

Advanced Probability Theory

Thierry Lévy *

2020-21

Contents

1	A reminder of elementary probability theory	2
1.1	Probability spaces	2
1.2	Independence	4
1.3	Random variables	5
1.4	Integration of random variables	7
1.5	Convergence of random variables	9
2	Conditional expectation	11
2.1	Definition and fundamental examples	11
2.2	Uniqueness and existence	15
2.3	Main properties	19
3	Martingales	22
3.1	Martingales and conserved quantities	22
3.2	Definition	24
3.3	The case where the filtration is generated by partitions	26
3.4	Stochastic integration	29
3.5	Almost sure convergence	31
3.6	Branching processes	36
3.7	Convergence in L^1	39
3.8	Stopping times	41
3.9	Convergence in L^p for $p > 1$	44
3.10	Square-integrable martingales	48
3.11	Uniform integrability	51
3.12	Backward martingales	55
4	Markov chains	58
4.1	Introduction	58
4.2	First definition and first properties	59
4.3	Construction of Markov chains	64
4.4	The Markov property	70
4.5	Recurrent and transient states	73
4.6	Invariant measures	80
4.7	Brief summary	86
4.8	The ergodic theorem	88

*Sorbonne Université – LPSM – 4, place Jussieu – F-75005 Paris
thierry.levy@sorbonne-universite.fr

1 A reminder of elementary probability theory

1.1 Probability spaces

In order to describe an experiment from the point of view of probability theory, one considers a *probability space*, that is, a triple $(\Omega, \mathcal{A}, \mathbb{P})$, where

- Ω is a set: the set of all possible outcomes of the experiment,
- \mathcal{A} is a σ -field on Ω : the set of subsets of Ω that it is possible to consider, to name, to think of,
- \mathbb{P} is a probability measure on (Ω, \mathcal{A}) , which to each $A \in \mathcal{A}$ associates the probability $\mathbb{P}(A)$ that the outcome of the experiment belongs to A .

By definition, the fact that \mathcal{A} is a σ -field means that \mathcal{A} is a subset of the set $\mathcal{P}(\Omega)$ of all subsets of Ω (in other words, the elements of \mathcal{A} are subsets of Ω) such that

- \mathcal{A} contains the empty set \emptyset ,
- for all $A \in \mathcal{A}$, the set $\Omega \setminus A$ also belongs to \mathcal{A} ,
- for all sequence $(A_n)_{n \geq 0}$ of elements of \mathcal{A} , the union $\bigcup_{n \geq 0} A_n$ also belongs to \mathcal{A} .

The fact that \mathbb{P} is a probability measure on (Ω, \mathcal{A}) means that \mathbb{P} is a function $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ which satisfies

- $\mathbb{P}(\Omega) = 1$,
- for all sequence $(A_n)_{n \geq 0}$ of elements of \mathcal{A} which are pairwise disjoint,

$$\mathbb{P} \left(\bigcup_{n \geq 0} A_n \right) = \sum_{n \geq 0} \mathbb{P}(A_n).$$

Exercise 1.1 Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Prove that $\mathbb{P}(\emptyset) = 0$.

The first fundamental example of probability space is $(\Omega, \mathcal{P}(\Omega), \mathbb{P})$, where Ω is a finite set, endowed with the σ -field of all subsets of Ω , and \mathbb{P} is the uniform measure defined, for all $A \subset \Omega$, by $\mathbb{P}(A) = \frac{|A|}{|\Omega|}$. Here, $|A|$ denotes the number of elements of A .

Another fundamental example is $([0, 1], \mathcal{B}_{[0,1]}, \lambda)$, where $\mathcal{B}_{[0,1]}$ is the σ -field of Borel subsets of the real interval $[0, 1]$, that is, the smallest σ -field on $[0, 1]$ which contains all open subsets of $[0, 1]$, and λ is the Lebesgue measure on $([0, 1], \mathcal{B}_{[0,1]})$, the unique measure on this measurable space which for all a, b such that $0 \leq a \leq b \leq 1$ satisfies $\lambda([a, b]) = b - a$.

Let us agree on the following question of terminology: some authors call a set countable if it has the cardinality of \mathbb{N} . For us, a countable set is a set which has the cardinality of a subset of \mathbb{N} . The difference is that for us, a finite set is countable.

Exercise 1.2 Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. An element $\omega \in \Omega$ is called an atom of \mathbb{P} if $\{\omega\} \in \mathcal{A}$ and $\mathbb{P}(\{\omega\}) > 0$. Can you give an upper bound to the number of atoms of \mathbb{P} of mass at least $\frac{1}{10}$? Prove that the set of atoms of \mathbb{P} is a countable subset of Ω . Prove that $\{\omega \in \Omega : \omega \text{ is not an atom of } \mathbb{P}\}$ belongs to \mathcal{A} .

Exercise 1.3 How many σ -fields does there exist on the set $\{1, 2, 3\}$? Define, for all $n \geq 0$, B_n as the number of σ -fields on a finite set with n elements. Prove that for all $n \geq 0$, the equality $B_{n+1} = \sum_{k=0}^n \binom{n}{k} B_k$ holds. The number B_n is called the n -th Bell number.

The most important technical property of a probability measure is the following. Consider a sequence $(A_n)_{n \geq 0}$ of elements of \mathcal{A} which is non-decreasing in the sense that $A_0 \subset A_1 \subset A_2 \subset \dots$. Then

$$\mathbb{P} \left(\bigcup_{n \geq 0} A_n \right) = \sup \{ \mathbb{P}(A_n) : n \geq 0 \} = \lim_{n \rightarrow \infty} \mathbb{P}(A_n). \quad (1)$$

Exercise 1.4 Why does the limit on the right-hand side exist? Prove the two equalities.

Exercise 1.5 Prove that for any non-increasing sequence $(A_n)_{n \geq 0}$ of elements of \mathcal{A} one has $\mathbb{P}(\bigcap_{n \geq 0} A_n) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n)$? Is this true on any measure space? Is (1) true on any measure space?

Another important property is this one: for all sequence $(A_n)_{n \geq 0}$ of elements of \mathcal{A} , the inequality

$$\mathbb{P} \left(\bigcup_{n \geq 0} A_n \right) \leq \sum_{n \geq 0} \mathbb{P}(A_n)$$

holds.

Yet another important result is formulated in the next exercise.

Exercise 1.6 (Borel-Cantelli lemma) Let $(A_n)_{n \geq 0}$ be a sequence of elements of \mathcal{A} . How would you describe in words the set $S = \bigcap_{p \geq 0} \bigcup_{n \geq p} A_n$? Prove that S belongs to \mathcal{A} . Assume now that $\sum_{n \geq 0} \mathbb{P}(A_n) < \infty$ and compute $\mathbb{P}(S)$.

Exercise 1.7 Prove that for any family $(A_i)_{i \in I}$ of elements of \mathcal{A} which are pairwise disjoint and such that $\bigcup_{i \in I} A_i$ belongs to \mathcal{A} , the inequality

$$\mathbb{P} \left(\bigcup_{i \in I} A_i \right) \geq \sum_{i \in I} \mathbb{P}(A_i)$$

holds, where the right-hand side is defined by

$$\sum_{i \in I} \mathbb{P}(A_i) = \sup \left\{ \sum_{i \in F} \mathbb{P}(A_i) : F \subset I, F \text{ finite} \right\}.$$

Note that the index set I may be uncountable. Can you think of an example where the inequality is strict?

Exercise 1.8 Let (Ω, \mathcal{A}) be a measurable space. Prove that a function $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ is a probability measure if and only if the following conditions hold:

- $\mathbb{P}(\emptyset) = 0$ and $\mathbb{P}(\Omega) = 1$,

- for all $A, B \in \mathcal{A}$, $\mathbb{P}(A) + \mathbb{P}(B) = \mathbb{P}(A \cup B) + \mathbb{P}(A \cap B)$,
- for all sequence $(A_n)_{n \geq 0}$ of elements of \mathcal{A} ,

$$\mathbb{P} \left(\bigcup_{p \geq 0} \bigcap_{n \geq p} A_n \right) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(A_n).$$

Is it possible to remove the assumption that $\mathbb{P}(\emptyset) = 0$ from this list ?

