk).

Let  $\theta \in [0, 1]$  be a parameter, and let  $X_1, X_2, \dots$  be a sequence of independent and identically distributed  $\{\pm 1\}$ -valued random variables with

$$\mathbb{P}[X_n = +1] = \theta \quad \text{and} \quad \mathbb{P}[X_n = -1] = 1 - \theta.$$

The argument in Example VI.4 shows that the infinite series  $\sum_{n=1}^{\infty} X_n$  diverges almost surely (we had  $\theta = \frac{1}{2}$ , but the same argument works for any  $\theta$ ). In this example we consider a different but related question.

We think of  $X_n$  as the  $n$ :th step of a random walker:  $X_n = +1$  is interpreted as a step forward and  $X_n = -1$  as a step backwards, and the parameter  $\theta$  gives the probability of a forward step. If the walker starts from the origin, then her position after the first  $s$  steps is

$$W_s := \sum_{n=1}^s X_n.$$

Let us consider the question of whether the random walker eventually advances arbitrarily much: this is described by the event

$$\left\{ \omega \in \Omega \mid \sup_s W_s(\omega) = +\infty \right\}$$

that the sequence  $W_1, W_2, \dots$  of walker's positions is not bounded from above. We leave it to the reader to check that this event  $\{\sup_s W_s = +\infty\}$  belongs to the tail  $\sigma$ -algebra  $\mathcal{T}_\infty$ .<sup>3</sup> Having assumed the independence of the steps  $X_1, X_2, \dots$ , it follows from Kolmogorov's 0-1 law that the probability of advancing arbitrarily much is either zero or one:

$$\mathbb{P}\left[\sup_s W_s = +\infty\right] \in \{0, 1\}.$$

In other words, either the walker almost surely advances arbitrarily much, or she almost surely does not. The theorem does not, however, tell which of these two options is true. The answer turns out to depend on the parameter  $\theta$ : one can show that we have

$$\mathbb{P}\left[\sup_s W_s = +\infty\right] = 1 \quad \text{if and only if} \quad \theta \geq \frac{1}{2}.$$

It is natural to refine the question of advancement of the walker slightly. By including also considerations of  $\limsup_s W_s$  and  $\liminf_s W_s$ , the reader can show that also

$$\left\{ \lim_{s \rightarrow \infty} W_s = +\infty \right\} \quad \text{and} \quad \left\{ \lim_{s \rightarrow \infty} W_s = -\infty \right\}$$

are tail events. The former event describes the walker escaping towards  $+\infty$ , and the latter describes the the walker escaping towards  $-\infty$ . Concerning these, Kolmogorov's 0-1 law says that either the walker almost surely escapes towards  $+\infty$  (resp.  $-\infty$ ), or the probability of such escape is zero. Which is the case again depends on  $\theta$ . One can show that

$$\mathbb{P}\left[\lim_{s \rightarrow \infty} W_s = +\infty\right] = 1 \quad \text{if and only if} \quad \theta > \frac{1}{2}.$$

whereas symmetrically

$$\mathbb{P}\left[\lim_{s \rightarrow \infty} W_s = -\infty\right] = 1 \quad \text{if and only if} \quad \theta < \frac{1}{2}.$$

This seems to leave just one question: where on earth does that walker go if  $\theta = \frac{1}{2}$ ?<sup>4</sup>

<sup>3</sup>It is important to note that we still talk about the tail  $\sigma$ -algebra of the sequence  $X_1, X_2, \dots$  of steps, not the tail  $\sigma$ -algebra of the sequence  $W_1, W_2, \dots$  of random walk positions.

<sup>4</sup>The answer is: all over the place! Namely, at  $\theta = \frac{1}{2}$  we in fact have

$$\mathbb{P}\left[\limsup_{s \rightarrow \infty} W_s = +\infty \text{ and } \liminf_{s \rightarrow \infty} W_s = -\infty\right] = 1,$$

so the walker almost surely advances arbitrarily far in both directions and meanwhile returns to the origin infinitely often as well.



## Lecture VII

### Integration theory

Let  $(S, \mathcal{S}, \mu)$  be a measure space. The goal of this lecture is to define, for all reasonable functions  $f: S \rightarrow \mathbb{R}$ , the integral

$$\int f \, d\mu$$

of the function  $f$  with respect to the measure  $\mu$  — denoted, when we want to emphasize the space  $S$  and the variable  $s \in S$ , also by

$$\int_S f(s) \, d\mu(s).$$

This general theory of integration was originally conceived of by Henri Lebesgue. The theory is much more flexible than integration in the sense of Riemann, and it moreover leads to some remarkably powerful tools.

One of our main motivations for integration in probability theory is to get a precise mathematical definition of the expected value, which we will address in more detail in the next lecture (Lecture VIII). There are, however, also other equally important uses of the construction — special cases include:

**Summation:** Integration with respect to the counting measure (Example II.10) is just summation,

$$\int_S f(s) \, d\mu_{\#}(s) = \sum_{s \in S} f(s).$$

From the general theory we develop for integrals, we will in particular get a general and precise definition of summation and many useful tools for calculating with sums.

**Riemann integral and generalization:** For real-valued functions  $f$  on the real line  $\mathbb{R}$ , integration with respect to the Lebesgue measure  $\Lambda$  (Example II.12) generalizes<sup>1</sup> the familiar Riemann integral. For this case, we therefore use also the very familiar notation

$$\int_{\mathbb{R}} f(x) \, d\Lambda(x) = \int_{-\infty}^{\infty} f(x) \, dx.$$

**Expected value:** The expected value  $E[X]$  of a real valued random variable  $X$  will be defined as the integral of the function  $X: \Omega \rightarrow \mathbb{R}$  with respect to the probability measure  $P$  on  $(\Omega, \mathcal{F})$ ,

$$E[X] := \int_{\Omega} X(\omega) \, dP(\omega).$$

---

<sup>1</sup>To be precise, for any Riemann-integrable function  $f: \mathbb{R} \rightarrow \mathbb{R}$  such that also  $|f|$  is Riemann-integrable, the Riemann integral of  $f$  agrees with the integral of  $f$  with respect to the Lebesgue measure  $\Lambda$ . In principle both integrals are defined for somewhat more general functions  $f$ , but the Lebesgue integral is without a doubt the more fruitful generalization.

The construction of the integral will proceed by increasing the complexity of the allowed functions step by step:

- (1) indicator functions
- (2) non-negative simple functions
- (3) non-negative measurable functions
- (4) all integrable functions.

Many fundamental results about integration are likewise proved step by step — this proof strategy is affectionately referred to as the “*standard machine*”.

For these steps, recall also the notations

$$\begin{aligned} \text{m}\mathcal{S} &= \text{the set of } \mathcal{S}\text{-measurable functions } S \rightarrow [-\infty, +\infty] \\ \text{m}\mathcal{S}^+ &= \text{the set of } \mathcal{S}\text{-measurable functions } S \rightarrow [0, +\infty] \\ \text{s}\mathcal{S} &= \text{the set of simple functions } S \rightarrow \mathbb{R} \\ \text{s}\mathcal{S}^+ &= \text{the set of simple functions } S \rightarrow [0, +\infty] \end{aligned}$$

In order to clearly distinguish between the (in principle) different definitions in each stage, during this lecture we use different notation for the integrals in the different stages as follows:

$$\begin{aligned} \int^\square &\text{ — the integral for non-negative simple functions (step 2)} \\ \int^+ &\text{ — the integral for non-negative measurable functions (step 3)} \\ \int &\text{ — the integral for all integrable functions (step 4).} \end{aligned}$$

It will be shown in Lemmas VII.5 and VII.13, however, that the more general definitions coincide with the earlier ones, so in later lectures there will be no need to make these distinctions, and we can safely only use the notation  $\int$ .

Besides the construction, we prove some basic properties of the integral. In particular we will show *linearity*<sup>2</sup>

$$\begin{aligned} \int (cf) \, d\mu &= c \int f \, d\mu \\ \int (f_1 + f_2) \, d\mu &= \int f_1 \, d\mu + \int f_2 \, d\mu \end{aligned}$$

and *monotonicity*

$$f \leq g \quad \implies \quad \int f \, d\mu \leq \int g \, d\mu,$$

which in fact must be verified separately at each step of the construction, i.e., for the integrals  $\int^\square$ ,  $\int^+$ , and  $\int$ .

Finally, in Section VII.4 we establish powerful general convergence theorems for integrals.

---

<sup>2</sup>In the end, the scalar  $c$  is allowed to be an arbitrary real number. During the steps which address integrals of non-negative functions, however, it will of course only be meaningful to allow non-negative scalars.

### The approximation lemma

Especially when going from step (2) of the construction to step (3), i.e., from non-negative simple functions to all non-negative measurable functions, it will be important to keep in mind the following approximation result from Lecture III.

**Lemma** (Lemma III.18). *For any non-negative measurable function*

$$f \in \mathfrak{m}\mathcal{S}^+$$

*there exists a sequence<sup>3</sup>*

$$f_1, f_2, f_3, \dots \in \mathfrak{s}\mathcal{S}^+$$

*of non-negative simple functions increasing pointwise to  $f$ , i.e.,*

$$f_n(s) \uparrow f(s) \quad \text{for all } s \in S.$$

## VII.1. Integral for non-negative simple functions

### Definition of integral of non-negative simple functions

For a subset  $A \subset S$ , the indicator function  $\mathbb{I}_A$  takes the value 1 on  $A$  and 0 elsewhere,

$$\mathbb{I}_A(s) = \begin{cases} 1 & \text{if } s \in A \\ 0 & \text{if } s \notin A. \end{cases}$$

If  $A \in \mathcal{S}$ , then we want to define the integral of such a function simply as the measure of the set  $A$ ,

$$\int^{\square} \mathbb{I}_A \, d\mu := \mu[A]. \quad (\text{VII.1})$$

This is step 1 in our definition of the integral. Note that already in this case it is possible for the integral to become infinite if  $\mu[A] = +\infty$ .

Any non-negative simple function  $h \in \mathfrak{s}\mathcal{S}^+$  can be written as

$$h = \sum_{j=1}^n a_j \mathbb{I}_{A_j},$$

where  $a_1, \dots, a_n$  are non-negative coefficients and  $A_1, \dots, A_n$  are measurable sets. To achieve linearity of the integral for simple functions, our only option is the following definition.

**Definition VII.1** (Integral of a non-negative simple function).

For  $h = \sum_{j=1}^n a_j \mathbb{I}_{A_j}$  with  $a_1, \dots, a_n \in [0, +\infty)$  and  $A_1, \dots, A_n \in \mathcal{S}$ , define

$$\int^{\square} h \, d\mu := \sum_{j=1}^n a_j \mu[A_j]. \quad (\text{VII.2})$$

<sup>3</sup> One concrete way of constructing the approximating sequence is to set  $f_n = \varsigma_n \circ f$ , where  $\varsigma_n: [0, +\infty] \rightarrow [0, n]$  is the  $n$ :th staircase function with steps of size  $2^{-n}$  and truncation at level  $n$  — see Equation (III.5) for definition and Figure III.2 for illustration.

A priori, this definition depends on the choices made in the expression of  $h$  as a linear combination of indicators, so we should verify well-definedness.

**Exercise VII.1** (Well-definedness of the integral of non-negative simple functions).

Show that if  $h \in \mathcal{S}^+$  can be written in two ways,

$$h = \sum_{j=1}^n a_j \mathbb{I}_{A_j} \quad \text{and} \quad h = \sum_{k=1}^{n'} a'_k \mathbb{I}_{A'_k},$$

as a non-negative linear combination of indicator functions, then we have

$$\sum_{j=1}^n a_j \mu[A_j] = \sum_{k=1}^{n'} a'_k \mu[A'_k].$$

Conclude that the integral of  $h$  is well-defined by (VII.2).

### Properties of integral of non-negative simple functions

Let us then verify the basic properties of the integral defined in (VII.3). These will be used already in the next step.

**Lemma VII.2** (Linearity of the integral for simple functions).

(a) If  $h \in \mathcal{S}^+$  and  $c \geq 0$ , then  $ch \in \mathcal{S}^+$  and

$$\int^{\square} (ch) \, d\mu = c \int^{\square} h \, d\mu.$$

(b) If  $h, g \in \mathcal{S}^+$ , then  $h + g \in \mathcal{S}^+$  and

$$\int^{\square} (h + g) \, d\mu = \int^{\square} h \, d\mu + \int^{\square} g \, d\mu.$$

*Proof.* Both assertions are derived from Definition VII.1 relying on its well-definedness (Exercise VII.1), which permits us to use whichever decomposition to linear combination of indicators we find the most convenient.

*proof of (a):* Write  $h \in \mathcal{S}^+$  as a linear combination of indicators with non-negative coefficients,  $h = \sum_{j=1}^n a_j \mathbb{I}_{A_j}$ , with  $a_j \geq 0$  and  $A_j \in \mathcal{S}$  for  $j = 1, \dots, n$ . Then we have

$$ch = c \sum_{j=1}^n a_j \mathbb{I}_{A_j} = \sum_{j=1}^n c a_j \mathbb{I}_{A_j}.$$

Therefore, according to (VII.2), the integral is

$$\int^{\square} (ch) \, d\mu = \sum_{j=1}^n c a_j \mu[A_j] = c \sum_{j=1}^n a_j \mu[A_j] = c \int^{\square} h \, d\mu.$$

This proves the assertion (a).

*proof of (b):* Write  $h \in \mathcal{S}^+$  and  $g \in \mathcal{S}^+$  as linear combinations of indicators with non-negative coefficients,  $h = \sum_{j=1}^n a_j \mathbb{I}_{A_j}$  and  $g = \sum_{k=1}^m b_k \mathbb{I}_{B_k}$ . Then we have

$$h + g = \sum_{j=1}^n a_j \mathbb{I}_{A_j} + \sum_{k=1}^m b_k \mathbb{I}_{B_k}.$$

Therefore, according to (VII.2), the integral is

$$\int^{\square} (h + g) \, d\mu = \sum_{j=1}^n a_j \mu[A_j] + \sum_{k=1}^m b_k \mu[B_k] = \int^{\square} h \, d\mu + \int^{\square} g \, d\mu.$$

This proves the assertion (b).  $\square$

**Lemma VII.3** (Monotonicity of the integral for simple functions).

Suppose that  $h, g \in s\mathcal{S}^+$  and  $h \leq g$  pointwise. Then we have

$$\int^{\square} h \, d\mu \leq \int^{\square} g \, d\mu.$$

*Proof.* We can first of all write

$$h = \sum_{j=1}^n a_j \mathbb{I}_{A_j}$$

with  $a_1, \dots, a_n$  the finitely many different values of  $h$  and  $A_j = h^{-1}(\{a_j\})$  the sets where this value is taken. Similarly we can write

$$g = \sum_{k=1}^m b_k \mathbb{I}_{B_k}$$

with  $b_1, \dots, b_m$  the finitely many different values of  $g$  and  $B_k = g^{-1}(\{b_k\})$ . The intersections  $A_j \cap B_k$  are the sets where simultaneously  $f$  takes the value  $a_j$  and  $g$  takes the value  $b_k$  — we can use them as refinements, and write the alternative expressions

$$h = \sum_{j=1}^n \sum_{k=1}^m a_j \mathbb{I}_{A_j \cap B_k} \quad \text{and} \quad g = \sum_{j=1}^n \sum_{k=1}^m b_k \mathbb{I}_{A_j \cap B_k},$$

which have the advantage that the same indicators appear in both. Observe that the assumption  $h \leq g$  implies that  $a_j \leq b_k$  whenever  $A_j \cap B_k \neq \emptyset$ . On the other hand, whenever  $A_j \cap B_k = \emptyset$  we of course have  $\mu[A_j \cap B_k] = 0$ , so the desired conclusion

$$\int^{\square} h \, d\mu = \sum_{j=1}^n \sum_{k=1}^m a_j \mu[A_j \cap B_k] \leq \sum_{j=1}^n \sum_{k=1}^m b_k \mu[A_j \cap B_k] = \int^{\square} g \, d\mu.$$

follows by comparing the non-zero terms in the sums.  $\square$

## VII.2. Integral for non-negative measurable functions

### Definition of integral of non-negative measurable functions

To preserve monotonicity, the integral of a nonnegative function  $f$  should be at least as large as the integral of any simple function  $h \leq f$  below it. The following definition is thus motivated by preserving monotonicity without unnecessarily giving up anything extra.

**Definition VII.4** (Integral of a non-negative measurable function).

For a non-negative measurable function  $f \in m\mathcal{S}^+$  we define

$$\int^+ f \, d\mu := \sup_{\substack{h \in s\mathcal{S}^+ \\ 0 \leq h \leq f}} \int^{\square} h \, d\mu. \quad (\text{VII.3})$$

Note that it is possible for the integral defined by the above supremum to become infinite — even for finite measures for which the integrals of all simple functions are finite.

We start by checking that the new more general definition (VII.3) of integral agrees with the earlier one (VII.2) whenever both are defined.

**Lemma VII.5** (No conflict between the two first definitions of integral).

*For any non-negative simple function  $f \in \mathfrak{s}\mathcal{S}^+$  we have*

$$\int^{\square} f \, d\mu = \int^{+} f \, d\mu.$$

*Proof.* Assume that  $f \in \mathfrak{s}\mathcal{S}^+$ .

For any  $h \in \mathfrak{s}\mathcal{S}^+$  such that  $h \leq f$  we have by monotonicity, Lemma VII.3,

$$\int^{\square} h \, d\mu \leq \int^{\square} f \, d\mu.$$

By taking the supremum over such  $h$  as in the definition (VII.3), we get

$$\int^{+} f \, d\mu = \sup_{\substack{h \in \mathfrak{s}\mathcal{S}^+ \\ 0 \leq h \leq f}} \int^{\square} h \, d\mu \leq \int^{\square} f \, d\mu.$$

To prove the converse inequality, just consider  $h = f$  in the supremum (VII.3), to get

$$\int^{+} f \, d\mu = \sup_{\substack{h \in \mathfrak{s}\mathcal{S}^+ \\ 0 \leq h \leq f}} \int^{\square} h \, d\mu \geq \int^{\square} f \, d\mu.$$

□

## Properties of integral of non-negative measurable functions

The key properties of integrals continue to hold. Let us first verify monotonicity.

**Proposition VII.6** (Monotonicity of the integral for non-negative functions).

*Suppose that  $f, g \in \mathfrak{m}\mathcal{S}^+$  and  $f \leq g$  pointwise. Then we have*

$$\int^{+} f \, d\mu \leq \int^{+} g \, d\mu.$$

*Proof.* Whenever  $h \in \mathfrak{s}\mathcal{S}^+$  is such that  $h \leq f$ , we also have  $h \leq g$ , so the supremum in the definition of the integral of  $g$  is over a larger collection than the supremum in the definition of the integral of  $f$ . We get

$$\int^{+} f \, d\mu = \sup_{\substack{h \in \mathfrak{s}\mathcal{S}^+ \\ 0 \leq h \leq f}} \int^{\square} h \, d\mu \leq \sup_{\substack{h \in \mathfrak{s}\mathcal{S}^+ \\ 0 \leq h \leq g}} \int^{\square} h \, d\mu = \int^{+} g \, d\mu.$$

□

As a consequence, we get that the only way that a positive function can have vanishing integral is for the function to be equal to zero almost everywhere (i.e. except in a set of measure zero).

**Proposition VII.7** (No nontrivial non-negative function has vanishing integral).

If  $f \in \mathfrak{m}\mathcal{S}^+$  and

$$\int^+ f \, d\mu = 0$$

then

$$\mu\left[\{s \in S \mid f(s) > 0\}\right] = 0.$$

*Proof.* For  $n \in \mathbb{N}$ , define

$$A_n = \left\{s \in S \mid f(s) \geq \frac{1}{n}\right\} = f^{-1}\left(\left[\frac{1}{n}, +\infty\right)\right) \in \mathcal{S}.$$

We can then express the set of interest as

$$A = \left\{s \in S \mid f(s) > 0\right\} = \bigcup_{n=1}^{\infty} A_n.$$

Consider the non-negative simple function  $h = \frac{1}{n}\mathbb{I}_{A_n}$ . By construction of  $A_n$ , we have  $h \leq f$ , and we can therefore estimate the supremum that defines the integral of  $f$  from below by

$$\int^+ f \, d\mu \geq \int^{\square} \frac{1}{n}\mathbb{I}_{A_n} \, d\mu = \frac{1}{n}\mu[A_n].$$

Using the assumption that the integral of  $f$  vanishes, we get

$$0 \leq \mu[A_n] \leq n \int^+ f \, d\mu = 0,$$

so  $\mu[A_n] = 0$ . Finally, by the countable subadditivity (II.10) of Lemma II.19 (i.e., “the union bound”) we get the conclusion

$$\mu[A] = \mu\left[\bigcup_{n=1}^{\infty} A_n\right] \leq \sum_{n=1}^{\infty} \mu[A_n] = \sum_{n=1}^{\infty} 0 = 0.$$

□

The converse to Proposition VII.7 also holds.

**Exercise VII.2** (Almost everywhere vanishing functions have vanishing integrals).

Check directly from Definition VII.4 that if  $f \in \mathfrak{m}\mathcal{S}^+$  and

$$\mu\left[\{s \in S \mid f(s) > 0\}\right] = 0,$$

then we have

$$\int^+ f \, d\mu = 0.$$

At this stage, yet another fundamentally important property of integration arises, concerning the behavior of integrals under pointwise monotone approximation in the integrands (i.e., the functions to be integrated). The result, *Monotone convergence theorem (MCT)*, is at the same time

- “really all there is to integration theory”<sup>4</sup>
- a very practical tool for calculations.

<sup>4</sup>This is a quote from David Williams, [Wil91, Section 5.3]. One aspect of Williams’ message is certainly to stress the importance of the MCT. Another aspect might be to point out that besides the MCT, all other steps in the construction of integral and derivation of its basic properties are rather intuitive and straightforward. For this latter reason we have chosen to postpone the somewhat lengthy proof of the Monotone convergence theorem to Appendix D, so that it does not interrupt the flow of the rest of the steps.

**Theorem VII.8** (Monotone convergence theorem).

If  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}^+$  and  $f_n \uparrow f$  as  $n \rightarrow \infty$ , then we have

$$\int^+ f_n \, d\mu \uparrow \int^+ f \, d\mu \quad \text{as } n \rightarrow \infty.$$

The proof is given in Appendix D.

Recall that by Lemma III.18, any non-negative measurable function  $f \in \mathfrak{m}\mathcal{S}^+$  can be approximated by simple functions  $f_1, f_2, \dots \in \mathfrak{s}\mathcal{S}^+$  in a pointwise increasing way,  $f_n \uparrow f$  as  $n \rightarrow \infty$ . Together with the observation of Lemma VII.5, Monotone convergence theorem thus gives the expression

$$\int^+ f \, d\mu = \lim_{n \rightarrow \infty} \int^{\square} f_n \, d\mu$$

for the integral of  $f$  as an increasing limit of integrals (VII.2) of simple functions. This should reassure the reader that integrals, despite their slightly abstract definition, are rather innocent constructions, after all.

It is now easy to verify that also linearity continues to hold.

**Proposition VII.9** (Linearity of the integral for non-negative functions).

(a) If  $f \in \mathfrak{m}\mathcal{S}^+$  and  $c \geq 0$ , then  $cf \in \mathfrak{m}\mathcal{S}^+$  and

$$\int^+ (cf) \, d\mu = c \int^+ f \, d\mu.$$

(b) If  $f, g \in \mathfrak{m}\mathcal{S}^+$ , then  $f + g \in \mathfrak{m}\mathcal{S}^+$  and

$$\int^+ (f + g) \, d\mu = \int^+ f \, d\mu + \int^+ g \, d\mu.$$

*Proof.* Both parts can be easily verified using the Monotone convergence theorem<sup>5</sup> and the corresponding properties of integrals of simple functions.

Let us only prove part (b) — the proof of part (a) is similar. So let  $f, g \in \mathfrak{m}\mathcal{S}^+$ . As in Lemma III.18, pick sequences  $f_1, f_2, \dots \in \mathfrak{s}\mathcal{S}^+$  and  $g_1, g_2, \dots \in \mathfrak{s}\mathcal{S}^+$  of simple functions such that  $f_n \uparrow f$  and  $g_n \uparrow g$  as  $n \rightarrow \infty$ . Then by linearity of limits, we have also  $f_n + g_n \uparrow f + g$  as  $n \rightarrow \infty$ . Now apply the Monotone convergence theorem (Theorem VII.8) to each of these three monotone approximations, and use also the linearity of integral for simple functions

---

<sup>5</sup>The careful reader wary of circular reasoning will now inspect that the proof of the Monotone convergence theorem in Appendix D does not rely on linearity of  $\int^+$ , but instead uses only properties of  $\int^{\square}$  and the definition of  $\int^+$ .



(Lemma VII.2), to get

$$\begin{aligned}
 & \int^+ (f + g) \, d\mu \\
 &= \lim_{n \rightarrow \infty} \int^{\square} (f_n + g_n) \, d\mu && \text{(by MCT and } f_n + g_n \uparrow f + g) \\
 &= \lim_{n \rightarrow \infty} \left( \int^{\square} f_n \, d\mu + \int^{\square} g_n \, d\mu \right) && \text{(by Lemma VII.2)} \\
 &= \lim_{n \rightarrow \infty} \int^{\square} f_n \, d\mu + \lim_{n \rightarrow \infty} \int^{\square} g_n \, d\mu && \text{(by linearity of limits)} \\
 &= \int^+ f \, d\mu + \int^+ g \, d\mu && \text{(by MCT and } f_n \uparrow f \text{ and } g_n \uparrow g).
 \end{aligned}$$

□

### VII.3. Integral for integrable functions

For an arbitrary measurable function  $f \in m\mathcal{S}$ , define two non-negative functions  $f_+, f_- : S \rightarrow [0, +\infty]$  by

$$f_+(s) = \max \{f(s), 0\} \quad \text{and} \quad f_-(s) = \max \{-f(s), 0\}. \quad (\text{VII.4})$$

Then we have  $f_+, f_- \in m\mathcal{S}^+$  by Proposition III.14, and we can write

$$f = f_+ - f_- \quad \text{and} \quad |f| = f_+ + f_-.$$

We call  $f_+$  the *positive part* and  $f_-$  the *negative part* of  $f$ .

**Definition VII.10** (Integrable function).

We say that a function  $f \in m\mathcal{S}$  is *integrable* with respect to the measure  $\mu$  and denote  $f \in \mathcal{L}^1(\mu)$  if

$$\int^+ |f| \, d\mu < +\infty.$$

**Remark VII.11** (Equivalent condition for integrability of a function).

Since  $|f| = f_+ + f_-$  and the integral of non-negative functions is additive,  $\int^+ |f| \, d\mu = \int^+ f_+ \, d\mu + \int^+ f_- \, d\mu$ , a measurable function is integrable if and only if

$$\int^+ f_+ \, d\mu < +\infty \quad \text{and} \quad \int^+ f_- \, d\mu < +\infty.$$

#### Definition of integral of integrable functions

In view of the decomposition  $f = f_+ - f_-$ , our final definition of the integral is unsurprising.

**Definition VII.12** (Integral of an integrable function).

For  $f \in \mathcal{L}^1(\mu)$  we define the integral as

$$\int f \, d\mu = \int^+ f_+ \, d\mu - \int^+ f_- \, d\mu, \quad (\text{VII.5})$$

where  $f_+$  and  $f_-$  are the positive and negative parts of the function  $f$ .

By Remark VII.11 above, the two terms in formula (VII.5) are finite, so the expression is well defined.

For non-negative functions  $f \in \mathfrak{m}\mathcal{S}^+$  the new definition of integral coincides with the earlier one.

**Lemma VII.13** (No conflict between the two last definitions of integral).

*For any non-negative integrable function  $f \in \mathfrak{m}\mathcal{S}^+ \cap \mathcal{L}^1(\mu)$  we have*

$$\int f \, d\mu = \int^+ f \, d\mu.$$

*Proof.* For a non-negative function  $f \in \mathfrak{m}\mathcal{S}^+$  we have  $f_+ = f$  and  $f_- = 0$ . Therefore the agreement of the two integrals is obvious from the definition (VII.5).  $\square$

### Properties of integral of integrable functions

Now that the functions involved may have both positive and negative signs, estimating integrals by the following triangle inequality becomes important.

**Theorem VII.14** (Triangle inequality for integrals).

*For any integrable function  $f \in \mathcal{L}^1(\mu)$  we have*

$$\left| \int f \, d\mu \right| \leq \int |f| \, d\mu.$$

*Proof.* Observe first that the right hand side is

$$\int |f| \, d\mu = \int^+ |f| \, d\mu = \int^+ (f_+ + f_-) \, d\mu = \int^+ f_+ \, d\mu + \int^+ f_- \, d\mu.$$

We will prove  $\int f \, d\mu \leq \int |f| \, d\mu$  and  $-\int f \, d\mu \leq \int |f| \, d\mu$ , which together imply the assertion  $|\int f \, d\mu| \leq \int |f| \, d\mu$ .

By (VII.5) and the non-negativity of integral  $\int f_- \, d\mu$  of the non-negative function  $f_-$ , we have

$$\begin{aligned} \int f \, d\mu &= \int^+ f_+ \, d\mu - \int^+ f_- \, d\mu \\ &\leq \int^+ f_+ \, d\mu + \int^+ f_- \, d\mu = \int |f| \, d\mu. \end{aligned}$$

Similarly, we have

$$\begin{aligned} -\int f \, d\mu &= -\int^+ f_+ \, d\mu + \int^+ f_- \, d\mu \\ &\leq \int^+ f_+ \, d\mu + \int^+ f_- \, d\mu = \int |f| \, d\mu. \end{aligned}$$

$\square$

In particular, in proving linearity, we use the triangle inequality to check that the sum of integrable functions is integrable in the first place.

**Theorem VII.15** (Linearity of the integral).

(a) For any  $f \in \mathcal{L}^1(\mu)$  and  $c \in \mathbb{R}$  we have  $cf \in \mathcal{L}^1(\mu)$  and

$$\int cf \, d\mu = c \int f \, d\mu.$$

(b) For any  $f, g \in \mathcal{L}^1(\mu)$  we have  $f + g \in \mathcal{L}^1(\mu)$  and

$$\int (f + g) \, d\mu = \int f \, d\mu + \int g \, d\mu.$$

*Proof of part (a).* If  $c \geq 0$  then we have  $(cf)_+ = cf_+$  and  $(cf)_- = cf_-$ . Integrability of  $cf$  follows easily by Remark VII.11, since

$$\int^+ (cf)_+ \, d\mu = \int^+ cf_+ \, d\mu = c \underbrace{\int^+ f_+ \, d\mu}_{<+\infty} < +\infty$$

and similarly  $\int^+ (cf)_- \, d\mu < +\infty$ . The asserted formula is also a direct consequence of the definition (VII.5) of  $\int$ , and linearity of  $\int^+$  (Proposition VII.9),

$$\int^+ cf \, d\mu = \int^+ cf_+ \, d\mu - \int^+ cf_- \, d\mu = c \left( \int^+ f_+ \, d\mu - \int^+ f_- \, d\mu \right) = c \int f \, d\mu.$$

If  $c < 0$  then we have  $(cf)_+ = -cf_-$  and  $(cf)_- = -cf_+$ , and one can proceed similarly.

*Proof of part (b).* Note that  $|f + g| \leq |f| + |g|$ . If  $f, g \in \mathcal{L}^1(\mu)$ , we thus have

$$\int^+ |f + g| \, d\mu \leq \int^+ (|f| + |g|) \, d\mu = \int^+ |f| \, d\mu + \int^+ |g| \, d\mu < +\infty,$$

and therefore  $f + g \in \mathcal{L}^1(\mu)$ . The formula for the integral of  $f + g$  can be proven by carefully considering which possibilities contribute to the positive and negative parts of  $f + g$ .  $\square$

Also monotonicity continues to hold.

**Theorem VII.16** (Monotonicity of the integral).

If  $f, g \in \mathcal{L}^1(\mu)$  and  $f \leq g$ , then we have

$$\int f \, d\mu \leq \int g \, d\mu.$$

*Proof.* If we have the pointwise inequality  $f \leq g$ , then the positive and negative parts of these functions satisfy  $f_+ \leq g_+$  and  $f_- \geq g_-$ . It therefore follows from the defining formula (VII.5) of  $\int$ , and monotonicity of  $\int^+$  (Proposition VII.6), that

$$\begin{aligned} \int f \, d\mu &= \int^+ f_+ \, d\mu - \int^+ f_- \, d\mu \\ &\leq \int^+ g_+ \, d\mu - \int^+ g_- \, d\mu = \int g \, d\mu. \end{aligned}$$

$\square$

## VII.4. Convergence theorems for integrals

Finally, we establish some of the powerful standard tools of measure theory, which address when it is legitimate to interchange the order of a limit and an integration. Recalling that integration contains as special cases also expected values and

summation, we thus ask under what conditions we have

$$\begin{aligned}\lim_{n \rightarrow \infty} \int f_n(s) \, d\mu(s) &\stackrel{?}{=} \int \left( \lim_{n \rightarrow \infty} f_n(s) \right) \, d\mu(s) \\ \lim_{n \rightarrow \infty} \sum_{k=0}^{\infty} a_k^{(n)} &\stackrel{?}{=} \sum_{k=0}^{\infty} \left( \lim_{n \rightarrow \infty} a_k^{(n)} \right) \\ \lim_{n \rightarrow \infty} \mathbb{E}[X_n] &\stackrel{?}{=} \mathbb{E} \left[ \lim_{n \rightarrow \infty} X_n \right].\end{aligned}$$

In particular, derivatives are limits of difference quotients, so the results can be applied to justifying the interchange of differentiation and integral, sum, or expected value, i.e.

$$\begin{aligned}\frac{d}{d\lambda} \int f_n(s; \lambda) \, d\mu(s) &\stackrel{?}{=} \int \frac{d}{d\lambda} f_n(s; \lambda) \, d\mu(s) \\ \frac{d}{d\lambda} \sum_{k=0}^{\infty} a_k(\lambda) &\stackrel{?}{=} \sum_{k=0}^{\infty} \left( \frac{d}{d\lambda} a_k(\lambda) \right) \\ \frac{d}{d\lambda} \mathbb{E}[f(X_n; \lambda)] &\stackrel{?}{=} \mathbb{E} \left[ \frac{d}{d\lambda} f(X_n; \lambda) \right].\end{aligned}$$

To fully appreciate the positive results we will derive, it is advisable to first think about the ways the conclusion could fail.

**Exercise VII.3** (Interchanging limit and integration is not always possible).

- (a) Consider the Lebesgue measure  $\Lambda$  on  $\mathbb{R}$ . For  $n \in \mathbb{N}$ , define  $f_n: \mathbb{R} \rightarrow \mathbb{R}$  by

$$f_n(x) = \begin{cases} n & \text{if } 0 < x \leq \frac{1}{n} \\ 0 & \text{otherwise.} \end{cases}$$

Calculate the following

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} f_n \, d\Lambda, \quad \int_{\mathbb{R}} \left( \lim_{n \rightarrow \infty} f_n \right) \, d\Lambda, \quad \lim_{n \rightarrow \infty} \int_{\mathbb{R}} f_n \, d\Lambda.$$

- (b) Repeat the calculations of part (a) for the functions  $\sqrt{f_n}$  and  $f_n^2$ , instead.  
(c) Consider  $\mathbb{Z}$  with the counting measure  $\mu_{\#}$ . For  $n \in \mathbb{N}$ , define  $g_n: \mathbb{Z} \rightarrow \mathbb{R}$  by

$$g_n(k) = \begin{cases} 1 & \text{if } k = n \\ 0 & \text{if } k \neq n. \end{cases}$$

Calculate the following

$$\lim_{n \rightarrow \infty} \int_{\mathbb{Z}} g_n \, d\mu_{\#}, \quad \int_{\mathbb{Z}} \left( \lim_{n \rightarrow \infty} g_n \right) \, d\mu_{\#}, \quad \lim_{n \rightarrow \infty} \int_{\mathbb{Z}} g_n \, d\mu_{\#}.$$

We will obtain several results about what conditions are sufficient to permit changing limits and integration. Besides

(MCT) Monotone convergence theorem (Theorem VII.8)

another such result of great importance is

(DCT) Dominated convergence theorem (Theorem VII.19)

For expected values, there is a very practical special case of the Dominated convergence theorem:

(BCT) Bounded convergence theorem (Corollary VII.21)

The Monotone convergence theorem was already introduced as a tool in step (3) of the construction. The proof of the Dominated convergence theorem relies on it, and another intermediate result that is occasionally useful on its own right as well:

Fatou's Lemma (Lemma VII.17)

### Fatou's lemma

Fatou's lemma is a general convergence result for non-negative integrals, which does not even require the limits to exist! It instead addresses lower limits which always exist. A minor drawback is that the conclusion is merely an inequality in one direction.

**Lemma VII.17** (Fatou's lemma).

*For any sequence  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}^+$  of non-negative measurable functions, we have*

$$\int \left( \liminf_n f_n \right) d\mu \leq \liminf_n \int f_n d\mu. \quad (\text{VII.6})$$

*Proof.* Define  $g = \liminf_n f_n$ . For  $k \in \mathbb{N}$ , define also

$$g_k := \inf_{n \geq k} f_n.$$

Then we have  $g_k \in \mathfrak{m}\mathcal{S}^+$  (by Proposition III.14), and  $g_k \uparrow g$  as  $k \rightarrow \infty$  by construction. By the Monotone convergence theorem, we therefore have

$$\int g_k d\mu \uparrow \int g d\mu \quad \text{as } k \rightarrow \infty.$$

We have  $f_n \geq g_k$  for any  $n \geq k$ , so monotonicity of integral gives  $\int f_n d\mu \geq \int g_k d\mu$  for  $n \geq k$ , and thus also

$$\inf_{n \geq k} \int f_n d\mu \geq \int g_k d\mu.$$

In this inequality, take the limit as  $k \rightarrow \infty$  to obtain

$$\liminf_n \int f_n d\mu \geq \lim_{k \rightarrow \infty} \int g_k d\mu = \int g d\mu,$$

which is the asserted inequality. □

**Exercise VII.4** (Strict inequality in Fatou's lemma).

Find an example of a sequence of non-negative functions for which there is a strict inequality in Fatou's lemma, (VII.6).

Under suitable additional conditions, we get a result about upper limits of non-negative integrals as well.

**Lemma VII.18** (Reverse Fatou's lemma).

*Suppose that  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}^+$  is a sequence of non-negative measurable functions such that there exists a non-negative measurable function  $g \in \mathfrak{m}\mathcal{S}^+$  which uniformly bounds the sequence,  $f_n \leq g$  for all  $n \in \mathbb{N}$ , and which is itself integrable,  $\int g d\mu < +\infty$ . Then we have*

$$\limsup_n \int f_n d\mu \leq \int \left( \limsup_n f_n \right) d\mu. \quad (\text{VII.7})$$

*Proof.* Apply Fatou's lemma to the sequence of functions  $g - f_n$ ,  $n \in \mathbb{N}$ , and cancel the finite number  $\int g \, d\mu$  from both sides.  $\square$

### Dominated convergence theorem

Arguably the most practical general convergence result is Lebesgue's dominated convergence theorem: it states that pointwise convergence and uniform boundedness of integrands by an integrable function are sufficient to permit taking the limit inside the integral.

**Theorem VII.19** (Dominated convergence theorem).

*Suppose that  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}$  is a sequence of measurable functions such that there exists a non-negative measurable function  $g \in \mathfrak{m}\mathcal{S}^+$  which uniformly bounds the absolute values of the sequence,  $|f_n| \leq g$  for all  $n \in \mathbb{N}$ , and which is itself integrable,  $\int g \, d\mu < +\infty$ . Then if the pointwise limit  $f = \lim_{n \rightarrow \infty} f_n$  exists, we have*

$$\lim_{n \rightarrow \infty} \int |f_n - f| \, d\mu = 0 \quad (\text{VII.8})$$

$$\text{and} \quad \lim_{n \rightarrow \infty} \int f_n \, d\mu = \int f \, d\mu. \quad (\text{VII.9})$$

*Proof.* Since  $|f_n| \leq g$  for all  $n \in \mathbb{N}$ , and  $f_n \rightarrow f$  as  $n \rightarrow \infty$ , we also have  $|f| \leq g$ . Therefore the triangle inequality gives

$$|f_n - f| \leq |f_n| + |f| \leq 2g.$$

We have  $\int 2g < +\infty$  by integrability of  $g$ , so reverse Fatou's lemma can be applied to the sequence of functions  $|f_n - f|$ ,  $n \in \mathbb{N}$ . Together with the assumed pointwise convergence  $\lim_{n \rightarrow \infty} |f_n - f| = 0$ , the reverse Fatou's lemma gives

$$\limsup_n \int |f_n - f| \, d\mu \leq \int (\limsup_n |f_n - f|) \, d\mu = \int 0 \, d\mu = 0,$$

which proves the first assertion (VII.8).

Then use linearity and triangle inequality for integrals to get

$$\left| \int f_n \, d\mu - \int f \, d\mu \right| = \left| \int (f_n - f) \, d\mu \right| \leq \int |f_n - f| \, d\mu.$$

The right hand side tends to zero as  $n \rightarrow \infty$  by the first part, (VII.8). This proves the second assertion (VII.9).  $\square$

### Bounded convergence theorem

For finite measures (see Definition II.9) — and thus probability measures, in particular — there is a very easy and practical special case of the Dominated convergence theorem, known as the Bounded convergence theorem. The underlying reason for it, as well as many other simplifications in finite measure spaces, is the following.

**Lemma VII.20** (Constant functions are integrable on finite measure spaces).

*Suppose that  $(S, \mathcal{S}, \mu)$  is a finite measure space. Let  $c \in \mathbb{R}$  be a constant, and let  $g: S \rightarrow \mathbb{R}$  be the constant function  $g(s) = c$  for all  $s \in S$ . Then  $g \in \mathcal{L}^1(\mu)$ .*

*Proof.* The constant function  $g(s) = c$  is in particular a simple function,  $g = c\mathbb{I}_S$ , so we have

$$\int^+ |g| d\mu = \int^\square |c| d\mu = |c| \underbrace{\mu[S]}_{<+\infty} < +\infty.$$

□

In particular, constant functions are legitimate choices for the dominating function  $g$  in the Dominated convergence theorem, Theorem VII.19. This immediately yields the following.

**Corollary VII.21** (Bounded convergence theorem).

*Suppose that  $(S, \mathcal{S}, \mu)$  is a finite measure space and  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}$  is a bounded sequence<sup>6</sup> of measurable functions such that the pointwise limit function  $f = \lim_{n \rightarrow \infty} f_n$  exists. Then we have*

$$\lim_{n \rightarrow \infty} \int f_n d\mu = \int f d\mu.$$

## VII.5. Integrals over subsets and restriction of measures

Let  $(S, \mathcal{S}, \mu)$  be a measure space.

We adopt the following notation for *integrals over subsets*: if  $f: S \rightarrow \mathbb{R}$  is a function and if  $R \subset S$  is a measurable subset,  $R \in \mathcal{S}$ , then the integral of  $f$  over  $R$  is defined as

$$\int_R f(s) d\mu(s) := \int_S \mathbb{I}_R(s) f(s) d\mu(s) \quad (\text{VII.10})$$

whenever either  $\mathbb{I}_R f \geq 0$  or  $\mathbb{I}_R f \in \mathcal{L}^1(\mu)$ .

A clear justification for the above convention is obtained in the following exercise.

**Exercise VII.5** (Restriction of a measure to a subset).

Let  $(S, \mathcal{S}, \mu)$  is a measure space and  $R \subset S$  a measurable subset,  $R \in \mathcal{S}$ .

(a) Let  $\mathcal{S}'$  consist of those  $A \subset R$  such that  $A = B \cap R$  for some  $B \in \mathcal{S}$ , that is,

$$\mathcal{S}' := \{B \cap R \mid B \in \mathcal{S}\}.$$

Show that  $\mathcal{S}'$  is a  $\sigma$ -algebra on  $R$ .

(b) When  $A = B \cap R$  for some  $B \in \mathcal{S}$ , set

$$\mu'[A] := \mu[B \cap R].$$

Show that  $\mu'$  is a measure on the measurable space  $(R, \mathcal{S}')$ .

(c) For any function  $f: S \rightarrow \mathbb{R}$ , denote by  $f|_R: R \rightarrow \mathbb{R}$  the restriction of  $f$  to  $R$ . Prove that we have

$$\int_R f|_R d\mu' = \int_S \mathbb{I}_R f d\mu,$$

for (1) indicator functions  $f$ , (2) non-negative simple functions  $f$ , (3) non-negative measurable functions  $f$ , and finally for (4) all measurable functions  $f$  such that  $f|_R \in \mathcal{L}^1(\mu')$ .

<sup>6</sup>This means that there exists some  $c \in [0, +\infty)$  such that  $|f_n(s)| \leq c$  for all  $s \in S$  and  $n \in \mathbb{N}$ .

### VII.6. Riemann integral vs. Lebesgue integral

The reader can now verify that the Riemann integral (familiar from undergraduate calculus) and Lebesgue integral coincide in the following setups.

**Exercise VII.6** (Riemann integral and Lebesgue integral on closed intervals).

Consider a closed interval  $[a, b] \subset \mathbb{R}$ . Suppose that  $f: [a, b] \rightarrow \mathbb{R}$  is a continuous function. Show that

$$\int_a^b f(x) dx = \int_{[a,b]} f(x) d\Lambda(x),$$

where the right hand side denotes the integral of  $f$  with respect to the Lebesgue measure  $\Lambda$  restricted to  $[a, b] \subset \mathbb{R}$ , and the left hand side denotes the Riemann integral of  $f$ .

**Hint:** The Riemann integral is defined by the condition that the upper and lower Riemann sums associated to subdivisions of the interval  $[a, b]$  tend to the same limit as the meshes of the subdivisions become finer. First show that these upper and lower Riemann sums associated to a subdivision are equal to the integrals of certain simple functions. Then use monotonicity to compare the Lebesgue integral of  $f$  to these.

**Exercise VII.7** (Improper Riemann integrals and Lebesgue integral).

- (a) Suppose that  $f: [a, b) \rightarrow \mathbb{R}$  is continuous. Show that if  $\int_{[a,b)} |f| d\Lambda < +\infty$ , then we have the following equality of the improper Riemann integral and the Lebesgue integral:

$$\int_a^b f(x) dx := \lim_{\varepsilon \downarrow 0} \int_a^{b-\varepsilon} f(x) dx = \int_{[a,b)} f(x) d\Lambda(x).$$

- (b) Suppose that  $f: \mathbb{R} \rightarrow \mathbb{R}$  is continuous. Show that if  $\int_{\mathbb{R}} |f| d\Lambda < +\infty$ , then we have the following equality of the improper Riemann integral and the Lebesgue integral:

$$\int_{-\infty}^{+\infty} f(x) dx := \lim_{r \uparrow +\infty} \int_{-r}^r f(x) dx = \int_{\mathbb{R}} f(x) d\Lambda(x).$$

**Hint:** Exercise VII.6 and Dominated convergence theorem.

**Remark VII.22** (On the relationship between Riemann-integral and Lebesgue integral).

- There are plenty of functions for which the Riemann integral does not exist, while the Lebesgue integral is well-defined. A simple example is the indicator function  $\mathbb{I}_{\mathbb{Q} \cap [a,b]}$  of the rational numbers on the interval. The Lebesgue integrals  $\int \mathbb{I}_{\mathbb{Q} \cap [a,b]} d\Lambda = 0$  and  $\int \mathbb{I}_{[a,b] \setminus \mathbb{Q}} d\Lambda = b - a$  reflect the fact that the countable set  $\mathbb{Q} \cap [a, b]$  is negligible, in a measure-theoretic sense, compared to its complement  $[a, b] \setminus \mathbb{Q}$  (which, in particular, is uncountable). The Riemann integral fails to account for this “evident” size difference.
- Pointwise limits of Riemann integrable functions need not be Riemann integrable, so one could not hope to have the general and powerful convergence theorems of Section VII.4 for the Riemann integral.<sup>7</sup>
- The coincidence of Riemann integral and Lebesgue integral for well-behaved functions allows us to use the familiar techniques such as the fundamental theorem of calculus, integration by parts, etc., without the need to prove them again.
- Given how common integration on the real axis is, it is convenient to have a concise notation for it. By virtue of the coincidence of the two notions of integration on  $\mathbb{R}$  for all well-behaved functions, with almost no risk of confusion we can use the notation

$$dx := d\Lambda(x)$$

for integration with respect to Lebesgue measure  $\Lambda$ .

<sup>7</sup>Besides the generality (at once we constructed, e.g., integrals in all dimensions, expected values, infinite sums, ...), arguably the most significant advantage of the integration theory developed in this chapter is indeed the convergence theorems.



## Lecture VIII

### Expected values

In this lecture we begin examining the role and use of expected values in probability theory.