1.2 Independence

Let us now turn to the concept of independence. Two events A and B are independent (with respect to \mathbb{P}) if $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$.

Exercise 1.9 *What can you say about an event which is independent of itself ?*

In order to define the independence of more than two events, we need to be more careful. We say that A_1, \dots, A_n are independent if for all $k \in \{2, \dots, n\}$ and all choice of $1 \leq i_1 < \dots < i_k \leq n$, the equality $\mathbb{P}(A_{i_1} \cap \dots \cap A_{i_k}) = \mathbb{P}(A_{i_1}) \dots \mathbb{P}(A_{i_k})$ holds.

Exercise 1.10 *Find three events A, B, C such that $\mathbb{P}(A \cap B \cap C) = \mathbb{P}(A)\mathbb{P}(B)\mathbb{P}(C)$ but A and B are not independent. Can you find an example where $\mathbb{P}(A)\mathbb{P}(B)\mathbb{P}(C) > 0$?*

It turns out that the best definition is in terms of σ -fields. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. We say that n sub- σ -fields $\mathcal{B}_1, \dots, \mathcal{B}_n$ of \mathcal{A} are independent if for all choice of events $B_1 \in \mathcal{B}_1, \dots, B_n \in \mathcal{B}_n$, the equality $\mathbb{P}(B_1 \cap \dots \cap B_n) = \mathbb{P}(B_1) \dots \mathbb{P}(B_n)$ holds. We say that an arbitrary family $(\mathcal{B}_i)_{i \in I}$ of sub- σ -fields of \mathcal{A} is independent if any finite sub-family $\mathcal{B}_{i_1}, \dots, \mathcal{B}_{i_n}$ is independent.

Exercise 1.11 *Let \mathcal{C} be a subset of $\mathcal{P}(\Omega)$. Prove that there exists a unique σ -field \mathcal{A} on Ω such that $\mathcal{C} \subset \mathcal{A}$ and, for all σ -field \mathcal{B} on Ω such that $\mathcal{C} \subset \mathcal{B}$, one has $\mathcal{A} \subset \mathcal{B}$. In words, \mathcal{A} is the smallest σ -field on Ω which contains \mathcal{C} . It is called the σ -field generated by \mathcal{C} and it is denoted by $\sigma(\mathcal{C})$.*

For example, the Borel σ -field of \mathbb{R} is the σ -field generated by the class of open subsets of \mathbb{R} .

Compute $\sigma(\emptyset)$ and, for all $A \in \mathcal{A}$, $\sigma(\{A\})$. Let $\Omega = A_1 \cup \dots \cup A_n$ be a partition of Ω . This means that the events A_1, \dots, A_n are non-empty and pairwise disjoint. Describe $\sigma(\{A_1, \dots, A_n\})$. What can you say if instead of a finite partition we consider a countable partition $\Omega = \bigcup_{n \geq 0} A_n$? An arbitrary partition $\Omega = \bigcup_{i \in I} A_i$? What is the σ -field on \mathbb{R} generated by $\{\{t\} : t \in \mathbb{R}\}$?

In the following exercise, you will show that our previous two definitions of independence are consistent.

Exercise 1.12 *Let A_1, \dots, A_n be n events. Prove that it is equivalent to say that the events A_1, \dots, A_n are independent or to say that the σ -fields $\sigma(\{A_1\}), \dots, \sigma(\{A_n\})$ are independent.*

The following example warns us against an easily made mistake.

Exercise 1.13 Let us toss a coin twice. Consider the events “the first coin gives tail”, “the second coin gives tail”, “the two coins give the same result”. Prove that any two of these three events are independent, but the three together are not independent.

Exercise 1.14 For $n \geq 2$, set $\Omega = \{\omega = (\omega_1, \dots, \omega_n) \in \{-1, 1\}^n : \omega_1 \dots \omega_n = 1\}$ and consider the uniform probability \mathbb{P} on $(\Omega, \mathcal{P}(\Omega))$. For each $i \in \{1, \dots, n\}$, define $A_i = \{\omega \in \Omega : \omega_i = 1\}$. Compute, for all $k \in \{1, \dots, n\}$ and all $1 \leq i_1 < \dots < i_k \leq n$, the probability $\mathbb{P}(A_{i_1} \cap \dots \cap A_{i_k})$. What is this an example of? (“Of a silly exercise” is not the expected answer.)

1.3 Random variables

A random variable is the mathematical notion that represents a quantity whose value depends on the outcome of the experiment which is performed. Formally, a (real-valued) random variable on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ is a measurable function $X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$. Measurable means by definition that for all $B \in \mathcal{B}_{\mathbb{R}}$, the set $X^{-1}(B) = \{\omega \in \Omega : X(\omega) \in B\}$ belongs to \mathcal{A} .

For example, the function $X : [0, 1] \rightarrow \mathbb{R}$ defined by $X(t) = t^2 - 1$ is a random variable on the probability space $([0, 1], \mathcal{B}_{[0,1]}, \lambda)$. Random variables are in fact rarely defined by formulas like this one, and this example is rather atypical.

Exercise 1.15 Consider a random variable $X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$. Prove that $\{X^{-1}(B) : B \in \mathcal{B}_{\mathbb{R}}\}$ is a sub- σ -field of \mathcal{A} . It is called the σ -field generated by X and it is denoted by $\sigma(X)$.

Exercise 1.16 Let $X : \Omega \rightarrow \mathbb{R}$ be a function. For all $x \in \mathbb{R}$, consider the subset

$$\{X \leq x\} = X^{-1}((-\infty, x]) = \{\omega \in \Omega : X(\omega) \leq x\}$$

of Ω . Prove that X is a random variable if and only if for all $a \in \mathbb{R}$, one has $\{X \leq a\} \in \mathcal{A}$.

It is sometimes convenient to allow infinite values for random variables. To do so, one considers the extended real line $\overline{\mathbb{R}} = [-\infty, +\infty]$. This set is endowed with the σ -field $\mathcal{B}_{\overline{\mathbb{R}}} = \sigma(\mathcal{B}_{\mathbb{R}} \cup \{-\infty\}, \{+\infty\})$.

Exercise 1.17 Prove that a subset A of $\overline{\mathbb{R}}$ belongs to $\mathcal{B}_{\overline{\mathbb{R}}}$ if and only if $A \cap \mathbb{R}$ belongs to $\mathcal{B}_{\mathbb{R}}$.

If you are familiar with the abstract notion of a topology on a set: is $\mathcal{B}_{\overline{\mathbb{R}}}$ the Borel σ -field of a topology on $\overline{\mathbb{R}}$, that is, the σ -field generated by the class of all open subsets of $\overline{\mathbb{R}}$?

Just as for events, we say that n random variables X_1, \dots, X_n defined on the same probability space are independent if the σ -fields $\sigma(X_1), \dots, \sigma(X_n)$ are independent. More generally, we say that a family $(X_i)_{i \in I}$ of random variables is independent if every finite sub-family X_{i_1}, \dots, X_{i_n} is independent.

Exercise 1.18 Let A be an event. Prove that the indicator function $\mathbb{1}_A$ defined by $\mathbb{1}_A(\omega) = 1$ if $\omega \in A$ and $\mathbb{1}_A(\omega) = 0$ otherwise, is a random variable. Prove that it is equivalent to say that the events A_1, \dots, A_n are independent or to say that the random variables $\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_n}$ are independent.

Exercise 1.19 What can you say about a random variable which is independent of itself?

Exercise 1.20 *Is it equivalent to say that a finite family of random variables are pairwise independent and to say that they are independent ?*

Let $X : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ be a random variable. The function $\mathbb{P} \circ X^{-1} : \mathcal{B}_{\mathbb{R}} \rightarrow [0, 1]$ defined by $\mathbb{P} \circ X^{-1}(B) = \mathbb{P}(X^{-1}(B))$ is a probability measure on $(\mathbb{R}, \mathcal{B}_{\mathbb{R}})$.

Exercise 1.21 *Prove this assertion.*

The probability measure $\mathbb{P} \circ X^{-1}$ is called the distribution, or the law, of X . It is also often denoted by \mathbb{P}_X . The line

$$\mathbb{P} \circ X^{-1}(B) = \mathbb{P}_X(B) = \mathbb{P}(\{\omega \in \Omega : X(\omega) \in B\}) = \mathbb{P}(\{X \in B\}) = \mathbb{P}(X \in B)$$

gives five common notations for the same quantity. Let us emphasize that $\{X \in B\}$ is a notation for $\{\omega \in \Omega : X(\omega) \in B\}$, that is, for $X^{-1}(B)$.

The distribution of a real-valued random variable is thus a probability measure on the real line. The following rough classification of these probability measures is often used. We say that the distribution of X is discrete, or atomic, if there exists a sequence x_1, x_2, \dots of atoms of \mathbb{P}_X , finite or infinite (but in any case countable, see Exercise 1.2), such that $\mathbb{P}_X(\{x_1\}) + \mathbb{P}_X(\{x_2\}) + \dots = 1$. In other words, there exists a countable subset $C \subset \mathbb{R}$ such that $\mathbb{P}(X \in C) = 1$. The typical cases are $C = \mathbb{N}$ and $C = \mathbb{Z}$, and we then speak of integer-valued random variables, but there are also other interesting cases of discrete random variables. In this course, we will define and use many random variables with values in $\mathbb{N} \cup \{\infty\}$.

The distribution of a discrete random variable is completely described by a countable set $C = \{x_1, x_2, \dots\}$ such that $\mathbb{P}(X \in C) = 1$ and the data, for each $i \geq 1$, of the probability $p_i = \mathbb{P}(X = x_i)$. It is the framework of elementary probability theory, especially when C is finite.

A random variable X such that \mathbb{P}_X has no atom, that is, such that $\mathbb{P}(X = x) = 0$ for all $x \in \mathbb{R}$, is said to be diffuse.

Exercise 1.22 *Find a diffuse real-valued random variable. Find a real-valued random variable which is neither diffuse nor discrete.*

Among diffuse random variables, there is the practically very important class of random variables for which \mathbb{P}_X is absolutely continuous with respect to the Lebesgue measure. This means that for all Borel subset N of \mathbb{R} such that $\lambda(N) = 0$, one has $\mathbb{P}(X \in N) = 0$. In this case, since we are working with probability measures, which are σ -finite, the Radon-Nikodym theorem can be applied to produce a non-negative measurable function $f : (\mathbb{R}, \mathcal{B}_{\mathbb{R}}) \rightarrow (\mathbb{R}_+, \mathcal{B}_{\mathbb{R}_+})$ which has the property that for all Borel subset B of \mathbb{R} , the equality

$$\mathbb{P}_X(B) = \mathbb{P}(X \in B) = \int_B f \, d\lambda = \int_{\mathbb{R}} f \mathbb{1}_B \, d\lambda$$

holds. This situation is extremely convenient, because the computation of probabilities is reduced to the computation of ordinary integrals with respect to the Lebesgue measure.

Our quick zoology of probability measures on \mathbb{R} is summarised by Figure 1 below.

Exercise 1.23 Let $f : \mathbb{R} \rightarrow \mathbb{R}_+$ be a measurable function. Prove that the function $\mu : \mathcal{A} \rightarrow [0, +\infty]$ defined by $\mu(A) = \int_A f d\lambda$ is a measure on $(\mathbb{R}, \mathcal{B}_{\mathbb{R}})$. We shall write $\mu = f\lambda$. Under which condition on f is μ a probability measure ?

On $(\mathbb{R}, \mathcal{B}_{\mathbb{R}})$, consider the Lebesgue measure λ and the counting measure κ , which by definition is such that $\kappa(A)$ is the number of elements of A if A is finite, and $\kappa(A) = +\infty$ if A is infinite. Prove that λ is absolutely continuous with respect to κ . Prove that there does not exist a measurable function $f : \mathbb{R} \rightarrow \mathbb{R}_+$ such that $\lambda = f\kappa$. Is this not in contradiction with the Radon-Nikodym theorem ?

Exercise 1.24 Prove that every probability measure on $(\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ can be written in a unique way as the sum of a discrete measure and a diffuse measure.

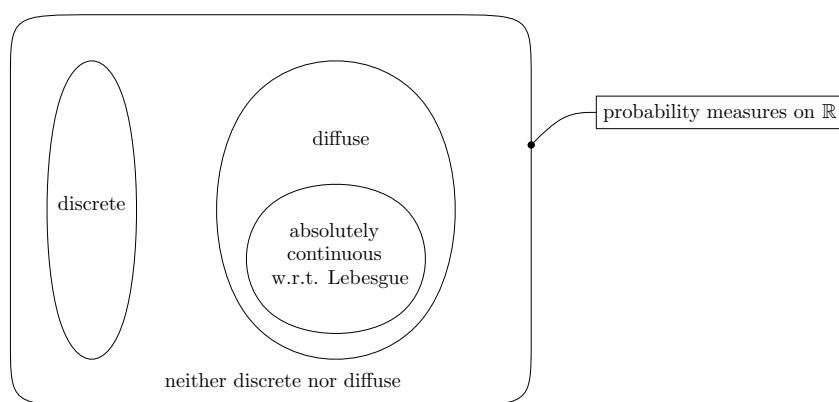


Figure 1: Do you agree with this partition ? Can you find a probability measure in each class ?

1.4 Integration of random variables

Consider a real-valued random variable X . If X admits an integral on $(\Omega, \mathcal{A}, \mathbb{P})$ in the sense of Lebesgue's integration theory, we say that X admits an expectation, which is defined by

$$\mathbb{E}[X] = \int_{\Omega} X d\mathbb{P}.$$

The following exercises help you to remember the key points of the theory of integration.

Exercise 1.25 (Definition of the integral) Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Recall the definition of the integral with respect to \mathbb{P} of a non-negative measurable function $X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}_+, \mathcal{B}_{\mathbb{R}_+})$, for example as a supremum of elementary integrals of simple functions. Consider an arbitrary measurable real-valued function $X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$. Define $X^+ = \max(X, 0)$ and $X^- = \max(-X, 0)$, so that $X = X^+ - X^-$. Explain how to define the integral of X if one at least of the two functions X^+ and X^- have a finite integral. What is the integral of X if both X^+ and X^- have infinite integrals ? When is it the case that $|X|$ has a finite integral ? (In this case, we say that X is integrable).

Let us emphasize that a non-negative random variable always admits an expectation, possibly equal to $+\infty$.

Exercise 1.26 (Convergence theorems) *Recall the statement of the monotone convergence theorem. Does the theorem also hold for a decreasing sequence of functions? Deduce Fatou's lemma from the monotone convergence theorem and then the dominated convergence theorem from Fatou's lemma. (You may find that this is not as difficult as you expect or remember. The monotone convergence theorem is the one most important convergence theorem of the theory of integration, and the other convergence theorems follow relatively easily from it.)*

A few inequalities are of crucial importance. The simplest one, but very useful, is the Markov inequality. It says that for a non-negative random variable X and a non-negative real a , one has

$$a\mathbb{P}(X \geq a) \leq \mathbb{E}[X].$$

Exercise 1.27 *Prove the Markov inequality.*

The Hölder inequality states that if X and Y are non-negative random variables and if $p, q > 1$ are two real numbers such that $\frac{1}{p} + \frac{1}{q} = 1$, then

$$\mathbb{E}[XY] \leq \mathbb{E}[X^p]^{\frac{1}{p}} \mathbb{E}[Y^q]^{\frac{1}{q}}.$$

The case where $p = q = 2$ is an instance of the Cauchy-Schwarz inequality.

One says that a random variable X admits a moment of order $p \geq 1$ if the random variable $|X|^p$ is integrable.