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. We will use the notation  $\mathbf{E}$  for expected values. The *expected value* of a real-valued random variable  $X: \Omega \rightarrow \mathbb{R}$  is just the integral with respect to the probability measure  $\mathbf{P}$ ,

$$\mathbf{E}[X] := \int_{\Omega} X(\omega) \, d\mathbf{P}(\omega).$$

Let us summarize how it is constructed step by step for increasingly general random variables  $X: \Omega \rightarrow \mathbb{R}$  as in the previous lecture (Lecture VII):

- (1) If  $X = \mathbb{I}_A$  for an event  $A \in \mathcal{F}$ , then we set  $\mathbf{E}[X] = \mathbf{P}[A]$ .
- (2) If  $X = \sum_{j=1}^n a_j \mathbb{I}_{A_j}$  is simple<sup>1</sup>, then we set  $\mathbf{E}[X] = \sum_{j=1}^n a_j \mathbf{P}[A_j]$ .
- (3) If  $X$  is non-negative, then we set

$$\mathbf{E}[X] = \sup_{\substack{H \in \mathfrak{s}\mathcal{F}^+ \\ 0 \leq H \leq X}} \mathbf{E}[H].$$

- (4) If  $X$  is integrable, i.e.,  $\mathbf{E}[X_+] < +\infty$  and  $\mathbf{E}[X_-] < +\infty$ , then we set

$$\mathbf{E}[X] = \mathbf{E}[X_+] - \mathbf{E}[X_-].$$

Since expected value is by definition a special case of integral, the basic properties of integrals lead to the corresponding basic properties of expected values:

**Linearity:** (Proposition VII.9 and Theorem VII.15)

For either  $X, Y \in \mathfrak{m}\mathcal{F}^+$  (two non-negative random variables), or for  $X, Y \in \mathcal{L}^1(\mathbf{P})$  (two integrable random variables), we have

$$\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y].$$

For either  $X \in \mathfrak{m}\mathcal{F}^+$  and  $c \geq 0$ , or for  $X \in \mathcal{L}^1(\mathbf{P})$  and  $c \in \mathbb{R}$ , we have

$$\mathbf{E}[cX] = c \mathbf{E}[X].$$

**Monotonicity:** (Proposition VII.6 and Theorem VII.16)

For either  $X, Y \in \mathfrak{m}\mathcal{F}^+$  (two non-negative random variables) or for  $X, Y \in \mathcal{L}^1(\mathbf{P})$  (two integrable random variables), we have

$$X \leq Y \quad \implies \quad \mathbf{E}[X] \leq \mathbf{E}[Y]$$

---

<sup>1</sup>In Section VII.1 we were careful to only integrate simple functions which are non-negative, in order to avoid undefined expressions of the form  $+\infty - \infty$ . For expected values, because of finiteness  $\mathbf{P}[\Omega] = 1 < +\infty$  of the total mass of a probability measure, this never causes problems, so formula (VII.2) can be directly taken as a definition of the expected value of any simple function (random variable).

**Triangle inequality:**

(Theorem VII.14)

For  $X \in \mathcal{L}^1(\mathbb{P})$ , we have

$$\left| \mathbb{E}[X] \right| \leq \mathbb{E}[|X|].$$

Besides these basic properties, also the powerful convergence theorems of Section VII.4 continue hold for expected values.

**VIII.1. Expected values in terms of laws**

Recall from Definition III.6 that the law of a random variable  $X: \Omega \rightarrow \mathbb{R}$  (or the distribution of  $X$ ) is the probability measure  $P_X$  on  $\mathbb{R}$  given by

$$P_X[B] = \mathbb{P}[X \in B] \quad \text{for } B \in \mathcal{B}.$$

We next show how the expected value of  $g(X)$  can be calculated as an integral with respect to the law  $P_X$ .

**Theorem VIII.1** (Expected values in terms of laws).

Let  $X: \Omega \rightarrow \mathbb{R}$  be a random variable with law  $P_X$ , and let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be a Borel function. Then we have the equivalence of the following integrability properties

$$g(X) \in \mathcal{L}^1(\mathbb{P}) \quad \iff \quad g \in \mathcal{L}^1(P_X).$$

If either (then both) of the above integrability properties holds, then we have

$$\mathbb{E}[g(X)] = \int_{\mathbb{R}} g(x) \, dP_X(x).$$

*Proof.* The proof follows the “standard machine”, i.e., we verify the statement for (1) indicator functions, (2) simple functions, (3) non-negative measurable functions, and finally (4) all measurable functions  $g$ .

*step 1:* Consider an indicator function  $g = \mathbb{I}_B$  of a Borel set  $B \in \mathcal{B}$ . Observe, as in (IV.1) and (IV.2), that  $\mathbb{I}_B(X(\omega)) = \mathbb{I}_{X^{-1}(B)}(\omega)$  for all  $\omega \in \Omega$ . Thus we have

$$\begin{aligned} \mathbb{E}[\mathbb{I}_B(X)] &= \mathbb{E}[\mathbb{I}_{X^{-1}(B)}] && \text{(because } \mathbb{I}_B \circ X = \mathbb{I}_{X^{-1}(B)}) \\ &= \mathbb{P}[X^{-1}(B)] && \text{(by step 1 of def. of } \mathbb{E}) \\ &= P_X[B] && \text{(by def. of } P_X) \\ &= \int \mathbb{I}_B \, dP_X. && \text{(by step 1 of def. of integral)} \end{aligned}$$

This shows that  $\mathbb{E}[g(X)] = \int g \, dP_X$  for  $g = \mathbb{I}_B$ .

*step 2:* Consider a simple Borel function  $g \in \mathcal{s}\mathcal{B}^+$ . Write  $g = \sum_{j=1}^n a_j \mathbb{I}_{B_j}$ . We have

$$\begin{aligned} \mathbb{E}\left[\sum_{j=1}^n a_j \mathbb{I}_{B_j}(X)\right] &= \sum_{j=1}^n a_j \mathbb{E}[\mathbb{I}_{B_j}(X)] && \text{(by linearity of } \mathbb{E}) \\ &= \sum_{j=1}^n a_j \int \mathbb{I}_{B_j} \, dP_X && \text{(by step 1)} \\ &= \int \left(\sum_{j=1}^n a_j \mathbb{I}_{B_j}\right) \, dP_X && \text{(by linearity of integral).} \end{aligned}$$

This shows that  $\mathbb{E}[g(X)] = \int g \, dP_X$  for  $g \in \mathcal{s}\mathcal{B}^+$ .

*step 3:* Consider a non-negative Borel function  $g \in \mathfrak{m}\mathcal{B}^+$ . Take a monotone increasing pointwise approximation  $g_n \uparrow g$  as  $n \rightarrow \infty$  by non-negative simple Borel functions  $g_n \in \mathfrak{s}\mathcal{B}^+$ . We then also have  $g_n(X) \uparrow g(X)$ , and  $g_n(X) \in \mathfrak{s}\mathcal{F}^+$  are non-negative simple random variables. Using Monotone convergence theorem (Theorem VII.8) both for  $\mathbf{P}$  and for  $P_X$ , we get

$$\begin{aligned} \mathbb{E}[g(X)] &= \lim_{n \rightarrow \infty} \mathbb{E}[g_n(X)] && \text{(by MCT for } \mathbf{P} \text{)} \\ &= \lim_{n \rightarrow \infty} \int g_n \, dP_X && \text{(by step 2)} \\ &= \int g \, dP_X. && \text{(by MCT for } P_X \text{).} \end{aligned}$$

*step 4:* Consider a Borel function  $g \in \mathfrak{m}\mathcal{B}$ , and let  $g_+, g_- \in \mathfrak{m}\mathcal{B}^+$  be its positive and negative parts. The positive and negative parts of the random variable  $g(X) = g \circ X$  are  $(g(X))_+ = g_+ \circ X \in \mathfrak{m}\mathcal{F}^+$  and  $(g(X))_- = g_- \circ X \in \mathfrak{m}\mathcal{F}^+$ . By step 3 we thus have

$$\begin{aligned} \mathbb{E}[(g(X))_+] &= \mathbb{E}[g_+(X)] = \int g_+ \, dP_X \\ \text{and } \mathbb{E}[(g(X))_-] &= \mathbb{E}[g_-(X)] = \int g_- \, dP_X. \end{aligned}$$

Therefore we first of all have the equivalent conditions

$$\begin{aligned} \mathbb{E}[(g(X))_+] < +\infty &\iff \int g_+ \, dP_X < +\infty \\ \text{and } \mathbb{E}[(g(X))_-] < +\infty &\iff \int g_- \, dP_X < +\infty, \end{aligned}$$

which by Remark VII.11 shows the equivalence  $g(X) \in \mathcal{L}^1(\mathbf{P}) \iff g \in \mathcal{L}^1(P_X)$ . Moreover, when this integrability holds, step 4 of the definition of the integrals gives

$$\begin{aligned} \mathbb{E}[g(X)] &= \mathbb{E}[(g(X))_+] - \mathbb{E}[(g(X))_-] \\ \text{and } \int g \, dP_X &= \int g_+ \, dP_X - \int g_- \, dP_X. \end{aligned}$$

These coincide by the above equalities, which finishes the proof.  $\square$

**Exercise VIII.1** (Discrete random numbers).

A random variable is *discrete* if its range  $A = X(\Omega)$  is finite or countably infinite. The probability mass function of a discrete random variable  $X$  is defined by  $p_X(x) = \mathbf{P}[X = x]$ . Prove that any discrete real-valued random variable satisfies:

- $\mathbb{E}[h(X)] = \sum_{x \in A} h(x) p_X(x)$  for all Borel functions  $h : \mathbb{R} \rightarrow [0, \infty)$ .
- $h(X) \in \mathcal{L}^1(\Omega, \mathcal{F}, \mathbf{P})$  if and only if  $\sum_{x \in A} |h(x)| p_X(x) < \infty$ .
- Explain why the formula in (a) is true for all  $h \in \mathcal{L}^1(\mathbb{R}, \mathcal{B}, P_X)$ .

**Hint:** Recall the argument in the proof of Theorem VIII.1.

**Exercise VIII.2** (Almost sure equality).

Let  $X$  and  $Y$  be real-valued random variables that are equal almost surely, that is,

$$\mathbf{P}[X = Y] = 1.$$

- In order to make sure that  $\mathbf{P}[X = Y]$  is a meaningful probability, explain first why the set  $\{\omega \in \Omega \mid X(\omega) = Y(\omega)\}$  is measurable.
- Prove that the laws of  $X$  and  $Y$  are the same, and conclude that we in particular have  $\mathbb{E}[X] = \mathbb{E}[Y]$  (whenever the expected values exist).

### Probability densities of continuous distributions

**Definition VIII.2** (Probability density of a continuous distribution).

Let  $X \in \mathfrak{m}\mathcal{F}$  be a real valued random variable. If there exists a Borel function

$$f_X : \mathbb{R} \rightarrow [0, +\infty)$$

such that<sup>2</sup>

$$P_X[B] := \mathbb{P}[X \in B] = \int_B f_X(x) dx \quad (\text{VIII.1})$$

for all  $B \in \mathcal{B}$ , then we say that  $X$  has a *continuous distribution* (or a *continuous law*), and we say that  $f_X$  is a *density function* of  $X$ .

**Example VIII.3** (Gaussian distribution).

Let  $\mathfrak{m} \in \mathbb{R}$  and  $\mathfrak{s} > 0$ . A random variable  $X$  is said to have a *gaussian distribution* with mean  $\mathfrak{m}$  and variance  $\mathfrak{s}^2$  if its distribution is continuous and

$$f_X(x) = \frac{1}{\sqrt{2\pi \mathfrak{s}^2}} \exp\left(-\frac{1}{2\mathfrak{s}^2}(x - \mathfrak{m})^2\right)$$

is a density function of  $X$ .

Gaussian distributions are also called *normal distributions*, and the particular case of zero mean  $\mathfrak{m} = 0$  and unit variance  $\mathfrak{s}^2 = 1$  is called the *standard normal distribution*.

**Example VIII.4** (Exponential distribution).

Let  $\lambda > 0$ . A random variable  $X$  is said to have an *exponential distribution* with parameter  $\lambda$  if its distribution is continuous and

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$$

is a density function of  $X$ .

**Exercise VIII.3** (Random numbers with continuous distribution).

Assume that  $X$  has a continuous law with a density function  $f_X$ .

- (a) Show that  $X$  is integrable if and only if  $\int_{\mathbb{R}} |x| f_X(x) dx < \infty$ .
- (b) If  $X$  is integrable, show that expectation of  $X$  is given by

$$\mathbb{E}[X] = \int_{\mathbb{R}} x f_X(x) dx.$$

- (c) If  $X$  is integrable with expected value  $\mathfrak{m} := \mathbb{E}[X]$ , show that the variance of  $X$  can be computed as

$$\text{Var}(X) := \mathbb{E}[(X - \mathbb{E}[X])^2] = \int_{\mathbb{R}} (x - \mathfrak{m})^2 f_X(x) dx.$$

- (d) Can a random variable with continuous law have more than one density function?

**Hint:** Remember the hint for Exercise VIII.1, and seek a unified argument leading to (a) and (b) — perhaps even (c).

**Exercise VIII.4** (Expected value and variance of an exponential random variable).

Assume that a random variable  $X$  has an exponential distribution with parameter  $\lambda$  (see Example VIII.4). Use the previous exercise to calculate the expected value  $\mathbb{E}[X]$  and the variance  $\text{Var}(X)$  of  $X$ .

<sup>2</sup>The right hand side of (VIII.1) is an integral over the subset  $B \subset \mathbb{R}$  in the sense of (VII.10),

$$\int_B f_X(x) dx := \int_{\mathbb{R}} f_X(x) \mathbb{I}_B(x) d\Lambda(x).$$

**Exercise VIII.5** (The absolute value of a random number).

Let  $Y = |X|$  where  $X$  is a real-valued random variable.

- Prove that  $Y$  is a random variable.
- Assume that we know the cumulative distribution function  $F_X(x) = \mathbb{P}[X \leq x]$  of  $X$ . What is the cumulative distribution function of  $Y$ ?
- Assume that  $X$  has continuous distribution with a density function  $f_X(x)$ . Does  $Y$  also have a continuous distribution in this case? If yes, write down an expression for a density function  $f_Y$  of  $Y$  in terms of  $f_X$ . If not, explain why not.

## VIII.2. Applications of convergence theorems for expected values

In this section we look at a few consequences of the convergence theorems in various aspects of probability theory.

### Expected values of random series with non-negative terms

Let us observe that a non-negative random variable can only have finite expected value if the random variable is almost surely finite.<sup>3</sup>

**Lemma VIII.5** (Finite expected value implies almost sure finiteness).

If  $X \in \mathfrak{m}\mathcal{F}^+$  is a non-negative random variable such that  $\mathbb{E}[X] < +\infty$ , then we have  $X < +\infty$  almost surely, i.e.,  $\mathbb{P}[X < +\infty] = 1$ .

*Proof.* Let  $A = \{\omega \in \Omega \mid X(\omega) = +\infty\}$  be the event that the random variable  $X$  takes an infinite value. Then for any  $n \in \mathbb{N}$  we obviously have  $X \geq n \mathbb{I}_A$ , so by monotonicity of expected values we get

$$n \mathbb{P}[A] = \mathbb{E}[n \mathbb{I}_A] \leq \mathbb{E}[X].$$

Dividing this by  $n$ , we get  $\mathbb{P}[A] \leq \frac{1}{n} \mathbb{E}[X]$ . By assumption  $X$  has finite expected value,  $\mathbb{E}[X] < +\infty$ , so by letting  $n \rightarrow \infty$  we get  $\mathbb{P}[A] \leq 0$ . This shows that  $\mathbb{P}[A] = 0$ , so passing to the complement we prove the claim

$$\mathbb{P}[X < +\infty] = \mathbb{P}[A^c] = 1 - \mathbb{P}[A] = 1 - 0 = 1.$$

□

Next we note that in a random sum with non-negative terms, we are allowed to interchange the order of summation and expected value.

**Lemma VIII.6** (Expected value of a series of non-negative random terms).

Suppose that  $X_1, X_2, \dots \in \mathfrak{m}\mathcal{F}^+$  is a sequence of non-negative random variables. Consider the random infinite series  $\sum_{k=1}^{\infty} X_k$ . Then we have

$$\mathbb{E}\left[\sum_{k=1}^{\infty} X_k\right] = \sum_{k=1}^{\infty} \mathbb{E}[X_k].$$

<sup>3</sup>A similar conclusion holds for integrals generally, not just expected values. We leave it to the reader to precisely formulate the statement and see that the proof carries through.

*Proof.* By definition, these infinite sums are the limits of their finite partial sums, in particular

$$\sum_{k=1}^{\infty} X_k(\omega) = \lim_{n \rightarrow \infty} \sum_{k=1}^n X_k(\omega) \quad \text{for all } \omega \in \Omega.$$

By non-negativity of the terms, the sequence of partial sums is increasing, i.e.,

$$\sum_{k=1}^n X_k(\omega) \uparrow \sum_{k=1}^{\infty} X_k(\omega) \quad \text{as } n \rightarrow \infty, \quad \text{for all } \omega \in \Omega.$$

In this pointwise increasing approximation, we can apply the Monotone convergence theorem (MCT), and get

$$\begin{aligned} \mathbb{E} \left[ \sum_{k=1}^{\infty} X_k \right] &= \mathbb{E} \left[ \lim_{n \rightarrow \infty} \sum_{k=1}^n X_k \right] && \text{(by definition of infinite sum)} \\ &= \lim_{n \rightarrow \infty} \mathbb{E} \left[ \sum_{k=1}^n X_k \right] && \text{(by MCT)} \\ &= \lim_{n \rightarrow \infty} \sum_{k=1}^n \mathbb{E} [X_k] && \text{(by linearity of expected value)} \\ &= \sum_{k=1}^{\infty} \mathbb{E} [X_k]. && \text{(by definition of infinite sum)} \end{aligned}$$

The assertion is thus proven.  $\square$

We can combine the above observations to obtain a sufficient condition for almost sure convergence of random series with non-negative terms.

**Proposition VIII.7** (Guaranteeing almost sure convergence of a random series).

*Suppose that  $X_1, X_2, \dots \in \mathfrak{m}\mathcal{F}^+$  is a sequence of non-negative random variables. Suppose moreover that we have*

$$\sum_{k=1}^{\infty} \mathbb{E} [X_k] < +\infty.$$

*Then we have almost surely  $\sum_{k=1}^{\infty} X_k < +\infty$  and almost surely  $X_k \rightarrow 0$  as  $k \rightarrow \infty$ , i.e.,*

$$\mathbb{P} \left[ \sum_{k=1}^{\infty} X_k < +\infty \right] = 1 \quad \text{and} \quad \mathbb{P} \left[ \lim_{k \rightarrow \infty} X_k = 0 \right] = 1.$$

*Proof.* Note that the convergence  $\sum_{k=1}^{\infty} X_k < +\infty$  of the series implies that its terms tend to zero,  $\lim_{k \rightarrow \infty} X_k = 0$ , so for the corresponding events we have

$$\left\{ \omega \in \Omega \mid \sum_{k=1}^{\infty} X_k(\omega) < +\infty \right\} \subset \left\{ \omega \in \Omega \mid \lim_{k \rightarrow \infty} X_k(\omega) = 0 \right\}$$

and for their probabilities therefore

$$\mathbb{P} \left[ \sum_{k=1}^{\infty} X_k < +\infty \right] \leq \mathbb{P} \left[ \lim_{k \rightarrow \infty} X_k = 0 \right].$$

It therefore suffices to show that the probability on the left hand side is equal to one. By Lemma VIII.6 and the assumption of convergence of the series of expected values, we get

$$\mathbb{E} \left[ \sum_{k=1}^{\infty} X_k \right] = \sum_{k=1}^{\infty} \mathbb{E} [X_k] < +\infty.$$

Lemma VIII.5 therefore implies what we claimed,

$$1 = \mathbb{P}\left[\sum_{k=1}^{\infty} X_k < +\infty\right] \leq \mathbb{P}\left[\lim_{k \rightarrow \infty} X_k = 0\right].$$

□

Indicators of events, in particular, are non-negative random variables. The first Borel – Cantelli lemma can thus be seen as consequence of the above observations.

**Lemma** (Borel–Cantelli lemma: convergence part, Lemma V.7).

Suppose that  $E_1, E_2, \dots \in \mathcal{F}$  are such that  $\sum_{k=1}^{\infty} \mathbb{P}[E_k] < +\infty$ . Then we have

$$\mathbb{P}\left[“E_k \text{ occurs infinitely often}”\right] = 0.$$

*Proof.* Consider the indicators  $\mathbb{I}_{E_k}$  of the events  $E_k$ ,  $k = 1, 2, \dots$ . Note that the random series

$$N := \sum_{k=1}^{\infty} \mathbb{I}_{E_k}$$

counts the (random) number  $N$  of events in the sequence  $E_1, E_2, \dots$ , which occur. By assumption, we have

$$\sum_{k=1}^{\infty} \mathbb{E}[\mathbb{I}_{E_k}] = \sum_{k=1}^{\infty} \mathbb{P}[E_k] < +\infty.$$

Therefore by Proposition VIII.7 we in particular have that  $\mathbb{P}[N < +\infty] = 1$ , i.e., almost surely only finitely many of the events in the sequence  $E_1, E_2, \dots$  occur. The complementary event is that infinitely many of the events occur, so we have

$$\mathbb{P}\left[“E_k \text{ occurs infinitely often}”\right] = \mathbb{P}[N = +\infty] = 1 - \mathbb{P}[N < +\infty] = 1 - 1 = 0.$$

This finishes our alternative proof of the first Borel–Cantelli lemma. □

## Differentiation inside expected values

As an application of the Dominated convergence theorem, the reader can verify that exchanging the order of expected values and differentiation is permitted for example in the following situation.

**Exercise VIII.6** (Differentiation inside expectation).

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and  $X: \Omega \rightarrow \mathbb{R}$  a random variable. Suppose that  $h: \mathbb{R} \times (a, b) \rightarrow \mathbb{R}$  is a continuous function,  $(x, \lambda) \mapsto h(x, \lambda)$ . Assume that the partial derivative with respect to the second variable,  $\frac{\partial}{\partial \lambda} h: \mathbb{R} \times (a, b) \rightarrow \mathbb{R}$ , is also continuous. Assume moreover, that for some integrable random variable  $Y \in \mathcal{L}^1(\Omega, \mathcal{F}, \mathbb{P})$  we have

$$\text{for all } \lambda \in (a, b) \text{ and } \omega \in \Omega : \quad \left| \left( \frac{\partial}{\partial \lambda} h \right) (X(\omega), \lambda) \right| \leq Y(\omega).$$

(a) Show that for any  $\lambda, \lambda' \in (a, b)$  with  $\lambda \neq \lambda'$  we have  $\left| \frac{h(X, \lambda') - h(X, \lambda)}{\lambda' - \lambda} \right| \leq Y$ .

**Hint:** Write  $h(x, \lambda_2) - h(x, \lambda_1) = \int_{\lambda_1}^{\lambda_2} \left( \frac{\partial}{\partial \lambda} h \right) (x, \lambda) d\lambda$ , and recall the triangle inequality for integrals.

(b) Show that

$$\frac{d}{d\lambda} \mathbb{E}[h(X, \lambda)] = \mathbb{E}\left[\left(\frac{\partial}{\partial \lambda} h\right)(X, \lambda)\right].$$

**Hint:** When  $\lambda_n \rightarrow \lambda$ , consider  $\frac{1}{\lambda_n - \lambda} \left( \mathbb{E}[h(X, \lambda_n)] - \mathbb{E}[h(X, \lambda)] \right)$ . Use part (a) and dominated convergence.

Differentiation inside expected value is in particular often used for various generating functions.

**Exercise VIII.7** (Differentiating the moment generating function).

Define the moment generating function of  $X: \Omega \rightarrow \mathbb{R}$  by  $M_X(\lambda) := \mathbf{E}[e^{\lambda X}]$ . Assume that  $\mathbf{E}[e^{\varepsilon|X|}] < \infty$  for some  $\varepsilon > 0$ .

(a) Show that  $M'_X(0) = \mathbf{E}[X]$ .

**Hint:** Find a way to apply Exercise VIII.6, perhaps by first showing that for  $|\lambda| < \varepsilon$  one has  $|xe^{\lambda x}| \leq Ce^{\varepsilon|x|}$ .

(d) Explain, without detailed calculations, why  $M''_X(0) = \mathbf{E}[X^2]$ . Find also a similar formula for  $\mathbf{E}[X^n]$  for all  $n \in \mathbb{N}$ .

### VIII.3. Space of $p$ -integrable random variables

Let  $X: \Omega \rightarrow \widehat{\mathbb{R}}$  be a random variable. Recall that we say that  $X$  is integrable and write  $X \in \mathcal{L}^1(\mathbf{P})$  if

$$\mathbf{E}[|X|] < +\infty.$$

The following generalization is encountered very often.

**Definition VIII.8** ( $p$ -integrability).

Let  $p > 0$ . We say that a random variable  $X$  is  $p$ -integrable and denote  $X \in \mathcal{L}^p(\mathbf{P})$  if

$$\mathbf{E}[|X|^p] < +\infty. \tag{VIII.2}$$

In probability spaces,  $p$ -integrability implies  $r$ -integrability for any  $r \leq p$ , as we will show in the next lemma. This is used in particular to conclude that if it is meaningful to talk about the *moment*  $\mathbf{E}[X^p]$  of order  $p$ , it is also meaningful to talk about all lower order *moments*  $\mathbf{E}[X^r]$ . Let us emphasize that in very common measure spaces of infinite total mass,  $p$ -integrability does not imply integrability of lower order, so the lemma is specific to probability theory.

**Lemma VIII.9** (Finiteness of lower order moments).

Let  $0 < r < p$ . Suppose that  $X \in \mathcal{L}^p(\mathbf{P})$ . Then we have  $X \in \mathcal{L}^r(\mathbf{P})$ , and moreover

$$\mathbf{E}[|X|^r] \leq 1 + \mathbf{E}[|X|^p].$$

*Proof.* Since  $0 < r < p$ , we have that  $|x|^r < |x|^p$  whenever  $|x| > 1$ , and  $|x|^r \leq 1$  whenever  $|x| \leq 1$ . Define the event

$$A = \left\{ \omega \in \Omega \mid |X(\omega)| > 1 \right\}.$$

Then we have the pointwise estimate

$$\begin{aligned} |X(\omega)|^r &\leq \mathbb{I}_{A^c}(\omega) + \mathbb{I}_A(\omega) |X(\omega)|^p \\ &\leq 1 + |X(\omega)|^p. \end{aligned}$$

By monotonicity and linearity of the expected value, this gives

$$\mathbf{E}[|X|^r] \leq \mathbf{E}[1 + |X|^p] \leq 1 + \mathbf{E}[|X|^p].$$

Both asserted results follow from this.  $\square$



The space  $\mathcal{L}^p(\mathbf{P})$  of all  $p$ -integrable random variables is a vector space.

**Lemma VIII.10** ( $p$ -integrable random variables form a vector space).

Let  $p \geq 0$ . Then we have the following.

- (a) If  $X \in \mathcal{L}^p(\mathbf{P})$  and  $a \in \mathbb{R}$ , then  $aX \in \mathcal{L}^p(\mathbf{P})$ .
- (b) If  $X, Y \in \mathcal{L}^p(\mathbf{P})$ , then  $X + Y \in \mathcal{L}^p(\mathbf{P})$  and

$$\mathbb{E}[|X + Y|^p] \leq 2^p \left( \mathbb{E}[|X|^p] + \mathbb{E}[|Y|^p] \right) \quad (\text{VIII.3})$$

*Proof.* Part (a) is clear, since  $\mathbb{E}[|aX|^p] = |a|^p \mathbb{E}[|X|^p]$ . It therefore only remains to prove (b).

For any  $x, y \in \mathbb{R}$  note that

$$|x + y| \leq |x| + |y| \leq 2 \max\{|x|, |y|\}.$$

The mapping  $t \mapsto t^p$  is increasing  $[0, \infty) \rightarrow [0, \infty)$ , so applying it to the above gives

$$|x + y|^p \leq 2^p \max\{|x|^p, |y|^p\}.$$

This of course implies also

$$|x + y|^p \leq 2^p (|x|^p + |y|^p).$$

Applying this inequality pointwise to the values of the random variables  $X$  and  $Y$  gives

$$|X(\omega) + Y(\omega)|^p \leq 2^p (|X(\omega)|^p + |Y(\omega)|^p) \quad \text{for all } \omega \in \Omega.$$

Taking expected values and using monotonicity and linearity, we get

$$\mathbb{E}[|X + Y|^p] \leq 2^p (\mathbb{E}[|X|^p] + \mathbb{E}[|Y|^p]),$$

and the assertions of part (b) follow.  $\square$

Since a constant random variable  $c$  is  $p$ -integrable for any  $p > 0$ , the vector space property in particular has the following consequence.

**Corollary VIII.11** (Adding a constant preserves  $p$ -integrability).

If  $X \in \mathcal{L}^p(\mathbf{P})$  and  $c \in \mathbb{R}$  is a constant, then also  $X + c \in \mathcal{L}^p(\mathbf{P})$ .

The estimates in Lemma VIII.9 were rather crude. The following very often useful inequality is one way to improve them.

**Exercise VIII.8** (Jensen's inequality).

Suppose that  $I \subset \mathbb{R}$  is an interval and  $\phi: I \rightarrow \mathbb{R}$  is a convex function, i.e., a function such that for any  $x, y \in I$  and  $\lambda \in (0, 1)$  we have

$$\phi(\lambda x + (1 - \lambda)y) \leq \lambda \phi(x) + (1 - \lambda) \phi(y).$$

Suppose moreover that  $X: \Omega \rightarrow I$  is a random variable with values on the interval  $I$ .

- (a) Show that for any  $z \in I$  there exists a number  $d \in \mathbb{R}$  such that  $\phi(x) \geq \phi(z) + (x - z)d$ .  
**Hint:** Consider the left and right derivatives of  $\phi$  at  $z$ .
- (b) Show that if  $X$  is a random variable with values in  $I \subset \mathbb{R}$ , and both  $X$  and  $\phi(X)$  are integrable, then we have

$$\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)].$$

**Hint:** Choose  $z = \mathbb{E}[X]$  and apply part (a) together with monotonicity and linearity of expected values.

**Exercise VIII.9** (The  $p$ -norm controls lower norms).

Suppose that  $X \in \mathcal{L}^p(\mathbf{P})$  for some  $p > 0$ . Let  $0 < r < p$ . Using Exercise VIII.8, show that

$$\left(\mathbb{E}[|X|^r]\right)^{1/r} \leq \left(\mathbb{E}[|X|^p]\right)^{1/p},$$

## Lecture IX

### Product spaces and Fubini's theorem

Imagine observing two random phenomena simultaneously. The set of possible outcomes then is a Cartesian product, it consists of pairs

$$\Omega_1 \times \Omega_2 = \{(\omega_1, \omega_2) \mid \omega_1 \in \Omega_1, \omega_2 \in \Omega_2\},$$

whose two components are outcomes of the two phenomena.

**Example IX.1** (Two coin tosses).

For one coin toss the sample space is  $\{\text{H}, \text{T}\}$ , representing the possible outcomes “heads” and “tails”. For two coin tosses the sample space is

$$\{\text{H}, \text{T}\} \times \{\text{H}, \text{T}\} = \{(\text{H}, \text{H}), (\text{H}, \text{T}), (\text{T}, \text{H}), (\text{T}, \text{T})\}.$$

In this lecture we treat measures on such Cartesian product spaces. We construct in particular:

- the product  $\sigma$ -algebra on a Cartesian product of measurable spaces
- the product measure on a Cartesian product of ( $\sigma$ -finite) measure spaces.

One of the results of great practical importance is Fubini's Theorem (Theorem IX.9): for two measure spaces  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  and a function  $f: S_1 \times S_2 \rightarrow \mathbb{R}$ , under reasonable conditions the change of order of integration formula

$$\int_{S_1} \left( \int_{S_2} f(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) = \int_{S_2} \left( \int_{S_1} f(s_1, s_2) \, d\mu_1(s_1) \right) d\mu_2(s_2)$$

holds. Recalling that expected values and summation are special cases of integration, this implies (under reasonable assumptions) formulas such as

$$\int \left( \sum_{k=1}^{\infty} f_k(s) \right) d\mu(s) = \sum_{k=1}^{\infty} \int f_k(s) \, d\mu(s),$$
$$\mathbb{E} \left[ \sum_{k=1}^{\infty} X_k \right] = \sum_{k=1}^{\infty} \mathbb{E}[X_k],$$

etc.

#### Key tool: Monotone class theorem

In this lecture, we will repeatedly use the Monotone class theorem (Theorem C.2 from Appendix C). Let us recall its statement, and introduce convenient notation for it.

We denote as before

$$\begin{aligned} \text{m}\mathcal{S} & \text{--- the set of } \mathcal{S}\text{-measurable functions } S \rightarrow [-\infty, +\infty] \\ \text{s}\mathcal{S} & \text{--- the set of simple functions } S \rightarrow \mathbb{R} \end{aligned}$$

and in addition the set of bounded measurable functions will be denoted by

$$\text{b}\mathcal{S} := \left\{ f \in \text{m}\mathcal{S} \mid |f| \leq C \text{ for some } C \in \mathbb{R} \right\}.$$

We also continue to use the superscript  $+$  for the sets of non-negative functions of the corresponding types,

$$\begin{aligned} \text{m}\mathcal{S}^+ & := \left\{ f \in \text{m}\mathcal{S} \mid f \geq 0 \right\} \\ \text{s}\mathcal{S}^+ & := \left\{ f \in \text{s}\mathcal{S} \mid f \geq 0 \right\} \\ \text{b}\mathcal{S}^+ & := \left\{ f \in \text{b}\mathcal{S} \mid f \geq 0 \right\}. \end{aligned}$$

According to Definition C.1, a collection  $\mathcal{H}$  of functions  $S \rightarrow \mathbb{R}$  is said to be a monotone class if:

- (MC-1) The constant function 1 belongs to  $\mathcal{H}$ .
- (MC- $\mathbb{R}$ ) The class  $\mathcal{H}$  is a vector space over  $\mathbb{R}$ .
- (MC- $\uparrow$ ) If  $f_1, f_2, \dots \in \mathcal{H}$  is an increasing sequence of non-negative functions in  $\mathcal{H}$  such that the pointwise limit  $f_n \uparrow f$  is a bounded function  $f$ , then  $f \in \mathcal{H}$ .

The statement of the Monotone Class Theorem is the following.

**Theorem** (Monotone class theorem, Theorem C.2).

*Let  $\mathcal{H}$  be a monotone class of bounded functions from  $S$  to  $\mathbb{R}$  and let  $\mathcal{J}$  be a  $\pi$ -system on  $S$  such that  $\sigma(\mathcal{J}) = \mathcal{S}$ . Suppose that*

$$\mathbb{1}_A \in \mathcal{H} \quad \text{for every } A \in \mathcal{J}.$$

*Then we have*

$$\text{b}\mathcal{S} \subset \mathcal{H}.$$

### IX.1. Product sigma algebra

Let  $(S_1, \mathcal{S}_1)$  and  $(S_2, \mathcal{S}_2)$  be two measurable spaces. As the first step, we have to equip the Cartesian product

$$S_1 \times S_2 = \{(s_1, s_2) \mid s_1 \in S_1, s_2 \in S_2\}$$

with a  $\sigma$ -algebra. We think of the coordinate projections

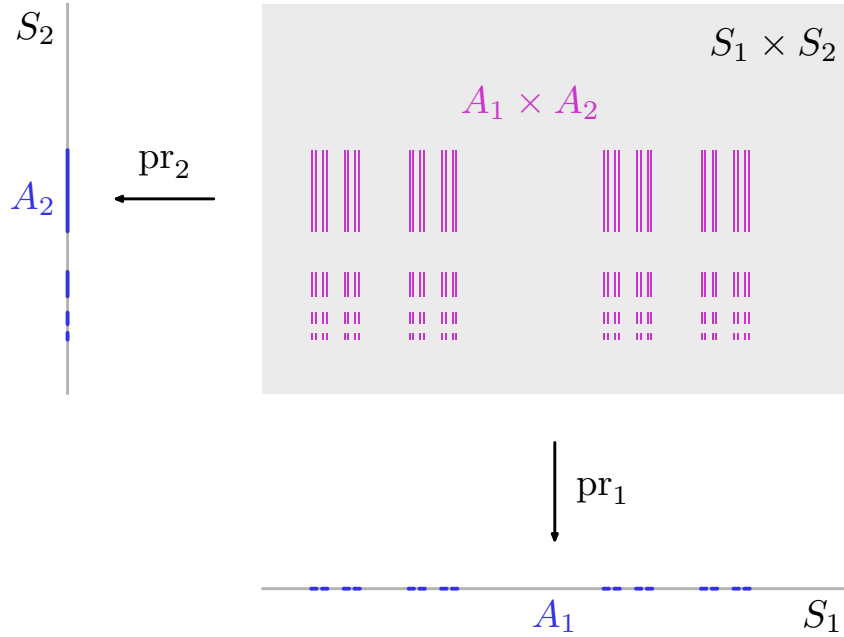
$$\begin{aligned} \text{pr}_1: S_1 \times S_2 & \rightarrow S_1 & \text{pr}_1(s_1, s_2) & = s_1 \\ \text{pr}_2: S_1 \times S_2 & \rightarrow S_2 & \text{pr}_2(s_1, s_2) & = s_2 \end{aligned}$$

as fundamental, and require that they must be measurable functions to  $(S_1, \mathcal{S}_1)$  and  $(S_2, \mathcal{S}_2)$ , respectively.

**Definition IX.2** (Product sigma algebra).

The *product  $\sigma$ -algebra*  $\mathcal{S}_1 \otimes \mathcal{S}_2$  on  $S_1 \times S_2$  is the  $\sigma$ -algebra generated by the functions  $\text{pr}_1$  and  $\text{pr}_2$ , i.e., the smallest  $\sigma$ -algebra on  $S_1 \times S_2$  with respect to which  $\text{pr}_1: S_1 \times S_2 \rightarrow S_1$  and  $\text{pr}_2: S_1 \times S_2 \rightarrow S_2$  are measurable.

The following lemma says that Cartesian products of measurable sets are  $\mathcal{S}_1 \otimes \mathcal{S}_2$ -measurable, and that  $\mathcal{S}_1 \otimes \mathcal{S}_2$  is the smallest  $\sigma$ -algebra with that property.



**Lemma IX.3** (A pi system for the product sigma algebra).

The collection

$$\mathcal{I} = \{A_1 \times A_2 \mid A_1 \in \mathcal{S}_1, A_2 \in \mathcal{S}_2\} \quad (\text{IX.1})$$

is a  $\pi$ -system on  $S_1 \times S_2$ , and we have  $\mathcal{S}_1 \otimes \mathcal{S}_2 = \sigma(\mathcal{I})$ .

**Exercise IX.1.** Prove Lemma IX.3.

The Euclidean plane is just the Cartesian product  $\mathbb{R}^2 = \mathbb{R} \times \mathbb{R}$  of two copies of the real axis. Does that mean that we have equipped the plane with two different  $\sigma$ -algebras: its Borel  $\sigma$ -algebra (generated by open subsets of the plane), and its product  $\sigma$ -algebra (generated by projections from the plane to the two axes)? In principle yes — but...

**Exercise IX.2** (The two natural  $\sigma$ -algebras on the plane are the same).

Let  $\mathcal{B}$  be the Borel  $\sigma$ -algebra on the real line  $\mathbb{R}$ , and  $\mathcal{B}(\mathbb{R}^2)$  the Borel  $\sigma$ -algebra on the plane  $\mathbb{R}^2$ . Show that

$$\mathcal{B} \otimes \mathcal{B} = \mathcal{B}(\mathbb{R}^2).$$

**Hint:** A set  $A \subset \mathbb{R}^n$  is open if and only if it can be written as a countable union of open boxes of the form  $(a_1, b_1) \times \cdots \times (a_n, b_n)$  with  $a_i, b_i$  being rational numbers.

### Freezing one coordinate

Our first use of the Monotone Class Theorem is to show that with any function  $f: S_1 \times S_2 \rightarrow \mathbb{R}$  which is measurable with respect to the product  $\sigma$ -algebra, keeping one of the two coordinates fixed and letting the other vary defines measurable functions  $S_1 \rightarrow \mathbb{R}$  and  $S_2 \rightarrow \mathbb{R}$ .

**Lemma IX.4** (Freezing a coordinate preserves measurability).

Let  $\mathcal{H}$  denote the class of functions  $f: S_1 \times S_2 \rightarrow \mathbb{R}$  such that  $f \in \mathfrak{b}(\mathcal{S}_1 \otimes \mathcal{S}_2)$  and

$$\begin{aligned} \forall s_1 \in S_1 \quad s_2 \mapsto f(s_1, s_2) \quad & \text{is } \mathcal{S}_2\text{-measurable } S_2 \rightarrow \mathbb{R} \\ \forall s_2 \in S_2 \quad s_1 \mapsto f(s_1, s_2) \quad & \text{is } \mathcal{S}_1\text{-measurable } S_1 \rightarrow \mathbb{R}. \end{aligned}$$

Then we have  $\mathcal{H} = \mathfrak{b}(\mathcal{S}_1 \otimes \mathcal{S}_2)$ .

*Proof.* By definition of  $\mathcal{H}$  we have  $\mathcal{H} \subset \mathfrak{b}(\mathcal{S}_1 \otimes \mathcal{S}_2)$ , so we must show  $\mathcal{H} \supset \mathfrak{b}(\mathcal{S}_1 \otimes \mathcal{S}_2)$ .

Clearly the class  $\mathcal{H}$  satisfies (MC-1), (MC- $\mathbb{R}$ ), and (MC- $\uparrow$ ). We will show that  $\mathcal{H}$  contains the indicator functions of all sets in the  $\pi$ -system  $\mathcal{I}$  of Lemma IX.3. Since  $\sigma(\mathcal{I}) = \mathcal{S}_1 \otimes \mathcal{S}_2$ , the Monotone Class Theorem then allows us to conclude that  $\mathcal{H} \supset \mathfrak{b}(\mathcal{S}_1 \otimes \mathcal{S}_2)$ .

So suppose that  $A_1 \in \mathcal{S}_1$  and  $A_2 \in \mathcal{S}_2$ , and consider the indicator function

$$\mathbb{I}_{A_1 \times A_2}: S_1 \times S_2 \rightarrow \mathbb{R}$$

of the set  $A_1 \times A_2 \subset S_1 \times S_2$ . We have

$$\mathbb{I}_{A_1 \times A_2}(s_1, s_2) = \mathbb{I}_{A_1}(s_1) \mathbb{I}_{A_2}(s_2),$$

so if  $s_1 \in A_1$ , the function  $s_2 \mapsto \mathbb{I}_{A_1 \times A_2}(s_1, s_2)$  is the indicator  $\mathbb{I}_{A_2}$  of the set  $A_2 \subset S_2$ , and if  $s_1 \notin A_1$ , it is zero. In either case,  $s_2 \mapsto \mathbb{I}_{A_1 \times A_2}(s_1, s_2)$  is  $\mathcal{S}_2$ -measurable  $S_2 \rightarrow \mathbb{R}$ . Similarly one shows that  $s_1 \mapsto \mathbb{I}_{A_1 \times A_2}(s_1, s_2)$  is  $\mathcal{S}_1$ -measurable  $S_1 \rightarrow \mathbb{R}$ , for any  $s_2 \in S_2$ . This shows that the indicator functions of all sets in the  $\pi$ -system  $\mathcal{I}$  of Lemma IX.3 are in the class  $\mathcal{H}$ . This finishes the proof.  $\square$

## IX.2. Product measure

Let  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  be two measure spaces. Our next goal is to construct a product measure, denoted by  $\mu_1 \otimes \mu_2$ , on the measurable space  $(S_1 \times S_2, \mathcal{S}_1 \otimes \mathcal{S}_2)$ .

First, in Section IX.2.1 we do that in the case when the measures are finite, and then in Section IX.2.2 we relax the assumption of finiteness to  $\sigma$ -finiteness.

### IX.2.1. Product of two finite measures

We now consider the case that  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  are two finite measure spaces, i.e., we assume that the total masses of the two measures are finite

$$\mu_1[S_1] < +\infty \quad \text{and} \quad \mu_2[S_2] < +\infty.$$

This assumption will be used frequently below to conclude that a function which is bounded and measurable is necessarily integrable (recall Lemma VII.20).

Suppose that  $f: S_1 \times S_2 \rightarrow \mathbb{R}$  is a bounded  $\mathcal{S}_1 \otimes \mathcal{S}_2$ -measurable function, i.e.,  $f \in \mathfrak{b}(\mathcal{S}_1 \otimes \mathcal{S}_2)$ . The functions  $s_2 \mapsto f(s_1, s_2)$  and  $s_1 \mapsto f(s_1, s_2)$  are clearly also

bounded, and they are measurable by Lemma IX.4. Thus they are integrable, and we can define the functions

$$\mathfrak{J}_1^f(s_1) := \int_{S_2} f(s_1, s_2) \, d\mu_2(s_2) \quad (\text{IX.2})$$

$$\mathfrak{J}_2^f(s_2) := \int_{S_1} f(s_1, s_2) \, d\mu_1(s_1). \quad (\text{IX.3})$$

The next lemma shows that these integrals over one of the variables define measurable functions. Moreover, it essentially establishes the conclusion of Fubini's theorem in the particular case of bounded measurable functions on finite measure spaces.

**Lemma IX.5** (Fubini's theorem for bounded functions).

Let  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  be two finite measure spaces, and let  $f \in b(\mathcal{S}_1 \otimes \mathcal{S}_2)$ . We then have the following:

- (i) The function  $s_1 \mapsto \mathfrak{J}_1^f(s_1)$  in (IX.2) is bounded and  $\mathcal{S}_1$ -measurable  $S_1 \rightarrow \mathbb{R}$ .
- (ii) The function  $s_2 \mapsto \mathfrak{J}_2^f(s_2)$  in (IX.3) is bounded and  $\mathcal{S}_2$ -measurable  $S_2 \rightarrow \mathbb{R}$ .
- (iii) We have the equality

$$\int_{S_1} \mathfrak{J}_1^f(s_1) \, d\mu_1(s_1) = \int_{S_2} \mathfrak{J}_2^f(s_2) \, d\mu_2(s_2).$$

*Proof.* Let  $\mathcal{H}$  be the collection of functions  $f \in b(\mathcal{S}_1 \otimes \mathcal{S}_2)$  for which (i), (ii), and (iii) hold. We must show  $\mathcal{H} = b(\mathcal{S}_1 \otimes \mathcal{S}_2)$ , and we will again do this using the Monotone Class Theorem.

We first claim that for any  $A_1 \in \mathcal{S}_1$  and  $A_2 \in \mathcal{S}_2$ , the indicator  $\mathbb{I}_{A_1 \times A_2}: S_1 \times S_2 \rightarrow \mathbb{R}$  is in the collection  $\mathcal{H}$ . Recall that  $\mathbb{I}_{A_1 \times A_2}(s_1, s_2) = \mathbb{I}_{A_1}(s_1)\mathbb{I}_{A_2}(s_2)$ , fix  $s_1 \in S_1$ , and integrate over the variable  $s_2$  to get

$$\begin{aligned} \mathfrak{J}_1^{\mathbb{I}_{A_1 \times A_2}}(s_1) &= \int_{S_2} \mathbb{I}_{A_1}(s_1)\mathbb{I}_{A_2}(s_2) \, d\mu_2(s_2) \\ &= \mathbb{I}_{A_1}(s_1) \int_{S_2} \mathbb{I}_{A_2}(s_2) \, d\mu_2(s_2) \\ &= \mu_2[A_2] \mathbb{I}_{A_1}(s_1). \end{aligned}$$

As a constant multiple of an indicator of the measurable set  $A_1$ , the function  $s_1 \mapsto \mathfrak{J}_1^{\mathbb{I}_{A_1 \times A_2}}(s_1)$  is  $\mathcal{S}_1$ -measurable and bounded, i.e., (i) holds for  $\mathbb{I}_{A_1 \times A_2}$ . Similarly one shows that

$$\mathfrak{J}_2^{\mathbb{I}_{A_1 \times A_2}}(s_2) = \mu_1[A_1] \mathbb{I}_{A_2}(s_2),$$

and concludes that  $s_2 \mapsto \mathfrak{J}_2^{\mathbb{I}_{A_1 \times A_2}}(s_2)$  is  $\mathcal{S}_2$ -measurable and bounded, i.e., (ii) holds for  $\mathbb{I}_{A_1 \times A_2}$ . For (iii), calculate

$$\begin{aligned} \int_{S_1} \mathfrak{J}_1^{\mathbb{I}_{A_1 \times A_2}}(s_1) \, d\mu_1(s_1) &= \int_{S_1} \mu_2[A_2] \mathbb{I}_{A_1}(s_1) \, d\mu_1(s_1) \\ &= \mu_2[A_2] \int_{S_1} \mathbb{I}_{A_1}(s_1) \, d\mu_1(s_1) \\ &= \mu_2[A_2] \mu_1[A_1], \end{aligned}$$

and similarly

$$\int_{S_2} \mathfrak{J}_2^{\mathbb{I}_{A_1 \times A_2}}(s_2) \, d\mu_2(s_2) = \mu_1[A_1] \mu_2[A_2].$$

These two integrals coincide, so we get that also (iii) holds for  $\mathbb{I}_{A_1 \times A_2}$ , and thus  $\mathbb{I}_{A_1 \times A_2} \in \mathcal{H}$ .

It remains to prove that  $\mathcal{H}$  is a monotone class. Properties (MC-1) and (MC- $\mathbb{R}$ ) are easy to verify. For property (MC- $\uparrow$ ), suppose that  $f_1, f_2, \dots \in \mathcal{H}$  are non-negative and  $f_n(s_1, s_2) \uparrow f(s_1, s_2)$  for all  $(s_1, s_2) \in S_1 \times S_2$ , where  $f: S_1 \times S_2 \rightarrow \mathbb{R}$  is bounded. Then for

any  $s_1 \in S_1$ , the functions  $s_2 \mapsto f_n(s_1, s_2)$  increase pointwise to  $s_2 \mapsto f(s_1, s_2)$  as  $n \rightarrow \infty$ , so by Monotone convergence theorem we get

$$\int_{S_2} f_n(s_1, s_2) \, d\mu_2(s_2) \uparrow \int_{S_2} f(s_1, s_2) \, d\mu_2(s_2),$$

i.e.,  $\mathfrak{J}_1^{f_n}(s_1) \uparrow \mathfrak{J}_1^f(s_1)$ . Since  $f_n \in \mathcal{H}$ , the functions  $s_1 \mapsto \mathfrak{J}_1^{f_n}(s_1)$  are  $\mathcal{S}_1$ -measurable, and as their pointwise limit also the function  $s_1 \mapsto \mathfrak{J}_1^f(s_1)$  must therefore be. Moreover, since  $f$  is bounded and the measure  $\mu_2$  is finite, the function  $s_1 \mapsto \mathfrak{J}_1^f(s_1)$  is bounded, so property (i) holds for  $f$ . Similarly one shows that  $\mathfrak{J}_2^{f_n}(s_2) \uparrow \mathfrak{J}_2^f(s_2)$  for any  $s_2 \in S_2$ , and property (ii) for  $f$  follows. Finally, use the Monotone convergence theorem twice (once for  $\mu_1$  and once for  $\mu_2$ ) to calculate

$$\begin{aligned} \int_{S_1} \mathfrak{J}_1^f(s_1) \, d\mu_1(s_1) &= \lim_{n \rightarrow \infty} \int_{S_1} \mathfrak{J}_1^{f_n}(s_1) \, d\mu_1(s_1) && (\mathfrak{J}_1^{f_n} \uparrow \mathfrak{J}_1^f \text{ and MCT for } \mu_1) \\ &= \lim_{n \rightarrow \infty} \int_{S_2} \mathfrak{J}_2^{f_n}(s_2) \, d\mu_2(s_2) && (\text{property (iii) for } f_n \in \mathcal{H}) \\ &= \int_{S_2} \mathfrak{J}_2^f(s_2) \, d\mu_2(s_2). && (\mathfrak{J}_2^{f_n} \uparrow \mathfrak{J}_2^f \text{ and MCT for } \mu_2) \end{aligned}$$

The equality above proves property (iii) for  $f$ , and shows that  $f \in \mathcal{H}$ , therefore establishing (MC $\uparrow$ ) and showing that  $\mathcal{H}$  is indeed a monotone class. This finishes the proof.  $\square$

Property (iii) of the previous lemma shows that the two formulas below are equal, and the following definition is therefore unambiguous.

**Definition IX.6** (Product measure).

Let  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  be two finite measure spaces. The *product measure*  $\mu_1 \otimes \mu_2$  on  $S_1 \times S_2$  is defined by

$$\begin{aligned} (\mu_1 \otimes \mu_2)[B] &= \int_{S_1} \left( \int_{S_2} \mathbb{I}_B(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) && \text{(IX.4)} \\ &= \int_{S_2} \left( \int_{S_1} \mathbb{I}_B(s_1, s_2) \, d\mu_1(s_1) \right) d\mu_2(s_2) \end{aligned}$$

for any  $B \in \mathcal{S}_1 \otimes \mathcal{S}_2$ .

This formula indeed gives rise to a measure on the product space, as we verify next.

**Lemma IX.7** (The product measure is a measure).

Let  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  be two finite measure spaces, and define  $\mu_1 \otimes \mu_2$  by formula (IX.4). Then  $\mu_1 \otimes \mu_2$  is a measure on  $(S_1 \times S_2, \mathcal{S}_1 \otimes \mathcal{S}_2)$ .

*Proof.* The defining formula (IX.4) only involves integrals of non-negative functions, so clearly  $(\mu_1 \otimes \mu_2)[B] \in [0, +\infty]$  for any  $B \in \mathcal{S}_1 \otimes \mathcal{S}_2$ . We must check the two defining properties of a measure.

The indicator of the empty set  $\emptyset \subset S_1 \times S_2$  is the zero function,  $\mathbb{I}_\emptyset(s_1, s_2) = 0$ . The integral of the zero function is zero, so from (IX.4) we directly get  $(\mu_1 \otimes \mu_2)[\emptyset] = 0$ .

Suppose that  $B_1, B_2, \dots \in \mathcal{S}_1 \otimes \mathcal{S}_2$  are disjoint, and let  $B := \bigcup_{k \in \mathbb{N}} B_k$ . Consider also the finite unions  $U_n = B_1 \cup \dots \cup B_n$  for  $n \in \mathbb{N}$ , to get an increasing sequence  $U_1 \subset U_2 \subset \dots$  of  $\mathcal{S}_1 \otimes \mathcal{S}_2$ -measurable sets increasing to  $\bigcup_{n \in \mathbb{N}} U_n = B$ . By disjointness we have  $\mathbb{I}_{U_n} = \sum_{k=1}^n \mathbb{I}_{B_k}$ . We first use this to calculate  $(\mu_1 \otimes \mu_2)[U_n]$  from (IX.4), using linearity of



integration,

$$\begin{aligned}
(\mu_1 \otimes \mu_2)[U_n] &= \int_{S_1} \left( \int_{S_2} \mathbb{I}_{U_n}(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \\
&= \int_{S_1} \left( \int_{S_2} \left( \sum_{k=1}^n \mathbb{I}_{B_k}(s_1, s_2) \right) d\mu_2(s_2) \right) d\mu_1(s_1) \\
&= \sum_{k=1}^n \left( \int_{S_1} \left( \int_{S_2} \mathbb{I}_{B_k}(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \right) = \sum_{k=1}^n (\mu_1 \otimes \mu_2)[B_k].
\end{aligned}$$

On the other hand, the indicators of the increasing sequence  $U_n$  of sets increase pointwise to the indicator of the limit set  $B$ , i.e.,  $\mathbb{I}_{U_n} \uparrow \mathbb{I}_B$ . We can use this to calculate  $(\mu_1 \otimes \mu_2)[B]$  from (IX.4), using Monotone convergence theorems (first for  $\mu_2$  and then for  $\mu_1$ ),

$$\begin{aligned}
(\mu_1 \otimes \mu_2)[B] &= \int_{S_1} \left( \int_{S_2} \mathbb{I}_B(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \\
&= \int_{S_1} \left( \int_{S_2} \lim_{n \rightarrow \infty} \mathbb{I}_{U_n}(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \\
&= \int_{S_1} \left( \lim_{n \rightarrow \infty} \int_{S_2} \mathbb{I}_{U_n}(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) && \text{(MCT for } \mu_2) \\
&= \lim_{n \rightarrow \infty} \int_{S_1} \left( \int_{S_2} \mathbb{I}_{U_n}(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) && \text{(MCT for } \mu_1) \\
&= \lim_{n \rightarrow \infty} (\mu_1 \otimes \mu_2)[U_n].
\end{aligned}$$

Combining these two calculations, we have shown that

$$\begin{aligned}
(\mu_1 \otimes \mu_2) \left[ \bigcup_{k \in \mathbb{N}} B_k \right] &= (\mu_1 \otimes \mu_2)[B] = \lim_{n \rightarrow \infty} (\mu_1 \otimes \mu_2)[U_n] \\
&= \lim_{n \rightarrow \infty} \sum_{k=1}^n (\mu_1 \otimes \mu_2)[B_k] = \sum_{k=1}^{\infty} (\mu_1 \otimes \mu_2)[B_k],
\end{aligned}$$

which is the countable additivity property for  $\mu_1 \otimes \mu_2$ .  $\square$

It took some work to construct the product measure, but to characterize it is very easy:

**Lemma IX.8** (A characterization of the product measure).

Let  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  be two finite measure spaces. Then the product measure  $\mu_1 \otimes \mu_2$  is the unique measure  $\nu$  on  $(S_1 \times S_2, \mathcal{S}_1 \otimes \mathcal{S}_2)$  such that for all  $A_1 \in \mathcal{S}_1$  and  $A_2 \in \mathcal{S}_2$  we have

$$\nu[A_1 \times A_2] = \mu_1[A_1] \mu_2[A_2].$$

*Proof.* An easy calculation (done already in the proof of Lemma IX.5) starting from the definition (IX.4) of the product measure gives

$$(\mu_1 \otimes \mu_2)[A_1 \times A_2] = \mu_1[A_1] \mu_2[A_2].$$

Thus  $\mu_1 \otimes \mu_2$  indeed satisfies the asserted formula.

Now if  $\nu$  is another measure on  $(S_1 \times S_2, \mathcal{S}_1 \otimes \mathcal{S}_2)$  for which the formula holds, then  $\nu$  and  $\mu_1 \otimes \mu_2$  coincide on all sets of the  $\pi$ -system (IX.1),  $\mathcal{S} = \{A_1 \times A_2 \mid A_1 \in \mathcal{S}_1, A_2 \in \mathcal{S}_2\}$ . In particular the total masses are equal,  $\nu[S_1 \times S_1] = (\mu_1 \otimes \mu_2)[S_1 \times S_1]$ . We can divide each by this finite total mass to get two probability measures, which agree on the  $\pi$ -system  $\mathcal{S}$ . Since  $\sigma(\mathcal{S}) = \mathcal{S}_1 \otimes \mathcal{S}_2$ , Dynkin's identification theorem then guarantees that these probability measures on  $(S_1 \times S_2, \mathcal{S}_1 \otimes \mathcal{S}_2)$  are equal. Multiplying the total mass back, we conclude that  $\nu = \mu_1 \otimes \mu_2$ . This shows the uniqueness of a measure satisfying the formula.  $\square$

Now we state and prove the main result about product measures.

**Theorem IX.9** (Fubini's theorem).

For a function  $f: S_1 \times S_2 \rightarrow [-\infty, +\infty]$ , consider the following three integrals:

$$\int_{S_1 \times S_2} f \, d(\mu_1 \otimes \mu_2) \quad (\text{IX.5})$$

$$\int_{S_1} \left( \int_{S_2} f(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \quad (\text{IX.6})$$

$$\int_{S_2} \left( \int_{S_1} f(s_1, s_2) \, d\mu_1(s_1) \right) d\mu_2(s_2). \quad (\text{IX.7})$$

We have:

- (a) If  $f$  is non-negative and measurable,  $f \in m(\mathcal{S}_1 \otimes \mathcal{S}_2)^+$ , then the integrals (IX.5), (IX.6), and (IX.7) are all in  $[0, +\infty]$ , and they are all equal.
- (b) If  $f$  is integrable,  $f \in \mathcal{L}^1(\mu_1 \otimes \mu_2)$ , then the integrals (IX.5), (IX.6), and (IX.7) are all in  $\mathbb{R}$ , and they are all equal.

*Proof of part (a):* We first claim that for all  $f \in b(\mathcal{S}_1 \otimes \mathcal{S}_2)$ , then the integrals (IX.5), (IX.6), and (IX.7) are all equal as real numbers. The equality of the last two was in fact shown in part (iii) of Lemma IX.5, and here one proceeds similarly. Consider the collection  $\mathcal{H}$  of functions  $f \in b(\mathcal{S}_1 \otimes \mathcal{S}_2)$  for which we have (i), (ii), and

- (iii') The integrals (IX.5), (IX.6), and (IX.7) are all equal as real numbers.

With minor modifications to the proof of Lemma IX.5, one can prove that this collection  $\mathcal{H}$  is a monotone class and contains indicators  $\mathbb{I}_{A_1 \times A_2}$  of sets  $A_1 \times A_2$ , and therefore  $\mathcal{H} = b(\mathcal{S}_1 \otimes \mathcal{S}_2)$ .

It in particular follows, because simple functions are necessarily bounded, that all non-negative simple functions satisfy (iii'). For a given non-negative measurable  $f \in m(\mathcal{S}_1 \otimes \mathcal{S}_2)^+$  we choose a pointwise increasing approximation  $f_n \uparrow f$  by non-negative simple functions  $f_n \in s(\mathcal{S}_1 \otimes \mathcal{S}_2)^+$ . By (iii'), for these approximating functions we have the equalities

$$\begin{aligned} \int_{S_1 \times S_2} f_n \, d(\mu_1 \otimes \mu_2) &= \int_{S_1} \left( \int_{S_2} f_n(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \\ &= \int_{S_2} \left( \int_{S_1} f_n(s_1, s_2) \, d\mu_1(s_1) \right) d\mu_2(s_2). \end{aligned}$$

We show that these three increase to (IX.5), (IX.6), and (IX.7) for  $f$ , respectively, and the assertion (a) then follows. For the first, just use Monotone convergence theorem for the product measure  $\mu_1 \otimes \mu_2$

$$\int_{S_1 \times S_2} f_n \, d(\mu_1 \otimes \mu_2) \uparrow \int_{S_1 \times S_2} f \, d(\mu_1 \otimes \mu_2).$$

Let us consider the second integral in the limit  $n \rightarrow \infty$ . Use first Monotone convergence theorem for  $\mu_2$  to get

$$\int_{S_2} f_n(s_1, s_2) \, d\mu_2(s_2) \uparrow \int_{S_2} f(s_1, s_2) \, d\mu_2(s_2),$$

which in particular shows that the inner integral in (IX.6) defines a non-negative measurable function of the variable  $s_1$ , as a pointwise limit of non-negative measurable functions. Use next the Monotone convergence theorem for  $\mu_1$  together with the previous result to get

$$\int_{S_1} \left( \int_{S_2} f_n(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1) \uparrow \int_{S_1} \left( \int_{S_2} f(s_1, s_2) \, d\mu_2(s_2) \right) d\mu_1(s_1).$$

The third one is shown similarly, with Monotone convergence theorem for  $\mu_1$  first and then for  $\mu_2$ . This finishes the proof of (a).

*Proof of part (b):* Suppose that  $f \in \mathcal{L}^1(\mu_1 \otimes \mu_2)$ , i.e., that  $f$  is  $\mathcal{S}_1 \otimes \mathcal{S}_2$ -measurable and

$$\int_{S_1 \times S_2} |f| \, d(\mu_1 \otimes \mu_2) < +\infty.$$

Write  $f = f_+ - f_-$ , where  $f_+, f_- \in m(\mathcal{S}_1 \otimes \mathcal{S}_2)^+$  are the positive and negative parts of  $f$ . From integrability of  $f$ , it follows that we have

$$\int_{S_1 \times S_2} f_+ \, d(\mu_1 \otimes \mu_2) < +\infty \quad \text{and} \quad \int_{S_1 \times S_2} f_- \, d(\mu_1 \otimes \mu_2) < +\infty.$$

We can then apply part (a) to  $f_+$  to get that all of the following integrals finite and equal

$$\begin{aligned} \int_{S_1 \times S_2} f_+ \, d(\mu_1 \otimes \mu_2) &= \int_{S_1} \left( \underbrace{\int_{S_2} f_+(s_1, s_2) \, d\mu_2(s_2)}_{=:\mathfrak{J}_1^{f_+}(s_1)} \right) d\mu_1(s_1) \\ &= \int_{S_2} \left( \underbrace{\int_{S_1} f_+(s_1, s_2) \, d\mu_1(s_1)}_{=:\mathfrak{J}_2^{f_+}(s_2)} \right) d\mu_2(s_2), \end{aligned}$$

and similarly for  $f_-$ . In particular this implies that for  $\mu_1$ -almost every  $s_1$  we have that  $\mathfrak{J}_1^{f_+}(s_1) < +\infty$  and  $\mathfrak{J}_1^{f_-}(s_1) < +\infty$ . Noting that the positive and negative parts of  $\mathfrak{J}_1^f$  are  $(\mathfrak{J}_1^f)_+ = \mathfrak{J}_1^{f_+}$  and  $(\mathfrak{J}_1^f)_- = \mathfrak{J}_1^{f_-}$ , we furthermore see that  $\mathfrak{J}_1^f \in \mathcal{L}^1(\mu_1)$ . Similarly we get  $\mathfrak{J}_2^f \in \mathcal{L}^1(\mu_2)$ . Then comparing the definition of the integral (IX.5)

$$\int_{S_1 \times S_2} f \, d(\mu_1 \otimes \mu_2) = \int_{S_1 \times S_2} f_+ \, d(\mu_1 \otimes \mu_2) - \int_{S_1 \times S_2} f_- \, d(\mu_1 \otimes \mu_2)$$

with the definitions of the double integrals (IX.6) and (IX.7) in terms of the positive and negative parts, the asserted equalities follow.  $\square$

**Remark IX.10.** In part (b) of Theorem IX.9 we assumed the integrability  $f \in \mathcal{L}^1(\mu_1 \otimes \mu_2)$  with respect to the product measure  $\mu_1 \otimes \mu_2$ . By part (a), however, this integrability follows if  $f$  is measurable and either one of the double integrals

$$\int_{S_1} \left( \int_{S_2} |f(s_1, s_2)| \, d\mu_2(s_2) \right) d\mu_1(s_1) \quad \text{or} \quad \int_{S_2} \left( \int_{S_1} |f(s_1, s_2)| \, d\mu_1(s_1) \right) d\mu_2(s_2)$$

of the absolute value are finite.

## IX.2.2. Product of two sigma finite measures

In Section IX.2.1 we constructed the product measure assuming that the two measures we started with were finite. Our next goal is to relax the assumption of finiteness. Let us, however, start by pointing out that it is necessary to make some assumptions in order to build a theory powerful enough to allow interchanging order of integrations. The following exercise contains a counterexample to the validity of such interchange property.

**Exercise IX.3** (It is not always possible to change the order of integrations).

Consider the following two measure spaces  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$ . The underlying set for both is the unit interval,

$$S_1 = [0, 1] \quad \text{and} \quad S_2 = [0, 1],$$

the  $\sigma$ -algebra for both is the Borel  $\sigma$ -algebra

$$\mathcal{S}_1 = \mathcal{B}([0, 1]) \quad \text{and} \quad \mathcal{S}_2 = \mathcal{B}([0, 1])$$

of the interval. The first measure is taken to be the restriction of the Lebesgue measure

$$\mu_1 = \Lambda \quad (\mu_1[(a, b)] = b - a)$$

to the interval, and the second measure is taken to be the counting measure

$$\mu_2 = \mu_{\#} \quad (\mu_2[B] = \#B)$$

on the interval. Define the function  $f: [0, 1] \times [0, 1] \rightarrow \mathbb{R}$  by

$$f(x, y) = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y. \end{cases}$$

- (a) Verify that that  $f$  is a non-negative function, which is  $\mathcal{S}_1 \otimes \mathcal{S}_2$ -measurable.  
 (b) Calculate both

$$\int_{[0,1]} \left( \int_{[0,1]} f(x, y) d\mu_{\#}(y) \right) d\Lambda(x) \quad \text{and} \quad \int_{[0,1]} \left( \int_{[0,1]} f(x, y) d\Lambda(x) \right) d\mu_{\#}(y)$$

and compare the results.

**Definition IX.11** (Sigma finite measures and measure spaces).

Let  $(S, \mathcal{S}, \mu)$  be a measure space. We say that the measure space is  $\sigma$ -finite or that the measure  $\mu$  is  $\sigma$ -finite if there exists a sequence  $A_1, A_2, \dots \in \mathcal{S}$  of measurable sets which together cover the entire space,  $S = \bigcup_{n=1}^{\infty} A_n$ , and each has finite measure,  $\mu[A_n] < +\infty$  for all  $n \in \mathbb{N}$ .

**Example IX.12** (The Lebesgue measure is sigma finite).

The real axis with the Lebesgue measure  $(\mathbb{R}, \mathcal{B}, \Lambda)$  is a  $\sigma$ -finite measure space. An example of a sequence of finite measure sets which covers the real axis is  $A_n = [-n, +n]$ , for  $n \in \mathbb{N}$ .

**Example IX.13** (The  $d$ -dimensional Lebesgue measure is sigma finite).

The  $d$ -dimensional Euclidean space with the  $d$ -dimensional Lebesgue measure  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d), \Lambda^d)$  is a  $\sigma$ -finite measure space. An example of a sequence of finite measure sets which covers the real axis is balls  $A_n = \{x \in \mathbb{R}^d \mid \|x\| \leq n\}$  of radii  $n$  centered at the origin, for  $n \in \mathbb{N}$ .

The following (counter-)example underlies the problem encountered in Exercise IX.3.

**Example IX.14** (The counting measure is sigma finite only on countable sets).

The counting measure  $\mu_{\#}$  on  $(S, \mathcal{S})$  is  $\sigma$ -finite if and only if  $S$  is a countable set. (Recall that the countable union of finite sets is countable!)

**Remark IX.15** (Sigma finiteness with disjoint pieces of finite measure).

If  $(S, \mathcal{S}, \mu)$  is a  $\sigma$ -finite measure space, then it is possible to choose disjoint measurable sets  $A'_1, A'_2, \dots \in \mathcal{S}$  which cover the entire space  $\bigcup_{n=1}^{\infty} A'_n = S$  and each has finite measure  $\mu[A'_n] < +\infty$  for all  $n \in \mathbb{N}$ . Indeed, if  $A_1, A_2, \dots$  form a sequence as in Definition IX.11, then it suffices to set  $A'_n := A_n \setminus (A_1 \cup \dots \cup A_{n-1})$  to achieve disjointness.

Suppose that  $(S_1, \mathcal{S}_1, \mu_1)$  and  $(S_2, \mathcal{S}_2, \mu_2)$  are two  $\sigma$ -finite measure spaces, and that  $A_1^{(1)}, A_2^{(1)}, \dots \in \mathcal{S}_1$  and  $A_1^{(2)}, A_2^{(2)}, \dots \in \mathcal{S}_2$  are their disjoint coverings by finite measure pieces,

$$\begin{aligned} S_1 &= \bigcup_{n=1}^{\infty} A_n^{(1)}, & S_2 &= \bigcup_{n=1}^{\infty} A_n^{(2)} \\ \mu_1[A_n^{(1)}] &< +\infty, & \mu_2[A_n^{(2)}] &< +\infty \quad \forall n \in \mathbb{N}. \end{aligned}$$

For any  $n, m \in \mathbb{N}$ , the truncation  $\mu_1^{(n)}$  of  $\mu_1$  to  $A_n^{(1)} \subset S_1$  and the truncation  $\mu_2^{(m)}$  of  $\mu_2$  to  $A_m^{(2)} \subset S_2$  are two finite measures (see Exercise II.1), so we know from

Section IX.2.1 how to construct on the product measure

$$\mu_1^{(n)} \otimes \mu_2^{(m)}.$$

The product of the two  $\sigma$ -finite measures can then be defined as the countable sum of these pieces,

$$(\mu_1 \otimes \mu_2)[B] := \sum_{n,m \in \mathbb{N}} (\mu_1^{(n)} \otimes \mu_2^{(m)})[B].$$

We leave it as an exercise to the reader to check that this definition does not depend on the chosen pieces  $A_n^{(i)}$ .<sup>1</sup> Fubini's theorem continues to hold for products of  $\sigma$ -finite measures. The proof can be done by splitting to the countably many pieces  $A_n^{(1)} \times A_m^{(2)} \subset S_1 \times S_2$  of finite measure.

---

<sup>1</sup>Hint: The formula  $(\mu_1 \otimes \mu_2)[A_1 \times A_2] = \mu_1[A_1] \mu_2[A_2]$  holds for any sets  $A_1 \in \mathcal{S}_1$ ,  $A_2 \in \mathcal{S}_2$  of finite measure, and Dynkin's identification theorem can then be used to characterize the restrictions of the product measure to any pieces of finite measure.



## Lecture X

### Probability on product spaces

In this lecture we look at some of the uses of product sigma algebras, product measures, and Fubini's theorem.

We will in particular consider random vectors and joint distributions of random variables.

#### X.1. Joint laws

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Recall from Definition III.6 that, in general, the law of a random variable  $X: \Omega \rightarrow S$  is the probability measure on  $(S, \mathcal{S})$  given by

$$P_X[B] = \mathbf{P}[X \in B].$$

As an example, for a real valued random variable  $X: \Omega \rightarrow \mathbb{R}$ , the law  $P_X$  is a probability measure on  $(\mathbb{R}, \mathcal{B})$ . To appreciate its role, recall from Theorem VIII.1 that for Borel functions  $h: \mathbb{R} \rightarrow \mathbb{R}$  we have

$$\mathbf{E}[h(X)] = \int_{\mathbb{R}} h(x) dP_X(x),$$

whenever  $h(X) \in \mathcal{L}^1(\mathbf{P})$  or equivalently  $h \in \mathcal{L}^1(P_X)$ . Below we generalize the considerations to several real valued random variables.

#### Definition of the joint law of two real random variables

Consider then two real valued random variables  $X$  and  $Y$  defined on the same probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ .

#### Definition X.1 (Joint law).

For  $X: \Omega \rightarrow \mathbb{R}$  and  $Y: \Omega \rightarrow \mathbb{R}$ , the *joint law* (or *joint distribution*) of  $X$  and  $Y$  is the probability measure on  $(\mathbb{R}^2, \mathcal{B}(\mathbb{R}^2))$  given by

$$P_{X,Y}[A] = \mathbf{P}[(X, Y) \in A].$$

#### Remark X.2 (Joint law as the law of a random vector).

The pair  $(X, Y)$  is a "random vector" in  $\mathbb{R}^2$  — or more precisely

$$\begin{aligned} Z: \Omega &\rightarrow \mathbb{R}^2 \\ \omega &\mapsto Z(\omega) = (X(\omega), Y(\omega)) \end{aligned}$$

is. Indeed, we check that  $Z: \Omega \rightarrow \mathbb{R}^2$  is  $\mathcal{F}/\mathcal{B}(\mathbb{R}^2)$ -measurable, and thus really a random variable with values in  $\mathbb{R}^2$ . To check this measurability, note first that for any  $B_1, B_2 \in \mathcal{B}$ ,

we have

$$\begin{aligned} Z^{-1}(B_1 \times B_2) &= \left\{ \omega \in \Omega \mid Z(\omega) \in B_1 \times B_2 \right\} \\ &= \left\{ \omega \in \Omega \mid X(\omega) \in B_1, Y(\omega) \in B_2 \right\} \\ &= X^{-1}(B_1) \cap Y^{-1}(B_2) \end{aligned}$$

which is an event in  $\mathcal{F}$  by the measurability of  $X$  and  $Y$ . The collection

$$\mathcal{J} = \{B_1 \times B_2 \mid B_1, B_2 \in \mathcal{B}\}$$

generates  $\mathcal{B} \otimes \mathcal{B} = \mathcal{B}(\mathbb{R}^2)$  by Lemma IX.3 and Exercise IX.2, so this observation is sufficient for the measurability of  $Z: \Omega \rightarrow \mathbb{R}^2$  according to Lemma III.9.

The joint law  $P_{X,Y}$  of  $X$  and  $Y$  is thus nothing but the law  $P_Z$  of the random vector  $Z = (X, Y)$ .

Analogously to Theorem VIII.1, for a Borel function  $h: \mathbb{R}^2 \rightarrow \mathbb{R}$  one can write the expected value of  $h(X, Y)$  in terms of the law  $P_{X,Y}$  as

$$\mathbb{E}[h(X, Y)] = \int_{\mathbb{R}^2} h(x, y) \, dP_{X,Y}(x, y), \quad (\text{X.1})$$

provided that  $h(X, Y) \in \mathcal{L}^1(\mathbb{P})$  or equivalently  $h \in \mathcal{L}^1(P_{X,Y})$ . We leave the detailed verification of this to the reader.<sup>1</sup>

**Lemma X.3** (Characterization of joint law).

*The law  $P_{X,Y}$  is uniquely characterized by the property that*

$$P_{X,Y}[B_1 \times B_2] = \mathbb{P}[(X, Y) \in B_1 \times B_2]$$

*for all  $B_1, B_2 \in \mathcal{B}$ .*

*Proof.* Recall that the collection  $\mathcal{J}$  of sets of the form  $B_1 \times B_2 \subset \mathbb{R}^2$  is a  $\pi$ -system which generates the product  $\sigma$ -algebra  $\mathcal{B} \otimes \mathcal{B}$  (Lemma IX.3). By Dynkin's identification theorem (Theorem II.26), the probability measure  $P_{X,Y}$  is uniquely characterized by its restriction to a generating  $\pi$ -system, so the assertion follows.  $\square$

The following example of joint laws is relevant to discrete time Markov processes on very general state spaces.

**Exercise X.1** (Transition probability kernels).

Let  $K$  be a transition probability kernel on  $(S, \mathcal{S})$ , i.e., a mapping  $S \times \mathcal{S} \rightarrow [0, +\infty)$  denoted by  $(s, A) \mapsto K_s[A]$  such that

- for any  $A \in \mathcal{S}$ , the mapping  $s \mapsto K_s[A]$  is  $\mathcal{S}$ -measurable  $S \rightarrow [0, +\infty)$
- for any  $s \in S$ , the mapping  $A \mapsto K_s[A]$  is a probability measure on  $(S, \mathcal{S})$ .

Let  $\mu$  be a probability measure on  $(S, \mathcal{S})$ .

(a) Define  $\mu K$  by

$$(\mu K)[A] = \int_S K_s[A] \, d\mu(s), \quad \text{for } A \in \mathcal{S}.$$

Show that  $\mu K$  is a probability measure on  $(S, \mathcal{S})$ .

<sup>1</sup>The idea is similar to Theorem VIII.1, but for this case, using the Monotone class theorem may be more convenient.



(b) Define, for  $\mathcal{S} \otimes \mathcal{S}$ -measurable subsets  $B \subset S \times S$

$$\nu[B] = \int_S \left( \int_S \mathbb{I}_B(s_1, s_2) dK_{s_1}(s_2) \right) d\mu(s_1).$$

Show that  $\nu$  is a probability measure on  $(S \times S, \mathcal{S} \otimes \mathcal{S})$ .

(c) Let  $X = (X_1, X_2)$  be a random “vector” in  $S \times S$  with distribution  $P_X = \nu$  given by (b). Show that the distributions  $P_{X_1}$  and  $P_{X_2}$  of its components  $X_1$  and  $X_2$  are  $\mu$  and  $\mu K$ , respectively.

**Note:**  $(X_1, X_2)$  can be viewed as the first two values of a discrete time Markov process with initial distribution  $\mu$  and transition kernel  $K$  (a measure valued generalization of transition matrix) on the state space  $S$ . The distribution of any number of first steps of the Markov process can be defined by generalizing the above construction. The existence of the whole Markov process can then be deduced from straightforward abstract extension theorems.

The notion of joint law of  $n$  real valued random variables is a straightforward generalization — and it is nothing but the law of the  $n$ -dimensional random vector whose components are the real valued random variables.

### Joint densities

Recall from (VIII.1) that a real valued random variable  $X: \Omega \rightarrow \mathbb{R}$  is said to have a *continuous distribution*, if there is a Borel-measurable function  $f_X: \mathbb{R} \rightarrow [0, +\infty]$  such that we have

$$P_X[B] = \mathbb{P}[X \in B] = \int_B f_X d\Lambda \quad \text{for all } B \in \mathcal{B}. \quad (\text{X.2})$$

In the above formula and below we use the convention (VII.10) for integrals over subsets. The function  $f_X$  is called the *probability density* of (the law of)  $X$ .

The Lebesgue measure  $\Lambda^2$  on  $\mathbb{R}^2$  is the product measure  $\Lambda^2 = \Lambda \otimes \Lambda$  of two Lebesgue measures on  $\mathbb{R}$ . It corresponds to the “area measure”, as it is for example determined by the measures

$$\Lambda^2[(a_1, b_1) \times (a_2, b_2)] = (b_1 - a_1)(b_2 - a_2)$$

of rectangles with  $a_1 \leq b_1$ ,  $a_2 \leq b_2$ .

Two random variables  $X: \Omega \rightarrow \mathbb{R}$  and  $Y: \Omega \rightarrow \mathbb{R}$  are said to have a *joint density*  $f_{X,Y}: \mathbb{R}^2 \rightarrow [0, +\infty]$  if we have

$$P_{X,Y}[A] = \mathbb{P}[(X, Y) \in A] = \int_A f_{X,Y} d\Lambda^2 \quad \text{for all } A \in \mathcal{B}(\mathbb{R}^2). \quad (\text{X.3})$$

The notion of joint density of more than two random variables is a straightforward generalization.

**Proposition X.4** (Marginal densities from joint density).

If  $X$  and  $Y$  have a joint density  $f_{X,Y}: \mathbb{R}^2 \rightarrow [0, +\infty]$ , then  $X$  has density  $f_X: \mathbb{R} \rightarrow [0, +\infty]$  given by

$$f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x, y) d\Lambda(y)$$

and  $Y$  has density  $f_Y: \mathbb{R} \rightarrow [0, +\infty]$  given by

$$f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x, y) \, d\Lambda(x).$$

*Proof.* We prove this using Fubini's theorem, Theorem IX.9.

Define  $f_X$  by the formula in the statement. We have to verify that it is a density for  $X$ . Let  $B \in \mathcal{B}$ . Write

$$P_X[B] = \mathbb{P}[X \in B] = \mathbb{P}[X \in B \text{ and } Y \in \mathbb{R}] = \mathbb{P}[(X, Y) \in B \times \mathbb{R}]$$

and then use the fact that  $f_{X,Y}$  is a joint density to write this as

$$P_X[B] = \int_{B \times \mathbb{R}} f_{X,Y} \, d\Lambda^2 = \int_{\mathbb{R}^2} \mathbb{I}_{B \times \mathbb{R}} f_{X,Y} \, d\Lambda^2.$$

Now observe that that  $\mathbb{I}_{B \times \mathbb{R}}(x, y) = \mathbb{I}_B(x)$  and recall that  $\Lambda^2 = \Lambda \otimes \Lambda$  is a product measure, and use Fubini's theorem to get

$$\begin{aligned} P_X[B] &= \int_{\mathbb{R}^2} \mathbb{I}_B(x) f_{X,Y}(x, y) \, d(\Lambda \otimes \Lambda)(x, y) \\ &= \int_{\mathbb{R}} \left( \int_{\mathbb{R}} \mathbb{I}_B(x) f_{X,Y}(x, y) \, d\Lambda(y) \right) d\Lambda(x). \end{aligned}$$

The inner integral here is

$$\int_{\mathbb{R}} \mathbb{I}_B(x) f_{X,Y}(x, y) \, d\Lambda(y) = \mathbb{I}_B(x) \int_{\mathbb{R}} f_{X,Y}(x, y) \, d\Lambda(y) = \mathbb{I}_B(x) f_X(x).$$

Therefore, we have obtained

$$P_X[B] = \int_{\mathbb{R}} \mathbb{I}_B(x) f_X(x) \, d\Lambda(x) = \int_B f_X(x) \, d\Lambda(x),$$

which shows that  $f_X$  is a density for  $X$ .

The claim about the density of  $Y$  is proven similarly.  $\square$

In summary, Proposition X.4 says that the existence of a joint density for a random vector guarantees the existence of probability densities for its components (called *marginal densities*). The converse does not hold, in general.

**Example X.5** (Marginal densities do not guarantee existence of joint density).

Let  $X$  be a real valued random variable with continuous distribution and density  $f_X: \mathbb{R} \rightarrow [0, +\infty]$ . Define a two-dimensional random vector  $Z = (X, X)$  whose two components are equal. Then both components have continuous distribution with density  $f_X$ . The joint law  $P_{X,X}$ , however, is supported on the line  $\{(x, y) \mid x = y\} \subset \mathbb{R}^2$ , which has measure zero (under the 2-dimensional Lebesgue measure  $\Lambda^2$ ). From this, it is easy to see that there can not exist a joint density for  $Z = (X, X)$ .

Of course not all probability distributions have densities. The following exercise concerns a very simple probability distribution in the  $d$ -dimensional Euclidean space.

**Exercise X.2** (Dirac measure on  $\mathbb{R}^d$ ).

The *Dirac measure* at a point  $a \in \mathbb{R}^d$  is defined by

$$\delta_a[A] = \begin{cases} 1 & \text{if } a \in A, \\ 0 & \text{otherwise.} \end{cases}$$

- Show that  $\delta_a$  is a probability measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ .
- Show that any Borel function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\delta_a$ -integrable (meaning that  $\int_{\mathbb{R}^d} |f(x)| \, d\delta_a(x) < +\infty$ ), and compute the integral  $\int_{\mathbb{R}^d} f(x) \, d\delta_a(x)$ .

Consider now the case  $d = 1$ , and fix  $a \in \mathbb{R}$ .

- (c) Does the measure  $\delta_a$  have a probability density function, i.e., a Borel function  $g: \mathbb{R}^n \rightarrow [0, +\infty]$  such that  $\delta_a[A] = \int_{\mathbb{R}} \mathbb{I}_A(x) g(x) dx$  for all  $A \in \mathcal{B}$ ? If yes, find out an expression for it. If not, explain why.

## X.2. Variances and covariances of square integrable random variables

Variances are important statistics of distributions of real valued random variables, and covariances are important statistics of the joint distributions of pairs of random variables. In order for these to be well-defined, we need the second moments of the random variables to exist.

The general notion of  $p$ -integrability was introduced in Definition VIII.8. For the present purposes, we are concerned with the particular case  $p = 2$ . A random variable  $X \in \mathcal{L}^2(\mathbf{P})$  is said to be *square integrable* — we recall that this means the following finiteness of second moment

$$\mathbf{E}[X^2] < +\infty.$$

By Lemma VIII.10 we have the vector space property of square integrable random variables: if  $X, Y \in \mathcal{L}^2(\mathbf{P})$  and  $a, b \in \mathbb{R}$ , then also  $aX + bY \in \mathcal{L}^2(\mathbf{P})$ . By Lemma VIII.9, square integrability implies in particular integrability, and thus  $\mathcal{L}^2(\mathbf{P}) \subset \mathcal{L}^1(\mathbf{P})$  is a vector subspace in the space of integrable random variables.

The following is a fundamentally important inequality for square integrable random variables.

**Theorem X.6** (Cauchy-Schwarz inequality).

*Suppose that  $X, Y \in \mathcal{L}^2(\mathbf{P})$  are two square integrable random variables. Then the product  $XY$  is integrable,  $XY \in \mathcal{L}^1(\mathbf{P})$ , and we have*

$$\left| \mathbf{E}[XY] \right| \leq \sqrt{\mathbf{E}[X^2] \mathbf{E}[Y^2]}. \quad (\text{X.4})$$

*Proof.* For any two real numbers  $x, y \in \mathbb{R}$  we have the following inequality

$$0 \leq (x - y)^2 = x^2 - 2xy + y^2.$$

By moving the cross-term to the other side and dividing by two we get  $xy \leq \frac{1}{2}x^2 + \frac{1}{2}y^2$ . By changing the sign of one of the numbers, we get also  $-xy \leq \frac{1}{2}x^2 + \frac{1}{2}y^2$ . Together these give

$$|xy| \leq \frac{1}{2}x^2 + \frac{1}{2}y^2.$$

Applying the above inequality to the values of the random variables  $X$  and  $Y$ , we get

$$|X(\omega)Y(\omega)| \leq \frac{1}{2}X(\omega)^2 + \frac{1}{2}Y(\omega)^2 \quad \text{for all } \omega \in \Omega.$$

Taking the expected values and using monotonicity and linearity, as well as the assumption  $X, Y \in \mathcal{L}^2(\mathbf{P})$ , we deduce

$$\mathbf{E}[|XY|] \leq \frac{1}{2}\mathbf{E}[X^2] + \frac{1}{2}\mathbf{E}[Y^2] < +\infty.$$

This shows that the product  $XY$  is indeed integrable,  $XY \in \mathcal{L}^1(\mathbf{P})$ .

Next observe that for any  $t \in \mathbb{R}$  and any  $\omega \in \Omega$  we have  $0 \leq (tX(\omega) + Y(\omega))^2$ . Taking expected values and using linearity and monotonicity we get that

$$0 \leq \mathbb{E}[(tX + Y)^2] = \mathbb{E}[t^2X^2 + 2tXY + Y^2] = t^2 \mathbb{E}[X^2] + 2t \mathbb{E}[XY] + \mathbb{E}[Y^2].$$

This is a quadratic polynomial in  $t$  which never becomes negative, so it can have at most one real root, and its discriminant is therefore either negative (if there are no roots) or zero (if there is one root), i.e.,

$$4 \mathbb{E}[XY]^2 - 4 \mathbb{E}[X^2] \mathbb{E}[Y^2] \leq 0.$$

Moving the second term to the other side and dividing by 4 yields

$$\mathbb{E}[XY]^2 \leq \mathbb{E}[X^2] \mathbb{E}[Y^2]. \quad (\text{X.5})$$

Taking the square roots then gives the asserted inequality.  $\square$

**Corollary X.7** (Square of expected value vs. expected value of square).

If  $X \in \mathcal{L}^2(\mathbb{P})$ , then we have also  $X \in \mathcal{L}^1(\mathbb{P})$  and

$$\mathbb{E}[X]^2 \leq \mathbb{E}[X^2]. \quad (\text{X.6})$$

*Proof.* Note that the constant random variable 1 is square integrable,  $1 \in \mathcal{L}^2(\mathbb{P})$ . The asserted inequality is derived by applying the squared Cauchy-Schwarz inequality (X.5) to  $X$  and 1 as follows

$$\mathbb{E}[X]^2 = \mathbb{E}[X \cdot 1]^2 \stackrel{(\text{X.5})}{\leq} \mathbb{E}[X^2] \underbrace{\mathbb{E}[1^2]}_{=1} = \mathbb{E}[X^2].$$

$\square$

Suppose now that  $X, Y \in \mathcal{L}^2(\mathbb{P})$  are square integrable random variables. Since we have  $\mathcal{L}^2(\mathbb{P}) \subset \mathcal{L}^1(\mathbb{P})$  (Lemma VIII.9) the expected values  $\mathbb{E}[X]$  and  $\mathbb{E}[Y]$  are well-defined and finite. Moreover, since adding constants does not affect  $\mathcal{L}^2(\mathbb{P})$  membership (Corollary VIII.11), we have also  $X - \mathbb{E}[X] \in \mathcal{L}^2(\mathbb{P})$  and  $Y - \mathbb{E}[Y] \in \mathcal{L}^2(\mathbb{P})$ . By Theorem X.6 then, we have  $(X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) \in \mathcal{L}^1(\mathbb{P})$ . Therefore the following definition makes sense.

**Definition X.8** (Variance and covariance).

The *variance* of a square integrable random variable  $X \in \mathcal{L}^2(\mathbb{P})$  is

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

The *covariance* of two square integrable random variables  $X, Y \in \mathcal{L}^2(\mathbb{P})$  is

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])].$$

**Proposition X.9** (Formulas for variance and covariance).

For  $X \in \mathcal{L}^2(\mathbb{P})$ , we have

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2. \quad (\text{X.7})$$

For  $X, Y \in \mathcal{L}^2(\mathbb{P})$ , we have

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y]. \quad (\text{X.8})$$

*Proof.* The first formula (X.7) is in fact a special case of the second (with  $Y = X$ ), so it is enough to prove (X.8). Denote  $\mathbf{m}_X = \mathbb{E}[X]$  and  $\mathbf{m}_Y = \mathbb{E}[Y]$ . Then calculate, by expanding and using linearity of expected values,

$$\begin{aligned} \text{Cov}(X, Y) &= \mathbb{E}\left[(X - \mathbf{m}_X)(Y - \mathbf{m}_Y)\right] \\ &= \mathbb{E}\left[XY - X\mathbf{m}_Y - \mathbf{m}_X Y + \mathbf{m}_X \mathbf{m}_Y\right] \\ &= \mathbb{E}[XY] - \underbrace{\mathbb{E}[X]}_{=\mathbf{m}_X} \mathbf{m}_Y - \mathbf{m}_X \underbrace{\mathbb{E}[Y]}_{=\mathbf{m}_Y} + \mathbf{m}_X \mathbf{m}_Y \\ &= \mathbb{E}[XY] - \mathbf{m}_X \mathbf{m}_Y. \end{aligned}$$

This is the asserted formula. □

### X.3. Independence and products

#### Equivalent conditions for independence of random numbers

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space, and  $X, Y: \Omega \rightarrow \mathbb{R}$  two real valued random variables on it. Denote the laws of  $X$  and  $Y$  by  $P_X$  and  $P_Y$ , respectively, and the joint law of  $X$  and  $Y$  by  $P_{X,Y}$ .

**Theorem X.10** (Independence in terms of laws).

*For two real valued random variables  $X, Y: \Omega \rightarrow \mathbb{R}$ , the following conditions are equivalent:*

- (i)  $X$  and  $Y$  are independent (i.e.,  $X \perp Y$ )
- (ii)  $P_{X,Y} = P_X \otimes P_Y$
- (iii) for all  $x, y \in \mathbb{R}$  we have  $\mathbb{P}[X \leq x, Y \leq y] = \mathbb{P}[X \leq x] \mathbb{P}[Y \leq y]$ .

*If moreover  $X$  and  $Y$  have a joint density  $f_{X,Y}: \mathbb{R}^2 \rightarrow [0, +\infty]$ , then (i), (ii), and (iii) are also equivalent to*

- (iv)  $X$  and  $Y$  have densities  $f_X$  and  $f_Y$ , respectively, such that  $f_{X,Y}(x, y) = f_X(x) f_Y(y)$  for  $\Lambda^2$ -almost all  $(x, y)$ .

In the proof, we repeatedly use the characterization of the product measure based on Lemma IX.8: if  $P$  and  $Q$  are two probability measures on  $\mathbb{R}$ , then  $P \otimes Q$  is the unique measure on  $\mathbb{R}^2$  such that

$$(P \otimes Q)[B_1 \times B_2] = P[B_1] Q[B_2]$$

for all  $B_1, B_2 \in \mathcal{B}$ .

*Proof of Theorem X.10:* The equivalence “(i)  $\Leftrightarrow$  (iii)” was already proven in Corollary V.6. We prove the implications “(i)  $\Rightarrow$  (ii)” and “(ii)  $\Rightarrow$  (iii)” to establish the first assertion about the equivalence of (i), (ii), and (iii). Then, under the additional hypothesis that  $X$  and  $Y$  have a joint density  $f_{X,Y}$ , we prove the implication “(ii)  $\Rightarrow$  (iv)” and leave it as an exercise to prove “(iv)  $\Rightarrow$  (iii)”. The second assertion follows.

*proof of “(i)  $\Leftrightarrow$  (iii)”:* Corollary V.6.

*proof of “(i)  $\Rightarrow$  (ii)”*: Suppose that we have independence  $X \perp\!\!\!\perp Y$ . Let  $B_1, B_2 \in \mathcal{B}$ . Then we have  $X^{-1}(B_1) \in \sigma(X)$  and  $Y^{-1}(B_2) \in \sigma(Y)$ , so independence can be used to calculate

$$\begin{aligned} P_{X,Y}[B_1 \times B_2] &= \mathbb{P}\left[X^{-1}(B_1) \cap Y^{-1}(B_2)\right] \\ &= \mathbb{P}\left[X^{-1}(B_1)\right] \mathbb{P}\left[Y^{-1}(B_2)\right] && \text{(because } X \perp\!\!\!\perp Y\text{)} \\ &= P_X[B_1] P_Y[B_2]. \end{aligned}$$

By Lemma IX.8, the product measure  $P_X \otimes P_Y$  is the unique measure for which this formula holds, so we conclude that indeed  $P_{X,Y} = P_X \otimes P_Y$ .

*proof of “(ii)  $\Rightarrow$  (iii)”*: Suppose that we have  $P_{X,Y} = P_X \otimes P_Y$ . Let  $x, y \in \mathbb{R}$ , and consider  $B_1 = (-\infty, x]$  and  $B_2 = (-\infty, y]$ . Then we have

$$\begin{aligned} \mathbb{P}\left[X \leq x, Y \leq y\right] &= P_{X,Y}[B_1 \times B_2] \\ &= P_X[B_1] P_Y[B_2] && \text{(since } P_{X,Y} = P_X \otimes P_Y\text{)} \\ &= \mathbb{P}[X \leq x] \mathbb{P}[Y \leq y], \end{aligned}$$

which is what we needed to show.

*proof of “(iv)  $\Rightarrow$  (iii)” assuming joint density*: Exercise.

*proof of “(ii)  $\Rightarrow$  (iv)” assuming joint density*: Suppose that  $P_{X,Y} = P_X \otimes P_Y$ , and assume moreover that there exists a joint density  $f_{X,Y}: \mathbb{R}^2 \rightarrow [0, +\infty)$  for  $X$  and  $Y$ . Recall from Proposition X.4 that  $X$  and  $Y$  then have densities

$$\begin{aligned} f_X(x) &= \int_{\mathbb{R}} f_{X,Y}(x, y) \, d\Lambda(y) \\ f_Y(y) &= \int_{\mathbb{R}} f_{X,Y}(x, y) \, d\Lambda(x). \end{aligned}$$

Define a new measure  $\mu$  on  $\mathbb{R}^2$  by the formula

$$\mu[A] := \int_A f_X(x) f_Y(y) \, d\Lambda^2(x, y) \quad \text{for all } A \in \mathcal{B}(\mathbb{R}^2) = \mathcal{B} \otimes \mathcal{B}.$$

By Fubini's theorem with the positive measurable function  $(x, y) \mapsto f_X(x) f_Y(y)$  in this integral, for any  $B_1, B_2 \in \mathcal{B}$  we get

$$\begin{aligned} \mu[B_1 \times B_2] &= \int_{B_1 \times B_2} f_X(x) f_Y(y) \, d\Lambda^2(x, y) \\ &= \int_{B_2} f_Y(y) \left( \underbrace{\int_{B_1} f_X(x) \, d\Lambda(x)}_{=P_X[B_1]} \right) d\Lambda(y) && \text{(Fubini)} \\ &= P_X[B_1] \int_{B_2} f_Y(y) \, d\Lambda(y) \\ &= P_X[B_1] P_Y[B_2]. \end{aligned}$$

Again, by Lemma IX.8 the product measure  $P_X \otimes P_Y$  is the only measure for which this formula holds, so we get  $\mu = P_X \otimes P_Y$ . But by assumption also  $P_{X,Y} = P_X \otimes P_Y$ . This shows that  $P_{X,Y} = \mu$ , which by the construction of  $\mu$  shows that  $(x, y) \mapsto f_X(x) f_Y(y)$  is a joint density of  $X$  and  $Y$ . Since also  $f_{X,Y}$  is a joint density of  $X$  and  $Y$ , we must have  $f_{X,Y}(x, y) = f_X(x) f_Y(y)$  for  $\Lambda^2$ -almost all points  $(x, y) \in \mathbb{R}^2$ .  $\square$

**Exercise X.3** (Disintegration of independent random variables).

Let  $X$  and  $Y$  be independent real valued random variables with laws  $P_X$  and  $P_Y$ . Let  $h: \mathbb{R}^2 \rightarrow \mathbb{R}$  be Borel function.