Exercise 1.28 *Prove that if $1 \leq p < p'$ and if a random variable X admits a moment of order p' , then it also admits a moment of order p . Find an explicit bound for $\mathbb{E}[|X|^p]$.*

Exercise 1.29 *Let X be a non-negative random variable. Prove that the function from $[1, \infty)$ to $[0, \infty]$ defined by $p \mapsto \log \mathbb{E}[X^p]$ is convex on $[1, \infty)$ (one says that the function $p \mapsto \mathbb{E}[X^p]$ is log-convex). Check that the statement and your proof of it make sense and are correct even if some or all moments of X are infinite.*

Exercise 1.30 *Let X be a non-negative random variable. With the usual agreement that $\inf \emptyset = +\infty$, define $r = \inf\{p \in (0, \infty) : \mathbb{E}[X^p] = \infty\}$. Prove by giving examples that r can be any element of $[0, \infty]$. Is it true that $r = \infty$ implies that X is bounded? When $0 < r < \infty$, is it always true, sometimes true and sometimes false, or never true that $\mathbb{E}[X^r] < \infty$?*

Another important inequality is Minkowski's inequality. It states that for all $p \geq 1$ and for any two non-negative random variables X and Y , the inequality

$$\mathbb{E}[(X + Y)^p]^{\frac{1}{p}} \leq \mathbb{E}[X^p]^{\frac{1}{p}} + \mathbb{E}[Y^p]^{\frac{1}{p}}$$

holds. In particular, the sum of two random variables which admit a moment of order p also admits a moment of order p .

This suggests to define, for all real $p \geq 1$, the set $L^p(\Omega, \mathcal{A}, \mathbb{P})$ as the set of real-valued random variables on Ω which admit a moment of order p , quotiented by the subspace formed by the random variables X such that $\mathbb{E}[|X|^p] = 0$.

Exercise 1.31 Check that $L^p(\Omega, \mathcal{A}, \mathbb{P})$ is a vector space and that the function $X \mapsto \|X\|_p = \mathbb{E}[|X|^p]^{\frac{1}{p}}$ is a norm on $L^p(\Omega, \mathcal{A}, \mathbb{P})$.

Exercise 1.32 Let X be a random variable. With the agreement that $\log 0 = -\infty$ and $\log(+\infty) = +\infty$, prove that

$$\lim_{p \rightarrow \infty} \frac{1}{p} \log \mathbb{E}[|X|^p] = \log \|X\|_\infty,$$

where the norm on the right-hand side is defined by

$$\|X\|_\infty = \inf\{M \in [0, +\infty] : \mathbb{P}(|X| \leq M) = 1\}.$$

Is it true that $\lim_{p \rightarrow \infty} \|X\|_p = \|X\|_\infty$?

It is a very important fact that the normed vector space $(L^p(\Omega, \mathcal{A}, \mathbb{P}), \|\cdot\|_p)$ is complete. In particular, $(L^2(\Omega, \mathcal{A}, \mathbb{P}), \|\cdot\|_2)$ is a Hilbert space, since the norm $\|\cdot\|_2$ is induced by the scalar product $\langle X, Y \rangle = \mathbb{E}[XY]$.

1.5 Convergence of random variables

Let us conclude this reminder of classical probability theory by defining three notions of convergence of random variables. Let $(X_n)_{n \geq 0}$ be a sequence of real-valued random variables defined on $(\Omega, \mathcal{A}, \mathbb{P})$. Let X be a random variable on the same probability space.

We say that the sequence $(X_n)_{n \geq 0}$ converges almost surely towards X if there exists an event $N \subset \Omega$ such that $\mathbb{P}(N) = 0$ and for all $\omega \notin N$, the sequence $(X_n(\omega))_{n \geq 0}$ converges to $X(\omega)$. This is the pointwise convergence outside a negligible event.

Let $p \geq 1$ be fixed. We say that the sequence $(X_n)_{n \geq 0}$ converges towards X in L^p if X admits a moment of order p and $\|X_n - X\|_p$ converges to 0. This is the norm convergence in the Banach space $L^p(\Omega, \mathcal{A}, \mathbb{P})$.

We say that the sequence $(X_n)_{n \geq 0}$ converges in probability towards X if for all $\varepsilon > 0$, the sequence of probabilities $\mathbb{P}(|X_n - X| > \varepsilon)$ converges to 0 as n tends to infinity.

The next exercises are perhaps more substantial than some of the earlier ones.

Exercise 1.33 (Modes of convergence) Prove that a sequence which converges almost surely or in L^p for some $p \geq 1$ also converges in probability towards the same limit. Prove that a sequence of random variables has at most one limit in any of the three modes of convergence which we have defined.

Prove that from a sequence which converges in probability one can extract a sub-sequence which converges almost surely.

Exercise 1.34 (L^p is complete) Let $(X_n)_{n \geq 1}$ be a Cauchy sequence in the normed vector space $(L^p(\Omega, \mathcal{A}, \mathbb{P}), \|\cdot\|_p)$.

Prove that there exists a sub-sequence $(X_{n_k})_{k \geq 1}$ of $(X_n)_{n \geq 1}$ such that for all $k \geq 1$ and for all $l \geq k$, one has $\|X_{n_k} - X_{n_l}\|_p^p \leq 2^{-k}$. Prove that for almost all $\omega \in \Omega$, the sequence $(X_{n_k}(\omega))_{k \geq 1}$ is a Cauchy sequence. Denote by $X(\omega)$ its limit.

Use Fatou's lemma to prove that X belongs to L^p and that the sequence $(X_{n_k})_{k \geq 1}$ converges in L^p to X . Prove that the sequence $(X_n)_{n \geq 1}$ converges in L^p to X .

The fact that L^p is complete allows one to use the machinery of functional analysis. In particular, L^2 is a Hilbert space as we already mentioned. This has at least two important consequences.

Firstly, given any closed linear subspace F of $L^2(\Omega, \mathcal{A}, \mathbb{P})$, the orthogonal subspace F^\perp is a closed subspace which satisfies $F \oplus F^\perp = L^2(\Omega, \mathcal{A}, \mathbb{P})$. In particular, there exists an orthogonal projection p_F whose image is exactly F .

Secondly, any continuous linear form on $L^2(\Omega, \mathcal{A}, \mathbb{P})$ is of the form $Y \mapsto \mathbb{E}[XY]$ for some (uniquely defined) $X \in L^2(\Omega, \mathcal{A}, \mathbb{P})$. This is called the Riesz representation theorem.

Exercise 1.35 (Radon-Nikodym theorem) *Let (Ω, \mathcal{A}) be a measurable space. Let \mathbb{P} and \mathbb{Q} be two probability measures on (Ω, \mathcal{A}) . Assume that $\mathbb{Q} \ll \mathbb{P}$, that is, \mathbb{Q} is absolutely continuous with respect to \mathbb{P} , which means that for all $A \in \mathcal{A}$, $\mathbb{P}(A) = 0$ implies $\mathbb{Q}(A) = 0$.*

Define the measure $\mu = \mathbb{P} + \mathbb{Q}$ on (Ω, \mathcal{A}) . Note that this is not a probability measure. Prove that $L^2(\Omega, \mathcal{A}, \mu)$ is a subspace of $L^1(\Omega, \mathcal{A}, \mathbb{Q})$ and that the mapping $\mathbb{E}_{\mathbb{Q}} : L^2(\Omega, \mathcal{A}, \mu) \rightarrow \mathbb{R}$ defined by $\mathbb{E}_{\mathbb{Q}}[X] = \int_{\Omega} X d\mathbb{Q}$ is a continuous linear form. Deduce that there exists $\delta \in L^2(\Omega, \mathcal{A}, \mu)$ such that for all $X \in L^2(\Omega, \mathcal{A}, \mu)$, one has $\mathbb{E}_{\mathbb{Q}}[X] = \int_{\Omega} X \delta d\mu$.

Prove that $\delta \in [0, 1)$ \mathbb{P} -almost surely. Set $D = \frac{\delta}{1-\delta}$. Prove that for all non-negative $X : \Omega \rightarrow \mathbb{R}_+$, one has $\int_{\Omega} X d\mathbb{Q} = \int_{\Omega} XD d\mathbb{P}$. Prove that the same equality holds for any $X \in L^1(\Omega, \mathcal{A}, \mathbb{Q})$.