(a) Prove that

$$\mathbb{E}[|h(X, Y)|] = \int_{\mathbb{R}} \mathbb{E}[|h(x, Y)|] \, dP_X(x) = \int_{\mathbb{R}} \mathbb{E}[|h(X, y)|] \, dP_Y(y).$$

(b) Prove that when  $\mathbb{E}[|h(X, Y)|] < \infty$ ,

$$\mathbb{E}[h(X, Y)] = \int_{\mathbb{R}} \mathbb{E}[h(x, Y)] dP_X(x) = \int_{\mathbb{R}} \mathbb{E}[h(X, y)] dP_Y(y).$$

**Hint:** You may start from (X.1). Then keep in mind product measure and Fubini's theorem.

**Exercise X.4** (Product of two uniformly distributed numbers).

Let  $U_1$  and  $U_2$  be independent uniformly distributed random variables on  $[0, 1]$ , so that both have the function  $\mathbb{I}_{[0,1]}$  as their probability density function. Define  $X = U_1 U_2$ .

- Calculate the cumulative distribution function  $F_X(x) = \mathbb{P}[X \leq x]$ .
- What is the distribution  $P_X[B] = \mathbb{P}[X \in B]$  of  $X$ ?
- Does  $X$  have a probability density function? If yes, find out an expression for it. If not, explain why.

## Independence and expected value of product

**Theorem X.11** (Expected value of product of independent random variables).

Suppose that  $X, Y: \Omega \rightarrow \mathbb{R}$  are two independent random variables which are integrable,  $X, Y \in \mathcal{L}^1(\mathbb{P})$ . Then also their product is integrable,  $XY \in \mathcal{L}^1(\mathbb{P})$ , and we have

$$\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y]. \quad (\text{X.9})$$

*Proof.* We prove the statement using the “standard machine”, i.e., successively for (1): indicators, (2): simple random variables, (3): non-negative random variables, and (4): all integrable random variables.

*step 1 (indicator random variables):* If  $X = \mathbb{I}_A$  and  $Y = \mathbb{I}_B$ , then  $XY = \mathbb{I}_{A \cap B}$ . In this case we have

$$\mathbb{E}[XY] = \mathbb{E}[\mathbb{I}_{A \cap B}] = \mathbb{P}[A \cap B].$$

and

$$\mathbb{E}[X] \mathbb{E}[Y] = \mathbb{E}[\mathbb{I}_A] \mathbb{E}[\mathbb{I}_B] = \mathbb{P}[A] \mathbb{P}[B].$$

Of course  $A = X^{-1}(\{1\})$  and  $B = Y^{-1}(\{1\})$ , so by the assumed independence of  $X$  and  $Y$ , we have  $\mathbb{P}[A \cap B] = \mathbb{P}[A] \mathbb{P}[B]$ . This shows that (X.9) holds for independent indicator random variables.

*step 2 (simple random variables):* Suppose that  $X, Y \in \mathfrak{s}\mathcal{F}$  are simple random variables. Write

$$X = \sum_{i=1}^n a_i \mathbb{I}_{A_i} \quad \text{and} \quad Y = \sum_{j=1}^m b_j \mathbb{I}_{B_j}$$

where  $a_1, \dots, a_n$  are the distinct non-zero values of  $X$ , and  $b_1, \dots, b_m$  are the distinct non-zero values of  $Y$ , and  $A_i = X^{-1}(\{a_i\})$  and  $B_j = Y^{-1}(\{b_j\})$ . Independence of  $X$  and  $Y$  shows that  $A_i \perp B_j$ , and thus also  $\mathbb{I}_{A_i} \perp \mathbb{I}_{B_j}$ . Now calculate, using linearity and the result

of step 0,

$$\begin{aligned}
\mathbb{E}[XY] &= \mathbb{E}\left[\sum_{i=1}^n \sum_{j=1}^m a_i b_j \mathbb{I}_{A_i} \mathbb{I}_{B_j}\right] \\
&= \sum_{i=1}^n \sum_{j=1}^m a_i b_j \mathbb{E}[\mathbb{I}_{A_i} \mathbb{I}_{B_j}] \\
&= \sum_{i=1}^n \sum_{j=1}^m a_i b_j \mathbb{E}[\mathbb{I}_{A_i}] \mathbb{E}[\mathbb{I}_{B_j}] && \text{(by step 1)} \\
&= \mathbb{E}\left[\sum_{i=1}^n a_i \mathbb{I}_{A_i}\right] \mathbb{E}\left[\sum_{j=1}^m b_j \mathbb{I}_{B_j}\right] = \mathbb{E}[X] \mathbb{E}[Y].
\end{aligned}$$

This shows that (X.9) holds for independent simple random variables.

*step 3 (non-negative random variables):* Assume that  $X$  and  $Y$  are non-negative random variables. Let  $\varsigma_n: [0, +\infty] \rightarrow [0, n]$  be the  $n$ :th staircase function, and define the simple functions  $X_n = \varsigma_n \circ X$  and  $Y_n = \varsigma_n \circ Y$ . Then independence of  $X$  and  $Y$  implies the independence of  $X_n$  and  $Y_n$  also, by Exercise V.2. Moreover, by construction we have the pointwise increasing approximations,  $X_n \uparrow X$  and  $Y_n \uparrow Y$ , and also  $X_n Y_n \uparrow XY$ . Monotone convergence theorem implies that as  $n \rightarrow \infty$ , we have  $\mathbb{E}[X_n] \uparrow \mathbb{E}[X]$ ,  $\mathbb{E}[Y_n] \uparrow \mathbb{E}[Y]$ , and  $\mathbb{E}[X_n Y_n] \uparrow \mathbb{E}[XY]$ . On the other hand, by step 2, we have

$$\mathbb{E}[X_n Y_n] = \mathbb{E}[X_n] \mathbb{E}[Y_n].$$

In the limit  $n \rightarrow \infty$  we now get (X.9) for independent non-negative random variables.

*step 4 (integrable random variables):* Consider finally general integrable random variables  $X, Y \in \mathcal{L}^1(\mathbb{P})$ . Split these to positive and negative parts,  $X = X_+ - X_-$  and  $Y = Y_+ - Y_-$ , with  $X_+, X_-, Y_+, Y_-$  non-negative. We have  $X_+ = \max\{X, 0\}$  etc., so by Exercise V.2, independence of  $X$  and  $Y$  implies also the independences  $X_+ \perp\!\!\!\perp Y_+$ ,  $X_+ \perp\!\!\!\perp Y_-$ ,  $X_- \perp\!\!\!\perp Y_+$ , and  $X_- \perp\!\!\!\perp Y_-$ . Then we calculate

$$\begin{aligned}
\mathbb{E}[XY] &= \mathbb{E}[(X_+ - X_-)(Y_+ - Y_-)] \\
&= \mathbb{E}[X_+ Y_+] - \mathbb{E}[X_+ Y_-] - \mathbb{E}[X_- Y_+] + \mathbb{E}[X_- Y_-] \\
&= \mathbb{E}[X_+] \mathbb{E}[Y_+] - \mathbb{E}[X_+] \mathbb{E}[Y_-] - \mathbb{E}[X_-] \mathbb{E}[Y_+] + \mathbb{E}[X_-] \mathbb{E}[Y_-] && \text{(by step 3)} \\
&= (\mathbb{E}[X_+] - \mathbb{E}[X_-]) (\mathbb{E}[Y_+] - \mathbb{E}[Y_-]) \\
&= \mathbb{E}[X] \mathbb{E}[Y].
\end{aligned}$$

This finishes the proof. □

**Proposition X.12** (Variance is additive for independent random variables).

If  $X, Y \in \mathcal{L}^2(\mathbb{P})$  are independent, then we have

$$\text{Cov}(X, Y) = 0 \tag{X.10}$$

and

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) \tag{X.11}$$

*Proof.* By independence and Theorem X.11, we have  $\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y]$ . Using formula (X.8) then immediately shows that the covariance vanishes,  $\text{Cov}(X, Y) = 0$ .

Then calculate the variance of  $X + Y$  using formula (X.7),

$$\begin{aligned}
\text{Var}(X + Y) &= \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2 \\
&= \mathbb{E}[X^2 + 2XY + Y^2] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 \\
&= \mathbb{E}[X^2] + 2 \mathbb{E}[XY] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - 2 \mathbb{E}[X] \mathbb{E}[Y] - \mathbb{E}[Y]^2.
\end{aligned}$$



Using again the consequence  $E[XY] = E[X]E[Y]$  of independence, we have a cancellation, and we get

$$\text{Var}(X + Y) = E[X^2] + E[Y^2] - E[X]^2 - E[Y]^2.$$

This leads to (X.11), if we just recall formula (X.7) for the variances of  $X$  and  $Y$ .  $\square$

**Exercise X.5** (Additivity of variance for independent random variables).

Suppose that  $X_1, \dots, X_n$  are independent. Show that then  $\text{Var}(\sum_{k=1}^n X_k) = \sum_{k=1}^n \text{Var}(X_k)$ . Show also, by finding a counterexample, that the formula does not hold in general without the assumption of independence.

#### X.4. A formula for the expected value

We finish this lecture by a straightforward, yet often practical use of the Fubini's theorem in a probabilistic setup. It gives a formula for the expected value of a non-negative random variable in terms of its cumulative distribution function.

**Proposition X.13** (Expected value in terms of c.d.f.).

Let  $X: \Omega \rightarrow [0, +\infty)$  be a non-negative random number. Let  $F_X: \mathbb{R} \rightarrow [0, 1]$  be its cumulative distribution function

$$F_X(x) = \mathbf{P}[X \leq x].$$

Then we have

$$E[X] = \int_0^\infty (1 - F_X(x)) \, dx.$$

*Proof.* Consider the product space

$$\Omega \times [0, +\infty),$$

equipped with the product  $\sigma$ -algebra  $\mathcal{F} \otimes \mathcal{B}([0, +\infty))$  and with the product measure  $\mathbf{P} \otimes \Lambda$  of the probability measure  $\mathbf{P}$  on the sample space  $\Omega$  and the Lebesgue measure  $\Lambda$  restricted to the non-negative real axis  $[0, +\infty)$ .

Define a subset

$$A := \{(\omega, t) \mid t < X(\omega)\} \subset \Omega \times [0, +\infty).$$

This subset is measurable<sup>2</sup>,  $A \in \mathcal{F} \otimes \mathcal{B}([0, +\infty))$ , so its indicator is a non-negative measurable function  $\mathbb{I}_A \in m(\mathcal{F} \otimes \mathcal{B}([0, +\infty)))^+$  concretely given by

$$\mathbb{I}_A(\omega, t) = \begin{cases} 1 & \text{if } t < X(\omega) \\ 0 & \text{if } t \geq X(\omega). \end{cases}$$

By Fubini's theorem, Theorem IX.9, it does not matter in which order we integrate the two variables of this function:

$$\int_0^\infty \left( \int_\Omega \mathbb{I}_A(\omega, t) dP(\omega) \right) dt = \int_\Omega \left( \int_0^\infty \mathbb{I}_A(\omega, t) dt \right) dP(\omega). \quad (\text{X.12})$$

The rest of the proof consists of calculating both of the above expressions separately.

Let us start from the left-hand side of (X.12). The inner integral is calculated as

$$\int_\Omega \mathbb{I}_A(\omega, t) dP(\omega) = \int_\Omega \mathbb{I}_{\{t < X\}} dP = P[t < X] = 1 - P[t \geq X] = 1 - F(t).$$

The double integral on the left-hand side of (X.12) is therefore

$$\int_0^\infty (1 - F(t)) dt. \quad (\text{X.13})$$

Consider then the right-hand side of (X.12). The inner integral is calculated as

$$\int_0^\infty \mathbb{I}_A(\omega, t) dt = \int_0^\infty \mathbb{I}_{[0, X(\omega))} dt = \Lambda[0, X(\omega)] = X(\omega).$$

The double integral on the right-hand side of (X.12) is therefore

$$\int_\Omega X(\omega) dP(\omega) = E[X]. \quad (\text{X.14})$$

By (X.12), the results (X.13) and (X.14) are equal, which proves the assertion.  $\square$

Variations of the idea above are very commonly used. As another example of the same idea, the following two exercises characterize the  $p$ -integrability of a non-negative random variable using its cumulative distribution function.

**Exercise X.6** (Moments with cumulative distribution function).

Let  $X$  be a non-negative random variable, and let  $F(x) = P[X \leq x]$  be its cumulative distribution function. Prove that for  $p \geq 1$  we have

$$E[X^p] = p \int_0^\infty t^{p-1} (1 - F(t)) dt$$

**Hint:** Try using as few of the following hints as possible:

- (1) Find a function  $f: [0, \infty) \rightarrow [0, \infty)$  such that for any  $x \geq 0$  one has  $x^p = \int_0^x f(t) dt$ .
- (2) Then we have  $X(\omega)^p = \int_0^{X(\omega)} f(t) dt$  for any  $\omega \in \Omega$ .

<sup>2</sup> Briefly, the measurability follows from writing  $A$  as the preimage  $A = g^{-1}((0, +\infty))$  of the Borel set  $(0, +\infty)$  under the function  $g: \Omega \times [0, +\infty) \rightarrow \mathbb{R}$  defined as  $g = X \circ \text{pr}_1 - \text{pr}_2$ , which is  $\mathcal{F} \otimes \mathcal{B}$ -measurable. Let us, however, spell out the details below.

Observe first that

$$(\omega, t) \mapsto \omega \quad \text{and} \quad (\omega, t) \mapsto t$$

are measurable with respect to the product  $\sigma$ -algebra  $\mathcal{F} \otimes \mathcal{B}$  directly according to Definition IX.2 — they are exactly the projections  $\text{pr}_1$  and  $\text{pr}_2$  from  $\Omega \times [0, +\infty)$  to the two factors. Further composing the former of these with the measurable function  $X: \Omega \rightarrow [0, +\infty)$  gives that  $(\omega, t) \mapsto X(\omega)$  is also  $\mathcal{F} \otimes \mathcal{B}$ -measurable. Then as a linear combination of these measurable functions, the function given by  $g(\omega, t) = X(\omega) - t$  is  $\mathcal{F} \otimes \mathcal{B}$ -measurable. Now we can recognize

$$A = \left\{ (\omega, t) \mid X(\omega) > t \right\} = \left\{ (\omega, t) \mid g(\omega, t) > 0 \right\} = g^{-1}((0, +\infty)).$$

(3) Now write the expected value of  $X^p$  as

$$\mathbb{E}[X^p] = \mathbb{E}\left[\int_0^{X(\omega)} f(t) dt\right],$$

and try to apply Fubini's theorem — but can you handle the random upper limit of the integral?

(4) Note that you could have alternatively written  $x^p = \int_0^x f(t) dt = \int_0^\infty \mathbb{I}_{[0,x)}(t) f(t) dt$ , and the upper limit would not have posed any problems.

**Exercise X.7** (Asymptotics of c.d.f. and  $p$ -integrability).

For a given  $p \geq 1$ , we say that a non-negative random variable  $X$  is  $p$ -integrable and write  $X \in \mathcal{L}^p$  if and only if the random variable  $|X|^p$  is integrable.

Assume that for some  $\alpha > 0$  the cumulative distribution function  $F$  of  $X$  satisfies

$$\lim_{x \rightarrow \infty} \left( x^\alpha (1 - F(x)) \right) = c > 0.$$

Prove that  $X \in \mathcal{L}^p$  if and only if  $p < \alpha$ .

**Hint:** Use the previous exercise.



## Probabilistic notions of convergence

Very often in stochastics, we want to assert that some sequence of random variables tends to a limit in a suitable probabilistic sense.

Some examples of such contexts could be:

- convergence of averages  $\frac{X_1 + \dots + X_n}{n}$  as  $n \rightarrow \infty$
- convergence of states  $X_t$  of a stochastic process as time  $t$  increases describes the long term behavior of the process
- limits of various random quantities as some parameter of the model tends to an idealized value, e.g.,
  - size of physical system  $\rightarrow \infty$  in thermodynamics and statistical physics
  - size of input data  $\rightarrow \infty$  in randomized algorithms
  - signal to noise ratio  $\rightarrow \infty$  in communications.

In this lecture and the next we will treat two very famous and important convergence results concerning sums of independent identically distributed real valued random variables  $X_1, X_2, \dots$ :

**(LLN):** “*Laws of large numbers*”: Under what assumptions and in which sense do the averages tend to the expected value

$$\frac{1}{n} \sum_{j=1}^n X_j \longrightarrow \mathbb{E}[X_i] ?$$

**(CLT):** “*Central limit theorems*”: Under what assumptions and in which sense do we have the normal approximation

$$\frac{1}{\sqrt{n}} \sum_{j=1}^n (X_j - \mathbb{E}[X_j]) \longrightarrow \mathcal{N}(0, \sigma^2) ?$$

### XI.1. Notions of convergence in stochastics

For the rest of this section, assume that

$X_1, X_2, X_3, \dots$  and  $X$  are real-valued random variables.

Let us then introduce two important probabilistic notions of convergence.

#### Convergence almost surely

Under any possible outcome  $\omega \in \Omega$ , the realized values  $X_1(\omega), X_2(\omega), \dots$  of the random variables form a sequence of real numbers. Therefore for a fixed  $\omega$ , the

meaning of the limit

$$\lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) \quad (\text{XI.1})$$

is familiar from undergraduate calculus. Now recall that we can form the event

$$E = \left\{ \omega \in \Omega \mid \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) \right\} \quad (\text{XI.2})$$

consisting of those outcomes  $\omega$  for which we have the limit (XI.1).

Pointwise convergence  $X_n \rightarrow X$  means that (XI.1) holds for all outcomes, i.e., that  $E = \Omega$ . As pointed out already in Example O.2, however, this is often too much to ask for. Instead, it is meaningful to ask about the probability  $\mathbb{P}[E]$  of the event (XI.2). If this event occurs almost surely, then we talk about almost sure limit.

**Definition XI.1** (Convergence almost surely).

We say that  $X_n$  tends to  $X$  almost surely as  $n \rightarrow \infty$ , if

$$\mathbb{P} \left[ \lim_{n \rightarrow \infty} X_n = X \right] = 1.$$

In this case we denote  $X_n \xrightarrow{\text{a.s.}} X$ .

This probabilistic notion of a limit should be intuitively easy to understand — we are giving up pointwise convergence only on an exceptional event  $E^c$  which has probability zero. Although we have thus relaxed the extremely stringent requirement of pointwise convergence, this is still a very strong notion of convergence.

### Convergence in probability

Occasionally the notion of almost sure convergence is still too much to hope for. The following notion of convergence is less restrictive.

**Definition XI.2** (Convergence in probability).

We say that  $X_n$  tends to  $X$  in probability as  $n \rightarrow \infty$ , if for all  $\varepsilon > 0$  we have

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ |X_n - X| < \varepsilon \right] = 1. \quad (\text{XI.3})$$

In this case we denote  $X_n \xrightarrow{\mathbb{P}} X$ .

**Remark XI.3.** The limit (XI.3) is equivalent to the following limit of probabilities of complementary events

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ |X_n - X| \geq \varepsilon \right] = 0. \quad (\text{XI.4})$$

Because there exist many techniques to upper bound probabilities, this latter formulation is used more often in practice.

### Comparison of the notions of convergence

To understand these notions and to appreciate their differences, it is instructive to unravel the definition of the limit (XI.1) and the limit in the above definition, and compare the notions:

**pointwise convergence:**  $X_n \rightarrow X$  as  $n \rightarrow \infty$  if

$$\forall \omega \in \Omega \quad \forall \varepsilon > 0 \quad \exists n_0(\omega, \varepsilon) \in \mathbb{N} \quad \forall n > n_0(\omega, \varepsilon) : \\ |X_n(\omega) - X(\omega)| < \varepsilon.$$

**convergence almost surely:**  $X_n \xrightarrow{\text{a.s.}} X$  as  $n \rightarrow \infty$  if

$$\mathbb{P} \left[ \left\{ \omega \in \Omega \mid \forall \varepsilon > 0 \quad \exists n_0(\omega, \varepsilon) \in \mathbb{N} \quad \forall n > n_0(\omega, \varepsilon) : \right. \right. \\ \left. \left. |X_n(\omega) - X(\omega)| < \varepsilon \right\} \right] = 1.$$

**convergence in probability:**  $X_n \xrightarrow{\mathbb{P}} X$  as  $n \rightarrow \infty$  if

$$\forall \varepsilon > 0 \quad \forall \delta > 0 \quad \exists n_0(\varepsilon, \delta) \in \mathbb{N} \quad \forall n > n_0(\varepsilon, \delta) : \\ \mathbb{P} \left[ \left\{ \omega \in \Omega \mid |X_n(\omega) - X(\omega)| < \varepsilon \right\} \right] > 1 - \delta.$$

We leave it as an exercise to the reader to verify that convergence almost surely is indeed stronger than convergence in probability.

**Exercise XI.1** (Convergence almost surely implies convergence in probability).

Let  $X_1, X_2, \dots$  be real valued random variables such that  $X_n \xrightarrow{\text{a.s.}} X$ . Show that  $X_n \xrightarrow{\mathbb{P}} X$ .

The converse is not true in general: convergence in probability does not imply almost sure convergence. However, it does imply almost sure convergence along some subsequence, as the following exercise shows.

**Exercise XI.2** (Convergence in probability implies convergence a.s. along a subsequence).

Assume that  $X_1, X_2, \dots$  are real valued random variables and  $X_n \xrightarrow{\mathbb{P}} X$ . Let  $(a_k)_{k \in \mathbb{N}}$  and  $(b_k)_{k \in \mathbb{N}}$  be two sequences of positive real numbers such that  $a_k \downarrow 0$  and  $\sum_{k=1}^{\infty} b_k < +\infty$  — for example  $a_k = \frac{1}{k}$  and  $b_k = 2^{-k}$ .

(a) Show that there exists  $(n_k)_{k \in \mathbb{N}}$  so that  $n_k \in \mathbb{N}$  for all  $k$  and  $n_1 < n_2 < \dots$  and

$$\mathbb{P} \left[ |X_{n_k} - X| \geq a_k \right] \leq b_k.$$

(b) With the sequence  $(n_k)_{k \in \mathbb{N}}$  chosen as in part (a), show that

$$\mathbb{P} \left[ |X_{n_k} - X| \geq a_k \text{ for infinitely many } k \right] = 0.$$

(c) With the sequence  $(n_k)_{k \in \mathbb{N}}$  chosen as in part (a), show that  $X_{n_k} \xrightarrow{\text{a.s.}} X$  as  $k \rightarrow \infty$ .

**Hint:** Recall a suitable Borel-Cantelli lemma.

Most notions of convergence behave well under continuous functions. The following exercise concerns the case of convergence in probability.

**Exercise XI.3** (Continuous transformations and convergence in probability).

Let  $X, X_1, X_2, \dots$  be real-valued random variables defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . Suppose that  $X_n \xrightarrow{\mathbb{P}} X$ .

(a) For any uniformly continuous<sup>1</sup>  $h: \mathbb{R} \rightarrow \mathbb{R}$ , show that we have  $h(X_n) \xrightarrow{\mathbb{P}} h(X)$ .

<sup>1</sup>The conclusion is actually valid for any continuous  $h: \mathbb{R} \rightarrow \mathbb{R}$ . Can you prove that, too?

(b) For any bounded<sup>2</sup> uniformly continuous  $h: \mathbb{R} \rightarrow \mathbb{R}$ , show that we have

$$\mathbb{E}[|h(X_n) - h(X)|] \rightarrow 0 \quad \text{and} \quad \mathbb{E}[h(X_n)] \rightarrow \mathbb{E}[h(X)].$$

## XI.2. Weak and strong laws of large numbers

Having introduced the notions of convergence of random variables, we now return to laws of large numbers. In the rest of this lecture we will prove two versions of laws of large numbers, which differ primarily by the notion of convergence used.

**Theorem XI.4** (Weak law of large numbers with bounded second moments).

Let  $X_1, X_2, \dots$  be independent real valued random variables. Assume that for some  $\mathbf{m} \in \mathbb{R}$  and  $K_2 < +\infty$  we have

$$\mathbb{E}[X_j] = \mathbf{m} \quad \text{and} \quad \mathbb{E}[X_j^2] \leq K_2 \quad \text{for all } j \in \mathbb{N}.$$

Then we have

$$\frac{1}{n}(X_1 + \dots + X_n) \xrightarrow{\mathbb{P}} \mathbf{m} \quad \text{as } n \rightarrow \infty.$$

**Theorem XI.5** (Strong law of large numbers with bounded fourth moments).

Let  $X_1, X_2, \dots$  be independent real valued random variables. Assume that for some  $\mathbf{m} \in \mathbb{R}$  and  $K_4 < +\infty$  we have

$$\mathbb{E}[X_j] = \mathbf{m} \quad \text{and} \quad \mathbb{E}[X_j^4] \leq K_4 \quad \text{for all } j \in \mathbb{N}.$$

Then we have

$$\frac{1}{n}(X_1 + \dots + X_n) \xrightarrow{\text{a.s.}} \mathbf{m} \quad \text{as } n \rightarrow \infty.$$

Theorem XI.4 is called a *weak law of large numbers*, because it asserts the convergence of the averages in the weaker sense of convergence in probability. Correspondingly, Theorem XI.5 is called a *strong law of large numbers*, because it asserts the convergence of the averages in the stronger sense of convergence almost surely. In both statements, we have chosen to make assumptions on boundedness of suitable moments, in order to simplify their proofs. Versions of weak and strong laws of large numbers could be made with various different assumptions. In particular, in Section XI.5 we present the statement of one of the most general formulations of such results: Kolmogorov's strong law of large numbers. Such improvements are occasionally important, but the main objectives for the present lecture are simply to appreciate the differences in the notions of convergence, and to get a flavor of techniques that are used in proving laws of large numbers.

The following exercise requires some calculations, but it illustrates that weak and strong laws of large numbers are valid under different assumptions.

**Exercise XI.4** (A subtle sequence of averages).

Let  $X_3, X_4, \dots$  be independent random variables such that for  $k = 3, 4, \dots$  we have

$$\mathbb{P}[X_k = 0] = 1 - \frac{1}{k \log(k)} \quad \text{and} \quad \mathbb{P}[X_k = +k] = \frac{1}{2k \log(k)} = \mathbb{P}[X_k = -k].$$

<sup>2</sup> The conclusions are not valid without the assumption of boundedness. Can you give a counterexample in that case?



- (a) Calculate the expected value and variance of  $X_k$ .
- (b) Show that we have

$$\sum_{j=3}^{\infty} \frac{1}{j \log(j)} = \infty \quad \text{and} \quad \frac{1}{n^2} \sum_{j=3}^n \frac{j}{\log(j)} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

**Hint:** Recall that the integral function of  $x \mapsto \frac{1}{x \log(x)}$  is  $x \mapsto \log(\log(x))$ . Consider breaking the second sum into two pieces,  $j \leq b_n$  and  $j > b_n$  with a judiciously chosen  $b_n$ , so that you can easily estimate the pieces.

For  $n \geq 3$ , define the average

$$A_n = \frac{1}{n-2} \sum_{k=3}^n X_k.$$

- (c) Does the sequence  $(A_n)_{n=3,4,\dots}$  of averages converge almost surely?
- (d) Does the sequence  $(A_n)_{n=3,4,\dots}$  of averages converge in probability?

**Hint:** For this part you can use Chebyshev's inequality presented below, in Corollary XI.7.

### XI.3. Proof of the weak law of large numbers

The weak law of large numbers, Theorem XI.4, follows easily from inequalities of Markov and Chebyshev, which we present below.

#### Markov's inequality and Chebyshev's inequality

While almost obvious, the following observation is extremely useful.

**Lemma XI.6** (Markov's inequality).

If  $X: \Omega \rightarrow \mathbb{R}$  is a random variable, then for any  $a > 0$  we have

$$\mathbb{P}[|X| \geq a] \leq \frac{1}{a} \mathbb{E}[|X|]. \tag{XI.5}$$

*Proof.* Let  $a > 0$  and define the event

$$E = \left\{ \omega \in \Omega \mid |X(\omega)| \geq a \right\} = X^{-1}((-\infty, -a] \cup [a, \infty)).$$

For all  $\omega \in \Omega$  we then have

$$|X(\omega)| \geq a \mathbb{I}_E(\omega),$$

so by monotonicity of expected value we get

$$\mathbb{E}[|X|] \geq \mathbb{E}[a \mathbb{I}_E] = a \mathbb{P}[E].$$

Dividing this by  $a$  gives the asserted inequality  $\mathbb{P}[E] \leq \frac{1}{a} \mathbb{E}[|X|]$ . □

Markov's inequality leads to the following upper bound for the probability of fluctuations of a random variable  $X$  from its mean by more than a multiple of its standard deviation.

**Corollary XI.7** (Chebyshev's inequality).

Suppose that  $X \in \mathcal{L}^2(\mathbb{P})$ . Denote  $\mathbf{m} := \mathbb{E}[X]$  and  $\mathbf{s}^2 := \text{Var}(X) = \mathbb{E}[(X - \mathbf{m})^2]$ .

Then for any  $c > 0$  we have

$$\mathbb{P}\left[|X - \mathbf{m}| \geq c\right] \leq \frac{\mathfrak{s}^2}{c^2}. \quad (\text{XI.6})$$

*Proof.* Let  $Y = (X - \mathbf{m})^2$ , so  $\mathbb{E}[|Y|] = \mathfrak{s}^2$ . Note that the event of interest can be written as

$$E = \left\{ \omega \in \Omega \mid |X(\omega) - \mathbf{m}| \geq c \right\} = \left\{ \omega \in \Omega \mid |Y(\omega)| \geq c^2 \right\}.$$

The assertion follows by applying Markov's inequality (XI.5) to  $Y$  with  $a = c^2$ :

$$\mathbb{P}\left[|Y| \geq c^2\right] \leq \frac{1}{c^2} \mathbb{E}[|Y|] = \frac{\mathfrak{s}^2}{c^2}.$$

□

### Proving the weak law of large numbers with Chebyshev's inequality

With Chebyshev's inequality, it is easy to prove a weak law of large numbers.

*Proof of Theorem XI.4.* Let  $S_n = \sum_{j=1}^n X_j$  and  $Y_n = \frac{1}{n} S_n$ . We want to show  $Y_n \xrightarrow{\mathbb{P}} \mathbf{m}$ . Note that by linearity of expected value we have

$$\mathbb{E}[S_n] = \sum_{j=1}^n \underbrace{\mathbb{E}[X_j]}_{=\mathbf{m}} = n\mathbf{m} \quad \text{and} \quad \mathbb{E}[Y_n] = \frac{1}{n} \mathbb{E}[S_n] = \frac{1}{n} n\mathbf{m} = \mathbf{m}.$$

Note also that by assumption  $\mathbb{E}[X_j^2] \leq K_2$  we get a bound on variance,

$$\text{Var}(X_j) = \mathbb{E}[X_j^2] - \mathbb{E}[X_j]^2 \leq \mathbb{E}[X_j^2] \leq K_2.$$

By independence we have (see Proposition X.12 and Exercise X.5)

$$\text{Var}(S_n) = \sum_{j=1}^n \underbrace{\text{Var}(X_j)}_{\leq K_2} \leq nK_2 \quad \text{and} \quad \text{Var}(Y_n) = \text{Var}\left(\frac{1}{n} S_n\right) = \frac{1}{n^2} \text{Var}(S_n) \leq \frac{K_2}{n}.$$

Therefore, applying Chebyshev's inequality to  $Y_n$ , we get for any  $\varepsilon > 0$

$$\mathbb{P}\left[|Y_n - \mathbf{m}| \geq \varepsilon\right] \leq \frac{\text{Var}(Y_n)}{\varepsilon^2} \leq \frac{K_2}{n\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0,$$

which shows the convergence in probability  $Y_n \xrightarrow{\mathbb{P}} \mathbf{m}$ . □

The ideas in the weak law of large numbers can be used to prove *Weierstrass' approximation theorem*, which states that polynomials are dense in the space of continuous functions on a compact interval, with respect to the uniform norm  $\|\cdot\|_\infty$  (see Example B.10). The following exercise has the proof broken down to six easy steps.

#### Exercise XI.5 (Weierstrass' approximation theorem).

Let  $U_1, U_2, \dots$  be independent random variables with uniform distribution on  $[0, 1]$ . For  $p \in [0, 1]$ , consider the events  $E_1^{(p)}, E_2^{(p)}, \dots$  defined by  $E_j^{(p)} = \{U_j \leq p\}$ .

- Define  $S_n^{(p)} = \sum_{j=1}^n \mathbb{I}_{E_j^{(p)}}$ . Calculate the expected value  $\mathbb{E}[S_n^{(p)}]$ , and show that the variance is  $\text{Var}(S_n^{(p)}) \leq \frac{n}{4}$ .
- Show that for any  $\delta > 0$  we have  $\mathbb{P}\left[\left|S_n^{(p)}/n - p\right| \geq \delta\right] \leq \frac{1}{4n\delta^2}$ .

Let  $f: [0, 1] \rightarrow \mathbb{R}$  be some function.

- Show that  $B_n(p) := \mathbb{E}\left[f(S_n^{(p)}/n)\right]$  is a polynomial in  $p$ .

(d) Show that  $|B_n(p) - f(p)| \leq \mathbb{E} \left[ |f(S_n^{(p)}/n) - f(p)| \right]$ .

Suppose that  $f: [0, 1] \rightarrow \mathbb{R}$  is continuous. By compactness of  $[0, 1]$ , we know the following. The function  $f$  is bounded: there exists a  $K < +\infty$  such that  $|f(p)| \leq K$  for all  $p \in [0, 1]$ . Also,  $f$  is uniformly continuous: for all  $\varepsilon > 0$  there exists a  $\delta > 0$  such that  $|f(p) - f(q)| < \varepsilon$  whenever  $|p - q| < \delta$ .

(e) Let  $\varepsilon > 0$  and choose  $\delta > 0$  by uniform continuity of  $f$  as above. Consider the event  $A_n^{(p)} = \left\{ |S_n^{(p)}/n - p| \geq \delta \right\}$ . Show that we have

$$|f(S_n^{(p)}/n) - f(p)| \leq 2K \mathbb{I}_{A_n^{(p)}} + \varepsilon \mathbb{I}_{\Omega \setminus A_n^{(p)}}.$$

(f) Prove the Weierstrass' approximation theorem: for any  $\varepsilon > 0$  there exists a polynomial  $B_n$  such that  $|f(p) - B_n(p)| < 2\varepsilon$  for all  $p \in [0, 1]$ .

### XI.4. Proof of the strong law of large numbers

*Proof of Theorem XI.5.* We first prove the statement assuming  $\mathbf{m} = 0$ , and then reduce the general case to this particular case by appropriately centering the random variables.

*the case  $\mathbf{m} = 0$ :* Assume  $\mathbf{m} = 0$ . Denote

$$S_n = \sum_{j=1}^n X_j.$$

The assertion of the strong law of large numbers is that the event

$$E := \left\{ \omega \in \Omega \mid \lim_{n \rightarrow \infty} \frac{S_n(\omega)}{n} = 0 \right\}$$

occurs almost surely. Instead of considering  $E$  directly, however, we will consider another event

$$E' := \left\{ \omega \in \Omega \mid \sum_{n=1}^{\infty} \left( \frac{S_n(\omega)}{n} \right)^4 < +\infty \right\},$$

and show that it is almost sure, and that it implies the original event  $E$ .

To see that the occurrence of  $E'$  implies the occurrence of  $E$ , note that by definition of  $E'$  we have

$$\sum_{n=1}^{\infty} \left( \frac{S_n(\omega)}{n} \right)^4 < +\infty \quad \text{for all } \omega \in E'.$$

In particular, the terms of this convergent sum must tend to zero,

$$\left( \frac{S_n(\omega)}{n} \right)^4 \xrightarrow{n \rightarrow \infty} 0 \quad \text{for all } \omega \in E'.$$

Applying furthermore the continuous function  $t \mapsto t^{1/4}$  to this limit, we conclude that

$$\frac{S_n(\omega)}{n} \xrightarrow{n \rightarrow \infty} 0 \quad \text{for all } \omega \in E'.$$

This proves that  $E' \subset E$ , establishing the desired implication.

To prove that  $E'$  occurs almost surely, the main body of work consists of showing that

$$\sum_{n=1}^{\infty} \mathbb{E} \left[ \left( \frac{S_n}{n} \right)^4 \right] < +\infty. \tag{XI.7}$$

Once this is done, Lemmas VIII.6 and VIII.5 imply that we have  $\sum_{n=1}^{\infty} \left( \frac{S_n}{n} \right)^4 < +\infty$  almost surely, i.e., that  $\mathbb{P}[E'] = 1$ . This will finish the proof, because we then have

$$1 = \mathbb{P}[E'] \leq \mathbb{P}[E] \quad \text{since } E' \subset E.$$

It remains to prove (XI.7). Note first of all that since  $X_j \in \mathcal{L}^4(\mathbf{P})$  for all  $j$  by assumption, Lemma VIII.10 shows that also  $S_n = X_1 + \cdots + X_n \in \mathcal{L}^4(\mathbf{P})$ . We will compute the fourth moment of  $S_n$  by expanding the multinomial

$$\begin{aligned} S_n^4 &= (X_1 + \cdots + X_n)^4 \\ &= \sum_{1 \leq i \leq n} X_i^4 + \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \frac{4!}{3!1!} X_i^3 X_j + \frac{1}{2!} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \frac{4!}{2!2!} X_i^2 X_j^2 \\ &\quad + \frac{1}{2!} \sum_{\substack{1 \leq i, j, k \leq n \\ i, j, k \text{ different}}} \frac{4!}{2!1!2!} X_i^2 X_j X_k + \frac{1}{4!} \sum_{\substack{1 \leq i, j, k, \ell \leq n \\ i, j, k, \ell \text{ different}}} \frac{4!}{1!4!} X_i X_j X_k X_\ell. \end{aligned} \quad (\text{XI.8})$$

We want to take expected values of this expansion, so let us make some observations regarding that. By assumption we have  $X_i \in \mathcal{L}^4(\mathbf{P})$  for any  $i \in \mathbb{N}$ , so from Lemma VIII.9 we get finiteness of lower order moments  $\mathbb{E}[|X_i|] < \infty$ ,  $\mathbb{E}[|X_i|^2] < \infty$ , and  $\mathbb{E}[|X_i|^3] < \infty$  as well. For any  $i \neq j$  the random variables  $X_i$  and  $X_j$  were assumed independent, so we have also the independence of  $X_i^3 \in \mathcal{L}^1(\mathbf{P})$  and  $X_j \in \mathcal{L}^1(\mathbf{P})$  (recall Exercise V.2). By Theorem X.11, we then get

$$\mathbb{E}[X_i^3 X_j] = \mathbb{E}[X_i^3] \underbrace{\mathbb{E}[X_j]}_{=0} = 0.$$

Similarly one argues that

$$\mathbb{E}[X_i^2 X_j X_k] = \mathbb{E}[X_i^2] \underbrace{\mathbb{E}[X_j]}_{=0} \underbrace{\mathbb{E}[X_k]}_{=0} = 0$$

and

$$\mathbb{E}[X_i X_j X_k X_\ell] = \underbrace{\mathbb{E}[X_i]}_{=0} \underbrace{\mathbb{E}[X_j]}_{=0} \underbrace{\mathbb{E}[X_k]}_{=0} \underbrace{\mathbb{E}[X_\ell]}_{=0} = 0.$$

Only two terms of the expansion (XI.8) of the fourth power thus contribute to the expected value,

$$\mathbb{E}[S_n^4] = \sum_{1 \leq i \leq n} \mathbb{E}[X_i^4] + \frac{1}{2!} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \frac{4!}{2!2!} \mathbb{E}[X_i^2 X_j^2].$$

For the first contribution we have the bound  $\mathbb{E}[X_i^4] \leq K_4$  by assumption, and for the second we use Cauchy-Schwarz inequality

$$\mathbb{E}[X_i^2 X_j^2] \leq \sqrt{\mathbb{E}[X_i^4] \mathbb{E}[X_j^4]} \leq \sqrt{K_4 K_4} = K_4.$$

This gives

$$\begin{aligned} \mathbb{E}[S_n^4] &\leq \sum_{1 \leq i \leq n} K_4 + \frac{1}{2!} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \frac{4!}{2!2!} K_4 = n K_4 + n(n-1) \frac{1}{2!} \frac{4!}{2!2!} K_4 \\ &= n K_4 + (3n^2 - 3n) K_4 \\ &\leq 3n^2 K_4. \end{aligned}$$

With this estimate we get

$$\sum_{n=1}^{\infty} \mathbb{E}\left[\left(\frac{S_n}{n}\right)^4\right] = \sum_{n=1}^{\infty} \frac{1}{n^4} \mathbb{E}[S_n^4] \leq \sum_{n=1}^{\infty} \frac{3n^2 K_4}{n^4} = 3K_4 \sum_{n=1}^{\infty} \frac{1}{n^2} < +\infty.$$

This is exactly (XI.7), and the proof is complete for the case  $\mathbf{m} = 0$ .

*the general case:* Finally, let us reduce the general case,  $\mathbb{E}[X_j] = \mathbf{m} \in \mathbb{R}$  for all  $j \in \mathbb{N}$ , to the case above. We define  $\tilde{X}_j = X_j - \mathbf{m}$ . This is a “centered” version of the original term  $X_j$ , i.e., it has vanishing expected value

$$\mathbb{E}[\tilde{X}_j] = 0.$$

If we can apply the previously proven centered case of the strong law of large numbers, we get

$$\frac{1}{n} \sum_{j=1}^n \tilde{X}_j \xrightarrow{\text{a.s.}} 0.$$

But since  $\tilde{X}_j = X_j - \mathbf{m}$ , this translates to

$$\frac{1}{n} \sum_{j=1}^n X_j - \mathbf{m} \xrightarrow{\text{a.s.}} 0,$$

which in turn is equivalent to the desired almost sure convergence

$$\frac{1}{n} \sum_{j=1}^n X_j \xrightarrow{\text{a.s.}} \mathbf{m}.$$

The proof will thus be complete if we are permitted to apply the previous case to the centered random variables  $\tilde{X}_j$ . The independence of  $\tilde{X}_1, \tilde{X}_2, \dots$  follows from the assumed independence of  $X_1, X_2, \dots$  (as in Exercise V.2(a)). To derive the boundedness of the fourth moments of  $\tilde{X}_1, \tilde{X}_2, \dots$  from the assumed boundedness of the fourth moments of  $X_1, X_2, \dots$ , we use (VIII.3):

$$\mathbb{E}[|\tilde{X}_j|^4] = \mathbb{E}[|X_j - \mathbf{m}|^4] \leq 2^4 \left( \mathbb{E}[|X|^4] + \mathbb{E}[\mathbf{m}^4] \right) \leq 16(K_4 + \mathbf{m}^4).$$

We see that there exists a constant  $\tilde{K}_4 < +\infty$  (given by the last expression above) such that  $\mathbb{E}[\tilde{X}_j^4] \leq \tilde{K}_4$  for all  $j \in \mathbb{N}$ . This finishes the proof.  $\square$

## XI.5. Kolmogorov's strong law of large numbers

Our laws of large numbers, Theorems XI.4 and XI.5, were formulated with assumptions of bounded moments of suitable order. The existence of expected values only requires that the random variables are integrable,  $X_1, X_2, \dots \in \mathcal{L}^1(\mathbf{P})$ . We finish by giving the statement of a result of Kolmogorov, which does not use any higher moment assumptions. The result is a strong law of large numbers, because it gives almost sure convergence, but it also guarantees yet another notion of convergence which we introduce before the statement.

### Convergence in $\mathcal{L}^1$

The following notion of convergence already appeared implicitly in the Dominated convergence theorem (Theorem VII.19).

**Definition XI.8** (Convergence in  $\mathcal{L}^1$ ).

Suppose that  $X_1, X_2, \dots \in \mathcal{L}^1(\mathbf{P})$  and  $X \in \mathcal{L}^1(\mathbf{P})$ . We say that  $X_n$  tends to  $X$  in  $\mathcal{L}^1$  as  $n \rightarrow \infty$ , if we have

$$\lim_{n \rightarrow \infty} \mathbb{E}[|X_n - X|] = 0.$$

In this case we denote  $X_n \xrightarrow{\mathcal{L}^1} X$ .

The notion is stronger than convergence in probability.

**Exercise XI.6** (Convergence in  $\mathcal{L}^1$  implies convergence in probability).

Assume that  $X_1, X_2, \dots \in \mathcal{L}^1(\mathbf{P})$  are integrable random variables and  $X_n \xrightarrow{\mathcal{L}^1} X$ . Show that  $X_n \xrightarrow{\mathbf{P}} X$ .

**Hint:** Apply Markov's inequality.

### Statement of Kolmogorov's strong law of large numbers

**Theorem** (Kolmogorov's strong law of large numbers).

Let  $X_1, X_2, \dots \in \mathcal{L}^1(\mathbf{P})$  be independent and identically distributed integrable random variables with  $\mathbf{E}[X_j] = \mathbf{m}$  for all  $j \in \mathbb{N}$ . Then we have

$$\frac{1}{n}(X_1 + \dots + X_n) \xrightarrow{\text{a.s.}} \mathbf{m}$$

and

$$\frac{1}{n}(X_1 + \dots + X_n) \xrightarrow{\mathcal{L}^1} \mathbf{m}$$

as  $n \rightarrow \infty$ .

The proof can be found in, e.g., [Wil91, Chapters 12 and 14].

For the expected value to make sense, we had to at least assume  $X_j \in \mathcal{L}^1(\mathbf{P})$  in the above statement. In Theorems XI.4 and XI.5 we assumed more — that  $X_j \in \mathcal{L}^2(\mathbf{P})$  and  $X_j \in \mathcal{L}^4(\mathbf{P})$ , respectively. Note, however, that while the above result is thus relaxing the moment assumptions, it is not strictly speaking a generalization of the strong law of Theorem XI.5, because it assumes that the sequence consists of identically distributed random variables. Both formulations are useful.

## Central limit theorem and convergence in distribution

Consider a sequence  $X_1, X_2, \dots$  of independent and identically distributed random numbers, and form the sums

$$S_n = X_1 + \dots + X_n$$

of the first  $n$  members of the sequence. We are interested in the behavior of the sums  $S_n$  with a large number  $n$  of terms.

If the random variables are integrable,  $X_j \in \mathcal{L}^1(\mathbf{P})$ , with expected values  $\mathbf{E}[X_j] = \mathbf{m}$ , then we have  $\mathbf{E}[S_n] = n\mathbf{m}$ . The law of large numbers<sup>1</sup> then says that the sum  $S_n$  concentrates around the value  $n\mathbf{m}$ , when we look at it in a scale proportional to  $n$ , or more precisely

$$\frac{S_n - n\mathbf{m}}{n} \longrightarrow 0 \quad (\text{almost surely, in probability, and in } \mathcal{L}^1(\mathbf{P})).$$

If the random variables are square integrable,  $X_j \in \mathcal{L}^2(\mathbf{P})$ , with expected values  $\mathbf{E}[X_j] = \mathbf{m}$  and variances  $\mathbf{Var}(X_j) = \mathbf{s}^2$ , then by the independence of the terms we have  $\mathbf{Var}(S_n) = n\mathbf{s}^2$ . Chebyshev's inequality then says that the fluctuations of the sum  $S_n$  around the value  $n\mathbf{m}$  do not exceed a scale proportional to  $\sqrt{n}$ , or more precisely

$$\mathbf{P} \left[ \frac{|S_n - n\mathbf{m}|}{\sqrt{n}} \geq c \right] \leq \frac{n\mathbf{s}^2}{(c\sqrt{n})^2} = \frac{\mathbf{s}^2}{c^2} \quad (\text{for any } c > 0).$$

To understand the behavior of the sums  $S_n$  for large  $n$  in detail, it is therefore meaningful to look at  $S_n - n\mathbf{m}$  on a scale proportional to  $\sqrt{n}$ . One can show that

$$\frac{S_n - n\mathbf{m}}{\sqrt{n}}$$

does not converge anywhere almost surely or in probability, so in order to describe the limiting behavior, we need another notion of convergence. The appropriate notion is *convergence in distribution* (also known as *convergence in law* or *weak convergence*). Addressing exactly this, the Central limit theorem asserts that as  $n$  increases, the distribution of  $\frac{S_n - n\mathbf{m}}{\sqrt{n}}$  approaches a Gaussian distribution (Example VIII.3). An elementary interpretation of this is that the cumulative distribution functions have the following limit

$$\mathbf{P} \left[ \frac{S_n - n\mathbf{m}}{\mathbf{s}\sqrt{n}} \leq x \right] \xrightarrow{n \rightarrow \infty} \Phi(x) := \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \quad (\text{for all } x \in \mathbb{R}).$$

The right hand side above is the cumulative distribution of the standard normal distribution (Gaussian distribution with mean zero and unit variance).

<sup>1</sup>In this formulation, we need Kolmogorov's strong law of large numbers.

### XII.1. Characteristic functions

An important tool in the proof of the central limit theorem and in various other places in probability theory is characteristic functions. The characteristic function of a real valued random variable<sup>2</sup> is a single function, which encodes its distribution — it is essentially the Fourier transform of the law of the random variable.

Describing distributions in terms of their characteristic functions thus unifies the treatment of discrete distributions (which can be described by probability mass functions<sup>3</sup>) and continuous distributions (which can be described by probability densities<sup>4</sup>), as well as distributions, which are neither discrete nor continuous!

#### Complex valued random variables

The characteristic functions are complex valued, so we begin with a few remarks about complex valued random variables.

The set of complex numbers is denoted by  $\mathbb{C}$ , and the imaginary unit by  $\mathbf{i}$  ( $\mathbf{i} \in \mathbb{C}$  is a square root of  $-1$ , i.e.,  $\mathbf{i}^2 = -1$ ). A complex number  $z \in \mathbb{C}$  can be written uniquely as  $z = x + \mathbf{i}y$  with  $x, y \in \mathbb{R}$ . We call  $x = \Re(z)$  the real part and  $y = \Im(z)$  the imaginary part of  $z$ . The absolute value (or modulus) of  $z$  is  $|z| := \sqrt{x^2 + y^2}$ . We identify  $\mathbb{C}$  with the Euclidean plane  $\mathbb{R}^2$  (the set of complex numbers forms the “complex plane”), so that  $x$  and  $y$  are the two coordinates of  $z = x + \mathbf{i}y$  in the plane. The absolute value  $|z|$  is exactly the Euclidean norm, and we equip  $\mathbb{C}$  with the topology of the plane, and correspondingly with the Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{C}) \cong \mathcal{B}(\mathbb{R}^2) = \mathcal{B} \otimes \mathcal{B}$ .

A *complex valued random variable*  $Z = X + \mathbf{i}Y$  is constructed out of a pair  $X, Y$  of real valued random variables<sup>5</sup> (the real and imaginary parts of  $Z$ ). We say that  $Z$  is integrable if both its real and imaginary parts are, i.e., if  $X, Y \in \mathcal{L}^1(\mathbf{P})$ . We define the expected value of an integrable complex random variable  $Z$  to be the complex number

$$\mathbf{E}[Z] := \mathbf{E}[X] + \mathbf{i}\mathbf{E}[Y]$$

whose real and imaginary parts are the expected values of the real and imaginary parts of  $Z$ . Note that for  $z = x + \mathbf{i}y$  we have  $|x + \mathbf{i}y| \leq |x| + |y|$ . From monotonicity, linearity, and triangle inequality we therefore get

$$\mathbf{E}[|Z|] \leq \mathbf{E}[|X| + |Y|] = \mathbf{E}[|X|] + \mathbf{E}[|Y|] < \infty,$$

if  $Z$  is integrable.

The exponential of a complex number  $z \in \mathbb{C}$  is defined by the convergent power series

$$e^z = \sum_{n=0}^{\infty} \frac{1}{n!} z^n.$$

<sup>2</sup>More generally, characteristic functions of vector valued random variables could be defined, and have properties parallel with what we show in the setup of real valued random variables.

<sup>3</sup>See Definition II.15 and Exercise VIII.1, in particular.

<sup>4</sup>See Definition VIII.2 and Exercise VIII.3, in particular.

<sup>5</sup>As with random vectors,  $Z: \Omega \rightarrow \mathbb{C}$  is  $\mathcal{F}/\mathcal{B}(\mathbb{C})$ -measurable if and only if  $X, Y: \Omega \rightarrow \mathbb{R}$  are  $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable — see Remark X.2.



As a particular case, for  $z = i\phi$  with  $\phi \in \mathbb{R}$  we have Euler's formula

$$e^{i\phi} = \cos(\phi) + i \sin(\phi).$$

Note that we have  $|e^{i\phi}| = \sqrt{\cos^2(\phi) + \sin^2(\phi)} = 1$ . It is often convenient to write a complex number  $z \in \mathbb{C}$  in polar coordinates as  $z = r e^{i\phi}$ , where  $r = |z|$  and  $\phi \in [0, 2\pi)$ .

Expected values of complex valued random variables have familiar properties.

**Proposition XII.1** (Properties of expected values of complex random variables).

**Linearity:** If  $c_1, c_2 \in \mathbb{C}$  are complex numbers and  $Z_1, Z_2$  are integrable  $\mathbb{C}$ -valued random variables, then also  $c_1 Z_1 + c_2 Z_2$  is an integrable  $\mathbb{C}$ -valued random variable and  $\mathbf{E}[c_1 Z_1 + c_2 Z_2] = c_1 \mathbf{E}[Z_1] + c_2 \mathbf{E}[Z_2]$ .

**Triangle inequality:** If  $Z$  is an integrable  $\mathbb{C}$ -valued random variable, then we have  $|\mathbf{E}[Z]| \leq \mathbf{E}[|Z|]$ .

**Dominated convergence:** Suppose that  $Z_1, Z_2, \dots$  are  $\mathbb{C}$ -valued random variables and  $X \in \mathcal{L}^1(\mathbf{P})$  is an integrable random variable which dominates the absolute values,  $|Z_n| \leq X$  for all  $n \in \mathbb{N}$ . Then if the pointwise limit  $\lim_{n \rightarrow \infty} Z_n$  exists, we have  $\mathbf{E}[\lim_{n \rightarrow \infty} Z_n] = \lim_{n \rightarrow \infty} \mathbf{E}[Z_n]$ .

*Proof.* Linearity is proved directly from the definition by splitting each of  $c_1, c_2, Z_1, Z_2$  to real and imaginary parts. We leave the details as an exercise.

Triangle inequality can be proved as follows. The expected value is a complex number, so we can write it in polar coordinates as  $\mathbf{E}[Z] = r e^{i\phi}$ , where  $r = |\mathbf{E}[Z]|$  and  $\phi \in [0, 2\pi)$ . Then by linearity we have

$$r = e^{-i\phi} \mathbf{E}[Z] = \mathbf{E}[e^{-i\phi} Z] = \mathbf{E}\left[\Re(e^{-i\phi} Z)\right] + i \mathbf{E}\left[\Im(e^{-i\phi} Z)\right].$$

Since  $r \in \mathbb{R}$  by construction, the second term in fact has to vanish:  $\mathbf{E}[\Im(e^{-i\phi} Z)] = 0$ . If we furthermore use the fact that  $\Re(z) \leq |z|$  and monotonicity of real expected values, we thus get

$$r = \mathbf{E}\left[\Re(e^{-i\phi} Z)\right] \leq \mathbf{E}\left[|e^{-i\phi} Z|\right] = \mathbf{E}\left[|Z|\right].$$

Remembering that  $r = |\mathbf{E}[Z]|$ , this gives the triangle inequality.

Dominated convergence follows by splitting  $Z_n$  to real and imaginary parts and applying dominated convergence theorem to these separately: the same integrable random variable  $X$  which dominates  $|Z_n|$  also dominates  $\Re(Z_n)$  and  $\Im(Z_n)$ .  $\square$

**Exercise XII.1** (Independent complex random variables).

Let  $Z_1 = X_1 + iY_1$  and  $Z_2 = X_2 + iY_2$  be two independent, integrable complex valued random variables. Show that the product  $Z_1 Z_2$  is integrable, and that we have

$$\mathbf{E}[Z_1 Z_2] = \mathbf{E}[Z_1] \mathbf{E}[Z_2].$$

### Definition and first properties of characteristic functions

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : \Omega \rightarrow \mathbb{R}$  a real valued random variable. Note that for any  $\theta \in \mathbb{R}$  and any  $\omega \in \Omega$ , we have

$$e^{i\theta X(\omega)} = \cos(\theta X(\omega)) + i \sin(\theta X(\omega)). \tag{XII.1}$$

This shows that the real and imaginary parts of  $e^{i\theta X}$  are bounded random variables, and thus in particular integrable. Therefore the following definition makes sense.

**Definition XII.2** (Characteristic function).

The *characteristic function* of  $X$  is the function  $\varphi_X: \mathbb{R} \rightarrow \mathbb{C}$  given by

$$\begin{aligned}\varphi_X(\theta) &= \mathbb{E}\left[e^{i\theta X}\right] \\ &= \mathbb{E}\left[\cos(\theta X)\right] + i \mathbb{E}\left[\sin(\theta X)\right].\end{aligned}\tag{XII.2}$$

**Remark XII.3.** The function  $x \mapsto e^{i\theta x}$  is continuous, and therefore a Borel function by Corollary III.10 (i.e. the real and imaginary parts  $x \mapsto \cos(\theta x)$  and  $x \mapsto \sin(\theta x)$  are). Therefore by Theorem VIII.1, the expected value in (XII.2) can be written using the distribution  $P_X$  of  $X$ ,

$$\varphi_X(\theta) = \mathbb{E}\left[e^{i\theta X}\right] = \int_{\mathbb{R}} e^{i\theta x} dP_X(x).$$

This shows that the characteristic function  $\varphi_X$  of  $X$  only depends on the distribution  $P_X$  of  $X$ . Soon we will show that  $\varphi_X$  in fact contains enough information to fully determine the distribution  $P_X$ .

Let us give a few examples of characteristic functions.

**Example XII.4** (Characteristic function of exponential distribution).

Suppose that  $X \sim \text{Exp}(\lambda)$  is exponentially distributed with parameter  $\lambda > 0$  (see Example VIII.4), i.e.,  $X$  has a probability density

$$f_X(x) = \lambda e^{-\lambda x} \mathbb{I}_{[0, +\infty)}(x).$$

Let us compute its characteristic function using the formula of Exercise VIII.3,

$$\begin{aligned}\varphi_X(\theta) &= \mathbb{E}\left[e^{i\theta X}\right] = \int_{\mathbb{R}} e^{i\theta x} f_X(x) dx \\ &= \int_0^{\infty} e^{i\theta x} \lambda e^{-\lambda x} dx = \lambda \int_0^{\infty} e^{x(-\lambda + i\theta)} dx = \lambda \frac{-1}{-\lambda + i\theta} = \frac{1}{1 + i\theta/\lambda}.\end{aligned}$$

**Example XII.5** (Characteristic function of Poisson distribution).

Suppose that  $X \sim \text{Poisson}(\lambda)$  is Poisson distributed with parameter  $\lambda > 0$  (see Example II.16), i.e.,  $X$  has a probability mass function

$$p_X(n) = \mathbb{P}[X = n] = e^{-\lambda} \frac{\lambda^n}{n!} \quad \text{for } n \in \mathbb{Z}_{\geq 0} = \{0, 1, 2, \dots\}.$$

Let us compute its characteristic function using the formula of Exercise VIII.1,

$$\begin{aligned}\varphi_X(\theta) &= \mathbb{E}\left[e^{i\theta X}\right] = \sum_{n=0}^{\infty} p_X(n) e^{i\theta n} \\ &= \sum_{n=0}^{\infty} e^{-\lambda} \frac{\lambda^n}{n!} e^{i\theta n} = e^{-\lambda} \sum_{n=0}^{\infty} \frac{1}{n!} (\lambda e^{i\theta})^n = e^{-\lambda} e^{\lambda e^{i\theta}} = \exp(\lambda(e^{i\theta} - 1)).\end{aligned}$$

**Exercise XII.2** (The characteristic function of a standard Gaussian random variable).

Suppose that  $X \sim N(0, 1)$  is a real valued random variable with standard normal distribution (see Example VIII.3), i.e., a continuous distribution with density function

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad \text{for } x \in \mathbb{R}.$$

**Hint:** You can consider it known that  $\int_{-\infty}^{+\infty} f_X(x) dx = 1$ , and that the exponential of any complex number  $z \in \mathbb{C}$  is given by the convergent series  $e^z = \sum_{n=0}^{\infty} \frac{1}{n!} z^n$ .

(a) Let  $t \in \mathbb{R}$ . Show that  $\mathbb{E}\left[e^{tX}\right] = e^{t^2/2}$ .

**Hint:** Express the expected value in terms of the density, and perform a suitable change of variables  $x' = x + c$ .

(b) For  $x, t \in \mathbb{R}$ , show that  $e^{|tx|} \leq e^{tx} + e^{-tx}$ . Using this, prove that for any  $t \in \mathbb{R}$  we have

$$\mathbb{E} \left[ \sum_{n=0}^{\infty} \frac{1}{n!} |tX|^n \right] < +\infty.$$

(c) Prove that for any  $t \in \mathbb{R}$  we have

$$\mathbb{E}[e^{tX}] = \sum_{n=0}^{\infty} \frac{1}{n!} t^n \mathbb{E}[X^n].$$

(d) By comparing (a) with (c), deduce that for  $n \in \mathbb{N}$  we have

$$\mathbb{E}[X^n] = \begin{cases} \prod_{j=1}^{n/2} (2j-1) & \text{if } n \text{ is even} \\ 0 & \text{if } n \text{ is odd.} \end{cases}$$

(e) Prove that

$$\varphi_X(\theta) = e^{-\frac{1}{2}\theta^2} \quad \text{for } \theta \in \mathbb{R}.$$

We now state properties of characteristic functions that hold in general. You may directly inspect that the characteristic functions in Examples XII.4, XII.5 and Exercise XII.2 indeed have the stated properties.

**Proposition XII.6** (Basic properties of characteristic functions).

*Characteristic functions have the following properties:*

- (a) We have  $\varphi_X(0) = 1$ .
- (b) We have  $|\varphi_X(\theta)| \leq 1$  for all  $\theta \in \mathbb{R}$ .
- (c) The function  $\varphi_X: \mathbb{R} \rightarrow \mathbb{C}$  is continuous.
- (d) For any  $a, b \in \mathbb{R}$  we have  $\varphi_{aX+b}(\theta) = e^{ib\theta} \varphi_X(a\theta)$  for all  $\theta \in \mathbb{R}$ .
- (e) We have  $\varphi_{-X}(\theta) = \overline{\varphi_X(\theta)}$  for all  $\theta \in \mathbb{R}$ .
- (f) We have  $\varphi_X(-\theta) = \overline{\varphi_X(\theta)}$  for all  $\theta \in \mathbb{R}$ .

*Proof.* At  $\theta = 0$  we of course have  $\theta X(\omega) = 0$  for all  $\omega \in \Omega$  and thus  $e^{i\theta X(\omega)} = 1$ . We directly get  $\varphi_X(0) = \mathbb{E}[1] = 1$ , which proves (a).

Part (b) follows from triangle inequality:  $|\varphi_X(\theta)| = |\mathbb{E}[e^{i\theta X}]| \leq \mathbb{E}[|e^{i\theta X}|] = \mathbb{E}[1] = 1$ .

Continuity is proved as follows. We must show that for any sequence  $\theta_1, \theta_2, \dots \in \mathbb{R}$  such that  $\theta_n \rightarrow \theta$ , we have  $\varphi_X(\theta_n) \rightarrow \varphi_X(\theta)$ . Since  $\theta_n \rightarrow \theta$ , we get pointwise for all  $\omega \in \Omega$ , that  $e^{i\theta_n X(\omega)} \rightarrow e^{i\theta X(\omega)}$ , using the continuity of the exponential function. But the random variables  $e^{i\theta_n X}$  are also bounded, so we can use the Bounded convergence theorem (both real and imaginary parts are bounded real random variables which converge pointwise):

$$\varphi_X(\theta_n) = \mathbb{E}[e^{i\theta_n X}] \longrightarrow \mathbb{E}[e^{i\theta X}] = \varphi_X(\theta).$$

This proves part (c), the continuity of  $\varphi_X$ .

For part (d), observe that

$$e^{i\theta(aX(\omega)+b)} = e^{i\theta aX(\omega)} e^{i\theta b}$$

and use linearity.

Parts (e) and (f) follow by noting that the complex conjugate of  $e^{i\theta X(\omega)}$  is  $e^{-i\theta X(\omega)}$ .  $\square$

We can now for instance reduce the calculation of the characteristic function of a general Gaussian random variable to that of a standard Gaussian.

**Exercise XII.3** (Characteristic function of a Gaussian random variable).

Let  $\mathbf{m} \in \mathbb{R}$  and  $\mathfrak{s} > 0$ , and let  $X \sim N(\mathbf{m}, \mathfrak{s}^2)$  (see Example VIII.3). Use Exercise XII.2 and Proposition XII.6(d) to show that

$$\varphi_X(\theta) = e^{i\mathbf{m}\theta - \frac{1}{2}\mathfrak{s}^2\theta^2} \quad \text{for } \theta \in \mathbb{R}.$$

Another fundamental property of characteristic functions is that the characteristic function of a sum of independent terms is the pointwise product of the characteristic functions.

**Exercise XII.4** (Characteristic function of a sum of independent terms).

Suppose that  $X$  and  $Y$  are independent real valued random variables. Using Exercise XII.1, show that the characteristic function of their sum is

$$\varphi_{X+Y}(\theta) = \varphi_X(\theta) \varphi_Y(\theta) \quad \text{for } \theta \in \mathbb{R}.$$

### Lévy's inversion theorem

A fundamental property of the characteristic function of a random variable is that it contains all the information about the distribution of the random variable. This fact is made explicit by Lévy's inversion theorem, below.

**Theorem XII.7** (Lévy's inversion theorem).

Let  $X \in \mathfrak{m}\mathcal{F}$  be a real-valued random variable,  $P_X$  its distribution (a Borel probability measure on  $\mathbb{R}$ ), and  $\varphi_X: \mathbb{R} \rightarrow \mathbb{C}$  its characteristic function. Then for any  $a, b \in \mathbb{R}$ ,  $a < b$ , we have

$$\begin{aligned} & \lim_{T \rightarrow +\infty} \frac{1}{2\pi} \int_{-T}^{+T} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \varphi_X(\theta) \, d\theta \\ &= P_X[(a, b)] + \frac{1}{2} P_X[\{a\}] + \frac{1}{2} P_X[\{b\}]. \end{aligned}$$

In particular,  $\varphi_X$  uniquely determines  $P_X$ .

Moreover, if  $\int_{\mathbb{R}} |\varphi_X(\theta)| \, d\theta < +\infty$ , then  $X$  has a continuous probability density function  $f_X$  given by

$$f_X(x) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{-i\theta x} \varphi_X(\theta) \, d\theta.$$

The proof is given in Appendix F.

**Exercise XII.5** (Sum of independent Gaussian random variables).

Suppose that  $X_1 \sim N(\mathbf{m}_1, \mathfrak{s}_1^2)$  and  $X_2 \sim N(\mathbf{m}_2, \mathfrak{s}_2^2)$  are Gaussian random variables which are independent. Show that  $X_1 + X_2 \sim N(\mathbf{m}_1 + \mathbf{m}_2, \mathfrak{s}_1^2 + \mathfrak{s}_2^2)$ .

**Hint:** Use Lévy's inversion theorem together with Exercises XII.3 and XII.4.

**Exercise XII.6** (Sum of i.i.d. Bernoulli random variables).

- Let  $p \in [0, 1]$ . Calculate the characteristic function  $\varphi_B(\theta) = \mathbb{E}[e^{i\theta B}]$  of a random variable  $B$  such that  $\mathbb{P}[B = 1] = p$  and  $\mathbb{P}[B = 0] = 1 - p$  (we denote  $B \sim \text{Bernoulli}(p)$ ).
- Let  $p \in [0, 1]$  and  $n \in \mathbb{N}$ . Calculate the characteristic function  $\varphi_X(\theta) = \mathbb{E}[e^{i\theta X}]$  of a random variable  $X$  such that  $\mathbb{P}[X = k] = \binom{n}{k} p^k (1-p)^{n-k}$  for all  $k \in \{0, 1, 2, \dots, n\}$  (we denote  $X \sim \text{Bin}(n, p)$ ).

- (c) Let  $B_1, \dots, B_n$  be independent and identically distributed, with  $P[B_j = 1] = p$  and  $P[B_j = 0] = 1 - p$ , for all  $j$ . Compute the characteristic function of  $S = B_1 + \dots + B_n$  using part (a) and Exercise XII.4. Compare with the result of part (b), and conclude that  $S \sim \text{Bin}(n, p)$ .

### Taylor expansion of a characteristic function

By Lévy's inversion theorem, the characteristic function  $\varphi_X$  of a random variable  $X$  contains all information about the distribution  $P_X$  of  $X$ . In particular, it should contain the information about the expected value, variance, etc. To see why this is at least formally true, write the power series expansion

$$e^{i\theta X(\omega)} = \sum_{n=0}^{\infty} \frac{1}{n!} (i\theta X(\omega))^n = 1 + i\theta X(\omega) - \frac{1}{2}\theta^2 X(\omega)^2 + \dots \quad \text{for all } \omega \in \Omega.$$

If the expected value could be taken term by term in this expansion, then we would get

$$\varphi_X(\theta) = \mathbb{E} \left[ e^{i\theta X} \right] \stackrel{?}{=} 1 + i\theta \mathbb{E}[X] - \frac{1}{2}\theta^2 \mathbb{E}[X^2] + \dots$$

Formally, therefore, the expected value  $\mathbb{E}[X]$  seems to be encoded in the first order term in the Taylor expansion of  $\varphi_X(\theta)$  around the point  $\theta = 0$ , the variance  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$  in the terms up to order two, and more generally moment  $\mathbb{E}[X^n]$  of order  $n$  in the coefficient of  $\theta^n$ . Of course, this can only be meaningful if the random variable has moments of the correct order, i.e.,  $X \in \mathcal{L}^p(\mathbb{P})$  for high enough  $p \geq 1$ .

The following lemma makes precise sense of the above formal observation for square integrable random variables.<sup>6</sup>

**Proposition XII.8** (Taylor expansion of characteristic function).

Let  $X \in \mathcal{L}^2(\mathbb{P})$  be a square integrable random variable and let  $\varphi_X: \mathbb{R} \rightarrow \mathbb{C}$  be its characteristic function. Then we have

$$\varphi_X(\theta) = 1 + i\theta \mathbb{E}[X] - \frac{1}{2}\theta^2 \mathbb{E}[X^2] + \epsilon(\theta), \tag{XII.3}$$

where the function  $\epsilon: \mathbb{R} \rightarrow \mathbb{C}$  is an error term of smaller order than  $\theta^2$  in the sense that

$$\frac{|\epsilon(\theta)|}{|\theta|^2} \longrightarrow 0 \quad \text{as } \theta \rightarrow 0.$$

*Proof.* The idea is to Taylor expand  $e^{i\theta X}$  up to order two, with a controlled error term.

Start by observing that we have  $\frac{d}{du} e^{i\theta u} = i\theta e^{i\theta u}$ , so for any  $x \in \mathbb{R}$  we have

$$\int_0^x i\theta e^{i\theta u} du = e^{i\theta x} - 1.$$

Let us solve for  $e^{i\theta x}$ , and get

$$e^{i\theta x} = 1 + i\theta \int_0^x e^{i\theta u} du.$$

---

<sup>6</sup>The reader should think about how to modify the assumptions, statement, and proof to see moments up to order  $n$  in the Taylor expansion of the characteristic function.

Apply the same observation again to the integrand  $e^{i\theta u}$ , to get

$$\begin{aligned} e^{i\theta x} &= 1 + i\theta \int_0^x \left( 1 + i\theta \int_0^u e^{i\theta v} dv \right) du \\ &= 1 + i\theta x - \theta^2 \int_0^x \left( \int_0^u e^{i\theta v} dv \right) du. \end{aligned}$$

In this expression, write still  $e^{i\theta v} = 1 + (e^{i\theta v} - 1)$ , and perform the integrations of the first term to get

$$e^{i\theta x} = 1 + i\theta x - \theta^2 \frac{x^2}{2} - \theta^2 \int_0^x \left( \int_0^u (e^{i\theta v} - 1) dv \right) du. \quad (\text{XII.4})$$

The first three terms without integrations are the ones we care about, so let us introduce the following notation for the remainder that we want to get rid of,

$$R(\theta, x) := \int_0^x \left( \int_0^v (e^{i\theta v} - 1) dv \right) du$$

To estimate the magnitude of this remainder, note first that  $|R(\theta, -x)| = |R(\theta, x)|$ , so it is enough to consider  $x \geq 0$ . Then use the triangle inequality for integrals and the observation  $e^{i\theta v} - 1 = e^{i\theta v/2}(e^{i\theta v/2} - e^{-i\theta v/2}) = 2i e^{i\theta v/2} \sin(\theta v/2)$  to get the upper bound

$$|R(\theta, x)| \leq \int_0^{|x|} \left( \int_0^v |e^{i\theta v} - 1| dv \right) du \leq 2 \int_0^{|x|} \left( \int_0^v |\sin(\theta v/2)| dv \right) du$$

If we estimate the integrand in the last expression by  $|\sin(\theta v/2)| \leq 1$ , we find after the two integrations

$$|R(\theta, x)| \leq |x|^2, \quad (\text{XII.5})$$

and if we estimate it by  $|\sin(\theta v/2)| \leq \frac{1}{2} |\theta| |v|$ , we find after integrations

$$|R(\theta, x)| \leq \frac{1}{6} |\theta| |x|^3,$$

which in particular shows that for any  $x \in \mathbb{R}$  we have

$$|R(\theta, x)| \rightarrow 0 \quad \text{as } \theta \rightarrow 0. \quad (\text{XII.6})$$

Let us now apply (XII.4) pointwise to the values of the random variable  $X$ , to get

$$e^{i\theta X(\omega)} = 1 + i\theta X(\omega) - \frac{1}{2}\theta^2 X(\omega)^2 - \theta^2 R(\theta, X(\omega)).$$

With this, we can write the error term  $\epsilon(\theta)$  in the approximation (XII.3) in a manageable form. Namely, by linearity of expectation we have

$$\begin{aligned} \epsilon(\theta) &:= \varphi_X(\theta) - \left( 1 + i\theta \mathbf{E}[X] - \frac{1}{2}\theta^2 \mathbf{E}[X^2] \right) \\ &= \mathbf{E} \left[ e^{i\theta X} - 1 - i\theta X + \frac{1}{2}\theta^2 X^2 \right] \\ &= -\theta^2 \mathbf{E}[R(\theta, X)]. \end{aligned}$$

Then use the triangle inequality for expected values to control the magnitude of the error term,

$$|\epsilon(\theta)| \leq |\theta|^2 \mathbf{E}[|R(\theta, X)|].$$

The estimate (XII.5) shows that  $|R(\theta, X)| \leq |X|^2$  for any  $\theta$ , so by the assumption  $X \in \mathcal{L}^2(\mathbf{P})$  we have an integrable upper bound and we can use the Dominated convergence theorem in

$$\lim_{\theta \rightarrow 0} \mathbf{E}[|R(\theta, X)|] = \mathbf{E} \left[ \lim_{\theta \rightarrow 0} |R(\theta, X)| \right] \stackrel{(\text{XII.6})}{=} \mathbf{E}[0] = 0.$$

We conclude that

$$\frac{|\epsilon(\theta)|}{|\theta|^2} \leq \mathbf{E}[|R(\theta, X)|] \rightarrow 0$$

as  $\theta \rightarrow 0$ , and the proof is complete.  $\square$

## XII.2. Convergence in distribution

The notion of convergence in distribution is different from notion of convergence in Lecture XI: it is not strictly speaking about convergence of random variables (functions on the sample space  $\Omega$ ), but rather the convergence of their distributions (measures on  $\mathbb{R}$ , in the case of real-valued random variables).<sup>7</sup>

**Theorem XII.9** (Equivalent conditions for convergence in distribution).

Let  $X_1, X_2, \dots$  and  $X$  be real-valued random variables. Let also  $F_1, F_2, \dots$  and  $F$  be their cumulative distribution functions, and let  $\varphi_1, \varphi_2, \dots$  and  $\varphi$  be their characteristic functions, respectively. Then the following conditions are equivalent:

- (i) For all bounded continuous functions  $f: \mathbb{R} \rightarrow \mathbb{R}$  we have  $\mathbf{E}[f(X_n)] \rightarrow \mathbf{E}[f(X)]$  as  $n \rightarrow \infty$ .
- (ii) We have  $F_n(x) \rightarrow F(x)$  as  $n \rightarrow \infty$  for all points  $x \in \mathbb{R}$  such that  $F$  is continuous at  $x$ .
- (iii) We have  $\varphi_n(\theta) \rightarrow \varphi(\theta)$  as  $n \rightarrow \infty$  for all  $\theta \in \mathbb{R}$ .

The proof is given in Appendix F.

**Definition XII.10** (Convergence in distribution).

Let  $X_1, X_2, \dots$  and  $X$  be real-valued random variables. We say that  $X_n$  tend to  $X$  in distribution (or in law) as  $n \rightarrow \infty$  and denote  $X_n \xrightarrow{\text{law}} X$ , if any (then all) of the equivalent conditions of Theorem XII.9 hold.

**Remark XII.11** (Convergence in distribution formulated in terms of distributions).

Note that a bounded continuous function  $f: \mathbb{R} \rightarrow \mathbb{R}$  is a Borel function (Corollary III.10), so the expected values  $\mathbf{E}[f(X_n)]$  and  $\mathbf{E}[f(X)]$  can be written using the distributions  $P_{X_n}$  and  $P_X$  (Theorem VIII.1). Therefore condition (i) of Theorem XII.9 can be written as

$$\int_{\mathbb{R}} f \, dP_{X_n} \longrightarrow \int_{\mathbb{R}} f \, dP_X \quad \text{as } n \rightarrow \infty$$

for all such  $f$ . This formulation explains the terminology *convergence in distribution*.

## XII.3. Central limit theorem

We are now ready to state and prove the Central limit theorem (CLT). It is one of the most central<sup>8</sup> theorems in probability and statistics. It is for instance the rigorous justification for various normal approximations that are commonly used.

<sup>7</sup>In fact, for convergence in distribution it is not even necessary that the random variables are defined on the same probability space: we could have  $X: \Omega \rightarrow \mathbb{R}$  defined on  $(\Omega, \mathcal{F}, \mathbf{P})$  but  $X_n: \Omega_n \rightarrow \mathbb{R}$  each defined on its own probability space  $(\Omega_n, \mathcal{F}_n, \mathbf{P}_n)$ . In any case, the distributions  $P_{X_n}$  are probability measures on the real line  $\mathbb{R}$ , so the formulation of Remark XII.11 below still makes perfect sense. For simplicity of presentation, however, we choose not to keep explicitly mentioning and writing all the (possibly) different probability spaces during this lecture.

<sup>8</sup>pun intended / hence the name

**Theorem XII.12** (Central limit theorem).

Let  $X_1, X_2, \dots \in \mathcal{L}^2(\mathbf{P})$  be independent and identically distributed square integrable random variables. Denote

$$\mathbf{m} := \mathbf{E}[X_j] \quad \text{and} \quad \mathfrak{s} := \sqrt{\text{Var}(X_j)}.$$

Assume that  $\mathfrak{s} > 0$ .<sup>9</sup> For all  $n \in \mathbb{N}$ , let  $S_n = \sum_{j=1}^n X_j$ . Then we have

$$\frac{S_n - n\mathbf{m}}{\mathfrak{s}\sqrt{n}} \xrightarrow{\text{law}} Z \quad \text{as } n \rightarrow \infty, \quad (\text{XII.7})$$

where  $Z$  is a random variable with standard normal distribution  $\mathbf{N}(0, 1)$ .

*Proof.* By considering  $\tilde{X}_j = X_j - \mathbf{m}$  if necessary, we may assume that  $\mathbf{m} = 0$ . Likewise, by considering  $\frac{X_j}{\mathfrak{s}}$  if necessary, we may assume that  $\mathfrak{s} = 1$ . The goal is then to show that  $\frac{S_n}{\sqrt{n}} \xrightarrow{\text{law}} Z$ . We will prove this by verifying condition (iii) of Theorem XII.9, i.e., the pointwise convergence of the characteristic functions of  $\frac{S_n}{\sqrt{n}}$ .

By the assumption of identical distributions, the characteristic functions of all  $X_j$ ,  $j \in \mathbb{N}$ , are equal, so let us denote them by

$$\varphi(\theta) := \varphi_{X_j}(\theta) = \mathbf{E}[e^{i\theta X_j}].$$

By Proposition XII.8 and assumptions  $\mathbf{m} = 0$  and  $\mathfrak{s} = 1$ , we have

$$\varphi(\theta) = 1 - \frac{1}{2}\theta^2 + \epsilon(\theta), \quad (\text{XII.8})$$

where  $\frac{\epsilon(\theta)}{|\theta|^2} \rightarrow 0$  as  $\theta \rightarrow 0$ .

Now calculate the characteristic function of the sum  $S_n = \sum_{j=1}^n X_j$  using independence,

$$\varphi_{S_n}(\theta) = \mathbf{E}\left[e^{i\theta \sum_{j=1}^n X_j}\right] = \mathbf{E}\left[\prod_{j=1}^n e^{i\theta X_j}\right] \stackrel{(\perp)}{=} \prod_{j=1}^n \mathbf{E}[e^{i\theta X_j}] = \varphi(\theta)^n.$$

The characteristic function of  $\frac{S_n}{\sqrt{n}}$  is then

$$\varphi_{S_n/\sqrt{n}}(\theta) = \mathbf{E}\left[e^{i\theta S_n/\sqrt{n}}\right] = \varphi_{S_n}\left(\frac{\theta}{\sqrt{n}}\right) = \left(\varphi\left(\frac{\theta}{\sqrt{n}}\right)\right)^n.$$

By (XII.8), we have

$$\varphi\left(\frac{\theta}{\sqrt{n}}\right) = 1 - \frac{1}{2}\left(\frac{\theta}{\sqrt{n}}\right)^2 + \epsilon\left(\frac{\theta}{\sqrt{n}}\right) = 1 - \frac{\theta^2}{2n} + r_n,$$

where  $\frac{r_n}{1/n} \rightarrow 0$  as  $n \rightarrow \infty$ . By substituting this in the expression for the characteristic function of  $\frac{S_n}{\sqrt{n}}$ , we get

$$\varphi_{S_n/\sqrt{n}}(\theta) = \left(1 - \frac{\theta^2}{2n} + r_n\right)^n.$$

The limit in Exercise XII.7 gives

$$\lim_{n \rightarrow \infty} \varphi_{S_n/\sqrt{n}}(\theta) = \lim_{n \rightarrow \infty} \left(1 - \frac{\theta^2}{2n} + r_n\right)^n = e^{-\frac{1}{2}\theta^2}.$$

Since this is the characteristic function  $\varphi_Z(\theta)$  of a standard normal distributed random variable  $Z \sim \mathbf{N}(0, 1)$  according to Exercise XII.2, the proof is complete.  $\square$

**Exercise XII.7** (Complex exponential as a limit).

Let  $z \in \mathbb{C}$  and suppose that  $r_1, r_2, \dots \in \mathbb{C}$  are such that  $n r_n \rightarrow 0$  as  $n \rightarrow \infty$ . Show that

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

<sup>9</sup>This is always true unless  $X_j$  are almost surely constants. The case of almost surely constant random variables is not probabilistically interesting.



## Appendix A

### Set theory preliminaries

This appendix reviews necessary background on set theory, in particular the notion of countability, which is crucial in probability theory and measure theory.

#### A.1. Intersections and unions of sets

Let  $A, B$  be two sets. The *intersection*  $A \cap B$  is defined as the set of those elements which belong to both  $A$  and  $B$ ,

$$A \cap B = \{x \mid x \in A \text{ and } x \in B\}.$$

The *union*  $A \cup B$  is defined as the set of those elements which belong to at least one of the sets  $A$  and  $B$ ,

$$A \cup B = \{x \mid x \in A \text{ or } x \in B\}.$$

More generally, let  $(A_i)_{i \in I}$  be a collection of sets  $A_i$  indexed by  $i \in I$ . The union of the collection is defined as

$$\bigcup_{i \in I} A_i = \{x \mid x \in A_i \text{ for some } i \in I\}.$$

The intersection of the collection is defined provided the collection is non-empty ( $I \neq \emptyset$ ) as

$$\bigcap_{i \in I} A_i = \{x \mid x \in A_i \text{ for all } i \in I\}.$$

A collection  $(A_i)_{i \in I}$  of sets is said to be *disjoint* if no two different members of the collection have common elements, i.e., for all  $i, j \in I$ ,  $i \neq j$ , we have  $A_i \cap A_j = \emptyset$ . If the collection  $(A_i)_{i \in I}$  of sets is disjoint, then we say that the union  $\bigcup_{i \in I} A_i$  is a *disjoint union*. Disjoint unions enjoy additivity properties in probability theory and measure theory, according to the axiomatic properties in Definitions II.4 and II.5.

#### A.2. Set differences and complements

Let  $A, B$  be two sets. The *set difference*  $A \setminus B$  is defined as the set of those elements which belong to  $A$  but do not belong to  $B$ ,

$$A \setminus B = \{x \mid x \in A \text{ and } x \notin B\}.$$

When it is clear from the context that we are considering subsets of a particular reference set  $S$  (often the sample space  $S = \Omega$ ), then the *complement* of a subset  $A \subset S$ , denoted by  $A^c$ , is the set of those elements (of  $S$ ) which do not belong to  $A$

$$A^c = S \setminus A = \{x \in S \mid x \notin A\}.$$

The following basic results of set theory tell how unions and intersections behave under complements. They are known as *De Morgan's laws*.

**Proposition A.1** (De Morgan's laws).

Let  $(A_i)_{i \in I}$  be a non-empty indexed collection of subsets  $A_i \subset S$  of a fixed set  $S$ . Then we have

$$\left(\bigcup_{i \in I} A_i\right)^c = \bigcap_{i \in I} A_i^c \quad \text{and} \quad \left(\bigcap_{i \in I} A_i\right)^c = \bigcup_{i \in I} A_i^c.$$

**Exercise A.1.** Prove the De Morgan's laws above.

### A.3. Images and preimages of sets under functions

Let  $S$  and  $S'$  be two sets and

$$f: S \rightarrow S'$$

a function from  $S$  to  $S'$ . For  $A \subset S$ , the *image* of  $A$  under  $f$  is the subset  $f(A) \subset S'$  consisting of all those elements  $s' \in S'$  such that  $s' = f(s)$  for some  $s \in A$ ,

$$f(A) = \{f(s) \mid s \in A\} \subset S'.$$

For  $A' \subset S'$ , the *preimage* of  $A'$  under  $f$  is the subset  $f^{-1}(A') \subset S$  consisting of all those elements  $s \in S$  whose image  $f(s)$  belongs to the subset  $A'$ ,

$$f^{-1}(A') = \{s \in S \mid f(s) \in A'\} \subset S.$$

**Exercise A.2** (Properties of preimages).

Show that

- (a)  $f^{-1}(A'^c) = (f^{-1}(A'))^c$
- (b)  $f^{-1}\left(\bigcup_i A'_i\right) = \bigcup_i f^{-1}(A'_i)$
- (c)  $f^{-1}\left(\bigcap_i A'_i\right) = \bigcap_i f^{-1}(A'_i)$ .

### A.4. Cartesian products

Let  $A, B$  be two sets. The *Cartesian product*  $A \times B$  is defined as the set of ordered pairs  $(a, b)$  whose first member  $a$  belongs to the set  $A$  and second member  $b$  belongs to the set  $B$ , i.e.,

$$A \times B = \{(a, b) \mid a \in A, b \in B\}.$$

More generally, when  $A_1, \dots, A_n$  are sets, the  $n$ -fold Cartesian product  $A_1 \times \dots \times A_n$  is the set of ordered  $n$ -tuples  $(a_1, \dots, a_n)$  such that  $a_k \in A_k$  for each  $k = 1, \dots, n$

$$A_1 \times \dots \times A_n = \{(a_1, \dots, a_n) \mid a_1 \in A_1, \dots, a_n \in A_n\}.$$

Even more generally, if  $(A_j)_{j \in J}$  is a collection of sets indexed by  $j \in J$ , the Cartesian product  $\prod_{j \in J} A_j$  is the set of indexed collections  $(a_j)_{j \in J}$  such that  $a_j \in A_j$  for

each  $j \in J$ ,

$$\prod_{j \in J} A_j = \{(a_j)_{j \in J} \mid \forall j \in J : a_j \in A_j\}.$$

The *axiom of choice* states that if all the factors are non-empty,  $A_j \neq \emptyset$  for all  $j \in J$ , then the Cartesian product is also non-empty  $\prod_{j \in J} A_j \neq \emptyset$ .

As a particular case of Cartesian products, if each factor is the same set,  $A_j = A$  for all  $j \in J$ , then the Cartesian product is alternatively denoted by  $A^J := \prod_{j \in J} A$ . In that case an element of  $A^J$  is an indexed collection  $(a_j)_{j \in J}$  of elements of  $A$ , which can be naturally identified with the function  $j \mapsto a_j$  from  $J$  to  $A$ . Therefore  $A^J$  is identified with the set of functions from  $J$  to  $A$ ,

$$A^J = \{f : J \rightarrow A \text{ function}\}.$$

### A.5. Power set

Given a set  $S$ , the set  $\mathcal{P}(S)$  of all subsets  $A \subset S$  of it is called the *power set* of  $S$ ,

$$\mathcal{P}(S) = \{A \mid A \subset S\}.$$

A subset  $A \subset S$  can be specified by indicating for each element  $s \in S$  whether it belongs to  $A$  or not, so it is natural to identify the power set  $\mathcal{P}(S)$  of  $S$  with the set  $\{0, 1\}^S$  of functions  $S \rightarrow \{0, 1\}$ . In particular if  $S$  is a finite set with  $\#S = n$  elements, then its power set is a finite set with  $\#\mathcal{P}(S) = 2^n$  elements.

It is good to keep in mind that the power set readily provides an easy first example of many notions introduced in probability theory. For instance, in view of Definitions I.1, II.23, and C.3, the collection  $\mathcal{P}(S)$  of all subsets of  $S$  is obviously a sigma algebra on  $S$ , a  $\pi$ -system on  $S$ , a d-system on  $S$ , etc.

### A.6. Sequences of sets

A sequence  $A_1, A_2, \dots$  of sets is said to be *increasing* if

$$A_1 \subset A_2 \subset A_3 \subset \dots.$$

In this case we denote

$$A_n \uparrow A,$$

where the *limit*  $A$  of the increasing sequence of sets is defined as the union

$$A = \bigcup_{n=1}^{\infty} A_n.$$

Likewise, a sequence  $A_1, A_2, \dots$  of sets is said to be *decreasing* if

$$A_1 \supset A_2 \supset A_3 \supset \dots.$$

In this case we denote

$$A_n \downarrow A,$$

where the *limit*  $A$  of the decreasing sequence of sets is defined as the intersection

$$A = \bigcap_{n=1}^{\infty} A_n.$$

Sequences of sets which are either increasing or decreasing are said to be *monotone*.

**Exercise A.3** (Characterization of limits of monotone sequences of sets).

Show that the limits of monotone sequences of sets can be characterized as follows.

(a) Suppose that  $A_n \uparrow A$ . Show that

$$\begin{aligned} x \in A &\iff \exists m \in \mathbb{N} \text{ such that } \forall n \geq m : x \in A_n \\ x \notin A &\iff \forall n \in \mathbb{N} : x \notin A_n. \end{aligned}$$

(b) Suppose that  $A_n \downarrow A$ . Show that

$$\begin{aligned} x \in A &\iff \forall n \in \mathbb{N} : x \in A_n \\ x \notin A &\iff \exists m \in \mathbb{N} \text{ such that } \forall n \geq m : x \notin A_n. \end{aligned}$$

For a sequence  $A_1, A_2, \dots$  of sets, we define its *upper limit* as

$$\limsup_n A_n := \bigcap_{m \in \mathbb{N}} \bigcup_{n \geq m} A_n.$$

Note that if we define  $C_m = \bigcup_{n \geq m} A_n$ , then the sequence  $C_1, C_2, \dots$  of sets is decreasing, and its limit  $\bigcap_{m \in \mathbb{N}} C_m$  is precisely  $\limsup_n A_n$ .

We also define the *lower limit* of the sequence of sets as

$$\liminf_n A_n := \bigcup_{m \in \mathbb{N}} \bigcap_{n \geq m} A_n.$$

Note that if we define  $D_m = \bigcap_{n \geq m} A_n$ , then the sequence  $D_1, D_2, \dots$  of sets is increasing, and  $\liminf_n A_n$  is its limit.

**Exercise A.4** (Characterization of upper and lower limits of sequences of sets).

Show that the limsup and liminf of a sequence  $A_1, A_2, \dots$  of sets can be characterized as follows:

$$\begin{aligned} \text{(a):} \quad & \limsup_n A_n = \{s \mid \forall m \in \mathbb{N} : \exists n \geq m : s \in A_n\} \\ \text{(b):} \quad & \liminf_n A_n = \{s \mid \exists m \in \mathbb{N} : \forall n \geq m : s \in A_n\}. \end{aligned}$$

## A.7. Countable and uncountable sets

In probability theory we need to distinguish between sets of different sizes: finite sets, countably infinite sets, and uncountably infinite sets. In fact, if one had to summarize probability theory (and measure theory) in a single phrase, it might be:

All countable operations in probability theory are defined to behave just as intuition dictates.

## Comparison of cardinalities

Cardinality is the set-theoretic notion of the size or “number of elements” of a set. The idea is that when comparing the sizes of two sets  $A$  and  $B$ , we attempt to match the elements of  $A$  to the elements of  $B$  by functions  $f: A \rightarrow B$ . Under a surjective function, each element of  $B$  has at least some element(s) of  $A$  matched to it, and we then would interpret that  $A$  has at least as many elements as  $B$ . Under an injective function  $f$ , all elements of  $A$  are matched to some elements of  $B$  without any two different elements ever being matched to the same element, and then we would interpret that  $A$  has at most as many elements as  $B$ . Comparison of cardinalities is done by asking about the existence of such functions.

### Definition A.2 (Comparison of cardinalities).

Let  $A$  and  $B$  be sets. We say that the *cardinality* of  $A$  is less than or equal to the cardinality of  $B$  if there exists an injective function  $f: A \rightarrow B$ .

### Remark A.3 (Comparison of cardinalities of finite sets).

Suppose that  $A$  and  $B$  are two finite sets. Let  $n = \#A$  be the number of elements in  $A$  and  $m = \#B$  be the number of elements in  $B$ . Then it is easy to see that there exists an injective function  $f: A \rightarrow B$  if and only if we have  $n \leq m$ . Thus for finite sets, the comparison of cardinalities amounts to just the comparison of the number of elements.

### Example A.4 (Subsets can not have larger cardinality).

If  $A$  is a subset of  $B$ ,  $A \subset B$ , then the cardinality of  $A$  is less than or equal to the cardinality of  $B$ , because the inclusion mapping  $\iota: A \rightarrow B$  defined by  $\iota(x) = x$  for all  $x \in A$  is injective.

### Example A.5 (Transitivity of comparison of cardinalities).

Let  $A, B, C$  be sets. Suppose that the cardinality of  $A$  is less than or equal to the cardinality of  $B$  and the cardinality of  $B$  is less than or equal to the cardinality of  $C$ . In that case there exists injective functions  $f: A \rightarrow B$  and  $\tilde{f}: B \rightarrow C$ . The composition  $\tilde{f} \circ f: A \rightarrow C$  is also injective, so we get that the cardinality of  $A$  is less than or equal to the cardinality of  $C$ . In other words, the comparison of cardinalities is transitive.

As suggested before Definition A.2, instead of requiring the existence of injective functions in the comparison of cardinalities, one can alternatively require the existence of surjective functions in the opposite direction. The following two exercises establish this alternative characterization. To solve these exercises, you are allowed to use the axiom of choice.

### Exercise A.5 (Comparison of cardinalities with surjective functions: necessity).

Show that if the cardinality of a non-empty set  $A \neq \emptyset$  is less than or equal to the cardinality of a set  $B$ , then there exists a surjective function  $g: B \rightarrow A$ .

### Exercise A.6 (Comparison of cardinalities with surjective functions: sufficiency).

Show that if there exists a surjective function  $g: B \rightarrow A$ , then the cardinality of the set  $A$  is less than or equal to the cardinality of the set  $B$ .

## Equal cardinalities

### Definition A.6 (Equal cardinalities).

We say that  $A$  and  $B$  have *equal cardinalities* if there exists an injective function  $f: A \rightarrow B$  and an injective function  $g: B \rightarrow A$ .

Clearly if there exists a bijective function  $f: A \rightarrow B$ , then  $A$  and  $B$  have equal cardinalities (we may then take  $g = f^{-1}$ ). The converse is also true, but it is not as obvious. The Schröder - Bernstein theorem states that if  $A$  and  $B$  have equal cardinalities, then there exists a bijective function  $f: A \rightarrow B$  (you could try to prove this as an exercise).

## Countable sets

For the purposes of probability theory and measure theory, countable cardinalities are the most crucial. We begin with the definition.

### Definition A.7 (Countability).

A set  $A$  is said to be *countable* if the cardinality of  $A$  is less than or equal to the cardinality of the set  $\mathbb{N} = \{1, 2, 3, \dots\}$  of natural numbers.

### Example A.8 (Subsets of natural numbers are countable).

From Example A.4 it follows that any subset  $S \subset \mathbb{N}$ , including the set  $\mathbb{N}$  of natural numbers itself, is countable.

Since countable sets are so important, we unravel the definition once more, and provide an alternative characterization and two useful sufficient conditions.

### Lemma A.9 (Criteria for countability).

- (a) *A set  $A$  is countable if and only if there exists an injective function  $f: A \rightarrow \mathbb{N}$ .*
- (b) *A non-empty set  $A \neq \emptyset$  is countable if and only if there exists a surjective function  $g: \mathbb{N} \rightarrow A$ .*
- (c) *If  $B$  is a countable set and there exists an injective function  $f: A \rightarrow B$ , then also the set  $A$  is countable.*
- (d) *If  $B$  is a countable set and there exists a surjective function  $g: B \rightarrow A$ , then also the set  $A$  is countable.*

*Proof.* Assertion (a) follows directly by combining Definitions A.2 and A.7.