2 Conditional expectation

There are three main parts in this chapter. The first is the definition of conditional expectation. We will make a lot of comments about it, in order to help you to relate this rather abstract definition to what you already know. The second part is the proof that conditional expectation exists and is unique. Uniqueness is easy, but in my opinion more instructive than many books indicate. Existence is not easy and is commonly proved in one of two different ways. We explain one, the other can for example be found in Durrett's book. The third part is a list of properties which will allow you to manipulate in practical computations the concept of conditional expectation. My own experience of mathematics is that learning a new notion involves repeatedly going from applying rules without questioning them too much and watching them at work on one hand, and thinking about the structure of the theory and going deeper into the understanding of the rules on the other hand.

2.1 Definition and fundamental examples

Definition 2.1 *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let \mathcal{B} be a sub- σ -field of \mathcal{A} , that is, a σ -field on Ω such that $\mathcal{B} \subset \mathcal{A}$. Let $X : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ be an integrable random variable.*

A conditional expectation of X given \mathcal{B} is an integrable real-valued random variable $Y : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ such that the following two conditions hold.

- 1. Y is measurable with respect to \mathcal{B} .*
- 2. For all $B \in \mathcal{B}$, one has $\int_B X d\mathbb{P} = \int_B Y d\mathbb{P}$.*

We shall prove that a conditional expectation of X given \mathcal{B} always exists and is unique almost surely, and we shall denote it by $\mathbb{E}[X|\mathcal{B}]$.

Exercise 2.1 *Assume that $\mathcal{B} = \{\emptyset, \Omega\}$. Prove that a conditional expectation of X given \mathcal{B} must be a constant random variable and find the unique possible value of this constant.*

Let us start by an example. Suppose that \mathcal{B} is generated by an event $C \in \mathcal{A}$, so that $\mathcal{B} = \{\emptyset, C, C^c, \Omega\}$. Let us assume that $\mathbb{P}(C)$ is neither 0 nor 1. If a conditional expectation of X given \mathcal{B} exists, it must be of the form

$$Y = \alpha \mathbb{1}_C + \beta \mathbb{1}_{C^c},$$

because this is (up to modification on a negligible set) the form of the most general \mathcal{B} -measurable random variable. The second condition which Y must satisfy applied with $B = C$ yields

$$\alpha \mathbb{P}(C) = \int_C Y d\mathbb{P} = \int_C X d\mathbb{P},$$

so that

$$\alpha = \frac{1}{\mathbb{P}(C)} \int_C X d\mathbb{P}.$$

Similarly,

$$\beta = \frac{1}{\mathbb{P}(C^c)} \int_{C^c} X d\mathbb{P}.$$

The second condition is satisfied for $B = \emptyset$ and a short computation shows that, with the choices of α and β above, it is also satisfied for $B = \Omega$. We have thus proved that a conditional expectation of X given \mathcal{B} exists and is unique. It is equal to

$$\mathbb{E}[X|\mathcal{B}] = \left(\frac{1}{\mathbb{P}(C)} \int_C X \, d\mathbb{P} \right) \mathbb{1}_C + \left(\frac{1}{\mathbb{P}(C^c)} \int_{C^c} X \, d\mathbb{P} \right) \mathbb{1}_{C^c}.$$

We see the connection with elementary conditional probabilities: the value of $\mathbb{E}[X|\mathcal{B}]$ on C is nothing but the expectation of X under the conditional probability $\mathbb{P}^C = \mathbb{P}(\cdot|C)$ defined on \mathcal{A} by

$$\mathbb{P}^C(A) = \frac{\mathbb{P}(A \cap C)}{\mathbb{P}(C)}.$$

Let us try to formulate this with words. We know what the expectation of X with respect to the conditional probability \mathbb{P}^C is, as well as its expectation with respect to the conditional probability \mathbb{P}^{C^c} . These are two real numbers. Now the conditional expectation $\mathbb{E}[X|\mathcal{B}]$ is a random variable on Ω , a real-valued function on Ω . Consider a point ω in Ω . Think of the σ -field \mathcal{B} as encoding the information to which we have access. What do we know about ω ? We know whether it belongs to C or to C^c , and no more. If it belongs to C , then the conditional expectation of X given \mathcal{B} evaluated at ω is equal to the expectation of X with respect to the conditional probability \mathbb{P}^C . If it belongs to C^c , the same conclusion holds with C replaced by C^c .

What we have done with two sets C and C^c can be generalised to an arbitrary finite or countably infinite number of sets. This is the object of the next exercise.

Exercise 2.2 Let $\{C_1, C_2, \dots\}$ be a partition of Ω into at least two (and possibly countably infinitely many) disjoint events with positive probability. Let $\mathcal{B} = \sigma(\{C_1, C_2, \dots\})$ be the σ -algebra generated by this partition. Describe \mathcal{B} . Prove that $\mathbb{E}[X|\mathcal{B}]$ exists and is unique, and is given by the formula

$$\mathbb{E}[X|\mathcal{B}] = \sum_{i \geq 1} \left(\frac{1}{\mathbb{P}(C_i)} \int_{C_i} X \, d\mathbb{P} \right) \mathbb{1}_{C_i}.$$

This first discussion revealed the relation between the definition of $\mathbb{E}[X|\mathcal{B}]$ and the elementary notion of conditional probability. However, it is not general enough for many purposes, and cannot be generalised to the case of an arbitrary sub- σ -field \mathcal{B} . Indeed, many interesting σ -fields are not generated by a partition.

Exercise 2.3 Consider on \mathbb{R} the Borel σ -field $\mathcal{B}_{\mathbb{R}}$ and define

$$\mathcal{C} = \{C \in \mathcal{B}_{\mathbb{R}} : C + 2\pi = C\},$$

where for any subset C of \mathbb{R} , we denote by $C + 2\pi$ the result of the translation of C by 2π , that is, $C + 2\pi = \{x + 2\pi : x \in C\}$. Prove that \mathcal{C} is a sub- σ -field of $\mathcal{B}_{\mathbb{R}}$. Is \mathcal{C} generated by a partition of \mathbb{R} ?

Exercise 2.4 Prove that every σ -field on a countable set is generated by a partition. Is it true that if a set Ω has the property that every σ -field on Ω is generated by a partition, then Ω is countable?

There is another situation where we can study by hand the existence of a conditional expectation. It is a case where, among other things that we will explain in a minute, we assume that \mathcal{B} is generated by a real random variable, say $Z : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$. Recall from Exercise 1.15 the definition of $\sigma(Z)$, the σ -field generated by Z .

Let us make the assumption that the random vector (X, Z) has a distribution which admits a density $f_{(X,Z)}$ with respect to the Lebesgue measure on \mathbb{R}^2 . This means that for all Borel subset D of \mathbb{R}^2 , we have

$$\mathbb{P}((X, Z) \in D) = \int_D f_{(X,Z)}(u, v) \, dudv.$$

In this situation, we know that the distribution of Z admits a density with respect to the Lebesgue measure on \mathbb{R} which is the function f_Z given (for almost every $v \in \mathbb{R}$) by

$$f_Z(v) = \int_{\mathbb{R}} f_{(X,Z)}(u, v) \, du.$$

Now let us look for $\mathbb{E}[X|\sigma(Z)]$, which is also denoted by $\mathbb{E}[X|Z]$. It is useful to understand what it means for a random variable to be measurable with respect to $\sigma(Z)$.

Exercise 2.5 *In this exercise, we work with sets and maps between sets, without σ -fields and without any notion of measurability. Let Ω be a set. Let $Y, Z : \Omega \rightarrow \mathbb{R}$ be two maps. Prove that Y is constant on every non-empty set of Ω of the form $Z^{-1}(\{t\})$, $t \in \mathbb{R}$, if and only if there exists a map $h : \mathbb{R} \rightarrow \mathbb{R}$ such that $Y = h \circ Z$.*

Proposition 2.2 *Let Y, Z be real-valued random variables on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. The random variable Y is measurable with respect to $\sigma(Z)$ if and only if there exists a Borel measurable function $h : \mathbb{R} \rightarrow \mathbb{R}$ such that $Y = h(Z) = h \circ Z$.*

The proof of this proposition is not essential for our purposes. I give it here for the sake of completeness and pleasure.