Assertion (b) follows by combining Definition A.7 with the characterization of Exercises A.5 and A.6.

Assertions (c) and (d) are similarly obtained using the transitivity in Example A.5.  $\square$

### Enumerations of countable sets

If  $A$  is countable and non-empty, then from Exercise A.5 it follows that there exists a surjective function  $g: \mathbb{N} \rightarrow A$ . We see that all elements of  $A$  are obtained in the following “list”

$$A = \{g(1), g(2), g(3), \dots\}.$$

Note, however, that repetitions are allowed in the above “list”, as  $g$  does not have to be injective. It is possible to remove repetitions and obtain an enumeration of the elements of  $A$ . To do this, one defines  $a_k \in A$  as the  $k$ :th value in the list above omitting repetitions. If the set  $A$  is finite, however, then there are only finitely many different values and the enumeration terminates at some point. Thus, for a finite set  $A$  with  $n$  elements, we have an enumeration

$$A = \{a_1, a_2, \dots, a_n\}.$$

An infinite set which is countable is said to be *countably infinite*, and for such a set  $A$ , we have an enumeration

$$A = \{a_1, a_2, a_3, \dots\}.$$

Note also that if the elements of  $A$  can be enumerated as above without repetition, then the mapping  $a_k \mapsto k$  is well defined and injective  $A \rightarrow \mathbb{N}$ . Therefore any set whose elements can be enumerated is countable.

### *Operations that preserve countability*

In probability theory and measure theory, countable operations work well. It is therefore crucial to understand clearly which set theoretic operations preserve countability.

Suppose that  $A_1$  and  $A_2$  are countable sets. By definition, there exists injective functions  $f_1: A_1 \rightarrow \mathbb{N}$  and  $f_2: A_2 \rightarrow \mathbb{N}$ . Consider the union  $A_1 \cup A_2$ , and note that it can be expressed as  $A_1 \cup A_2 = A_1 \cup (A_2 \setminus A_1)$ , where the latter is a disjoint union. The function  $f: A_1 \cup A_2 \rightarrow \mathbb{N}$  defined “piecewise” by

$$f(x) = \begin{cases} 2f_1(x) & \text{if } x \in A_1 \\ 2f_2(x) + 1 & \text{if } x \in A_2 \setminus A_1 \end{cases}$$

is clearly injective: it maps elements of  $A_1$  injectively to even natural numbers and the remaining elements injectively to odd natural numbers. From the existence of such an injective function we conclude that the union  $A_1 \cup A_2$  is countable. Using this argument inductively, we get that finite unions of countable sets remain countable.

**Lemma A.10** (Finite unions of countable sets are countable).

Let  $A_1, \dots, A_n$  be countable sets. Then the union

$$A_1 \cup \dots \cup A_n = \bigcup_{j=1}^n A_j$$

is also countable.

**Example A.11** (The set of integers is countable).

Consider the three sets:

$$\begin{aligned} A_1 &= \{1, 2, 3, \dots\} \\ A_2 &= \{0\} \\ A_3 &= \{-1, -2, -3, \dots\}. \end{aligned}$$

Each of them is countable: the set  $A_1 = \mathbb{N}$  is countable by Example A.8, the set  $A_2$  is countable because it is finite, and the set  $A_3$  is countable because it is in bijection with  $\mathbb{N}$  via  $x \mapsto -x$ . The set of all integers is the union of these three

$$\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\} = A_1 \cup A_2 \cup A_3$$

and as such  $\mathbb{Z}$  is itself countable by Lemma A.10

Consider now the set  $\mathbb{N} \times \mathbb{N}$ , the Cartesian product of the set of natural numbers with itself. We claim that  $\mathbb{N} \times \mathbb{N}$  is countable. To see this, note that the elements can be enumerated

$$\begin{aligned} \mathbb{N} \times \mathbb{N} = \{(n, m) \mid n, m \in \mathbb{N}\} = & \left\{ (1, 1), \right. \\ & (2, 1), (1, 2), \\ & (3, 1), (2, 2), (1, 3), \\ & (4, 1), (3, 2), (2, 3), (1, 4), \\ & \dots \left. \right\}. \end{aligned}$$

The enumeration shows that  $\mathbb{N} \times \mathbb{N}$  is also countable, as it gives rise to an injective function  $h: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ .

Suppose now that  $A_1$  and  $A_2$  are countable sets. Then there exists injective functions  $f_1: A_1 \rightarrow \mathbb{N}$  and  $f_2: A_2 \rightarrow \mathbb{N}$ . Now define  $f: A_1 \times A_2 \rightarrow \mathbb{N}$  by

$$f(x_1, x_2) = h(f_1(x_1), f_2(x_2)) \quad \text{for } x_1 \in A_1, x_2 \in A_2,$$

where  $h: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  is the injective function given by the above enumeration of  $\mathbb{N} \times \mathbb{N}$ . The function  $f: A_1 \times A_2 \rightarrow \mathbb{N}$  is injective, and we thus see that the Cartesian product  $A_1 \times A_2$  of countable sets  $A_1$  and  $A_2$  is again countable. Using this observation inductively, we get that Cartesian products of finitely many countable sets remain countable.

**Proposition A.12** (Finite Cartesian products of countable sets are countable).

*Let  $A_1, \dots, A_n$  be countable sets. Then the Cartesian product*

$$A_1 \times \dots \times A_n$$

*is also countable.*

**Example A.13** (The  $d$ -dimensional integer lattice  $\mathbb{Z}^d$  is countable).

Let  $d \in \mathbb{N}$ . The set

$$\mathbb{Z}^d = \underbrace{\mathbb{Z} \times \dots \times \mathbb{Z}}_{d \text{ times}} = \left\{ (x_1, \dots, x_d) \in \mathbb{R}^d \mid x_1, \dots, x_d \in \mathbb{Z} \right\}$$

of points with integer coordinates in the  $d$ -dimensional Euclidean space  $\mathbb{R}^d$  is countable, since it is the Cartesian product of  $d$  copies of the countable set  $\mathbb{Z}$ .

**Example A.14** (The set of rational numbers is countable).

Consider the set  $\mathbb{Q} \subset \mathbb{R}$  of rational numbers. The mapping  $g: \mathbb{Z} \times \mathbb{N} \rightarrow \mathbb{Q}$  defined by  $g(n, m) = \frac{n}{m}$  is surjective onto  $\mathbb{Q}$ . As a Cartesian product of the countable sets  $\mathbb{Z}$  and  $\mathbb{N}$ , the set  $\mathbb{Z} \times \mathbb{N}$  is countable by Proposition A.12. From Lemma A.9(d) and the existence of the surjective function  $g$  we get that also the set  $\mathbb{Q}$  of rational numbers is countable.

Using the above proposition about Cartesian products of countable sets we can strengthen our earlier observation about unions of countable sets: it turns out that countable unions of countable sets remain countable.



**Proposition A.15** (Countable unions of countable sets are countable).

Let  $A_1, A_2, A_3, \dots$  be countable sets. Then the union

$$\bigcup_{j=1}^{\infty} A_j = A_1 \cup A_2 \cup A_3 \cup \dots$$

is also countable.

*Proof.* We may assume that the sets  $A_j$  are non-empty,  $A_j \neq \emptyset$ , for all  $j \in \mathbb{N}$  (empty terms can be omitted from the union). Then, for each  $j \in \mathbb{N}$  there exists a surjective function

$$g_j: \mathbb{N} \rightarrow A_j.$$

We now define a function  $g$  on  $\mathbb{N} \times \mathbb{N}$  by

$$g(j, k) = g_j(k) \quad \text{for } j, k \in \mathbb{N},$$

and observe that this function

$$g: \mathbb{N} \times \mathbb{N} \rightarrow \bigcup_{j=1}^{\infty} A_j$$

is surjective onto the union  $\bigcup_{j=1}^{\infty} A_j$ . This shows that the cardinality of  $\bigcup_{j=1}^{\infty} A_j$  is less than or equal to the cardinality of  $\mathbb{N} \times \mathbb{N}$ . Since  $\mathbb{N} \times \mathbb{N}$  is countable, this shows that the union  $\bigcup_{j=1}^{\infty} A_j$  is also countable.  $\square$

## Uncountable sets

A set which is not countable is said to be *uncountable*. Since all finite sets are countable, an uncountable set is necessarily infinite.

In the previous section we saw that some rather nontrivial set theoretic operations preserve countability. We now provide examples of uncountable sets by a useful standard argument known as Cantor's diagonal extraction. The argument shows that countable Cartesian products of countable sets (or even of finite sets) are generally not countable.

**Example A.16** (The set of binary sequences is uncountable).

For each  $j \in \mathbb{N}$ , let  $A_j = \{0, 1\}$ . Consider the Cartesian product

$$B = \prod_{j=1}^{\infty} A_j = \{0, 1\}^{\mathbb{N}} = \left\{ (b_1, b_2, \dots) \mid b_1, b_2, \dots \in \{0, 1\} \right\}$$

of the sets  $A_1, A_2, \dots$ , which is most concretely interpreted as the set of infinite binary sequences  $b = (b_1, b_2, \dots)$  of zeroes and ones. The set  $B$  is a countable Cartesian product of finite sets. We claim that  $B$  itself is uncountable.

The diagonal argument proceeds by supposing, on the contrary, that  $B$  is countable. If this were the case, then we could find an enumeration

$$B = \{b^{(1)}, b^{(2)}, b^{(3)}, \dots\},$$

where the  $m$ :th element  $b^{(m)}$  is a binary sequence

$$b^{(m)} = (b_1^{(m)}, b_2^{(m)}, b_3^{(m)}, \dots).$$

Now define a binary sequence  $b' = (b'_1, b'_2, b'_3, \dots)$  by choosing for each  $j \in \mathbb{N}$  the  $j$ :th "digit"  $b'_j \in \{0, 1\}$  to be different from  $b_j^{(j)}$ , the  $j$ :th "digit" of the  $j$ :th element  $b^{(j)}$  in the enumeration. This construction of  $b' \in B$  is known as diagonal extraction. Now for any  $m \in \mathbb{N}$ , the binary sequence  $b'$  differs from  $b^{(m)}$  at least in the  $m$ :th digit, so  $b' \neq b^{(m)}$ . But the element  $b' \in B$  should appear in the enumeration of  $B$ , so we have derived a contradiction. We conclude that  $B$  can not be enumerated. Therefore  $B$  is in fact uncountable.

**Example A.17** (The set of real numbers is uncountable).

Let  $B = \{0, 1\}^{\mathbb{N}}$  be the set of binary sequences as in the previous example. Consider also the subset  $B' \subset B$  of those binary sequences which contain infinitely many zeroes,

$$B' = \left\{ (b_1, b_2, \dots) \in \{0, 1\}^{\mathbb{N}} \mid \forall m \in \mathbb{N} : \exists n \geq m : b_n = 0 \right\}.$$

We first claim that  $B'$  is also uncountable.

The complement  $B \setminus B'$  is the set of binary sequences which end with repeated ones,

$$B \setminus B' = \left\{ (b_1, b_2, \dots) \in \{0, 1\}^{\mathbb{N}} \mid \exists m \in \mathbb{N} : \forall n \geq m : b_n = 1 \right\}.$$

Let

$$R_m = \left\{ (b_1, b_2, \dots) \in \{0, 1\}^{\mathbb{N}} \mid \forall n \geq m : b_n = 1 \right\}.$$

denote the set of sequences where a repetition of ones has started by the  $m$ :th “digit”. Note that  $R_m$  is a finite set,  $\#R_m = 2^{m-1}$ , since we are only free to choose the values of the first  $m - 1$  “digits”. As a countable union of these finite sets, the complement

$$B \setminus B' = \bigcup_{m=1}^{\infty} R_m.$$

is countable by Proposition A.15. Now if  $B'$  would be countable, then the union  $B = B' \cup (B \setminus B')$  would also be countable, which is a contradiction with the conclusion of Example A.16. We thus conclude that  $B'$  is uncountable.

To prove that the set of real numbers is uncountable, we note that any real number has a binary expansion. More specifically, any number  $x \in [0, 1)$  has a binary expansion with its “digit sequence” in  $B'$ . Indeed, define a function  $f$  on  $[0, 1)$  by

$$f(x) = (b_1, b_2, \dots) \quad \text{where } b_j = \lfloor 2^j x \rfloor.$$

The sequence  $f(x) = (b_1, b_2, \dots)$  is a binary expansion of  $x$ ,

$$x = \sum_{j=1}^{\infty} b_j 2^{-j}.$$

It is easy to see that  $f(x) \in B'$  and that  $f: [0, 1) \rightarrow B'$  is surjective onto  $B'$  (for  $b \in B'$  and  $x = \sum_{j=1}^{\infty} b_j 2^{-j}$  we indeed have  $f(x) = b$ ). Therefore we conclude that the cardinality of  $B'$  is less than or equal to the cardinality of  $[0, 1)$ . But since  $B'$  is uncountable, also the set  $[0, 1)$  is uncountable.

Since  $[0, 1) \subset \mathbb{R}$  is a subset and  $[0, 1)$  is uncountable, also the set  $\mathbb{R}$  of real numbers is uncountable.

## Appendix B

### Topological preliminaries

This appendix reviews necessary background on topological notions, in particular properties of real numbers.

#### B.1. Topological properties of the real line

The set  $\mathbb{R}$  of real number inherits its topology from the natural notion of distance of points on the line: the distance  $\varrho(x, y)$  of two numbers  $x, y \in \mathbb{R}$  is the absolute value of their difference

$$\varrho(x, y) := |x - y|.$$

The distance function  $\varrho: \mathbb{R} \times \mathbb{R} \rightarrow [0, +\infty)$  is called the metric on  $\mathbb{R}$ , see Section B.2 for a summary of metric space topology more generally. Here we first remind the reader of some fundamental properties of the topology of the real line specifically.

#### Extended real line

Frequently during the present course it is convenient to extend the real line  $\mathbb{R}$  by two symbols,  $-\infty$  and  $+\infty$ , and consider the *extended real line*

$$\widehat{\mathbb{R}} = \mathbb{R} \cup \{-\infty\} \cup \{+\infty\}. \quad (\text{B.1})$$

We also alternatively denote the extended real line by  $\widehat{\mathbb{R}} = [-\infty, +\infty]$ , because it has the topology of a closed interval with  $-\infty$  and  $+\infty$  as its endpoints, as will be detailed later in Example B.19. Likewise, we denote

$$\begin{aligned} [-\infty, +\infty) &= \mathbb{R} \cup \{-\infty\} \\ (-\infty, +\infty] &= \mathbb{R} \cup \{+\infty\} \\ (-\infty, +\infty) &= \mathbb{R}. \end{aligned}$$

#### Supremum and infimum

One of the key defining properties of real numbers is the completeness property that non-empty bounded subsets have least upper bounds and greatest lower bounds. In the present text, supremum and infimum are used for these notions generalized to the setup of the extended real line, and to a setup where we do not even require non-emptiness and boundedness of the subset.

**Definition B.1** (Supremum and infimum).

Let  $A \subset \widehat{\mathbb{R}}$  be a subset of the extended real line.

The *supremum*,  $\sup(A)$ , is the least upper bound of  $A$ , i.e., the smallest  $M \in [-\infty, \infty]$  such that  $x \leq M$  for all  $x \in A$ .

The *infimum*,  $\inf(A)$ , is the greatest lower bound of  $A$ , i.e., the greatest  $M \in [-\infty, \infty]$  such that  $x \geq M$  for all  $x \in A$ .

For non-empty subsets  $\emptyset \neq A \subset \mathbb{R}$  of the real line, we have  $\sup(A) \in (-\infty, +\infty]$ , whereas for the empty set  $\emptyset$  we have  $\sup(\emptyset) = -\infty$ . Likewise, for non-empty subsets  $\emptyset \neq A \subset \mathbb{R}$  of the real line, we have  $\inf(A) \in [-\infty, +\infty)$ , whereas for the empty set  $\emptyset$  we have  $\inf(\emptyset) = +\infty$ .

For indexed collections  $(x_j)_{j \in J}$ , we denote the supremum and infimum also by

$$\sup_{j \in J} x_j := \sup \left( \{x_j \mid j \in J\} \right) \quad \text{and} \quad \inf_{j \in J} x_j := \inf \left( \{x_j \mid j \in J\} \right).$$

## Sequences of numbers

In probability theory, we frequently encounter sequences of real numbers. Crucial notions about sequences include in particular convergence (limits), monotonicity, and upper and lower limits (limsup and liminf).

### *Convergence of sequences of numbers*

The usual notion of limit of a sequence of real numbers is the following.

**Definition B.2** (Convergence of real number sequences).

Let  $x_1, x_2, x_3, \dots$  be a sequence of real numbers.

The sequence is said to *converge to* (or *tend to*) a *limit*  $x \in \mathbb{R}$  if for all  $\varepsilon > 0$  there exists  $n_0 = n_0(\varepsilon) \in \mathbb{N}$  such that for all  $n \geq n_0$  we have  $|x_n - x| < \varepsilon$ . We then denote

$$\lim_{n \rightarrow \infty} x_n = x \quad \text{or} \quad x_n \xrightarrow[n \rightarrow \infty]{} x.$$

In addition to the above usual notion of limit of the sequence inside  $\mathbb{R}$ , we consider also limits  $+\infty$  and  $-\infty$ .

**Definition B.3** (Convergence of real number sequences towards infinities).

Let  $x_1, x_2, x_3, \dots$  be a sequence of real numbers.

The sequence is said to *converge to* (or *tend to*)  $+\infty$  if for all  $M > 0$  there exists  $n_0 = n_0(M) \in \mathbb{N}$  such that for all  $n \geq n_0$  we have  $x_n > M$ . We then denote  $\lim_{n \rightarrow \infty} x_n = +\infty$  or  $x_n \rightarrow +\infty$ .

The sequence is said to *converge to* (or *tend to*)  $-\infty$  if for all  $M < 0$  there exists  $n_0 = n_0(M) \in \mathbb{N}$  such that for all  $n \geq n_0$  we have  $x_n < M$ . We then denote  $\lim_{n \rightarrow \infty} x_n = -\infty$  or  $x_n \rightarrow -\infty$ .

*Monotone sequences of numbers*

A sequence  $x_1, x_2, x_3, \dots$  of real numbers is said to be *increasing* if

$$x_1 \leq x_2 \leq x_3 \leq \dots .$$

If the sequence  $x_1, x_2, x_3, \dots$  is increasing, then it has a limit  $x \in (-\infty, +\infty]$ , and we denote

$$x_n \uparrow x.$$

It is easy to see that the limit of an increasing sequence is its supremum

$$x = \sup_{n \in \mathbb{N}} x_n.$$

Likewise, a sequence  $x_1, x_2, x_3, \dots$  of real numbers is said to be *decreasing* if

$$x_1 \geq x_2 \geq x_3 \geq \dots .$$

If the sequence  $x_1, x_2, x_3, \dots$  is decreasing, then it has a limit  $x \in [-\infty, +\infty)$ , and we denote

$$x_n \downarrow x.$$

It is easy to see that the limit of a decreasing sequence is its infimum

$$x = \inf_{n \in \mathbb{N}} x_n.$$

Sequences of numbers which are either increasing or decreasing are said to be *monotone*.

*Upper and lower limits of sequences*

For a sequence  $x_1, x_2, x_3, \dots$  of real numbers, we define the *upper limit*

$$\limsup_n x_n := \inf_{m \in \mathbb{N}} \left( \sup_{n \geq m} x_n \right)$$

Note that if we define  $c_m = \sup_{n \geq m} x_n$ , then the sequence  $c_1, c_2, \dots$  of numbers in  $(-\infty, +\infty]$  is decreasing, and its limit is  $\lim_{m \rightarrow \infty} c_m = \limsup_n x_n$ .

We also define the *lower limit*

$$\liminf_n x_n := \sup_{m \in \mathbb{N}} \left( \inf_{n \geq m} x_n \right)$$

Note that if we define  $d_m = \inf_{n \geq m} x_n$ , then the sequence  $d_1, d_2, \dots$  of numbers in  $[-\infty, +\infty)$  is increasing, and its limit is  $\lim_{m \rightarrow \infty} d_m = \liminf_n x_n$ .

**Proposition B.4** (Limit in terms of upper and lower limits).

*For any sequence  $x_1, x_2, x_3, \dots$  of real numbers we have*

$$\liminf_n x_n \leq \limsup_n x_n.$$

*The equality above holds if and only if the sequence is convergent, and in this case  $\liminf_n x_n = \limsup_n x_n = \lim_{n \rightarrow \infty} x_n$ .*

**Exercise B.1.** Prove Proposition B.4.

## Countability properties on the real line

We make use of the following facts on a few occasions.

**Exercise B.2** (Open subsets of the real line are countable unions of open intervals).

Show that any open set  $V \subset \mathbb{R}$  is the union of at most countably many open intervals.

**Hint:** Show that every point  $x \in V$  is contained in some interval  $(a, b) \subset V$  with rational endpoints  $a, b \in \mathbb{Q}$ .

**Proposition B.5** (Open sets as countable unions of disjoint open intervals).

Any open set  $V \subset \mathbb{R}$  is the union of countably many disjoint open intervals.

**Exercise B.3.** Give a proof of the above proposition, using the following steps.

(a) Show that any open set  $V \subset \mathbb{R}$  is the union of countably many open intervals.

**Hint:** Show that every point  $x \in V$  is contained in some interval  $(a, b) \subset V$  with rational endpoints  $a, b \in \mathbb{Q}$ .

(b) Show that any open set  $V \subset \mathbb{R}$  is the union of countably many disjoint open intervals.

**Hint:** Show that every point  $x \in V$  is contained in a unique maximal interval  $(a, b)$  within the set  $V$ . Use part (a) to show that there are at most countably many different such maximal intervals.

**Proposition B.6** (Monotone functions have countably many discontinuities).

A monotone function  $f: \mathbb{R} \rightarrow \mathbb{R}$  has countably many points of discontinuity.

**Exercise B.4.** Prove the above proposition.

**Hint:** Consider  $f$  restricted to an interval  $[k, k + 1]$ . For a given  $m \in \mathbb{N}$ , how many jumps of size at least  $\frac{1}{m}$  can  $f$  have on such an interval?

## B.2. Metric space topology

### Basic concepts of metric space topology

Recall that a metric space is a set  $\mathfrak{X}$  equipped with a metric, i.e., a function  $\varrho: \mathfrak{X} \times \mathfrak{X} \rightarrow [0, \infty)$  such that

$$\varrho(x, y) = 0 \iff x = y \quad (\varrho\text{-Sep})$$

$$\varrho(x, y) = \varrho(y, x) \quad \forall x, y \in \mathfrak{X} \quad (\varrho\text{-Sym})$$

$$\varrho(x, y) \leq \varrho(x, z) + \varrho(z, y) \quad \forall x, y, z \in \mathfrak{X}. \quad (\varrho\text{-Tri})$$

**Example B.7** (Real line as a metric space).

The set of real numbers  $\mathbb{R}$  equipped with the usual metric  $\varrho(x, y) = |x - y|$  is a metric space.

**Example B.8** (Euclidean space as a metric space).

Consider the  $d$ -dimensional real vector space  $\mathbb{R}^d$ . The *Euclidean norm* of a vector  $x = (x_1, \dots, x_d) \in \mathbb{R}^d$  is

$$\|x\| = \sqrt{x_1^2 + \dots + x_d^2}.$$

The space  $\mathbb{R}^d$  equipped with the metric  $\varrho(x, y) = \|x - y\|$  is a metric space.

**Example B.9** (Discrete metric).

Let  $S$  be any set. If we think of the points of  $S$  as discrete, then we may define the *discrete metric* on  $S$  by

$$\varrho_{\text{dis}}(x, y) = \begin{cases} 1 & \text{if } x \neq y \\ 0 & \text{if } x = y. \end{cases}$$

The set  $S$  equipped with the metric  $\varrho_{\text{dis}}$  is a metric space.

**Example B.10** (Uniform norm).

Consider the space of continuous real-valued functions on a closed interval  $[a, b] \subset \mathbb{R}$

$$\mathcal{C}([a, b]) := \{f: [a, b] \rightarrow \mathbb{R} \text{ continuous}\}.$$

The *supremum norm* (or *uniform norm*) of a function  $f \in \mathcal{C}([a, b])$  is

$$\|f\|_{\infty} = \sup_{x \in [a, b]} |f(x)|.$$

The space  $\mathcal{C}([a, b])$  equipped with the metric  $\varrho(f, g) = \|f - g\|_{\infty}$  is a metric space.

**Exercise B.5.** Verify in each of the examples above that the given metric indeed satisfies the axioms ( $\varrho$ -Sep), ( $\varrho$ -Sym), and ( $\varrho$ -Tri).

We will use the following topological notions in metric spaces.

**Definition B.11** (Open balls).

Let  $x \in \mathfrak{X}$  and  $r > 0$ . The (open) *ball* of radius  $r$  centered at  $x$  is the subset

$$\mathcal{B}_r(x) = \{y \in \mathfrak{X} \mid \varrho(x, y) < r\}.$$

**Definition B.12** (Open sets).

A subset  $A \subset \mathfrak{X}$  is *open*, if for all its points some ball centered at that point is contained in the set  $A$  (i.e.,  $\forall x \in A \exists r > 0 : \mathcal{B}_r(x) \subset A$ ).

**Exercise B.6** (Open balls are open sets).

Prove that any open ball  $\mathcal{B}_r(x) \subset \mathfrak{X}$  is an open set.

**Definition B.13** (Closed sets).

A subset  $A \subset \mathfrak{X}$  is *closed*, if its complement  $\mathfrak{X} \setminus A$  is open.

Note that for example the empty set  $\emptyset \subset \mathfrak{X}$  and the whole space  $\mathfrak{X}$  are both open and closed. There are also sets which are neither open nor closed.

**Definition B.14** (Limit of a sequence).

A sequence  $(x_n)_{n \in \mathbb{N}}$  of points  $x_n \in \mathfrak{X}$  *converges* to  $x \in \mathfrak{X}$  if  $\varrho(x_n, x) \rightarrow 0$  as  $n \rightarrow \infty$ . We then call  $x$  the *limit* of the sequence and denote

$$\lim_{n \rightarrow \infty} x_n = x \quad \text{or} \quad x_n \xrightarrow[n \rightarrow \infty]{} x.$$

**Exercise B.7** (The limit of a real number sequence in terms of metric).

Verify that the usual notion of a limit of a sequence of real numbers given in Definition B.2 coincides with Definition B.14 in the special case when the metric space is  $\mathbb{R}$  (Example B.7).

**Definition B.15** (Continuous functions between metric spaces).

If  $(\mathfrak{X}^{(1)}, \varrho^{(1)})$  and  $(\mathfrak{X}^{(2)}, \varrho^{(2)})$  are two metric spaces, then a function  $f: \mathfrak{X}^{(1)} \rightarrow \mathfrak{X}^{(2)}$  is *continuous* if for any convergent sequence  $(x_n)_{n \in \mathbb{N}}$  of points  $x_n \in \mathfrak{X}^{(1)}$  the sequence  $(f(x_n))_{n \in \mathbb{N}}$  converges in  $\mathfrak{X}^{(2)}$  and

$$\lim_{n \rightarrow \infty} f(x_n) = f\left(\lim_{n \rightarrow \infty} x_n\right).$$

The notions of convergence of sequences and continuity of functions can be formulated in purely topological terms, without reference to metric, only using the notion of open sets.

**Proposition B.16** (Topological characterization of continuity).

A function  $f: \mathfrak{X}^{(1)} \rightarrow \mathfrak{X}^{(2)}$  between two metric spaces  $(\mathfrak{X}^{(1)}, \varrho^{(1)})$  and  $(\mathfrak{X}^{(2)}, \varrho^{(2)})$  is continuous if and only if for every open set  $V$  in  $\mathfrak{X}^{(2)}$ , the preimage  $f^{-1}(V) = \{x \in \mathfrak{X}^{(1)} \mid f(x) \in V\}$  is open in  $\mathfrak{X}^{(1)}$ .

**Proposition B.17** (Topological characterization of limits).

A sequence  $(x_n)_{n \in \mathbb{N}}$  of points in a metric space  $(\mathfrak{X}, \varrho)$  converges to  $x \in \mathfrak{X}$  if and only if for every open set  $U \subset \mathfrak{X}$  containing the point  $x$ , there exists  $n_0 = n_0(U)$  such that for all  $n \geq n_0$  we have  $x_n \in U$ .

**Exercise B.8.** Verify that the convergence of a sequence and continuity of a function can be equivalently defined in terms of open sets as stated in Propositions B.17 and B.16 above.

**Definition B.18** (Homeomorphism).

If  $(\mathfrak{X}^{(1)}, \varrho^{(1)})$  and  $(\mathfrak{X}^{(2)}, \varrho^{(2)})$  are two metric spaces, then a function  $f: \mathfrak{X}^{(1)} \rightarrow \mathfrak{X}^{(2)}$  is a *homeomorphism* if  $f$  is bijective and both  $f: \mathfrak{X}^{(1)} \rightarrow \mathfrak{X}^{(2)}$  and its inverse  $f^{-1}: \mathfrak{X}^{(2)} \rightarrow \mathfrak{X}^{(1)}$  are continuous. Two spaces are *homeomorphic* if there exists a homeomorphism between them.

According to Proposition B.16, for a bijective function  $f$  to be a homeomorphism, a characterizing property is that a subset  $U \subset \mathfrak{X}^{(1)}$  is open if and only if its image  $f(U) \subset \mathfrak{X}^{(2)}$  is open. Since all topological properties can be formulated using the notion of open sets, homeomorphisms are exactly the mappings which preserve all topological properties.

**Example B.19** (The topology of the extended real line).

Consider the function  $h: [-\frac{\pi}{2}, +\frac{\pi}{2}] \rightarrow [-\infty, +\infty]$  defined by

$$h(s) = \begin{cases} -\infty & \text{if } s = -\frac{\pi}{2} \\ \tan(s) & \text{if } -\frac{\pi}{2} < s < +\frac{\pi}{2} \\ +\infty & \text{if } s = +\frac{\pi}{2}. \end{cases}$$

It is easy to see that  $h$  is bijective.

The restriction of the function  $h$  to the open interval  $(-\frac{\pi}{2}, +\frac{\pi}{2})$  is continuous (it is the trigonometric function  $\tan: (-\frac{\pi}{2}, +\frac{\pi}{2}) \rightarrow \mathbb{R}$ ) and has continuous inverse (the inverse is  $\arctan: \mathbb{R} \rightarrow (-\frac{\pi}{2}, +\frac{\pi}{2})$ ). It is therefore provides a homeomorphism from the open interval  $(-\frac{\pi}{2}, +\frac{\pi}{2})$  to the real line  $\mathbb{R}$ .

We can define a topology on the extended real line  $\widehat{\mathbb{R}} = [-\infty, +\infty]$  by requiring that  $h: [-\frac{\pi}{2}, +\frac{\pi}{2}] \rightarrow \widehat{\mathbb{R}}$  is a homeomorphism. With this definition, in particular, a subset  $A \subset \widehat{\mathbb{R}}$



is open if and only if  $h^{-1}(A) \subset [-\frac{\pi}{2}, +\frac{\pi}{2}]$  is open. The topology of  $\widehat{\mathbb{R}}$  can be obtained for example using the metric  $\varrho(x, y) = |h^{-1}(x) - h^{-1}(y)|$ , but since this choice of metric is not canonical and in particular does not agree with the usual metric on the subset  $\mathbb{R} \subset \widehat{\mathbb{R}}$ , we usually prefer not to use an explicitly chosen metric on  $\widehat{\mathbb{R}}$ .

**Exercise B.9** (Convergence on the extended real line).

Verify that our definition of convergence of a real sequence  $(x_n)_{n \in \mathbb{N}}$  to  $+\infty$  (respectively to  $-\infty$ ) in Definition B.3 is equivalent to the convergence of that sequence in the topological space  $\widehat{\mathbb{R}}$  to the point  $+\infty \in \widehat{\mathbb{R}}$  (resp. to  $-\infty \in \widehat{\mathbb{R}}$ ).



## Dynkin's identification and monotone class theorem

In this appendix we give the proof of Dynkin's identification theorem (Theorem II.26) and we state and prove a related result, the Monotone class theorem (Theorem C.2 below), which was used in Lectures IV and IX.

It is possible to study Sections C.2 and C.3 immediately after Lecture II, where Dynkin's identification theorem was stated.

Another option is to study this entire appendix after Lecture IV, where the Monotone class theorem is first used. This latter approach may be more convenient, since the techniques of proofs of both results are very closely related.

### C.1. Monotone class theorem

**Definition C.1** (Monotone class).

A collection  $\mathcal{H}$  of bounded functions from  $S$  to  $\mathbb{R}$  is said to be a *monotone class* if it satisfies the following conditions:

- (MC-1) The constant function 1 belongs to  $\mathcal{H}$ .
- (MC- $\mathbb{R}$ ) The class  $\mathcal{H}$  is a vector space<sup>1</sup> over  $\mathbb{R}$ .
- (MC- $\uparrow$ ) If  $f_1, f_2, \dots \in \mathcal{H}$  is an increasing sequence of non-negative functions in  $\mathcal{H}$  such that the pointwise limit  $f = \lim_{n \rightarrow \infty} f_n$  is a bounded function, then  $f \in \mathcal{H}$ .

**Theorem C.2** (Monotone class theorem).

Suppose that  $\mathcal{H}$  is a monotone class of bounded functions from  $S$  to  $\mathbb{R}$ . Let  $\mathcal{J}$  be a  $\pi$ -system on  $S$ . Then if  $\mathcal{H}$  contains the indicator function  $\mathbb{I}_A$  of every set  $A \in \mathcal{J}$  in the  $\pi$ -system, then  $\mathcal{H}$  contains all bounded  $\sigma(\mathcal{J})/\mathcal{B}$ -measurable functions.

### C.2. Auxiliary results

In the proof of both Dynkin's identification theorem and Monotone class theorem, we use the following definitions and auxiliary observations.

**Definition C.3** (D-system).

A collection  $\mathcal{D}$  of subsets of  $S$  is said to be a *d-system* on  $S$  if it satisfies the following conditions:

---

<sup>1</sup>That  $\mathcal{H}$  is a vector space means that it is stable under taking linear combinations of functions: if we have  $f_1, f_2 \in \mathcal{H}$  and  $c_1, c_2 \in \mathbb{R}$ , then we also have  $c_1 f_1 + c_2 f_2 \in \mathcal{H}$ .

- (D-1)  $S \in \mathcal{D}$   
 (D-d) if  $A, B \in \mathcal{D}$  and  $A \subset B$ , then  $B \setminus A \in \mathcal{D}$   
 (D- $\uparrow$ ) if  $A_1, A_2, \dots \in \mathcal{D}$  is an increasing sequence of subsets and  $A = \bigcup_{n \in \mathbb{N}} A_n$  is its limit, then  $A \in \mathcal{D}$ .

**Proposition C.4** (A characterization of sigma algebras).

*A collection  $\mathcal{C}$  of subsets of  $S$  is a  $\sigma$ -algebra if and only if it is both a d-system and a  $\pi$ -system.*

*Proof.* The “if” direction is obvious: a  $\sigma$ -algebra clearly satisfies both the conditions of a d-system and of a  $\pi$ -system. It thus suffices to prove the “only if” part.

Suppose that  $\mathcal{C}$  is both a  $\pi$ -system and a d-system. We must show that  $\mathcal{C}$  is a  $\sigma$ -algebra.

By property (D-1) we have  $S \in \mathcal{C}$ , so  $\mathcal{C}$  satisfies condition ( $\Sigma$ -1) of  $\sigma$ -algebras. If  $A \in \mathcal{C}$ , then we have  $A^c = S \setminus A \in \mathcal{C}$  by properties (D-1) and (D-d), so  $\mathcal{C}$  satisfies condition ( $\Sigma$ -c) of  $\sigma$ -algebras. It remains to show condition ( $\Sigma$ - $\cup$ ), i.e., that  $\mathcal{C}$  is stable under countable unions.

Consider first the union of just two sets  $A_1, A_2 \in \mathcal{C}$  in the collection. We have  $A_1^c, A_2^c \in \mathcal{C}$  by d-system properties and then  $A_1^c \cap A_2^c \in \mathcal{C}$  by  $\pi$ -system properties. Using De Morgan's law we then observe that

$$A_1 \cup A_2 = S \setminus (A_1^c \cap A_2^c) \in \mathcal{C},$$

again by d-system properties. From this, by induction one gets that  $\mathcal{C}$  is stable under finite unions.

Now consider a countable sequence of sets  $A_1, A_2, \dots \in \mathcal{C}$ . Denote  $G_n = A_1 \cup \dots \cup A_n$ . The induction above allows to conclude that  $G_n \in \mathcal{C}$ . Then  $G_1 \subset G_2 \subset \dots$  is an increasing sequence of subsets belonging to the collection  $\mathcal{C}$ , so by d-system property (D- $\uparrow$ ) we get that also the limit  $G = \bigcup_{n \in \mathbb{N}} G_n$  belongs to the collection  $\mathcal{C}$ . But by construction  $\bigcup_{n \in \mathbb{N}} G_n = \bigcup_{n \in \mathbb{N}} A_n$ , so we conclude that  $\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{C}$ . This shows that  $\mathcal{C}$  also satisfies condition ( $\Sigma$ - $\cup$ ) of  $\sigma$ -algebras, and finishes the proof.  $\square$

**Definition C.5** (D-system generated by a collection of subsets).

The d-system generated by a collection  $\mathcal{I}$  of subsets of  $S$  is the smallest d-system  $\mathcal{D}$  which contains the collection  $\mathcal{I}$ . We denote the d-system generated by the collection  $\mathcal{I}$  by  $d(\mathcal{I})$ .

**Remark C.6** (Well-definedness of generated d-systems).

The definition makes sense again essentially because the intersection of d-systems is a d-system. The smallest d-system with the property that they contain the collection  $\mathcal{I}$  is the intersection of all such d-systems (the intersection is over a non-empty collection since at the very least the power set  $\mathcal{P}(S)$  is such a d-system).

**Lemma C.7** (Dynkin's lemma).

*Suppose that  $\mathcal{I}$  is a  $\pi$ -system. Then the d-system  $d(\mathcal{I})$  generated by  $\mathcal{I}$  and the  $\sigma$ -algebra  $\sigma(\mathcal{I})$  generated by  $\mathcal{I}$  coincide,  $d(\mathcal{I}) = \sigma(\mathcal{I})$ .*

*Proof.* Since any sigma-algebra is a d-system, we clearly have  $d(\mathcal{I}) \subset \sigma(\mathcal{I})$ . By Proposition C.4 it thus suffices to show that  $d(\mathcal{I})$  is also a  $\pi$ -system. We show this in two steps.

In the first step our goal is to show that whenever  $B \in d(\mathcal{I})$  and  $C \in \mathcal{I}$ , we have  $B \cap C \in d(\mathcal{I})$ . Define therefore the collection

$$\mathcal{D}_1 = \left\{ B \in d(\mathcal{I}) \mid B \cap C \in d(\mathcal{I}) \text{ for all } C \in \mathcal{I} \right\}$$

of sets  $B$  with this property. Rephrasing our goal, we wish to show that this collection is simply  $\mathcal{D}_1 = d(\mathcal{J})$ . By construction we have  $\mathcal{D}_1 \subset d(\mathcal{J})$ . Moreover, since  $\mathcal{J}$  is a  $\pi$ -system, we have  $\mathcal{J} \subset \mathcal{D}_1$ . Since  $d(\mathcal{J})$  is the smallest d-system containing  $\mathcal{J}$ , showing that  $\mathcal{D}_1$  is a d-system will thus achieve our goal:  $\mathcal{D}_1 = d(\mathcal{J})$ .

For any  $C \in \mathcal{J}$  we have  $S \cap C = C \in \mathcal{J} \subset d(\mathcal{J})$ , so we get that  $S \in \mathcal{D}_1$ , i.e., property (D-1) holds for  $\mathcal{D}_1$ . If  $A, B \in \mathcal{D}_1$  and  $A \subset B$  and  $C \in \mathcal{J}$ , then we have

$$(B \setminus A) \cap C = \underbrace{(B \cap C)}_{\in d(\mathcal{J})} \setminus \underbrace{(A \cap C)}_{\in d(\mathcal{J})} \in d(\mathcal{J})$$

by property (D-d) of d-system  $d(\mathcal{J})$ . We get  $B \setminus A \in \mathcal{D}_1$ , and we thus see that property (D-d) holds for  $\mathcal{D}_1$ . If  $A_1, A_2, \dots \in \mathcal{D}_1$  is an increasing sequence of subsets and  $C \in \mathcal{J}$ , then we have

$$\left( \bigcup_{n \in \mathbb{N}} A_n \right) \cap C = \bigcup_{n \in \mathbb{N}} \underbrace{(A_n \cap C)}_{\in d(\mathcal{J})} \in d(\mathcal{J})$$

by property (D- $\uparrow$ ) of d-system  $d(\mathcal{J})$ . We get  $\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{D}_1$ , and we thus see that property (D- $\uparrow$ ) holds for  $\mathcal{D}_1$ . We conclude that  $\mathcal{D}_1$  is a d-system, and therefore that  $\mathcal{D}_1 = d(\mathcal{J})$ , which was our first goal.

In the second step our goal is to show that whenever  $A, B \in d(\mathcal{J})$ , we have  $A \cap B \in d(\mathcal{J})$ . Define therefore the collection

$$\mathcal{D}_2 = \left\{ A \in d(\mathcal{J}) \mid A \cap B \in d(\mathcal{J}) \text{ for all } B \in d(\mathcal{J}) \right\}$$

of sets  $A$  with this property. Rephrasing our second goal, we wish to show that this collection is simply  $\mathcal{D}_2 = d(\mathcal{J})$ . From the first step we got that  $\mathcal{J} \subset \mathcal{D}_2$  and by construction we have  $\mathcal{D}_2 \subset d(\mathcal{J})$ . Since  $d(\mathcal{J})$  is the smallest d-system containing  $\mathcal{J}$ , showing that  $\mathcal{D}_2$  is a d-system will thus achieve our second goal:  $\mathcal{D}_2 = d(\mathcal{J})$ . The proof that  $\mathcal{D}_2$  is a d-system is exactly parallel to the similar argument in the first step.

The conclusion of the second step precisely says that  $d(\mathcal{J})$  is a  $\pi$ -system, and the proof is thus complete.  $\square$

### C.3. Proof of Dynkin's identification theorem

Recall the statement of the nontrivial direction of Dynkin's identification theorem (Theorem II.26): we assume that

- $P_1$  and  $P_2$  are two probability measures on a measurable space  $(\Omega, \mathcal{F})$
- $\mathcal{J}$  is a  $\pi$ -system on  $\Omega$  such that  $\sigma(\mathcal{J}) = \mathcal{F}$
- for all  $E \in \mathcal{J}$  we have  $P_1[E] = P_2[E]$ .

Then the claim is that the two probability measures are in fact identical,

$$P_1 = P_2,$$

i.e., that the equalities  $P_1[E] = P_2[E]$  hold for all events  $E \in \mathcal{F}$ .

*Proof of Dynkin's identification theorem (Theorem II.26).* Let

$$\mathcal{D} = \left\{ E \in \mathcal{F} \mid P_1[E] = P_2[E] \right\}$$

be the collection of those  $E$  for which the desired equality  $P_1[E] = P_2[E]$  holds. To show that  $P_1 = P_2$  we must show that this collection contains all events,  $\mathcal{D} = \mathcal{F}$ .

We claim that  $\mathcal{D}$  is a d-system on  $\Omega$ . The defining properties are checked as follows:

- We have  $P_1[\Omega] = 1 = P_2[\Omega]$  by definition of probability measures. Therefore we see that  $\Omega \in \mathcal{D}$ , which shows property (D-1) for  $\mathcal{D}$ .

- Suppose that  $A, B \in \mathcal{D}$  and  $A \subset B$ . Write  $B = A \cup (B \setminus A)$ , which is a disjoint union, so by finite additivity of probabilities for disjoint sets (Lemma II.19) we get

$$P_1[B] = P_1[A] + P_1[B \setminus A]$$

and similarly for  $P_2$ . We solve this for the probability of  $B \setminus A$  and get

$$\begin{aligned} P_1[B \setminus A] &= P_1[B] - P_1[A] \\ &= P_2[B] - P_2[A] \quad (\text{because } A, B \in \mathcal{D}) \\ &= P_2[B \setminus A]. \end{aligned}$$

This shows that  $B \setminus A \in \mathcal{D}$ , and establishes property (D-d) for  $\mathcal{D}$ .

- Suppose that  $A_1, A_2, \dots \in \mathcal{D}$  and  $A_n \uparrow A$ . By monotone convergence of probabilities (Theorem II.22) we get  $P_1[A_n] \uparrow P_1[A]$  and similarly for  $P_2$ . Therefore we have

$$\begin{aligned} P_1[A] &= \lim_{n \rightarrow \infty} P_1[A_n] \\ &= \lim_{n \rightarrow \infty} P_2[A_n] \quad (\text{because } A_n \in \mathcal{D}) \\ &= P_2[A]. \end{aligned}$$

This shows that  $A \in \mathcal{D}$ , and establishes property (D- $\uparrow$ ) for  $\mathcal{D}$ .

We have shown that  $\mathcal{D}$  is a d-system. By assumption,  $\mathcal{I}$  is contained in it,  $\mathcal{I} \subset \mathcal{D}$ , so also the d-system generated by  $\mathcal{I}$  must be contained in it,  $d(\mathcal{I}) \subset \mathcal{D}$ . Since  $\mathcal{I}$  is a  $\pi$ -system, we get from Dynkin's lemma (Lemma C.7) that  $d(\mathcal{I}) = \sigma(\mathcal{I})$ , and by assumption we have  $\sigma(\mathcal{I}) = \mathcal{F}$ . Therefore we conclude  $\mathcal{F} \subset \mathcal{D}$ . By definition of  $\mathcal{D}$  this means that  $P_1[E] = P_2[E]$  hold for all events  $E \in \mathcal{F}$ .  $\square$

#### C.4. Proof of Monotone class theorem

With the preparations in Section C.2 we can prove also the Monotone class theorem. It is worth observing that the proof steps are exactly those of the “standard machine” of integration theory (see Lecture VII).

*Proof of Theorem C.2 (Monotone class theorem).* Let  $\mathcal{H}$  be a monotone class of bounded functions from  $S$  to  $\mathbb{R}$ . Define  $\mathcal{D}$  as the collection of subsets  $A \subset S$  whose indicator belongs to the monotone class,  $\mathbb{I}_A \in \mathcal{H}$ . Properties (MC-1), (MC- $\mathbb{R}$ ), and (MC- $\uparrow$ ) of  $\mathcal{H}$  imply that  $\mathcal{D}$  has properties (D-1), (D-d), and (D- $\uparrow$ ), i.e.,  $\mathcal{D}$  is a d-system.

Now assume, as in the statement, that  $\mathcal{I}$  is a  $\pi$ -system such that  $\mathcal{H}$  contains the indicator function of each member of  $\mathcal{I}$ . Then we have  $\mathcal{I} \subset \mathcal{D}$ . Since  $\mathcal{D}$  is a d-system, also  $d(\mathcal{I}) \subset \mathcal{D}$ . Furthermore  $d(\mathcal{I}) = \sigma(\mathcal{I})$  by Dynkin's lemma, so we actually have  $\sigma(\mathcal{I}) \subset \mathcal{D}$ . In other words, all indicator functions of sets in  $\sigma(\mathcal{I})$  are in the monotone class  $\mathcal{H}$ .

Any bounded simple  $\sigma(\mathcal{I})/\mathcal{B}$ -measurable function is a finite linear combination of indicator functions of sets in  $\sigma(\mathcal{I})$ , so by the above observation  $\sigma(\mathcal{I}) \subset \mathcal{D}$  and the vector space property (MC- $\mathbb{R}$ ) we have that the monotone class  $\mathcal{H}$  contains such bounded simple functions.

If  $f: S \rightarrow \mathbb{R}$  is a non-negative bounded  $\sigma(\mathcal{I})/\mathcal{B}$ -measurable function, by the approximation lemma (Lemma III.18) we may find a sequence  $f_1, f_2, \dots$  of non-negative simple  $\sigma(\mathcal{I})/\mathcal{B}$ -measurable functions such that  $f_n \uparrow f$ . We observed that the simple functions are in the monotone class,  $f_n \in \mathcal{H}$ . Therefore, by property (MC- $\uparrow$ ), their limit is as well,  $f \in \mathcal{H}$ .

We have shown that  $\mathcal{H}$  contains all non-negative bounded  $\sigma(\mathcal{I})/\mathcal{B}$ -measurable functions. For a general bounded  $\sigma(\mathcal{I})/\mathcal{B}$ -measurable function  $f$ , write  $f = f_+ - f_-$  where  $f_+ = \max(f, 0)$  and  $f_- = \max(-f, 0)$  are non-negative, bounded and  $\sigma(\mathcal{I})/\mathcal{B}$ -measurable. As this linear combination,  $f$  itself belongs to the monotone class,  $f \in \mathcal{H}$ . This finishes the proof.  $\square$

## Monotone convergence theorem

This appendix is devoted to the proof of the Monotone convergence theorem from Lecture VII. Let us recall its statement.

**Theorem** (Monotone convergence theorem, Theorem VII.8).

*If  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}^+$  and  $f_n \uparrow f$  as  $n \rightarrow \infty$ , then we have*

$$\int^+ f_n \, d\mu \uparrow \int^+ f \, d\mu \quad \text{as } n \rightarrow \infty.$$

We will prove this in a number of steps, gradually increasing the generality of the non-negative limit function  $f$  as well as the non-negative approximating functions  $f_n$ .