Proof. If $Y = h(Z)$ for some measurable function h , then for all $B \in \mathcal{B}_{\mathbb{R}}$, we have $Y^{-1}(B) = (h \circ Z)^{-1}(B) = Z^{-1}(h^{-1}(B)) \in \sigma(Z)$, so that Y is measurable with respect to $\sigma(Z)$.

The most interesting part is the converse. Assume that Y is measurable with respect to $\sigma(Z)$. Let us also assume, for a start, that Y is non-negative.

We shall use the binary expansion of real numbers:

$$14.5 \text{ (decimal)} = 1110.1 \text{ (binary)}.$$

Let us observe two things. Firstly, in the binary expansion of a real x , the digit just to the left of the dot, the 0 in the example above, is the parity of the integer part of x . Let us write it $[x] \bmod 2$. Secondly, for all $n \in \mathbb{Z}$, the n -th digit of the binary expansion of x , that is, the digit which is the coefficient of 2^n , is the digit located just to the left of the dot in the expansion of $2^{-n}x$. Putting these remarks together, we find that

$$x = \sum_{n \in \mathbb{Z}} 2^{-n} ((2^n x) \bmod 2).$$

We can write the same formula for a non-negative random variable

$$Y = \sum_{n \in \mathbb{Z}} 2^{-n} ((2^n Y) \bmod 2) = \lim_{N \rightarrow +\infty} \sum_{n=-N}^N 2^{-n} ((2^n Y) \bmod 2).$$

We use now the assumption. Since Y is $\sigma(Z)$ measurable, the set $C_n = \{(2^n Y) \bmod 2 = 1\}$ belongs to $\sigma(Z)$ for all $n \in \mathbb{Z}$. There exists thus a Borel subset B_n of \mathbb{R} such that $C_n = Z^{-1}(B_n)$. Define now a function h on \mathbb{R} by setting

$$h = \sum_{n \in \mathbb{Z}} 2^{-n} \mathbb{1}_{B_n}.$$

Then for all $\omega \in \Omega$,

$$h(Z(\omega)) = \sum_{n \in \mathbb{Z}} 2^{-n} \mathbb{1}_{B_n}(Z(\omega)) = \sum_{n \in \mathbb{Z}} 2^{-n} \mathbb{1}_{C_n}(\omega) = \sum_{n \in \mathbb{Z}} 2^{-n} ((2^n Y(\omega)) \bmod 2) = Y(\omega).$$

Observe that h may take the value $+\infty$, but on the set $\{h = +\infty\}$ we can give it the value 0 without changing the last computation, because Y takes real values.

Finally, we have expressed Y explicitly as a function of Z . The case where Y can take negative values is treated as usual by decomposing Y into its non-negative and non-positive parts. \square

Exercise 2.6 *Is the σ -field \mathcal{C} defined in Exercise 2.3 of the form $\sigma(Z)$ for some measurable function $Z : \mathbb{R} \rightarrow \mathbb{R}$?*

Let us come back to our original problem of computing $\mathbb{E}[X|Z]$. Thanks to the proposition which we have just proved, we know that we are looking for a random variable of the form $h(Z)$ for some function h . This function must satisfy the second defining property of conditional expectation for all $B \in \sigma(Z)$. By definition of $\sigma(Z)$, each $B \in \sigma(Z)$ is of the form $B = Z^{-1}(E)$ for some Borel subset E of \mathbb{R} . The condition to be satisfied is

$$\int_B h(Z) dP = \int_B X dP,$$

which, since $B = Z^{-1}(E)$, can be written as

$$\int_E h(v) f_Z(v) dv = \int_{\mathbb{R} \times E} u f_{(X,Z)}(u, v) dudv.$$

The last line was deduced from the previous one by a short computation that you might want to carefully check. Now the right-hand side can be written as

$$\int_E \left(\int_{\mathbb{R}} u f_{(X,Z)}(u, v) du \right) dv.$$

Provided the denominator is not zero, we see that the function h given by

$$h(v) = \frac{\int_{\mathbb{R}} u f_{(X,Z)}(u, v) du}{f_Z(v)} = \frac{\int_{\mathbb{R}} u f_{(X,Z)}(u, v) du}{\int_{\mathbb{R}} f_{(X,Z)}(u, v) du}$$

satisfies the desired relation. Hence, $h(Z)$ is a conditional expectation of X given $\sigma(Z)$.

To what extent is h unique? The integral of hf_Z with respect to the Lebesgue measure over each Borel subset of \mathbb{R} is prescribed. This means that the integral of h with respect to \mathbb{P}_Z over each Borel subset of \mathbb{R} is prescribed. If h' has the same integral over each Borel subset, then h and h' coincide \mathbb{P}_Z -almost surely. Let us state this as an exercise.

Exercise 2.7 *On a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, let X and X' be two random variables such that for all $A \in \mathcal{A}$ one has $\int_A X d\mathbb{P} = \int_A X' d\mathbb{P}$. Then $\mathbb{P}(X = X') = 1$.*

2.2 Uniqueness and existence

The result of Exercise 2.7 is elementary and we may not pay great attention to it, but in the context of conditional expectation, it is very important. Let us see why. Random variables are equivalence classes of measurable functions and therefore cannot be evaluated at a point of Ω . Indeed, two measurable functions which are equal almost everywhere may differ at any point which is not an atom of the underlying probability measure. In contrast, they can be integrated on any event, and this result tells us that this is all there is: a random variable is not meant to be evaluated, but it is certainly made to be integrated over events, and it is characterised by the collection of its integrals over all events.

This is important precisely because the conditional expectation $\mathbb{E}[X|\mathcal{B}]$ is defined not by its values at the points of Ω , but by the values of its integrals over all events of \mathcal{B} . It is an instance where a random variable is truly seen as an equivalence class of measurable functions, with no preferred choice of a specific measurable function in this class. This is in a sense the most intrinsic possible definition of a random variable.

These remarks hopefully make the structure of the definition of the conditional expectation $\mathbb{E}[X|\mathcal{B}]$ more natural. The first condition says that it is a \mathcal{B} -measurable random variable, and the second condition specifies its integral over every single event of \mathcal{B} .

By now, it should have become almost obvious that the conditional expectation is unique.

Lemma 2.3 *Recall the notation of Definition 2.1. If Y and Y' are conditional expectations of X given \mathcal{B} , then $Y = Y'$ almost surely.*

Proof. Indeed, Y and Y' are both measurable with respect to \mathcal{B} and they have the same integral over every event belonging to \mathcal{B} . □

We have settled the question of the uniqueness, but not yet that of the existence of the conditional expectation. In order to do this, we will introduce a new and quite different point of view on the conditional expectation. This point of view is usually explained in terms of “the best prediction about X which can be made using the information available in \mathcal{B} ”. This is certainly correct, but there is something about this phrasing that I have always found strange, and I will try to explain it in a convincing way.

Let us start by a comment on the usual expectation. Apart from its definition, we can understand the expectation of a random variable through the strong law of large numbers. This law tells us that the arithmetic mean of a sufficiently large sample of independent copies of our random variable will be arbitrarily close to its expectation. However, this can hardly be called a *prediction*, let alone a *best prediction*. In which sense does the expectation of a random variable constitute a prediction about this variable? If we think about a dice that we roll, there is little

hope that it will actually give the value 3.5. A similar comment could be made about a uniform random variable on the interval $[0, 1]$, and about many other random variables. The words *best prediction* should not be heard in the naive sense of *likeliest outcome*.

There is however a precise and simple sense in which the expectation of a random variable constitutes a best prediction of its outcome.

Exercise 2.8 *Let X be a square-integrable random variable. Prove that for all real number c , the following equality holds:*

$$\mathbb{E}[(X - c)^2] = \text{Var}(X) + (c - \mathbb{E}[X])^2.$$

Deduce that there is, among all constant random variables, one which is closer than any other to X in the normed vector space $(L^2(\Omega, \mathcal{A}, \mathbb{P}), \|\cdot\|_2)$.