### D.1. Monotone convergence theorem for simple functions

First, notice that the Monotone convergence theorem for integrals is certainly closely related to the following monotone convergence of measures, which we proved in Lecture II.

**Proposition** (Part (II.9) of Lemma II.19).

*If  $A_1, A_2, \dots \in \mathcal{S}$  and  $A_n \uparrow A$  as  $n \rightarrow \infty$ , then we have  $\mu[A_n] \uparrow \mu[A]$ .*

With this initial observation, our first step is the following monotone convergence result for simple functions approximating an indicator function.

**Lemma D.1** (Monotone convergence for simple approximations of an indicator).

*If we have  $A \in \mathcal{S}$ , and if  $h_1, h_2, \dots \in \mathfrak{s}\mathcal{S}^+$  are such that  $h_n \uparrow \mathbb{I}_A$  as  $n \rightarrow \infty$ , then we have*

$$\int^{\square} h_n \, d\mu \uparrow \mu[A] \quad \text{as } n \rightarrow \infty.$$

*Proof.* Since we have  $h_n \leq \mathbb{I}_A$ , the monotonicity of the integral  $\int^{\square}$  (Lemma VII.3) yields

$$\int^{\square} h_n \, d\mu \leq \int^{\square} \mathbb{I}_A \, d\mu = \mu[A],$$

so it suffices to prove that

$$\liminf_n \int^{\square} h_n \, d\mu \geq \mu[A].$$

Let  $0 < \varepsilon < 1$ . Define  $A_n = \{s \in S \mid h_n(s) > 1 - \varepsilon\}$ . Then by the assumption  $h_n \uparrow \mathbb{I}_A$ , we have  $A_n \uparrow A$ . Monotone convergence of measures (Lemma II.19) thus gives  $\mu[A_n] \uparrow \mu[A]$ .

Note also that  $h_n \geq (1-\varepsilon)\mathbb{I}_{A_n}$  by construction of  $A_n$ , so the monotonicity of the integral  $\int^\square$  gives

$$\int^\square h_n \, d\mu \geq (1-\varepsilon)\mu[A_n].$$

Taking the lower limit as  $n \rightarrow \infty$  of this inequality and recalling  $\mu[A_n] \uparrow \mu[A]$ , we get

$$\liminf_n \int^\square h_n \, d\mu \geq (1-\varepsilon)\mu[A].$$

Since  $\varepsilon > 0$  can be taken arbitrarily small, we conclude that

$$\liminf_n \int^\square h_n \, d\mu \geq \mu[A],$$

which finishes the proof.  $\square$

Next we prove the monotone convergence theorem for simple functions.

**Lemma D.2** (Monotone convergence theorem for simple functions).

If  $h_1, h_2, \dots \in s\mathcal{S}^+$  and  $h_n \uparrow h \in s\mathcal{S}^+$  as  $n \rightarrow \infty$ , then we have

$$\int^\square h_n \, d\mu \uparrow \int^\square h \, d\mu \quad \text{as } n \rightarrow \infty.$$

*Proof.* Write  $h = \sum_{k=1}^m a_k \mathbb{I}_{A_k}$ , with  $a_1, \dots, a_m > 0$  and  $A_1, \dots, A_m \in \mathcal{S}$  disjoint. Then we have that  $\frac{1}{a_k} \mathbb{I}_{A_k} h_n \uparrow \mathbb{I}_{A_k}$  as  $n \rightarrow \infty$  by assumption  $h_n \uparrow h$ . Now the assertion follows from Lemma D.1 and linearity of the integral.  $\square$

## D.2. Monotone convergence theorem for general non-negative functions

Above we established the Monotone convergence theorem for simple functions. The next steps first relax the assumption that the limit function is simple, and then relax the assumption that the approximating functions are simple.

Let us start by verifying that for any non-negative measurable function there exists at least some sequence of simple functions for which the conclusion of the Monotone convergence theorem holds.

**Lemma D.3** (Monotone convergence for some simple approximating sequence).

For any  $f \in m\mathcal{S}^+$  there exists a sequence  $g_1, g_2, \dots \in s\mathcal{S}^+$  such that as  $n \rightarrow \infty$ , we have

$$g_n \uparrow f \quad \text{and} \quad \int^\square g_n \, d\mu \uparrow \int^+ f \, d\mu.$$

*Proof.* By Definition VII.4 we have

$$\int^+ f \, d\mu := \sup_{\substack{h \in s\mathcal{S}^+ \\ 0 \leq h \leq f}} \int^\square h \, d\mu,$$

so there exists some sequence  $h_1, h_2, \dots \in s\mathcal{S}^+$  such that the integral of  $f$  is approximated as

$$\int^\square h_n \, d\mu \uparrow \int^+ f \, d\mu \quad \text{as } n \rightarrow \infty. \quad (\text{D.1})$$



In addition, by the approximation lemma (Lemma III.18), there exists a sequence  $f_1, f_2, \dots \in \mathcal{S}^+$  such that the function  $f$  is approximated as

$$f_n \uparrow f \quad \text{as } n \rightarrow \infty. \quad (\text{D.2})$$

Define a new sequence  $g_1, g_2, \dots \in \mathcal{S}^+$  in terms of the two sequences above as

$$g_n := \max \{f_n, h_1, h_2, \dots, h_n\}. \quad (\text{D.3})$$

By construction this sequence is pointwise increasing,  $g_1 \leq g_2 \leq \dots$ . The constructed functions are also simple,  $g_n \in \mathcal{S}^+$  (the possible values of the maximum of finitely many simple functions are the finitely many values of these simple functions together).

Moreover, we clearly have  $f_n \leq g_n \leq f$ , so from (D.2) we also get

$$g_n \uparrow f \quad \text{as } n \rightarrow \infty.$$

Likewise, we clearly have  $h_n \leq g_n \leq f$ , so by monotonicity of the integral  $\int^\square$  and definition of the integral  $\int^+$  we get the inequalities

$$\int^\square h_n \, d\mu \leq \int^\square g_n \, d\mu \leq \int^+ f \, d\mu.$$

These inequalities, together with (D.1), give

$$\int^\square g_n \, d\mu \uparrow \int^+ f \, d\mu \quad \text{as } n \rightarrow \infty.$$

This finishes the proof.  $\square$

In the remaining steps of the proof of the Monotone convergence theorem, we use the following auxiliary result about monotone increasing arrays twice.

**Lemma D.4** (Monotone arrays).

Let  $(t_n^{(r)})_{n \in \mathbb{N}, r \in \mathbb{N}}$  be an array of numbers  $t_n^{(r)} \in [0, +\infty]$ , which is increasing in both indices:

$$t_1^{(r)} \leq t_2^{(r)} \leq t_3^{(r)} \leq \dots \quad \text{for all } r \in \mathbb{N} \quad (\text{D.4})$$

$$t_n^{(1)} \leq t_n^{(2)} \leq t_n^{(3)} \leq \dots \quad \text{for all } n \in \mathbb{N}. \quad (\text{D.5})$$

For any  $r \in \mathbb{N}$ , denote the limit of the increasing sequence (D.4) by

$$t^{(r)} := \lim_{n \rightarrow \infty} t_n^{(r)},$$

and for any  $n \in \mathbb{N}$ , denote the limit of the increasing sequence (D.5) by

$$t_n := \lim_{r \rightarrow \infty} t_n^{(r)}.$$

Then the sequences  $t^{(1)}, t^{(2)}, \dots$  and  $t_1, t_2, \dots$  are both increasing, and their limits

$$t^{(\infty)} := \lim_{r \rightarrow \infty} t^{(r)} \quad \text{and} \quad t_\infty := \lim_{n \rightarrow \infty} t_n$$

coincide,

$$t^{(\infty)} = t_\infty.$$

*Proof.* We may assume that the array of numbers is uniformly bounded in the sense that for some  $M < +\infty$  we have  $t_n^{(r)} \leq M$  for all  $n \in \mathbb{N}$ ,  $r \in \mathbb{N}$  (consider for example  $\arctan(t_n^{(r)})$  if the original table is not uniformly bounded).

Recall the definition  $t_n = \lim_{r \rightarrow \infty} t_n^{(r)}$ , for all  $n \in \mathbb{N}$ . Since for any  $n$  and  $r$  we have

$$t_n^{(r)} \leq t_{n+1}^{(r)} \leq M,$$

taking the limits as  $r \rightarrow \infty$  we get

$$t_n \leq t_{n+1} \leq M.$$

This shows that the sequence  $t_1, t_2, \dots$  is increasing and bounded by  $M$ . It therefore has a limit, which we denote by  $t_\infty = \lim_{n \rightarrow \infty} t_n$ .

Similarly one shows that the sequence  $t^{(1)}, t^{(2)}, \dots$  is increasing and bounded. Denote its limit by  $t^{(\infty)} = \lim_{r \rightarrow \infty} t^{(r)}$ .

Let  $\varepsilon > 0$ . Since  $t_n \uparrow t_\infty$  as  $n \rightarrow \infty$ , we can choose some  $n_0$  such that  $t_{n_0} > t_\infty - \frac{1}{2}\varepsilon$ . Then, since  $t_{n_0}^{(r)} \uparrow t_{n_0}$  as  $r \rightarrow \infty$ , we can choose some  $r_0$  such that  $t_{n_0}^{(r_0)} > t_{n_0} - \frac{1}{2}\varepsilon$ . We now see, because the sequences are increasing, that

$$t^{(\infty)} \geq t^{(r_0)} \geq t_{n_0}^{(r_0)} > t_{n_0} - \frac{1}{2}\varepsilon > t_\infty - \varepsilon.$$

Since  $\varepsilon > 0$  was arbitrary, this shows  $t^{(\infty)} \geq t_\infty$ . Completely symmetrically one obtains the opposite inequality  $t_\infty \geq t^{(\infty)}$ . We conclude the equality  $t_\infty = t^{(\infty)}$ , and the lemma is proven.  $\square$

We now improve the result of Lemma D.3, and show that the specific choice of the approximating sequence of simple functions did not matter.

**Lemma D.5** (Monotone convergence for simple approximating sequences).

*For any  $f \in \mathfrak{m}\mathcal{S}^+$  and any sequence  $h_1, h_2, \dots \in \mathfrak{s}\mathcal{S}^+$  such that  $h_n \uparrow f$  as  $n \rightarrow \infty$ , we have*

$$\int^\square h_n \, d\mu \uparrow \int^+ f \, d\mu.$$

*Proof.* Given  $f \in \mathfrak{m}\mathcal{S}^+$ , use Lemma D.3 to get a sequence  $g^{(1)}, g^{(2)}, \dots \in \mathfrak{s}\mathcal{S}^+$  such that

$$g^{(r)} \uparrow f \quad \text{and} \quad \int^\square g^{(r)} \, d\mu \uparrow \int^+ f \, d\mu \quad \text{as } r \rightarrow \infty.$$

Suppose now that  $h_1, h_2, \dots \in \mathfrak{s}\mathcal{S}^+$  is another sequence such that  $h_n \uparrow f$  as  $n \rightarrow \infty$ . Define, for all  $n \in \mathbb{N}$  and  $r \in \mathbb{N}$  the function

$$f_n^{(r)} := \min \{h_n, g^{(r)}\},$$

which is simple,  $f_n^{(r)} \in \mathfrak{s}\mathcal{S}^+$ . Observe that in this situation

$$\begin{aligned} f_n^{(r)} \uparrow h_n & \quad \text{as } r \rightarrow \infty \\ \text{and } f_n^{(r)} \uparrow g^{(r)} & \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Also denote the integrals of these functions by

$$t_n^{(r)} := \int^\square f_n^{(r)} \, d\mu.$$

Consider first a fixed  $r \in \mathbb{N}$ . Since  $f_n^{(r)} \uparrow g^{(r)}$  as  $n \rightarrow \infty$ , we can apply the already proven Monotone convergence theorem for simple functions, Lemma D.2, to get

$$t_n^{(r)} = \int^\square f_n^{(r)} \, d\mu \uparrow \int^\square g^{(r)} \, d\mu =: t^{(r)} \quad \text{as } n \rightarrow \infty.$$

Now letting  $r \rightarrow \infty$ , we have by the choice of the sequence  $g^{(1)}, g^{(2)}, \dots$  that

$$t^{(r)} = \int^\square g^{(r)} \, d\mu \uparrow \int^+ f \, d\mu =: t^{(\infty)} \quad \text{as } r \rightarrow \infty.$$

Consider then a fixed  $n \in \mathbb{N}$ . Since  $f_n^{(r)} \uparrow h_n$  as  $r \rightarrow \infty$ , we can again apply the already proven Monotone convergence theorem for simple functions, Lemma D.2, to get

$$t_n^{(r)} = \int^{\square} f_n^{(r)} \, d\mu \uparrow \int^{\square} h_n \, d\mu =: t_n \quad \text{as } r \rightarrow \infty.$$

The integrals  $t_n$  of the functions  $h_n$  form an increasing sequence (the sequence of functions is increasing and the integral is monotone), so they have a limit, which we denote by  $t_\infty$ :

$$t_n = \int^{\square} h_n \, d\mu \uparrow t_\infty \quad \text{as } n \rightarrow \infty.$$

It now follows from Lemma D.4 on monotone increasing arrays that  $t_\infty = t^{(\infty)}$ . Recalling the definition of  $t^{(\infty)}$ , we have obtained that as  $n \rightarrow \infty$ ,

$$\int^{\square} h_n \, d\mu = t_n \uparrow t_\infty = t^{(\infty)} = \int^+ f \, d\mu$$

and the proof is complete.  $\square$

Now, finally, we are ready to prove the Monotone convergence theorem in full generality.

*Proof of the Monotone convergence theorem (Theorem VII.8).* Let  $f_1, f_2, \dots \in \mathfrak{m}\mathcal{S}^+$  be a sequence of non-negative measurable functions, which is pointwise increasing:

$$0 \leq f_1(s) \leq f_2(s) \leq \dots \quad \text{for all } s \in S.$$

Let  $f: S \rightarrow [0, +\infty]$  be the pointwise limit of this increasing sequence,

$$f(s) := \lim_{n \rightarrow \infty} f_n(s) \quad \text{for } s \in S,$$

so that we have  $f_n \uparrow f$ . By Proposition III.14 this limit function is measurable,  $f \in \mathfrak{m}\mathcal{S}^+$ .

For each  $n \in \mathbb{N}$ , use the staircase function construction of Lemma III.18 to get an increasing approximation of  $f_n$  by simple functions: define  $f_n^{(r)} := \varsigma_r \circ f_n$  so that

$$f_n^{(r)} \uparrow f_n \quad \text{as } r \rightarrow \infty.$$

Note also that for a fixed  $r \in \mathbb{N}$ , the  $r$ :th staircase function  $\varsigma_r: [0, +\infty] \rightarrow [0, r]$  is monotone, so when applied to the increasing sequence  $f_1 \leq f_2 \leq \dots$  it produces an increasing sequence  $f_1^{(r)} \leq f_2^{(r)} \leq \dots$ . By left-continuity of  $\varsigma_r$ , the limit of this increasing sequence is  $f^{(r)} := \varsigma_r \circ f$  (see Remark III.20), i.e.,

$$f_n^{(r)} \uparrow f^{(r)} \quad \text{as } n \rightarrow \infty.$$

In particular, we thus see that the array of simple functions  $(f_n^{(r)})_{n \in \mathbb{N}, r \in \mathbb{N}}$  is pointwise increasing in both of its indices,  $n$  and  $r$ .

Now for  $n \in \mathbb{N}$  and  $r \in \mathbb{N}$ , denote

$$t_n^{(r)} := \int^{\square} f_n^{(r)} \, d\mu.$$

Since the array of functions is increasing in both indices, by the monotonicity of the integral  $\int^{\square}$ , also the array  $(t_n^{(r)})_{n \in \mathbb{N}, r \in \mathbb{N}}$  of their integrals is increasing in both indices.

Consider first a fixed  $r \in \mathbb{N}$ . Then we have  $f_n^{(r)} \uparrow f^{(r)}$  as  $n \rightarrow \infty$ , where  $f^{(r)} = \varsigma_r \circ f$ . The functions here are all simple, so we can apply the already proven Monotone convergence theorem for simple functions, Lemma D.2, to get that

$$t_n^{(r)} := \int^{\square} f_n^{(r)} \, d\mu \uparrow \int^{\square} f^{(r)} \, d\mu =: t^{(r)} \quad \text{as } n \rightarrow \infty.$$

The functions  $f^{(r)} = \varsigma_r \circ f$  themselves constitute an increasing approximation of  $f$  by simple functions,

$$f^{(r)} \uparrow f \quad \text{as } r \rightarrow \infty.$$

Therefore we can use the already proven Monotone convergence theorem with simple approximating functions, Lemma D.5, to get

$$t^{(r)} := \int^{\square} f^{(r)} \, d\mu \uparrow \int^+ f \, d\mu =: t^{(\infty)} \quad \text{as } r \rightarrow \infty.$$

Consider then a fixed  $n \in \mathbb{N}$ . Then we have by construction  $f_n^{(r)} \uparrow f_n$  as  $r \rightarrow \infty$ . The approximating functions here are all simple, so we can apply the already proven Monotone convergence theorem with simple approximating functions, Lemma D.5, to get that

$$t_n^{(r)} := \int^{\square} f_n^{(r)} \, d\mu \uparrow \int^+ f_n \, d\mu =: t_n \quad \text{as } r \rightarrow \infty.$$

At this stage we apply again Lemma D.4: it says that as  $n \rightarrow \infty$  we have  $t_n \uparrow t_{\infty} = t^{(\infty)}$ . Recalling what  $t_n$  and  $t^{(\infty)}$  are, we have obtained

$$\int^+ f_n \, d\mu = t_n \uparrow t_{\infty} = t^{(\infty)} = \int^+ f \, d\mu.$$

This is exactly the assertion of the Monotone convergence theorem. □

## Orthogonal projections and conditional expected values

This appendix concerns orthogonal projections in the space  $\mathcal{L}^2(\mathbf{P})$  of square integrable random variables, and conditional expected values with respect to  $\sigma$ -algebras.

### E.1. Geometry of the space of square integrable random variables

The Cauchy-Schwarz inequality, Equation (X.4) in Theorem X.6,

$$|\mathbf{E}[XY]| \leq \sqrt{\mathbf{E}[X^2] \mathbf{E}[Y^2]} \quad \text{for } X, Y \in \mathcal{L}^2(\mathbf{P})$$

underlies a lot of familiar geometry in the space  $\mathcal{L}^2(\mathbf{P})$  of square integrable random variables. We begin by defining inner products, norm, and distances, and establishing results of familiar geometric flavor about them.

#### Inner product, norm, and distance

**Definition E.1** (Inner product and norm for square integrable random variables).

For  $X, Y \in \mathcal{L}^2(\mathbf{P})$ , we denote

$$\langle X, Y \rangle := \mathbf{E}[XY] \tag{E.1}$$

We call this the *inner product* of the random variables  $X$  and  $Y$ . If the inner product vanishes,  $\langle X, Y \rangle = 0$ , then we say that  $X$  and  $Y$  are *orthogonal* and we denote  $X \perp Y$ .

For  $X \in \mathcal{L}^2(\mathbf{P})$ , we denote

$$\|X\| := \sqrt{\langle X, X \rangle} = \sqrt{\mathbf{E}[X^2]} \tag{E.2}$$

We call this the *norm* (or more specifically *2-norm*) of the random variable  $X$ .

In this notation, the Cauchy-Schwarz inequality reads

$$|\langle X, Y \rangle| \leq \|X\| \|Y\|.$$

Corollary X.7, in turn, amounts to the following bound

$$\mathbf{E}[|X|] \leq \|X\|$$

for expected values in terms of the norm.

The norm leads to a notion of distance in  $\mathcal{L}^2(\mathbf{P})$ : for two square integrable random variables  $X, Y \in \mathcal{L}^2(\mathbf{P})$ , we interpret the norm of their difference

$$\|X - Y\|$$

as the distance between the two.<sup>1</sup>

Note that the distance satisfies the usual triangle inequality.

**Lemma E.2** (Triangle inequality).

For any  $X, Y \in \mathcal{L}^2(\mathbf{P})$ , we have

$$\|X + Y\| \leq \|X\| + \|Y\|. \quad (\text{E.3})$$

*Proof.* We have, by bilinearity of the inner product and Cauchy-Schwarz inequality

$$\begin{aligned} \|X + Y\|^2 &= \langle X + Y, X + Y \rangle = \langle X, X \rangle + \underbrace{2\langle X, Y \rangle}_{\leq 2\|X\|\|Y\|} + \langle Y, Y \rangle \\ &\leq \|X\|^2 + 2\|X\|\|Y\| + \|Y\|^2 = (\|X\| + \|Y\|)^2. \end{aligned}$$

The assertion follows by taking square roots.  $\square$

Regarding orthogonality, we have the following familiar formula.

**Lemma E.3** (Pythagoras' theorem).

If  $X, Y \in \mathcal{L}^2$  are orthogonal,  $X \perp Y$ , then we have

$$\|X + Y\|^2 = \|X\|^2 + \|Y\|^2.$$

*Proof.* By direct calculation using bilinearity and symmetry of the inner product, we get

$$\begin{aligned} \|X + Y\|^2 &= \langle X + Y, X + Y \rangle \\ &= \langle X, X \rangle + 2\underbrace{\langle X, Y \rangle}_{=0} + \langle Y, Y \rangle = \|X\|^2 + \|Y\|^2. \end{aligned}$$

$\square$

A similar calculation yields another useful formula.

**Lemma E.4** (Parallelogram law).

If  $X, Y \in \mathcal{L}^2$ , then we have

$$\|X + Y\|^2 + \|X - Y\|^2 = 2\|X\|^2 + 2\|Y\|^2.$$

*Proof.* By direct calculation using bilinearity and symmetry of the inner product, we get

$$\begin{aligned} \|X + Y\|^2 + \|X - Y\|^2 &= \langle X + Y, X + Y \rangle + \langle X - Y, X - Y \rangle \\ &= 2\langle X, X \rangle + 0\langle X, Y \rangle + 2\langle Y, Y \rangle = 2\|X\|^2 + 2\|Y\|^2. \end{aligned}$$

$\square$

<sup>1</sup>It is common to quotient the space  $\mathcal{L}^p(\mathbf{P})$  by the equivalence relation

$$X \stackrel{\text{a.s.}}{=} Y \iff \mathbf{P}\left[\left\{\omega \in \Omega \mid X(\omega) = Y(\omega)\right\}\right] = 1$$

of almost sure equality. This is quite natural, because we have  $\mathbf{E}[|X|^p] = \mathbf{E}[|Y|^p]$  whenever  $X \stackrel{\text{a.s.}}{=} Y$ , and moreover we have  $\mathbf{E}[|X|^p] = 0$  if and only if  $X \stackrel{\text{a.s.}}{=} 0$ . The quotient space  $L^p(\mathbf{P}) = \mathcal{L}^p(\mathbf{P}) / \stackrel{\text{a.s.}}{=}$  is a vector space, and the formula  $\|X\|_p = (\mathbf{E}[|X|^p])^{1/p}$  defines a norm in it. The statements of Proposition E.6 and Exercise E.1 then assert that  $L^2(\mathbf{P})$  and  $L^p(\mathbf{P})$  are Banach spaces, i.e., a complete normed vector spaces.

There would be certain advantages in identifying random variables which are almost surely equal. We, however, choose not to use this quotient space — agreeing with [Wil91], we find it preferable that random variables are functions on  $\Omega$  instead of equivalence classes of such functions.

The notion of distance, in turn, leads to a notion of convergence of sequences. Compare this notion of convergence in the space  $\mathcal{L}^2(\mathbb{P})$  of square integrable random variables with Definition XI.8 on convergence in the space  $\mathcal{L}^1(\mathbb{P})$  of integrable random variables.<sup>2</sup>

**Definition E.5** (Convergence in  $\mathcal{L}^2$ ).

Suppose that  $X_1, X_2, \dots \in \mathcal{L}^2(\mathbb{P})$  and  $X \in \mathcal{L}^2(\mathbb{P})$ . We say that  $X_n$  tends to  $X$  in  $\mathcal{L}^2$  as  $n \rightarrow \infty$  and denote  $X_n \xrightarrow{\mathcal{L}^2} X$ , if we have  $\|X_n - X\| \rightarrow 0$  or equivalently  $\mathbb{E}[(X_n - X)^2] \rightarrow 0$  as  $n \rightarrow \infty$ .

### Completeness of square integrable random variables

The space of square integrable random variables has the following completeness property.

**Proposition E.6** (The space of square integrable random variables is complete).

Suppose that  $X_1, X_2, \dots \in \mathcal{L}^2(\mathbb{P})$  is a sequence of square integrable random variables which is Cauchy in the sense that

$$\lim_{m \rightarrow \infty} \sup_{n, n' \geq m} \|X_n - X_{n'}\| = 0. \tag{E.4}$$

Then there exists a square integrable random variable  $X \in \mathcal{L}^2(\mathbb{P})$  such that we have  $X_n \xrightarrow{\mathcal{L}^2} X$ .

*Proof.* Assuming (E.4), we can choose  $m_1 < m_2 < \dots$  such that

$$\|X_n - X_{n'}\| \leq 2^{-k} \quad \text{whenever } n, n' \geq m_k. \tag{E.5}$$

By Corollary X.7, we then have also

$$\mathbb{E}[|X_n - X_{n'}|] \leq \sqrt{\mathbb{E}[(X_n - X_{n'})^2]} = \|X_n - X_{n'}\| \leq 2^{-k} \quad \text{for } n, n' \geq m_k.$$

In particular we get  $\mathbb{E}[|X_{m_{k+1}} - X_{m_k}|] \leq 2^{-k}$  and thus

$$\sum_{k=1}^{\infty} \mathbb{E}[|X_{m_{k+1}} - X_{m_k}|] < +\infty.$$

According to Lemma VIII.6, it follows that almost surely the series

$$\sum_{k=1}^{\infty} (X_{m_{k+1}} - X_{m_k})$$

converges absolutely. Let  $Y$  denote the sum of this almost surely convergent series,

$$\begin{aligned} Y &= \sum_{k=1}^{\infty} (X_{m_{k+1}} - X_{m_k}) = \lim_{\ell \rightarrow \infty} \sum_{k=1}^{\ell} (X_{m_{k+1}} - X_{m_k}) \\ &= \lim_{\ell \rightarrow \infty} (X_{m_{\ell+1}} - X_{m_\ell} + X_{m_\ell} - X_{m_{\ell-1}} + \dots + X_{m_3} - X_{m_2} + X_{m_2} - X_{m_1}) \\ &= \lim_{\ell \rightarrow \infty} (X_{m_{\ell+1}} - X_{m_1}). \end{aligned}$$

<sup>2</sup>The limit of a sequence converging in  $\mathcal{L}^2(\mathbb{P})$  (or in  $\mathcal{L}^1(\mathbb{P})$ , for that matter) is not strictly speaking uniquely defined. Rather, if  $X, X' \in \mathcal{L}^2(\mathbb{P})$  are  $\mathcal{L}^2$ -limits of the same sequence, then one can only conclude that they are almost surely equal,  $X \stackrel{\text{a.s.}}{=} X'$ . We have chosen to embrace this slight non-uniqueness of limits. The alternative approach would again be to quotient by the equivalence relation of almost sure equality,  $\stackrel{\text{a.s.}}{=}$ .

By setting  $X = Y + X_{m_1}$ , we conclude that almost surely

$$\lim_{\ell \rightarrow \infty} X_{m_\ell} = X.$$

For any  $n \geq m_k$ , as a consequence of (E.5) and Fatou's lemma, we obtain

$$\begin{aligned} \mathbb{E}[(X_n - X)^2] &= \mathbb{E}\left[\lim_{\ell \rightarrow \infty} (X_n - X_{m_\ell})^2\right] \\ &\leq \liminf_{\ell} \mathbb{E}[(X_n - X_{m_\ell})^2] \\ &= \liminf_{\ell} \|X_n - X_{m_\ell}\|^2 \leq 4^{-k}. \end{aligned}$$

This shows first of all that  $X_n - X \in \mathcal{L}^2(\mathbb{P})$ , so by the vector space property we get that  $X = X_n - (X_n - X)$  is square integrable as well. Moreover, it shows that

$$\lim_{n \rightarrow \infty} \mathbb{E}[(X_n - X)^2] = 0,$$

which completes the proof.  $\square$

**Exercise E.1** (Completeness of  $\mathcal{L}^p$ ).

Prove an analogous statement for  $\mathcal{L}^p(\mathbb{P})$ , using Exercise VIII.9 instead of Corollary X.7.

## Orthogonal projections to closed subspaces

**Definition E.7** (Closed subspaces of square integrable random variables).

A vector subspace  $\mathcal{V} \subset \mathcal{L}^2(\mathbb{P})$  in the space of square integrable random variables is said to be *closed*, if for any sequence  $X_1, X_2, \dots \in \mathcal{V}$  which converges in  $\mathcal{L}^2(\mathbb{P})$ , a limit can be found also within the subspace, i.e.,  $X_n \xrightarrow{\mathcal{L}^2} X \in \mathcal{V}$ .

We say that  $Y \in \mathcal{L}^2(\mathbb{P})$  is *orthogonal* to the subspace and denote  $Y \perp \mathcal{V}$ , if for all  $X \in \mathcal{V}$  we have  $Y \perp X$ .

Given a random variable and a subspace, the natural notion of projecting the random variable to the subspace is obtained by finding the nearest point within the subspace. The existence and almost uniqueness of such a nearest point is guaranteed by the following proposition, which also explains why this is an orthogonal projection.

**Proposition E.8** (Orthogonal projection).

Let  $\mathcal{V} \subset \mathcal{L}^2(\mathbb{P})$  be a closed subspace of square integrable random variables and let  $Y \in \mathcal{L}^2(\mathbb{P})$  be a square integrable random variable. Define the distance of  $Y$  to the subspace  $\mathcal{V}$  by

$$\Delta := \inf_{X \in \mathcal{V}} \|Y - X\|. \quad (\text{E.6})$$

Then for a random variable  $Z \in \mathcal{V}$  the following conditions are equivalent:

$$(i): \|Y - Z\| = \Delta \qquad (ii): Y - Z \perp \mathcal{V}.$$

Furthermore, there exists a random variable  $Z \in \mathcal{V}$  with these properties, and if  $\tilde{Z} \in \mathcal{V}$  is another such random variable, then we have  $Z = \tilde{Z}$  almost surely.

*Proof.* We first show the equivalence of the conditions (i) and (ii), then prove the existence of  $Z$  satisfying (ii), and finally prove the almost uniqueness.



*proof of (ii)  $\Rightarrow$  (i):* Suppose that  $Z \in \mathcal{V}$  satisfies (ii). Let  $Z' \in \mathcal{V}$  be any other point in the subspace. Because we have  $Z - Z' \in \mathcal{V}$ , property (ii) implies  $Y - Z \perp Z - Z'$ . Therefore Pythagoras' theorem gives

$$\begin{aligned} \|Y - Z'\|^2 &= \|Y - Z + Z - Z'\|^2 \\ &= \|Y - Z\|^2 + \|Z - Z'\|^2 \geq \|Y - Z\|^2. \end{aligned}$$

Therefore  $Z$  minimizes the distance to  $Y$  within the subspace,  $\|Y - Z\| = \Delta$ , i.e., (i) holds.

*proof of (i)  $\Rightarrow$  (ii):* Suppose that  $Z \in \mathcal{V}$  satisfies (i), i.e.,  $\|Y - Z\| = \Delta$ . Let  $V \in \mathcal{V}$ , and for  $t \in \mathbb{R}$  consider the vector  $Z + tV \in \mathcal{V}$ . Then by definition (E.6) we have

$$\begin{aligned} 0 &\leq \|Y - Z - tV\|^2 - \Delta^2 \\ &= \underbrace{\|Y - Z\|^2}_{=\Delta^2} - 2t \langle Y - Z, V \rangle + t^2 \|V\|^2 - \Delta^2 \\ &= -2t \langle Y - Z, V \rangle + t^2 \|V\|^2. \end{aligned}$$

If  $\langle Y - Z, V \rangle \neq 0$ , then this polynomial would obtain negative values for small positive or small negative  $t$ , which is a contradiction. Therefore we must have  $\langle Y - Z, V \rangle = 0$ , i.e., (ii) holds.

*proof of existence of minimizer:* By definition (E.6), we can find a sequence  $Z_1, Z_2, \dots \in \mathcal{V}$  such that

$$\|Y - Z_n\|^2 \leq \Delta^2 + \frac{1}{n}.$$

By the parallelogram law, Lemma E.4, for any  $n, n'$  we have

$$2\|Y - Z_n\|^2 + 2\|Y - Z_{n'}\|^2 = \|2Y - Z_n - Z_{n'}\|^2 + \|Z_n - Z_{n'}\|^2.$$

Note here, that we have  $\|2Y - Z_n - Z_{n'}\|^2 = 4\|Y - \frac{Z_n + Z_{n'}}{2}\|^2 \geq 4\Delta^2$ , since also  $\frac{Z_n + Z_{n'}}{2} \in \mathcal{V}$ . We thus conclude that

$$\begin{aligned} \|Z_n - Z_{n'}\|^2 &= 2\|Y - Z_n\|^2 + 2\|Y - Z_{n'}\|^2 - \|2Y - Z_n - Z_{n'}\|^2 \\ &\leq 2\left(\Delta^2 + \frac{1}{n}\right) + 2\left(\Delta^2 + \frac{1}{n'}\right) - 4\Delta^2 \\ &\leq \frac{2}{n} + \frac{2}{n'}. \end{aligned}$$

This shows that the sequence  $Z_1, Z_2, \dots$  is Cauchy. Therefore by Proposition E.6, the sequence  $Z_1, Z_2, \dots$  converges in  $\mathcal{L}^2(\mathbf{P})$ . Since  $\mathcal{V}$  is a closed subspace, a limit remains in the subspace,  $Z_n \xrightarrow{\mathcal{L}^2} Z \in \mathcal{V}$ .

By triangle inequality, we have

$$\Delta \leq \|Y - Z\| = \|Y - Z_n + Z_n - Z\| \leq \underbrace{\|Y - Z_n\|}_{\rightarrow \Delta} + \underbrace{\|Z_n - Z\|}_{\rightarrow 0} \xrightarrow{n \rightarrow \infty} \Delta,$$

so we conclude that  $\|Y - Z\| = \Delta$ . This shows the existence of a random variable  $Z \in \mathcal{V}$  satisfying (i), and therefore also (ii).

*proof of almost uniqueness:* Suppose that  $Z, \tilde{Z} \in \mathcal{V}$  are two random variables satisfying (i) and (ii). Then  $Z - \tilde{Z} \in \mathcal{V}$  and thus by (ii) we have  $Z - \tilde{Z} \perp Y - Z$ . Therefore property (i) and Pythagoras theorem again lead to

$$\begin{aligned} \Delta^2 &= \|Y - \tilde{Z}\|^2 = \|Y - Z + Z - \tilde{Z}\|^2 \\ &= \|Y - Z\|^2 + \|Z - \tilde{Z}\|^2 = \Delta^2 + \|Z - \tilde{Z}\|^2. \end{aligned}$$

This shows that  $\|Z - \tilde{Z}\|^2 = 0$ , which implies  $Z = \tilde{Z}$  almost surely.  $\square$

Given  $Y \in \mathcal{L}^2(\mathbf{P})$  and a closed subspace  $\mathcal{V} \subset \mathcal{L}^2(\mathbf{P})$ , a random variable  $Z$  which satisfies (i) and (ii) is called (a version of) the orthogonal projection of  $Y$  to  $\mathcal{V}$ .

## E.2. Conditional expected values

Throughout, let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space. Given an integrable random variable  $Y$  and a sub- $\sigma$ -algebra  $\mathcal{G} \subset \mathcal{F}$ , we will define and study the *conditional expected value*  $\mathbb{E}[Y|\mathcal{G}]$  of  $Y$  given  $\mathcal{G}$ . This should be interpreted as the best estimate of  $Y$  that can be made based on the information  $\mathcal{G}$ .<sup>3</sup> We start by discussing the case when  $Y$  is square-integrable.

### Orthogonal projections as best estimates

Observe that if  $\mathcal{G} \subset \mathcal{F}$  is a sub- $\sigma$ -algebra of  $\mathcal{F}$ , then the space

$$\mathcal{L}^2(\mathbb{P}) \cap \mathfrak{m}\mathcal{G}$$

of square integrable  $\mathcal{G}$ -measurable random variables is a closed subspace of  $\mathcal{L}^2(\mathbb{P})$ : it is clearly a vector subspace, and satisfies the completeness property of Proposition E.6 — therefore any sequence in it which converges has a limit in it (convergent sequences are necessarily Cauchy).

Suppose now that the  $\sigma$ -algebra  $\mathcal{G}$  represents information available to us, and our task is to form an estimate  $\widehat{Y}$  of a random quantity  $Y \in \mathcal{L}^2(\mathbb{P})$  based on that information. To formalize the property that the estimate is made based on information  $\mathcal{G}$ , we require that  $\widehat{Y}$  is a  $\mathcal{G}$ -measurable random variable.<sup>4</sup> The difference  $Y - \widehat{Y}$  between the actual value and our estimate is the error we commit, and thus the norm  $\|Y - \widehat{Y}\|$  tells us about the magnitude of the error.<sup>5</sup> In this sense, the best possible estimate is obtained by choosing  $\widehat{Y} \in \mathcal{L}^2(\mathbb{P}) \cap \mathfrak{m}\mathcal{G}$  which minimizes  $\|Y - \widehat{Y}\|$ . According to Proposition E.8, this is achieved by letting  $\widehat{Y}$  be the orthogonal projection of  $Y$  onto the subspace  $\mathcal{V} = \mathcal{L}^2(\mathbb{P}) \cap \mathfrak{m}\mathcal{G}$ .

We then observe a key property that the best estimate  $\widehat{Y}$  satisfies.

**Lemma E.9** (Orthogonal projection to  $\mathcal{G}$ -measurable random variables).

*For  $\mathcal{G} \subset \mathcal{F}$  a sub- $\sigma$ -algebra and  $Y \in \mathcal{L}^2$ , let  $\widehat{Y}$  be (a version of) the orthogonal projection of  $Y$  to the closed subspace  $\mathcal{L}^2(\mathbb{P}) \cap \mathfrak{m}\mathcal{G}$  of  $\mathcal{G}$ -measurable square integrable random variables. Then for any  $G \in \mathcal{G}$  we have*

$$\mathbb{E}[\mathbb{I}_G \widehat{Y}] = \mathbb{E}[\mathbb{I}_G Y].$$

*Proof.* Obviously  $\mathbb{I}_G$  is  $\mathcal{G}$ -measurable and square integrable, so  $\mathbb{I}_G \in \mathcal{V} := \mathcal{L}^2(\mathbb{P}) \cap \mathfrak{m}\mathcal{G}$ . By definition,  $\widehat{Y}$  satisfies property (ii) of Proposition E.8:  $Y - \widehat{Y} \perp \mathcal{V}$ . Thus we in particular have  $Y - \widehat{Y} \perp \mathbb{I}_G$ . This can be explicitly written as

$$0 = \langle Y - \widehat{Y}, \mathbb{I}_G \rangle = \langle Y, \mathbb{I}_G \rangle - \langle \widehat{Y}, \mathbb{I}_G \rangle = \mathbb{E}[Y \mathbb{I}_G] - \mathbb{E}[\widehat{Y} \mathbb{I}_G].$$

The asserted equality follows. □

This observation is the motivating idea, which underlies the abstract definition of conditional expected values.

<sup>3</sup>Recall from Lecture IV the interpretation of  $\sigma$ -algebras as information.

<sup>4</sup>Morally, the value  $\widehat{Y}$  of the estimate should be a deterministic function of the known information (see Doob's representation theorem, Theorem IV.5).

<sup>5</sup>In particular,  $\widehat{Y}$  must be square integrable for the magnitude of error in this sense to be finite.

**Definition of conditional expected value**

Let now  $Y \in \mathcal{L}^1(\mathbf{P})$  and let  $\mathcal{G} \subset \mathcal{F}$  be a sub- $\sigma$ -algebra.

**Definition E.10** (Conditional expected value).

A random variable  $\widehat{Y} \in \mathcal{L}^1(\mathbf{P})$  is said to be (a version of) the *conditional expected value* (denoted  $\mathbf{E}[Y|\mathcal{G}]$ ) of  $Y$  given  $\mathcal{G}$ , if  $\widehat{Y} \in \mathfrak{m}\mathcal{G}$  and

$$\mathbf{E}[\mathbb{I}_G \widehat{Y}] = \mathbf{E}[\mathbb{I}_G Y] \quad \text{for all } G \in \mathcal{G}. \quad (\text{E.7})$$

In particular, if  $Y \in \mathcal{L}^2(\mathbf{P})$ , then by Lemma E.9 the orthogonal projection of  $Y$  to  $\mathcal{L}^2(\mathbf{P}) \cap \mathfrak{m}\mathcal{G}$  is (a version of) this conditional expected value. In the general case of  $Y \in \mathcal{L}^1(\mathbf{P})$ , we still have to show that such conditional expected values exist. We start, however, by first addressing their uniqueness (up to the usual amount of ambiguity).

**Lemma E.11** (Almost uniqueness of conditional expected values).

Suppose that both  $\widehat{Y}$  and  $\widehat{Y}'$  are conditional expected values of  $Y$  given  $\mathcal{G}$ . Then the two are almost surely equal,  $\widehat{Y} \stackrel{\text{a.s.}}{=} \widehat{Y}'$ .

*Proof.* For  $n \in \mathbb{N}$ , let

$$G_n := \left\{ \omega \in \Omega \mid \widehat{Y}(\omega) - \widehat{Y}'(\omega) \geq \frac{1}{n} \right\}.$$

Then we have  $G_n \in \mathcal{G}$ , since  $\widehat{Y}$  and  $\widehat{Y}'$  are  $\mathcal{G}$ -measurable. By Markov's inequality, we get

$$\begin{aligned} \frac{1}{n} \mathbf{P}[G_n] &\leq \mathbf{E}[\mathbb{I}_{G_n} (\widehat{Y} - \widehat{Y}')] = \mathbf{E}[\mathbb{I}_{G_n} \widehat{Y}] - \mathbf{E}[\mathbb{I}_{G_n} \widehat{Y}'] \\ &= \mathbf{E}[\mathbb{I}_{G_n} Y] - \mathbf{E}[\mathbb{I}_{G_n} Y] = 0, \end{aligned}$$

where we used the definition of conditional expected values for both  $\widehat{Y}$  and  $\widehat{Y}'$ . We conclude that  $\mathbf{P}[G_n] = 0$ , and then using the union bound also

$$\mathbf{P}[\widehat{Y} > \widehat{Y}'] = \mathbf{P}\left[ \bigcup_{n \in \mathbb{N}} G_n \right] \leq \sum_{n \in \mathbb{N}} \mathbf{P}[G_n] = \sum_{n \in \mathbb{N}} 0 = 0.$$

By changing the roles of  $\widehat{Y}$  and  $\widehat{Y}'$ , one similarly derives  $\mathbf{P}[\widehat{Y} < \widehat{Y}'] = 0$ , and therefore

$$\mathbf{P}[\widehat{Y} \neq \widehat{Y}'] = 0.$$

By passing to the complementary events, this concludes the proof of  $\widehat{Y} \stackrel{\text{a.s.}}{=} \widehat{Y}'$ .  $\square$

We will denote the conditional expected values of  $Y$  given  $\mathcal{G}$  by

$$\mathbf{E}[Y|\mathcal{G}] \in \mathcal{L}^1(\mathbf{P}) \cap \mathfrak{m}\mathcal{G},$$

and not worry too much about the possibility that different choices for it could be made, since any two choices are anyway almost surely equal.

Admitting that conditional expected values exist, the reader can now for example verify the following.

**Exercise E.2** (Conditional expected value preserves non-negativity).

Show that if  $Y \geq 0$  (almost surely) then also  $\mathbf{E}[Y|\mathcal{G}] \geq 0$  (almost surely).

The remaining task is to show that the conditional expected value  $\mathbf{E}[Y|\mathcal{G}]$  exists not only when  $Y \in \mathcal{L}^2(\mathbf{P})$ , but generally for any  $Y \in \mathcal{L}^1(\mathbf{P})$ . We first do this by

assuming non-negativity of  $Y$ , and then it is routine to deal with the general case by splitting to positive and negative parts.

### Conditional expected value for integrable random variables

Let  $Y$  be a non-negative integrable random variable. Denote by  $Y \wedge n$  this random variable truncated at level  $n$ :

$$(Y \wedge n)(\omega) = \min \{Y(\omega), n\} \quad \text{for } \omega \in \Omega.$$

Then  $Y \wedge n$  is bounded and in particular square integrable,  $Y \wedge n \in \mathcal{L}^2(\mathbf{P})$ . Moreover, we have  $Y \wedge n \uparrow Y$  as  $n \rightarrow \infty$ .

We use this approximation by square integrable random variables to construct the conditional expected value of non-negative integrable  $Y$ .

**Lemma E.12** (Truncation approximation to conditional expected values).

*Let  $Y$  be a non-negative integrable random variable, and let  $Z_n$  be the orthogonal projection of  $Y \wedge n$  to the subspace  $\mathcal{L}^2(\mathbf{P}) \cap \mathfrak{m}\mathcal{G}$ . Then there exists an integrable  $\mathfrak{m}\mathcal{G}$ -measurable random variable  $Z$  such that  $Z_n \uparrow Z$  as  $n \rightarrow \infty$  (almost surely), and we have*

$$\mathbb{E}[\mathbb{I}_G Z] = \mathbb{E}[\mathbb{I}_G Y] \quad \text{for any } G \in \mathcal{G}.$$

*Proof.* Since  $Y \wedge (n+1) \geq Y \wedge n$  for any  $n \in \mathbb{N}$ , it follows from linearity of projection and Exercise E.2 that  $Z_{n+1} \geq Z_n$  (almost surely). Therefore the sequence  $Z_1, Z_2, \dots$  of (almost surely) non-negative random variables is (almost surely) increasing, and thus has a limit  $Z$ . By the  $\mathfrak{m}\mathcal{G}$ -measurability of each  $Z_n$ , the limit  $Z$  is also  $\mathfrak{m}\mathcal{G}$ -measurable. Now let  $G \in \mathcal{G}$ . Then using Lemma E.9 once and the Monotone convergence theorem twice, we calculate

$$\begin{aligned} \mathbb{E}[\mathbb{I}_G Z] &= \mathbb{E}\left[\mathbb{I}_G \left(\lim_{n \rightarrow \infty} Z_n\right)\right] = \lim_{n \rightarrow \infty} \mathbb{E}[\mathbb{I}_G Z_n] \\ &= \lim_{n \rightarrow \infty} \mathbb{E}[\mathbb{I}_G (Y \wedge n)] \\ &= \mathbb{E}\left[\mathbb{I}_G \left(\lim_{n \rightarrow \infty} (Y \wedge n)\right)\right] = \mathbb{E}[\mathbb{I}_G Y]. \end{aligned}$$

This is the asserted property of the limit  $Z$ , and taking  $G = \Omega$  in particular shows that indeed  $Z \in \mathcal{L}^1(\mathbf{P})$ .  $\square$

**Proposition E.13** (Existence of conditional expected values).

*For any  $Y \in \mathcal{L}^1(\mathbf{P})$ , (a version of) the conditional expected value  $\mathbb{E}[Y|\mathcal{G}]$  exists.*

*Proof.* Decompose  $Y$  to its positive and negative parts,  $Y = Y_+ - Y_-$ . Lemma E.12 can be used to construct conditional expected values  $\widehat{Y}_+$  and  $\widehat{Y}_-$  of the non-negative integrable random variables  $Y_+$  and  $Y_-$  given  $\mathcal{G}$ . Then  $\mathbb{E}[Y|\mathcal{G}] := \widehat{Y}_+ - \widehat{Y}_-$  clearly satisfies the defining property of conditional expected value of  $Y$  given  $\mathcal{G}$ .  $\square$

### Properties of conditional expected value

In the following we summarize a number of important properties of conditional expected values.

**Theorem E.14** (Properties of conditional expected values).

*Conditional expected values satisfy the following properties (interpreted in the almost sure sense), when  $Y$  and  $Y_1, Y_2, \dots$  are integrable random variables.*

- (i) *If  $Y \in \mathfrak{m}\mathcal{G}$ , then we have  $\mathbb{E}[Y|\mathcal{G}] = Y$ .*
- (ii) *We have  $\mathbb{E}[\mathbb{E}[Y|\mathcal{G}]] = \mathbb{E}[Y]$ .*
- (iii) *For  $c_1, c_2 \in \mathbb{R}$ , we have  $\mathbb{E}[c_1 Y_1 + c_2 Y_2 | \mathcal{G}] = c_1 \mathbb{E}[Y_1|\mathcal{G}] + c_2 \mathbb{E}[Y_2|\mathcal{G}]$ .*
- (iv) *If  $\mathcal{H} \subset \mathcal{G} \subset \mathcal{F}$  are  $\sigma$ -algebras, then we have  $\mathbb{E}[\mathbb{E}[Y|\mathcal{G}] | \mathcal{H}] = \mathbb{E}[Y|\mathcal{H}]$ .*
- (v) *If  $Z \in \mathfrak{m}\mathcal{G}$  and  $ZY \in \mathcal{L}^1(\mathbb{P})$ , then we have  $\mathbb{E}[ZY|\mathcal{G}] = Z \mathbb{E}[Y|\mathcal{G}]$ .*
- (vi) *If  $\mathcal{G} \perp \sigma(Y)$ , then we have  $\mathbb{E}[Y|\mathcal{G}] = \mathbb{E}[Y]$ .*

**Remark E.15** (Interpretations of the properties of conditional expected values).

The properties in the theorem above have rather intuitive interpretations.