If we think of $\mathbb{E}[X]$ not as a number but as a constant random variable (and this is very much in the spirit of our study of $\mathbb{E}[X|\mathcal{B}]$ in the case where $\mathcal{B} = \{\emptyset, C, C^c, \Omega\}$), then it is, among all the constant random variables, the one which is closest to X in the L^2 norm. Recalling Exercise 2.1, and introducing $\mathcal{B}_0 = \{\emptyset, \Omega\}$, we can reformulate this as follows: $\mathbb{E}[X|\mathcal{B}_0]$ is, among all random variables which are measurable with respect to \mathcal{B}_0 , the one which is closest to X in the L^2 distance.

The key to the existence of the conditional expectation is that this is true not just for the σ -field \mathcal{B}_0 , but for any other sub- σ -field of \mathcal{A} .

Theorem 2.4 *Let $X : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ be a square-integrable random variable. Let \mathcal{B} be a sub- σ -field of \mathcal{A} .*

There exists, among all square-integrable random variables which are measurable with respect to \mathcal{B} , one which is closer to X in the L^2 distance than any other. This random variable is moreover a conditional expectation of X given \mathcal{B} .

Proof. The space of square-integrable random variables which are measurable with respect to \mathcal{B} is the subspace $L^2(\Omega, \mathcal{B}, \mathbb{P})$ of $L^2(\Omega, \mathcal{A}, \mathbb{P})$. Since the L^2 distance on $L^2(\Omega, \mathcal{B}, \mathbb{P})$ is the distance induced by the ambient space $L^2(\Omega, \mathcal{A}, \mathbb{P})$ and since $L^2(\Omega, \mathcal{B}, \mathbb{P})$ endowed with this distance is complete, it is a closed subspace of $L^2(\Omega, \mathcal{A}, \mathbb{P})$.

Let $p : L^2(\Omega, \mathcal{A}, \mathbb{P}) \rightarrow L^2(\Omega, \mathcal{B}, \mathbb{P})$ be the orthogonal projection, which exists precisely because $L^2(\Omega, \mathcal{B}, \mathbb{P})$ is a closed linear subspace. Then from the general theory of the geometry of Hilbert spaces, we know that $p(X)$ is the element of $L^2(\Omega, \mathcal{B}, \mathbb{P})$ which is closest to X in the L^2 distance.

Let us now prove that $p(X)$ is a conditional expectation of X given \mathcal{B} . To start with, $p(X)$ is measurable with respect to \mathcal{B} . Then, let us consider an event $B \in \mathcal{B}$. We have

$$\int_B p(X) d\mathbb{P} - \int_B X d\mathbb{P} = \int_{\Omega} (p(X) - X) \mathbb{1}_B d\mathbb{P} = \langle p(X) - X, \mathbb{1}_B \rangle_{L^2(\Omega, \mathcal{A}, \mathbb{P})}.$$

But on one hand, $\mathbb{1}_B$ is an element of $L^2(\Omega, \mathcal{B}, \mathbb{P})$. On the other hand, since p is an orthogonal projection, $p(X) - X$ is orthogonal to $L^2(\Omega, \mathcal{B}, \mathbb{P})$. Hence, the scalar product of the right-hand side is equal to 0, so that

$$\int_B p(X) d\mathbb{P} = \int_B X d\mathbb{P},$$

proving that $p(X)$ is the conditional expectation of X given \mathcal{B} . □

We have proved the existence for all square-integrable random variables. It remains to go from square-integrable to integrable random variables (recall that $L^2(\Omega, \mathcal{A}, \mathbb{P}) \subset L^1(\Omega, \mathcal{A}, \mathbb{P})$).

Theorem 2.5 *Let $X : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ be an integrable random variable. Let \mathcal{B} be a sub- σ -field of \mathcal{A} . There exists a conditional expectation of X given \mathcal{B} .*

We first need to prove a fundamental property of the conditional expectation, which is its positivity. It is more easily done after proving its linearity. We do both in the following lemma.

Lemma 2.6 *Let $X_1, X_2 : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ be two random variables. Let \mathcal{B} be a sub- σ -field of \mathcal{A} . Assume that X_1 and X_2 admit conditional expectations with respect to \mathcal{B} .*

1. *For all reals a, b , the random variable $aX_1 + bX_2$ admits a conditional expectation given \mathcal{B} and $\mathbb{E}[aX_1 + bX_2|\mathcal{B}] = a\mathbb{E}[X_1|\mathcal{B}] + b\mathbb{E}[X_2|\mathcal{B}]$.*

2. *If $0 \leq X_1 \leq X_2$ almost surely, then $0 \leq \mathbb{E}[X_1|\mathcal{B}] \leq \mathbb{E}[X_2|\mathcal{B}]$ almost surely.*

Proof. 1. The random variable $a\mathbb{E}[X_1|\mathcal{B}] + b\mathbb{E}[X_2|\mathcal{B}]$ is integrable and \mathcal{B} -measurable. The linearity of the integral implies immediately that it is a conditional expectation of $aX_1 + bX_2$ given \mathcal{B} .

2. Thanks to the linearity, it suffices to prove that $0 \leq X$ almost surely implies $0 \leq \mathbb{E}[X|\mathcal{B}]$ almost surely. Consider an integer $k \geq 1$ and the event $B_k = \{\mathbb{E}[X|\mathcal{B}] \leq -\frac{1}{k}\}$. It belongs to \mathcal{B} . Hence, we have the inequalities

$$0 \leq \int_{B_k} X \, d\mathbb{P} = \int_{B_k} \mathbb{E}[X|\mathcal{B}] \, d\mathbb{P} \leq -\frac{1}{k}\mathbb{P}(B_k).$$

This implies $\mathbb{P}(B_k) = 0$. Hence, the event $B = \{\mathbb{E}[X|\mathcal{B}] < 0\}$ satisfies

$$\mathbb{P}(B) = \mathbb{P}\left(\bigcup_{k \geq 1} B_k\right) \leq \sum_{k \geq 1} \mathbb{P}(B_k) = 0.$$

Finally, $\mathbb{E}[X|\mathcal{B}]$ is non-negative almost surely. □

Let us now turn to the proof of the theorem.

Proof. Let us write $X = X^+ - X^-$, where $X^+ = \max(X, 0)$ and $X^- = (-X)^+$. Let us treat the case of X^+ first.

For each $n \geq 1$, let us consider the random variable $\min(X^+, n)$. It is non-negative and bounded by n . It is in particular square-integrable, so that it admits a conditional expectation given \mathcal{B} .

For each $m \leq n$, we have $0 \leq \min(X^+, m) \leq \min(X^+, n)$, so that, by the lemma, $0 \leq \mathbb{E}[\min(X^+, m)|\mathcal{B}] \leq \mathbb{E}[\min(X^+, n)|\mathcal{B}]$. Hence, the sequence of random variables $(\mathbb{E}[\min(X^+, n)|\mathcal{B}])_{n \geq 1}$ is non-decreasing. In particular, it has a limit Y_+ towards which it converges almost surely.

On the other hand, the sequence $(\min(X^+, n))_{n \geq 1}$ is also non-decreasing and converges to X . Hence, for all $B \in \mathcal{B}$, the monotone convergence theorem applied to the equality

$$\int_B \mathbb{E}[\min(X^+, n)|\mathcal{B}] \, d\mathbb{P} = \int_B \min(X^+, n) \, d\mathbb{P}$$

yields

$$\int_B Y_+ d\mathbb{P} = \int_B X^+ d\mathbb{P}.$$

In particular, since Y_+ is non-negative, $\int_\Omega Y_+ d\mathbb{P} = \int_\Omega X^+ d\mathbb{P} < +\infty$ and Y_+ is integrable. Since Y_+ is the almost sure limit of a sequence of \mathcal{B} -measurable random variables, it is also \mathcal{B} -measurable. Finally, the equality above proves that it is the conditional expectation of X^+ given \mathcal{B} .

The same argument proves that X^- admits a conditional expectation given \mathcal{B} , which we denote by Y_- .