- (i) The best estimate of a known quantity  $Y \in \mathfrak{m}\mathcal{G}$  is the quantity  $Y$  itself.
- (ii) The best estimate of a quantity  $Y$  is unbiased, in the sense that it has the same expected value as the quantity  $Y$  itself.
- (iii) The best estimate of a linear combination of quantities is the corresponding linear combination of the best estimates.
- (iv) Suppose that a person  $H$  possesses less information than a person  $G$ . If  $H$  tries to form an estimate about the best estimate that  $G$  makes about some quantity  $Y$ , then the best she can do is to use her own best estimate of the quantity  $Y$ .
- (v) Known quantities can be treated like constants when forming best estimates.
- (vi) Any information that is independent of  $Y$  can not be used to estimate  $Y$  any better than the expected value  $\mathbb{E}[Y]$  of  $Y$ .

*Proof of Theorem E.14.* Property (i) is immediate from the defining equation (E.7) and the (almost) uniqueness of conditional expected value (Lemma E.11).

Property (ii) follows by taking  $G = \Omega \in \mathcal{G}$  in the defining equation (E.7) of conditional expected value.

Property (iii) is a consequence of the linearity of the defining equation (E.7) and the (almost) uniqueness of conditional expected value (Lemma E.11).

To prove property (iv), note first that both  $\mathbb{E}[Y|\mathcal{H}]$  and  $\mathbb{E}[\mathbb{E}[Y|\mathcal{G}] | \mathcal{H}]$  are by construction  $\mathcal{H}$ -measurable and integrable, so it only remains to verify the defining equation (E.7). For  $H \in \mathcal{H} \subset \mathcal{G} \subset \mathcal{F}$ , using the defining property of conditional expected values, we get

$$\mathbb{E}[\mathbb{I}_H \mathbb{E}[\mathbb{E}[Y|\mathcal{G}] | \mathcal{H}]] = \mathbb{E}[\mathbb{I}_H \mathbb{E}[Y|\mathcal{G}]] = \mathbb{E}[\mathbb{I}_H Y].$$

Since conditional expected values are (almost) uniquely defined, this proves that

$$\mathbb{E}[\mathbb{E}[Y|\mathcal{G}] | \mathcal{H}] = \mathbb{E}[Y|\mathcal{H}].$$

We leave it as an exercise to the reader to prove property (v) by the “standard machine”, i.e., by verifying the claim successively when the random variable  $Z$  is a  $\mathcal{G}$ -measurable indicator, a simple random variable, a non-negative random variable, and finally in the full generality of the assertion.

To prove property (vi), note that the assumed independence  $\mathcal{G} \perp \sigma(Y)$  implies that  $\mathbb{I}_G \perp Y$  for any  $G \in \mathcal{G}$ . Therefore we get

$$\mathbb{E}[\mathbb{I}_G Y] = \mathbb{E}[\mathbb{I}_G] \mathbb{E}[Y] = \mathbb{E}[\mathbb{I}_G \mathbb{E}[Y]],$$

which by (almost) uniqueness of conditional expected value shows that  $\mathbb{E}[Y|\mathcal{G}] = \mathbb{E}[Y]$ .  $\square$

Also for example Monotone convergence theorem, Dominated convergence theorem, Fatou’s lemma, Jensen’s inequality, etc. hold in the appropriate form for conditional expected values, and their proofs are straightforward modifications of the corresponding ones for usual expected values.



## Characteristic functions

This appendix is devoted to the proof of two results from Lecture XII: Lévy's inversion theorem (Theorem XII.7), by which the distribution of a random variable is recovered from its characteristic function, and Theorem XII.9 about equivalent conditions for convergence in distribution formulated in terms of cumulative distribution functions or characteristic functions.

### F.1. Lévy's inversion theorem

Let

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

be a real valued random variable. In this section, we use the following notational conventions regarding its distribution:

- $\nu = P_X$ , the distribution of  $X$ , a probability measure on  $\mathbb{R}$  such that

$$\nu[B] = \mathbf{P}[X \in B] \quad \text{for Borel sets } B \in \mathcal{B}.$$

- $\varphi = \varphi_X$ , the characteristic function of  $X$ , a function  $\mathbb{R} \rightarrow \mathbb{C}$  such that

$$\varphi(\theta) = \mathbf{E}[e^{i\theta X}] = \int_{\mathbb{R}} e^{i\theta x} d\nu(x) \quad \text{for } \theta \in \mathbb{R}.$$

- $F = F_X$ , the cumulative distribution function of  $X$ , a function  $\mathbb{R} \rightarrow [0, 1]$  such that

$$F(x) = \mathbf{P}[X \leq x] = \nu[(-\infty, x]] \quad \text{for } x \in \mathbb{R}.$$

The statement of Lévy's inversion theorem is the following.

**Theorem** (Lévy's inversion theorem, Theorem XII.7).

*For any  $a, b \in \mathbb{R}$ ,  $a < b$ , we have*

$$\lim_{T \rightarrow +\infty} \frac{1}{2\pi} \int_{-T}^{+T} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \varphi(\theta) d\theta \tag{F.1}$$

$$= \nu[(a, b)] + \frac{1}{2} \nu[\{a\}] + \frac{1}{2} \nu[\{b\}] \tag{F.2}$$

*Moreover, if  $\int_{\mathbb{R}} |\varphi(\theta)| d\theta < +\infty$ , then  $X$  has a continuous probability density function  $f_X$  given by*

$$f_X(x) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{-i\theta x} \varphi(\theta) d\theta. \tag{F.3}$$

**Remark F.1.** Formula (F.2) is occasionally written in terms of the cumulative distribution function as well. Recall from Proposition II.30 that the cumulative distribution function  $F$  is

increasing and right-continuous: if  $x_n \downarrow x$  then  $F(x_n) \downarrow F(x)$ . Since it is increasing and bounded, the left limits also exist: they are defined by

$$F(x^-) := \lim_{x' \uparrow x} F(x').$$

By monotone convergence of probability measures (Theorem II.22), the left limits can be expressed in terms of the distribution  $\nu$  as the following measures of open semi-infinite intervals

$$\begin{aligned} F(x^-) &= \lim_{n \rightarrow \infty} F\left(x - \frac{1}{n}\right) \\ &= \lim_{n \rightarrow \infty} \nu\left[\left(-\infty, x - \frac{1}{n}\right]\right] \\ &= \nu\left[\bigcup_{n=1}^{\infty} \left(-\infty, x - \frac{1}{n}\right]\right] = \nu[(-\infty, x)]. \end{aligned}$$

In particular, if there is a discontinuity in the cumulative distribution function at a point  $x$ , then the size of the jump  $F(x) - F(x^-)$  is the probability mass located at the single point  $x$ ,

$$F(x) - F(x^-) = \nu[(-\infty, x]] - \nu[(-\infty, x)] = \nu[\{x\}]$$

(at continuity points  $x$  of  $F$  there is no probability mass,  $\nu[\{x\}] = 0$ ).

For any  $a < b$ , the increment from  $a$  to  $b$  of the cumulative distribution function is

$$F(b) - F(a) = \nu[(-\infty, b]] - \nu[(-\infty, a]] = \nu[(a, b]].$$

If we replace  $F(x)$  by the average  $\frac{1}{2}(F(x) + F(x^-))$  of the left and right limits, then the corresponding increment becomes

$$\frac{1}{2}(F(b) + F(b^-)) - \frac{1}{2}(F(a) + F(a^-)) = \nu[(a, b)] + \frac{1}{2}\nu[\{a\}] + \frac{1}{2}\nu[\{b\}].$$

This allows to alternatively write (F.2) in terms of the cumulative distribution function  $F$ .

In the proof we use the following auxiliary calculation.

**Lemma F.2** (An auxiliary integral).

For  $r \in \mathbb{R}$  define

$$S(r) := \int_0^r \frac{\sin(\theta)}{\theta} d\theta.$$

Then the limits as  $r \uparrow +\infty$  and as  $r \downarrow -\infty$  of this integral are

$$\lim_{r \uparrow +\infty} S(r) = \frac{\pi}{2} \quad \text{and} \quad \lim_{r \downarrow -\infty} S(r) = -\frac{\pi}{2}.$$

Moreover, for any  $c \in \mathbb{R}$ , we have

$$\int_0^r \frac{\sin(c\theta)}{\theta} d\theta = S(cr).$$

*Proof.* Let us start from the last part: performing the change of variables  $\theta' = c\theta$  we obtain the asserted formula

$$\int_0^r \frac{\sin(c\theta)}{\theta} d\theta = \int_0^{cr} \frac{\sin(\theta')}{\theta'} d\theta' = S(cr).$$

In particular taking  $c = -1$  we see that  $S(-r) = S(r)$ , so it suffices to consider  $r > 0$ .

The improper Riemann integral  $\lim_{r \uparrow +\infty} S(r)$  is one which is routine to evaluate with complex analysis and residue calculus. Indeed, the principal value integral of  $\frac{e^{iz}}{z}$  over the real line evaluates to  $i\pi$  times the residue at  $z = 0$ ,

$$\text{P.V.} \int_{-\infty}^{\infty} \frac{e^{iz}}{z} dz = i\pi \operatorname{Res}_{z=0} \left( \frac{e^{iz}}{z} \right) = i\pi.$$



The result follows from this by taking the imaginary part and using parity. □

*Proof of Theorem XII.7.* We first prove the equality of (F.1) and (F.2). Then we prove the assertion about the probability density assuming integrability of the characteristic function.

*proof of (F.1)=(F.2):* For any  $u, v \in \mathbb{R}$ , we have  $|e^{iu} - e^{iv}| \leq |u - v|$ . We deduce that

$$\int_{\mathbb{R}} \left( \int_{-T}^{+T} \underbrace{\left| \frac{e^{i\theta(x-a)} - e^{i\theta(x-b)}}{i\theta} \right|}_{\leq \frac{|i\theta b - i\theta a|}{|\theta|} \leq |b-a|} d\theta \right) d\nu(x) \leq 2T|b-a| < +\infty.$$

This finiteness of the double integral of absolute value justifies the use of Fubini's theorem to exchange the order of integrations in expression (F.1)

$$\begin{aligned} \int_{-T}^{+T} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \varphi(\theta) d\theta &= \int_{-T}^{+T} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \left( \int_{\mathbb{R}} e^{i\theta x} d\nu(x) \right) d\theta \\ &= \int_{\mathbb{R}} \left( \int_{-T}^{+T} \frac{e^{i\theta(x-a)} - e^{i\theta(x-b)}}{i\theta} d\theta \right) d\nu(x). \end{aligned}$$

In the last expression, the inner integral over  $\theta$  is such that the imaginary part of the integrand is an odd function and the real part of the integrand is an even function, so only the real part survives integration over the symmetric interval  $[-T, T]$ . We are able to express the result in terms of the integrals  $S(r)$  of Lemma F.2:

$$\begin{aligned} \int_{-T}^{+T} \frac{e^{i\theta(x-a)} - e^{i\theta(x-b)}}{i\theta} d\theta &= \int_{-T}^{+T} \frac{\sin(\theta(x-a)) - \sin(\theta(x-b))}{\theta} d\theta \\ &= 2 \int_0^{+T} \frac{\sin(\theta(x-a))}{\theta} d\theta - 2 \int_0^{+T} \frac{\sin(\theta(x-b))}{\theta} d\theta \\ &= 2S((x-a)T) - 2S((x-b)T). \end{aligned}$$

As  $T \rightarrow +\infty$ , it follows from Lemma F.2 that we have

$$\int_{-T}^{+T} \frac{e^{i\theta(x-a)} - e^{i\theta(x-b)}}{i\theta} d\theta \rightarrow \begin{cases} 2\pi & \text{if } a < x < b \\ \pi & \text{if } x = a \text{ or } x = b \\ 0 & \text{if } x < a \text{ or } x > b, \end{cases}$$

and the left hand side expression is uniformly bounded in  $T$  (as an upper bound for its absolute value we may use for example  $4 \sup_{r>0} |S(r)| < +\infty$ ). The bounded convergence theorem therefore says that the limit  $T \rightarrow +\infty$  can be interchanged with the integration over the probability measure  $\nu$ :

$$\int_{\mathbb{R}} \left( \int_{-T}^{+T} \frac{e^{i\theta(x-a)} - e^{i\theta(x-b)}}{i\theta} d\theta \right) d\nu(x) \xrightarrow{T \rightarrow \infty} 2\pi \nu[(a, b)] + \pi \nu[\{a\}] + \pi \nu[\{b\}].$$

This concludes the proof of the first assertion, the equality of (F.1) and (F.2).

*proof of probability density part:* Suppose now that  $\int_{\mathbb{R}} |\varphi(x)| dx < +\infty$ . Note that the function (F.3)

$$f_X(x) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{-i\theta x} \varphi(\theta) d\theta.$$

is continuous: if  $x_n \rightarrow x$  then we have the pointwise limit  $e^{-i\theta x_n} \varphi(\theta) \rightarrow e^{-i\theta x} \varphi(\theta)$  by continuity of the exponential function, and domination  $|e^{-i\theta x_n} \varphi(\theta)| \leq |\varphi(\theta)|$  by an integrable function, so it follows from dominated convergence theorem that  $f_X(x_n) \rightarrow f_X(x)$ . It therefore remains to show that this  $f_X$  is a probability density for the measure  $\nu$ .

Next note that

$$\left| \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \varphi(\theta) \right| \leq |a - b| |\varphi(\theta)|.$$

Therefore by the Dominated convergence theorem (domination by constant multiple of the integrable function  $|\varphi(x)|$ ), the left hand side (F.1) can be written simply as

$$\frac{1}{2\pi} \int_{\mathbb{R}} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \varphi(\theta) \, d\theta.$$

Also by dominated convergence (domination again by multiple of  $|\varphi(x)|$ ), if  $a_n \rightarrow a$  and  $b_n \rightarrow b$ , we get

$$\frac{1}{2\pi} \int_{\mathbb{R}} \frac{e^{-i\theta a_n} - e^{-i\theta b_n}}{i\theta} \varphi(\theta) \, d\theta \rightarrow \frac{1}{2\pi} \int_{\mathbb{R}} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta} \varphi(\theta) \, d\theta.$$

Using the result of the first part and rewriting (F.2) in terms of the cumulative distribution function  $F$ , we get

$$\frac{F(b_n) + F(b_n^-)}{2} - \frac{F(a_n) + F(a_n^-)}{2} \rightarrow \frac{F(b) + F(b^-)}{2} - \frac{F(a) + F(a^-)}{2}.$$

This shows that the cumulative distribution function  $F$  is continuous.

Again, since we have

$$\left| \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta(b-a)} \varphi(x) \right| \leq |\varphi(x)|,$$

we may calculate the derivative of the cumulative distribution function  $F$  by dominated convergence as follows:

$$\begin{aligned} F'(a) &= \lim_{b \rightarrow a} \frac{F(b) - F(a)}{b - a} = \lim_{b \rightarrow a} \frac{1}{2\pi} \int_{\mathbb{R}} \frac{e^{-i\theta a} - e^{-i\theta b}}{i\theta(b-a)} \varphi(\theta) \, d\theta \\ &= \frac{1}{2\pi} \int_{\mathbb{R}} e^{i\theta a} \varphi(\theta) \, d\theta. \end{aligned}$$

Thus the derivative is given by the expression in (F.3),  $F'(x) = f_X(x)$ . It follows that  $f_X$  is a probability density, because for any interval  $(a, b) \subset \mathbb{R}$  we have

$$\nu[(a, b)] = F(b) - F(a) = \int_a^b F'(x) \, dx = \int_a^b f_X(x) \, dx$$

and intervals form a  $\pi$ -system which uniquely determine the probability measure  $\nu$ .  $\square$

Let us at this point also prove a lemma which allows us to control the amount of probability mass that is outside a large interval  $[-r, +r]$  in terms of the behavior of the characteristic function  $\varphi(\theta)$  near  $\theta = 0$ .

**Lemma F.3** (A bound on the probability mass outside an interval).

Let  $\nu$  be a probability measure on  $(\mathbb{R}, \mathcal{B})$ , and let  $\varphi(\theta) = \int_{\mathbb{R}} e^{i\theta x} \, d\nu(s)$  be its characteristic function. Then for any  $r > 0$  we have

$$\nu[\mathbb{R} \setminus [-r, +r]] \leq \frac{r}{2} \int_{-2/r}^{2/r} (1 - \varphi(\theta)) \, d\theta.$$

*Proof.* Let us denote  $u = 2/r$  for convenience. Start with the following calculation of an integral

$$\int_{-u}^u (1 - e^{i\theta x}) \, d\theta = 2u - \frac{1}{ix} (e^{iux} - e^{-iux}) = 2u - \frac{2 \sin(ux)}{x} = 2u \left(1 - \frac{\sin(ux)}{ux}\right).$$

Let us then integrate both sides of this equation over the variable  $x$  with respect to the measure  $\nu$ . On the left hand side we get, using Fubini's theorem,

$$\begin{aligned} \int_{\mathbb{R}} \left( \int_{-u}^u (1 - e^{i\theta x}) \, d\theta \right) \, d\nu(x) &= \int_{-u}^u \left( \int_{\mathbb{R}} (1 - e^{i\theta x}) \, d\nu(x) \right) \, d\theta \\ &= \int_{-u}^u (1 - \varphi(\theta)) \, d\theta. \end{aligned}$$

Dividing by  $u$  and equating with what we get on the right hand side gives

$$\frac{1}{u} \int_{-u}^u (1 - \varphi(\theta)) \, d\theta = 2 \int_{\mathbb{R}} \left(1 - \frac{\sin(ux)}{ux}\right) \, d\nu(x).$$

Because  $\frac{\sin(\xi)}{\xi} \leq 1$  for any  $\xi \in \mathbb{R}$ , we see that the integrand on the right hand side here is non-negative,  $1 - \frac{\sin(ux)}{ux} \geq 0$ . Therefore omitting from the integral over  $\mathbb{R}$  the part over  $[-r, r] \subset \mathbb{R}$  yields a lower bound

$$\frac{1}{u} \int_{-u}^u (1 - \varphi(\theta)) \, d\theta \geq 2 \int_{\mathbb{R} \setminus [-r, r]} \left(1 - \frac{\sin(ux)}{ux}\right) \, d\nu(x).$$

Then observe that for  $|x| > r = 2/u$  we have  $|\frac{\sin(ux)}{ux}| \leq \frac{1}{u|x|} \leq \frac{1}{2}$ , so the integrand on the remaining part satisfies  $1 - \frac{\sin(ux)}{ux} \geq \frac{1}{2}$ , which yields

$$\frac{1}{u} \int_{-u}^u (1 - \varphi(\theta)) \, d\theta \geq 2 \int_{\mathbb{R} \setminus [-r, r]} \frac{1}{2} \, d\nu(x) = \nu[\mathbb{R} \setminus [-r, +r]].$$

Recalling that  $u = 2/r$ , this is exactly the asserted inequality. □

### F.2. Equivalent conditions for convergence in distribution

**Theorem** (Theorem XII.9).

Let  $X_1, X_2, \dots$  and  $X$  be real-valued random variables. Let also  $\nu_1, \nu_2, \dots$  and  $\nu$  be their laws, and let  $F_1, F_2, \dots$  and  $F$  be their cumulative distribution functions, and let  $\varphi_1, \varphi_2, \dots$  and  $\varphi$  be their characteristic functions, respectively. Then the following conditions are equivalent:

(i) For all bounded continuous functions  $f: \mathbb{R} \rightarrow \mathbb{R}$  we have

$$\int_{\mathbb{R}} f(x) \, d\nu_n(x) \longrightarrow \int_{\mathbb{R}} f(x) \, d\nu(x) \quad \text{as } n \rightarrow \infty.$$

(ii) We have  $F_n(x) \rightarrow F(x)$  as  $n \rightarrow \infty$  for all points  $x \in \mathbb{R}$  such that  $F$  is continuous at  $x$ .

(iii) We have  $\varphi_n(\theta) \rightarrow \varphi(\theta)$  as  $n \rightarrow \infty$  for all  $\theta \in \mathbb{R}$ .

*Proof.* We first show the equivalence of (i) and (ii) by proving both implications “(i)  $\Rightarrow$  (ii)” and “(ii)  $\Rightarrow$  (i)”.

Then we show that “(i)  $\Rightarrow$  (iii)”. Finally, we show that “(iii)  $\Rightarrow$  (ii)”, making use of the previously established implications. The equivalence of all three conditions then follows.

*proof of “(i)  $\Rightarrow$  (ii)”:* Assume (i), i.e., that for all bounded continuous functions  $f: \mathbb{R} \rightarrow \mathbb{R}$  we have

$$\int_{\mathbb{R}} f(x) \, d\nu_n(x) \xrightarrow{n \rightarrow \infty} \int_{\mathbb{R}} f(x) \, d\nu(x).$$

Let  $x_0 \in \mathbb{R}$  be a continuity point of  $F$ . Fix  $\varepsilon > 0$ . Then by continuity of  $F$  at  $x_0$ , for some  $\delta > 0$  we have

$$F(x_0 - \delta) > F(x_0) - \varepsilon \quad \text{and} \quad F(x_0 + \delta) < F(x_0) + \varepsilon. \tag{F.4}$$

Now define two piecewise linear functions  $f^-, f^+ : \mathbb{R} \rightarrow \mathbb{R}$  by

$$f^-(x) = \begin{cases} 1 & \text{if } x \leq x_0 - \delta \\ \frac{x_0 - x}{\delta} & \text{if } x_0 - \delta < x < x_0, \\ 0 & \text{if } x_0 \leq x \end{cases}, \quad f^+(x) = \begin{cases} 1 & \text{if } x \leq x_0 \\ \frac{x_0 + \delta - x}{\delta} & \text{if } x_0 < x < x_0 + \delta. \\ 0 & \text{if } x_0 + \delta \leq x \end{cases}$$

These functions are chosen so that we have the pointwise inequalities

$$\mathbb{I}_{(-\infty, x_0 - \delta]}(x) \leq f^-(x) \leq \mathbb{I}_{(-\infty, x_0]}(x) \leq f^+(x) \leq \mathbb{I}_{(-\infty, x_0 + \delta]}(x). \quad (\text{F.5})$$

If we integrate (F.5) over  $x$  with respect to the measure  $\nu$ , the indicators integrate to values of the cumulative distribution function  $F$ , and by monotonicity of integral, we get

$$F(x_0 - \delta) \leq \int f^- d\nu \leq F(x_0) \leq \int f^+ d\nu \leq F(x_0 + \delta).$$

Combining with inequalities (F.4) and rearranging terms, we derive the estimates

$$\int f^+ d\nu - \varepsilon < F(x_0) < \int f^- d\nu + \varepsilon.$$

Similarly we can integrate (F.5) with respect to  $\nu_n$  for any  $n \in \mathbb{N}$  to get

$$\int f^- d\nu_n \leq F_n(x_0) \leq \int f^+ d\nu_n.$$

To estimate  $F_n(x_0) - F(x_0)$ , we can combine these with the previous inequalities and rearrange to the form

$$\int f^- d\nu_n - \int f^- d\nu - \varepsilon < F_n(x_0) - F(x_0) < \int f^+ d\nu_n - \int f^+ d\nu + \varepsilon.$$

Now since the functions  $f^-$  and  $f^+$  are bounded and continuous, by our assumption (i) there exists some  $N$  such that for  $n \geq N$  we have

$$\left| \int f^- d\nu_n - \int f^- d\nu \right| < \varepsilon \quad \text{and} \quad \left| \int f^+ d\nu_n - \int f^+ d\nu \right| < \varepsilon.$$

For  $n \geq N$ , our estimate on  $F_n(x_0) - F(x_0)$  therefore yields

$$-2\varepsilon < F_n(x_0) - F(x_0) < +2\varepsilon.$$

Since  $\varepsilon > 0$  was arbitrary, we can conclude  $\lim_{n \rightarrow \infty} F_n(x_0) = F(x_0)$ , which establishes (ii).

*proof of “(ii)  $\Rightarrow$  (i)”:* Assume (ii), i.e., that  $\lim_{n \rightarrow \infty} F_n(x) = F(x)$  for all  $x \in D$ , where  $D \subset \mathbb{R}$  is the set of continuity points of  $F$ . The increasing function  $F$  can only have countably many points of discontinuity, so the complement  $\mathbb{R} \setminus D$  is countable, and thus  $D \subset \mathbb{R}$  is dense.<sup>1</sup>

Let  $\varepsilon > 0$ . Since  $F(x) \downarrow 0$  as  $x \downarrow -\infty$  and  $F(x) \uparrow 1$  as  $x \uparrow +\infty$  (Proposition II.30), and since  $D \subset \mathbb{R}$  is dense, we can choose  $a, b \in D$ ,  $a < b$ , such that  $F(b) - F(a) > 1 - \varepsilon$ . Moreover, since  $\lim_{n \rightarrow \infty} F_n(a) = F(a)$  and  $\lim_{n \rightarrow \infty} F_n(b) = F(b)$ , there exists some  $N_1$  such that we have

$$F_n(b) - F_n(a) > 1 - 2\varepsilon \quad \text{for all } n \geq N_1.$$

Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be bounded and continuous. On the compact interval  $[a, b] \subset \mathbb{R}$ , the function  $f$  is uniformly continuous, so for some  $\delta > 0$  we have  $|f(x) - f(y)| < \varepsilon$  whenever  $|x - y| < \delta$  and  $x, y \in [a, b]$ . Now choose points  $a = c_0 < c_1 < \dots < c_{k-1} < c_k = b$  such that  $c_j - c_{j-1} < \delta$  and  $c_j \in D$  for all  $j = 1, \dots, k$ . Then for any  $j = 1, \dots, k$ , we get

$$|f(x) - f(c_j)| < \varepsilon \quad \text{for } x \in [c_{j-1}, c_j].$$

Define the simple function  $h: \mathbb{R} \rightarrow \mathbb{R}$  by

$$h(x) = \sum_{j=1}^k f(c_j) \mathbb{I}_{(c_{j-1}, c_j]}(x)$$

The above estimate shows that  $|f(x) - h(x)| < \varepsilon$  for all  $x \in [a, b]$ . By boundedness of  $f$ , there exists a constant  $K > 0$  such that  $|f(x)| \leq K$  for all  $x \in \mathbb{R}$ . Since  $h$  vanishes outside  $(a, b]$ , the triangle inequality for integral with respect to  $\nu_n$  gives

$$\left| \int_{\mathbb{R}} f d\nu_n - \int_{\mathbb{R}} h d\nu_n \right| \leq \underbrace{\int_{(a,b]} |f - h| d\nu_n}_{\leq \varepsilon} + \underbrace{\int_{\mathbb{R} \setminus (a,b]} |f| d\nu_n}_{\leq K \nu_n[\mathbb{R} \setminus (a,b)]}.$$

<sup>1</sup>In fact, the proof of this implication only relies on having pointwise convergence of the cumulative distribution functions in some dense set.

When  $n \geq N_1$ , we have  $\nu_n[\mathbb{R} \setminus (a, b)] = 1 - \nu_n[(a, b)] = 1 - (F_n(b) - F_n(a)) < 2\varepsilon$ , and thus the triangle inequality implies

$$\left| \int_{\mathbb{R}} f \, d\nu_n - \int_{\mathbb{R}} h \, d\nu_n \right| \leq \varepsilon + K 2\varepsilon = (1 + 2K)\varepsilon.$$

Similarly, integrating now with respect to  $\nu$  instead, one shows that

$$\left| \int_{\mathbb{R}} f \, d\nu - \int_{\mathbb{R}} h \, d\nu \right| \leq (1 + K)\varepsilon.$$

It remains to consider the integrals of the function  $h$  with respect to both  $\nu_n$  and  $\nu$ . These integrals are expressible in terms of the cumulative distribution functions,

$$\begin{aligned} \int_{\mathbb{R}} h \, d\nu_n &= \sum_{j=1}^k f(c_j) \nu_n[(c_{j-1}, c_j]] \\ &= \sum_{j=1}^k f(c_j) (F_n(c_j) - F_n(c_{j-1})) \end{aligned}$$

and similarly

$$\int_{\mathbb{R}} h \, d\nu = \sum_{j=1}^k f(c_j) (F(c_j) - F(c_{j-1})).$$

The difference of the integrals of  $h$  with respect to these two can therefore be estimated as

$$\begin{aligned} \left| \int_{\mathbb{R}} h \, d\nu - \int_{\mathbb{R}} h \, d\nu_n \right| &= \left| \sum_{j=1}^k f(c_j) (F(c_j) - F_n(c_j) - F(c_{j-1}) + F_n(c_{j-1})) \right| \\ &\leq \sum_{j=1}^k |f(c_j)| \left( |F(c_j) - F_n(c_j)| + |F(c_{j-1}) - F_n(c_{j-1})| \right) \\ &\leq 2kK \max_{j=1, \dots, k} |F(c_j) - F_n(c_j)|. \end{aligned}$$

By our assumption (ii), we have  $\lim_{n \rightarrow \infty} F_n(c_j) = F(c_j)$  for each  $j = 1, \dots, k$ , so there exists  $N_2$  such that for  $n \geq N_2$  we have  $\max_{j=1, \dots, k} |F(c_j) - F_n(c_j)| < \frac{\varepsilon}{k}$ , and thus

$$\left| \int_{\mathbb{R}} h \, d\nu - \int_{\mathbb{R}} h \, d\nu_n \right| \leq 2K\varepsilon.$$

Combining the estimates we have obtained, for  $n \geq \max(N_1, N_2)$ , we have

$$\begin{aligned} &\left| \int_{\mathbb{R}} f \, d\nu - \int_{\mathbb{R}} f \, d\nu_n \right| \\ &\leq \underbrace{\left| \int_{\mathbb{R}} f \, d\nu - \int_{\mathbb{R}} h \, d\nu \right|}_{\leq (1+K)\varepsilon} + \underbrace{\left| \int_{\mathbb{R}} h \, d\nu - \int_{\mathbb{R}} h \, d\nu_n \right|}_{\leq 2K\varepsilon} + \underbrace{\left| \int_{\mathbb{R}} h \, d\nu_n - \int_{\mathbb{R}} f \, d\nu_n \right|}_{\leq (1+2K)\varepsilon} \\ &\leq (2 + 5K)\varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  was arbitrary, this shows that  $\int f \, d\nu_n \rightarrow \int f \, d\nu$  as  $n \rightarrow \infty$ , so we have established (i).

*proof of “(i)  $\Rightarrow$  (iii)”:* Assume that we have  $\mathbb{E}[f(X_n)] \rightarrow \mathbb{E}[f(X)]$  as  $n \rightarrow \infty$  for all bounded continuous functions  $f: \mathbb{R} \rightarrow \mathbb{R}$ . In particular, for any fixed  $\theta \in \mathbb{R}$  the functions  $f(x) = \cos(\theta x)$  and  $f(x) = \sin(\theta x)$  are bounded and continuous, so we get

$$\begin{aligned} \varphi_n(\theta) &= \mathbb{E} \left[ \exp(i\theta X_n) \right] = \mathbb{E} \left[ \cos(\theta X_n) \right] + i \mathbb{E} \left[ \sin(\theta X_n) \right] \\ &\xrightarrow{n \rightarrow \infty} \mathbb{E} \left[ \cos(\theta X) \right] + i \mathbb{E} \left[ \sin(\theta X) \right] = \varphi(\theta). \end{aligned}$$

*proof of “(iii)  $\Rightarrow$  (ii)”:* Assume (iii), i.e., for all  $\theta \in \mathbb{R}$  we have  $\varphi_n(\theta) \rightarrow \varphi(\theta)$  as  $n \rightarrow \infty$ . Let us denote by  $D \subset \mathbb{R}$  the set of continuity points of  $F$ . Our proof of (ii) consists of proving two parts:

- (1) *Existence of subsequential limits:* Given any subsequence indexed by  $(n_k)_{k \in \mathbb{N}}$ , we can find a further subsequence indexed by  $(n_{k_\ell})_{\ell \in \mathbb{N}}$  and a cumulative distribution function  $\tilde{F}$  such that  $F_{n_{k_\ell}}(x) \rightarrow \tilde{F}(x)$  as  $\ell \rightarrow \infty$  at all continuity points  $x$  of  $\tilde{F}$ .
- (2) *Uniqueness of subsequential limits:* If a cumulative distribution function  $\tilde{F}$  is the limit of  $F_{n_k}$  along some subsequence, pointwise in the set of continuity points of  $\tilde{F}$ , then we must have  $\tilde{F} = F$ .

From these two together we derive that  $F_n(x) \rightarrow F(x)$  for all  $x \in D$ . Indeed, suppose by contradiction that there exists  $x \in D$  such that  $F_n(x)$  does not tend to  $F(x)$ . Then there exists some  $\varepsilon > 0$  and a subsequence  $(n_k)_{k \in \mathbb{N}}$  such that  $|F_{n_k}(x) - F(x)| \geq \varepsilon$  for all  $k \in \mathbb{N}$ . Applying (1) to this subsequence, we find a further subsequence indexed by  $(n_{k_\ell})_{\ell \in \mathbb{N}}$  and  $\tilde{F}$  such that  $F_{n_{k_\ell}} \rightarrow \tilde{F}$ , pointwise in the set of continuity points of  $\tilde{F}$ . Then by (2) we should have  $\tilde{F} = F$ , but this is a contradiction, since  $|F_{n_{k_\ell}}(x) - F(x)| \geq \varepsilon$  holds in the subsequence. This shows that (1) and (2) indeed imply (ii).

Proving uniqueness of subsequential limits (2) is straightforward using the already established implications “(ii)  $\Rightarrow$  (i)” and “(i)  $\Rightarrow$  (iii)”. Indeed, suppose that  $F_{n_k}$  converges to  $\tilde{F}$ , pointwise in the set of continuity points of  $\tilde{F}$ . By the already established “(ii)  $\Rightarrow$  (i)”, this implies the convergence  $\int f d\nu_{n_k} \rightarrow \int f d\tilde{\nu}$  for all bounded continuous  $f$ , where the probability measure  $\tilde{\nu}$  has  $\tilde{F}$  as its cumulative distribution function. This, in turn, by the already established “(ii)  $\Rightarrow$  (i)”, implies that  $\varphi_{n_k}(\theta) \rightarrow \tilde{\varphi}(\theta)$  for all  $\theta \in \mathbb{R}$ , where  $\tilde{\varphi}$  is the characteristic function of  $\tilde{\nu}$ . But since the entire sequence  $\varphi_1, \varphi_2, \dots$  of characteristic functions has limit  $\varphi$ , the limit of the subsequence must be the same,  $\tilde{\varphi} = \varphi$ . Therefore we also get  $\tilde{\nu} = \nu$  (by Lévy’s inversion theorem) and  $\tilde{F} = F$ , and the uniqueness part (2) follows.

It remains to prove the existence of subsequential limits (1). The basic idea is to use Cantor’s diagonal extraction. In order to keep the notation simpler, let us show how to extract a convergent subsequence from the entire sequence  $(F_n)_{n \in \mathbb{N}}$  — the same argument works for extracting a convergent subsequence from any given subsequence. The set  $\mathbb{Q}$  of rational numbers is countable, so let  $\mathbb{Q} = \{q^{(1)}, q^{(2)}, \dots\}$  be an enumeration of it. Observe that for any rational point  $q^{(j)} \in \mathbb{Q}$ , the sequence  $(F_n(q^{(j)}))_{n \in \mathbb{N}}$  of values at this point is a bounded sequence:  $0 \leq F_n(q^{(j)}) \leq 1$  for all  $n \in \mathbb{N}$ . In particular, considering  $q^{(1)} \in \mathbb{Q}$  first, we can find a subsequence indexed by  $(n_k^{(1)})_{k \in \mathbb{N}}$  such that

$$\lim_{k \rightarrow \infty} F_{n_k^{(1)}}(q^{(1)})$$

exists. Then considering  $q^{(2)} \in \mathbb{Q}$ , we can find a subsequence of this already chosen subsequence, indexed now by  $(n_k^{(2)})_{k \in \mathbb{N}}$  such that

$$\lim_{k \rightarrow \infty} F_{n_k^{(2)}}(q^{(2)})$$

exists, and since this is a subsequence of the previous one, also

$$\lim_{k \rightarrow \infty} F_{n_k^{(2)}}(q^{(1)}) = \lim_{k \rightarrow \infty} F_{n_k^{(1)}}(q^{(1)}).$$

Continuing inductively, we find subsequences  $(n_k^{(\ell)})_{k \in \mathbb{N}}$ ,  $\ell = 1, 2, 3, \dots$ , such that

$$\text{the limit } \lim_{k \rightarrow \infty} F_{n_k^{(\ell)}}(q^{(j)}) \text{ exists for all } j \leq \ell.$$

Then if we consider the subsequence of diagonally selected indices  $n_\ell := n_\ell^{(\ell)}$ , we get that

$$\text{the limit } \lim_{\ell \rightarrow \infty} F_{n_\ell}(q^{(j)}) \text{ exists for all } j \in \mathbb{N},$$

since apart from finitely many first members, the index sequence  $(n_\ell)_{\ell \in \mathbb{N}}$  is a subsequence of  $(n_k^{(j)})_{k \in \mathbb{N}}$ . The limits

$$G(q^{(j)}) := \lim_{\ell \rightarrow \infty} F_{n_\ell}(q^{(j)})$$

define a function on the set of rational numbers

$$G: \mathbb{Q} \rightarrow [0, 1].$$

As the limit of the increasing functions  $F_{n_\ell}$  (restricted to  $\mathbb{Q}$ ), this function is also increasing:

$$G(q) \leq G(q') \quad \text{when } q, q' \in \mathbb{Q} \text{ and } q < q'.$$

We now claim that  $G(q) \downarrow 0$  as  $q \downarrow -\infty$  and  $G(q) \uparrow 1$  as  $q \uparrow -\infty$  — this is the place where the assumption (iii) is used. Let  $\varepsilon > 0$ . Then, since  $\varphi(0) = 1$  and  $\theta \mapsto \varphi(\theta)$  is continuous at  $\theta = 0$  (by Proposition XII.6), there exists  $\delta > 0$  such that

$$|1 - \varphi(\theta)| < \varepsilon \quad \text{when } |\theta| < \delta.$$

In particular we have

$$\frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi(\theta)) \, d\theta < \varepsilon.$$

By assumption (iii) we have  $\varphi_n(\theta) \rightarrow \varphi(\theta)$  for all  $\theta$ . Since we have  $|1 - \varphi_n(\theta)| \leq 1 + |\varphi_n(\theta)| \leq 2$ , the Bounded convergence theorem (Corollary VII.21) implies

$$\frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi_n(\theta)) \, d\theta \xrightarrow{n \rightarrow \infty} \frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi(\theta)) \, d\theta,$$

and we can conclude that there exists an  $N_\varepsilon$  such that for  $n \geq N_\varepsilon$  we have

$$\frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi_n(\theta)) \, d\theta < 2\varepsilon.$$

In view of Lemma F.3 this gives for  $n \geq N_\varepsilon$

$$\nu_n[\mathbb{R} \setminus [-R, R]] < 2\varepsilon,$$

where  $R = 2/\delta$ . For the corresponding cumulative distribution functions this implies

$$\begin{aligned} F_n(x) &< 2\varepsilon && \text{for } x < -R \\ F_n(x) &> 1 - 2\varepsilon && \text{for } x > -R. \end{aligned}$$

The inequalities inherited by the subsequential limit  $G(q) = \lim_{\ell \rightarrow \infty} F_{n_\ell}(q)$  are then

$$\begin{aligned} G(q) &\leq 2\varepsilon && \text{for } q < -R \\ G(q) &\geq 1 - 2\varepsilon && \text{for } q > -R. \end{aligned}$$

Since  $\varepsilon > 0$  was arbitrary, this shows that  $G(q) \downarrow 0$  as  $q \downarrow -\infty$  and  $G(q) \uparrow 1$  as  $q \uparrow -\infty$ .

The function  $G$  is not yet the cumulative distribution function we are looking for: it is only defined on the rational numbers, and it is not necessarily right-continuous. But we can use it to define

$$\tilde{F}(x) = \inf_{\substack{q > x \\ q \in \mathbb{Q}}} G(q)$$

(this in fact defines the smallest right continuous function above  $G$ ). From the definition it is clear that if  $q', q'' \in \mathbb{Q}$  and  $q' < x < q''$ , then

$$G(q') \leq \tilde{F}(x) \leq G(q'').$$

We leave it to the reader to check that this  $\tilde{F}$  is right continuous, increasing, and has limits 0 and 1 at  $-\infty$  and  $+\infty$ , respectively. Therefore  $\tilde{F}$  is a cumulative distribution function of some probability measure  $\tilde{\nu}$  on  $\mathbb{R}$ .

The only remaining claim is that the subsequence  $(F_{n_\ell})_{\ell \in \mathbb{N}}$  converges to  $\tilde{F}$  pointwise in the set  $\tilde{D} \subset \mathbb{R}$  of continuity points of  $\tilde{F}$ . Consider  $x \in \tilde{D}$  and let  $\varepsilon > 0$ . Then by continuity of  $\tilde{F}$  at  $x$  we have for some  $\delta > 0$

$$\tilde{F}(x - \delta) > \tilde{F}(x) - \varepsilon \quad \text{and} \quad \tilde{F}(x + \delta) < \tilde{F}(x) + \varepsilon.$$

Choose  $q' \in \mathbb{Q} \cap (x - \delta, x)$ . Then by the definition  $G(q') = \lim_{\ell \rightarrow \infty} F_{n_\ell}(q')$  there exists some  $L_\varepsilon$  such that for all  $\ell \geq L_\varepsilon$  we have

$$F_{n_\ell}(q') > G(q') - \varepsilon,$$

which yields

$$F_{n_\ell}(x) \geq F_{n_\ell}(q') > G(q') - \varepsilon \geq \tilde{F}(x - \delta) - \varepsilon > \tilde{F}(x) - 2\varepsilon.$$

Since  $\varepsilon > 0$  was arbitrary, we get

$$\liminf_{\ell} F_{n_{\ell}}(x) \geq \tilde{F}(x).$$

By choosing  $q'' \in \mathbb{Q} \cap (x, x + \delta)$ , one can similarly argue that

$$\limsup_{\ell} F_{n_{\ell}}(x) \leq \tilde{F}(x),$$

and these two combined imply

$$\lim_{\ell \rightarrow \infty} F_{n_{\ell}}(x) = \tilde{F}(x).$$

This establishes the existence of subsequential limits (1), and finishes the proof.  $\square$



## Index

- $\pi$ -system, *see also* pi-system
- $\sigma$ -algebra, *see also* sigma algebra
- $p$ -integrable, 78
  
- addition of functions, *see also* pointwise sum
  - of functions
- additivity
  - countable, 8
  - finite, 13
- almost sure, *see also* almost surely
  - event, 9
  - limit, 108
- almost surely, 9
- axiom of choice, 129
  
- ball, 141
- binary sequence, 135
- binomial distribution, 13
- Borel sigma algebra, 5
- Borel-measurable function, 22
  
- cardinality, 131
- Cartesian product, 81, 128
- Cauchy sequence
  - of square integrable random variables, 155
- Cauchy-Schwarz inequality, 97
- central limit theorem, 107
- characteristic function, 120, 159
- Chebyshev's inequality, 111
- closed set, 141
- closed subspace, 156
- complement, 127
- complex valued random variable, 118
- conditioned probability measure, 11
- continuous distribution, 74, 95
- continuous function, 141
- convergence
  - almost surely, 108
  - in  $\mathcal{L}^1$ , 115, 154
  - in probability, 108
  - of a sequence of numbers, 138
- convergence in distribution, 117, 125
- convergence in law, 117, 125
- convergent sequence, 141
- countable set, 132
- countably infinite, 132
- counting measure, 9, 90
  
- covariance, 98
- cumulative distribution function, 18, 103
  
- d-system, 143
- De Morgan's laws, 128
- decreasing sequence
  - of numbers, 139
  - of sets, 129
- density function, 74
  - joint, 95
  - marginal, 96
- difference
  - of sets, 127
- Dirac measure, 96
- discrete metric, 141
- disjoint, 127
- disjoint union, 127
- distribution
  - joint, 93
  - of a random variable, 23, 72
- dominated convergence theorem, 68
  
- Euclidean norm, 140
- ev., *see also* eventually
- event, ix, 1, 2
- eventually, 44
- expected value, ix, 55, 71
- exponential distribution, 74
- extended real line, 27, 137
  
- Fatou's lemma, 67
  - reverse, 67
- finite measure, 9
- finite measure space, 9
  
- gaussian distribution, 74
- geometric distribution, 13
- Goddess of Chance, ix
  
- homeomorphism, 142
  
- i.o., *see also* infinitely often
- image of a set under function, 128
- improper Riemann integral, 70
- increasing sequence
  - of numbers, 139
  - of sets, 129
- independence, 41

- of  $\sigma$ -algebras, 39
  - of events, 40
  - of random variables, 40
- indicator random variable, 24
- inf, *see also* infimum
- infimum, 137, 138
- infinitely often, 44
- inner product
  - of square integrable random variables, 153
- integrable
  - complex random variable, 118
- integrable function, 63
- integrable random variable, 63, 78
- integral, 55
  - of a non-negative function, 59
  - of a simple function, 57
  - of an integrable function, 63
  - over subset, 69
- intersection of sets, 127
  
- Jensen's inequality, 79
- joint density, 95
- joint law, 93
  
- Kolmogorov's strong law of large numbers, 116
  
- law, *see also* distribution
  - of a random variable, 23, 72
- law of large numbers, 107
- Lebesgue integral, 70
- Lebesgue measure, 10, 90
  - $d$ -dimensional, 10, 90
- liminf
  - of events, *see also* eventually
  - of sequence of numbers, 139
  - of sequence of sets, 130
- limit, 141
  - of a decreasing sequence of sets, 129
  - of a sequence of numbers, 138
  - of an increasing sequence of sets, 129
- limsup
  - of events, *see also* infinitely often
  - of sequence of numbers, 139
  - of sequence of sets, 130
- linearity
  - of integral, 56
- lower limit
  - of sequence of numbers, 139
  - of sequence of sets, 130
  
- marginal density, *see also* density function,
  - marginal, 96
- Markov process, 94
- Markov's inequality, 111
- MCT, *see also* Monotone convergence theorem
- measurable
  - set, 7
    - space, 7, 8
- measurable function, 21, 22
- measurable space, 7
- measure, 8
- measure space, 9
- moment, 78
- monotone class, 38, 82, 143
- Monotone class theorem, 82
- monotone convergence
  - of integrals, 62
  - of measures, 13
- Monotone convergence theorem, 61, 147
- monotone sequence
  - of numbers, 139
  - of sets, 130
- monotonicity
  - of integral, 56
  - of measures, 13
- multiplication of functions, *see also*
  - pointwise product of functions
  
- negative part
  - of a function, 63
- norm
  - of square integrable random variable, 153
  
- open set, 141
- orthogonality
  - of square integrable random variables, 153, 156
- outcome, ix
  
- pi-system ( $\pi$ -system), 17
- pointwise product of functions, 27
- pointwise scalar multiple of functions, 27
- pointwise sum of functions, 27
- Poisson distribution, 12
- positive part
  - of a function, 63
- power set, 129
- preimage of a set under function, 128
- probability, ix
- probability density, 95
- probability mass function, 12
- probability measure, 9
- probability space, 9
- product measure, 86
- product of functions, *see also* pointwise
  - product of functions
- product sigma algebra, 83
  
- random number, 22
- random variable, ix, 21, 22
- random walk, 54
- Riemann integral, 70
  
- sample space, ix
- scalar multiplication of functions, *see also*
  - pointwise scalar multiple of functions

- sequence
  - decreasing, 129, 139
  - increasing, 129, 139
  - monotone, 130, 139
  - of sets, 129
- sigma algebra
  - generated by collection of events, 3
  - generated by random variables, 36
- sigma algebra ( $\sigma$ -algebra), 2
- sigma finite, 90
- simple function, 30, 57
- square integrable random variable, 97
- staircase function, 31
- standard machine, 56, 146
- standard normal distribution, 74
- strong law of large numbers, 110
- subadditivity
  - countable, 14
  - finite, 13
- sum of functions, *see also* pointwise sum of functions
- sup, *see also* supremum
- supremum, 137, 138
- supremum norm, 141
- sure event, 9
  
- tail event, 48
- tail sigma algebra, 48
- tend
  - seeconverge, 138
- total mass
  - of a measure, 9
- truncated measure, 11
  
- uncountable set, 135
- uniform norm, 141
- uniform probability measure
  - continuous, 11
  - discrete, 10
  - on a finite set, 10
- union of sets, 127
- upper limit
  - of sequence of numbers, 139
  - of sequence of sets, 130
  
- variance, 98
  
- weak convergence, 117
- weak law of large numbers, 110
- Weierstrass' approximation theorem, 112



## References

### Bibliography

- [Dur10] Richard Durrett. *Probability: Theory and Examples*. Cambridge University Press, 4th edition, 2010.
- [JP04] Jean Jacod and Philip Protter. *Probability Essentials*. Cambridge Univ. Press, 2nd edition, 2004.
- [Kin16] Juha Kinnunen. *Measure and integral*. Aalto University, 2016. [https://math.aalto.fi/~jkkinnun/files/measure\\_and\\_integral.pdf](https://math.aalto.fi/~jkkinnun/files/measure_and_integral.pdf).
- [Wil91] David Williams. *Probability with Martingales*. Cambridge Univ. Press, 1991.