Finally, we claim that $Y = Y_+ - Y_-$ is the conditional expectation of X given \mathcal{B} . Indeed, Y is integrable, \mathcal{B} -measurable, and for all $B \in \mathcal{B}$, we have

$$\int_B Y d\mathbb{P} = \int_B Y_+ - Y_- d\mathbb{P} = \int_B Y_+ d\mathbb{P} - \int_B Y_- d\mathbb{P} = \int_B X^+ d\mathbb{P} - \int_B X^- d\mathbb{P} = \int_B X d\mathbb{P}.$$

This concludes the proof. \square

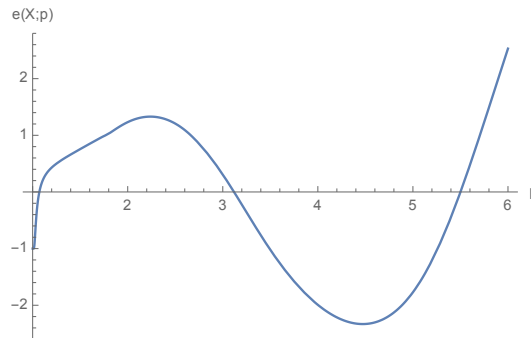
Let us summarise our approach. We analysed the definition of the conditional expectation in a few simple cases and this led us to an understanding of the uniqueness of the conditional expectation. We then proved its existence by using a strategy of proof which is not uncommon when one is dealing with L^1 spaces: we first proved the existence for L^2 random variables, using the rich geometry of Hilbert spaces, and then extended our result to L^1 by an argument of approximation.

Exercise 2.9 *Is it true that among all constant random variables, $\mathbb{E}[X]$ is closer to X than any other in the L^1 distance? Could we have characterised $\mathbb{E}[X|\mathcal{B}]$ for an integrable random variable by a property similar to the one we used for square-integrable random variables?*

Exercise 2.10 *Let X be a random variable that is uniformly distributed on the finite set $\{-2, -1, 1, 10\}$. Compute the expectation of X and draw the graph of the function $c \mapsto \mathbb{E}[|X - c|]$.*

For a general integrable random variable, describe the set of reals c that minimise $\mathbb{E}[|X - c|]$.

Exercise 2.11 *Let X be a bounded random variable. Prove that for all $p \in (1, \infty)$, the function $c \mapsto \|X - c\|_{L^p}$ attains its minimum for a unique value of c , which we denote by $e(X; p)$. What can you say about the function $p \mapsto e(X; p)$? Is it continuous? Does it have a limit as p tends to infinity? As p tends to 1? Here is a part of the graph of this function for a random variable X whose distribution is a slightly perturbed version of the distribution of the previous exercise:*



2.3 Main properties

We are now ready for the third part of this section, where we gather the useful properties of the conditional expectation.

Theorem 2.7 *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let X, Y, X_1, X_2, \dots be integrable random variables on this space. Let a, b be real numbers. Let \mathcal{B}, \mathcal{C} be sub- σ -fields of \mathcal{A} .*

1. $\mathbb{E}[aX + bY|\mathcal{B}] = a\mathbb{E}[X|\mathcal{B}] + b\mathbb{E}[Y|\mathcal{B}]$.
2. If X is \mathcal{B} -measurable and XY is integrable, then $\mathbb{E}[XY|\mathcal{B}] = X\mathbb{E}[Y|\mathcal{B}]$.
3. If X is \mathcal{B} -measurable, then $\mathbb{E}[X|\mathcal{B}] = X$.
4. If X is independent of \mathcal{B} , then $\mathbb{E}[X|\mathcal{B}] = \mathbb{E}[X]$ almost surely.
5. If $\mathcal{C} \subset \mathcal{B}$, then $\mathbb{E}[\mathbb{E}[X|\mathcal{B}]|\mathcal{C}] = \mathbb{E}[X|\mathcal{C}]$.
6. $\mathbb{E}[\mathbb{E}[X|\mathcal{B}]] = \mathbb{E}[X]$.
7. If $X \geq 0$ almost surely, then $\mathbb{E}[X|\mathcal{B}] \geq 0$ almost surely.
8. If $(X_n)_{n \geq 1}$ is a non-decreasing sequence of non-negative random variables, converging to an integrable random variable X , then $\mathbb{E}[X_n|\mathcal{B}] \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \mathbb{E}[X|\mathcal{B}]$.
9. If $(X_n)_{n \geq 1}$ is a sequence of non-negative random variables, then $\mathbb{E}[\liminf X_n|\mathcal{B}] \leq \liminf \mathbb{E}[X_n|\mathcal{B}]$.
10. If $(X_n)_{n \geq 1}$ is a sequence of non-negative random variables converging almost surely to a random variable X , and if there exists an integrable random variable Y such that for all $n \geq 1$ one has $X_n \leq Y$, then $\mathbb{E}[X_n|\mathcal{B}] \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \mathbb{E}[X|\mathcal{B}]$.
11. If $\phi: \mathbb{R} \rightarrow \mathbb{R}$ is a convex function and $\phi(X)$ is integrable, then $\phi(\mathbb{E}[X|\mathcal{B}]) \leq \mathbb{E}[\phi(X)|\mathcal{B}]$.

To prove this theorem yourself is one of the best exercises that you can do to at this stage. None of the ten first properties are difficult to prove. Among these, the least simple is the second, in that it is the least directly deduced from the definition. You should prove it after you have proved the eighth point, starting by the case where Y is an indicator function, then a simple function, and finally using an approximation argument.

For the eleventh property, use the fact that a convex function is the supremum of all affine functions which are inferior to it:

$$\phi(x) = \sup_{\substack{a, b \in \mathbb{R} \\ \forall y \in \mathbb{R}, ay + b \leq \phi(y)}} ax + b.$$

Here are some more exercises about conditional expectation. When nothing is specified, X, Y, \dots are integrable random variables on a probability space.

Exercise 2.12 *Let X be a standard Gaussian random variable. Compute $\mathbb{E}[X|X^2]$.*

Exercise 2.13 *Assume that $\mathbb{E}[X|\mathcal{B}]$ is a constant random variable. Prove that this constant is $\mathbb{E}[X]$. Is X necessarily independent of \mathcal{B} ?*

Exercise 2.14 *Let a be a real number. Let X be a random variable whose distribution admits the density*

$$f_X(x) = \frac{(x+a)^2 e^{-\frac{x^2}{2}}}{1+a^2 \sqrt{2\pi}}$$

with respect to the Lebesgue measure (check that f_X is indeed a density function). Try to understand before computing it that the sign of $\mathbb{E}[X|X^2]$ is almost surely constant, and to guess how it depends on a . Then compute $\mathbb{E}[X|X^2]$.

For all real number t , compute

$$\frac{tf_X(t) - tf_X(-t)}{f_X(t) + f_X(-t)}.$$

Do you see a relation between the two computations ?

Exercise 2.15 Which relations can you find between $\mathbb{E}[X|Y^2]$ and $\mathbb{E}[X||Y]$?

Exercise 2.16 Last year a student told me that $\mathbb{E}[X|X^2] = 0$ if and only if X and $-X$ have the same distribution. Do you think he was right ?

Exercise 2.17 Choose $p \in (0, 1)$. Let X and Y be integer-valued random variables such that for all $k, l \in \mathbb{N}$, one has

$$\mathbb{P}(X = k, Y = l) = (1 - p) \left(\frac{p}{e}\right)^k \frac{k^l}{l!}.$$

Compute $\mathbb{E}[Y|X]$.

Exercise 2.18 Let (X, Y) be a two-dimensional random vector whose distribution admits the density

$$f_{(X,Y)}(s, t) = e^{-s} \mathbb{1}_{[0,s]}(t)$$

with respect to the Lebesgue measure on \mathbb{R}^2 . Compute $\mathbb{E}[X|Y]$ and $\mathbb{E}[Y|X]$.

Exercise 2.19 (Do this exercise only if you know the definition and main properties of Gaussian random vectors.) Let (X, Y) be a two-dimensional Gaussian random vector. Prove that there exists a real number a such that $X - aY$ is independent of Y . Prove that there exists a real number b such that $\mathbb{E}[X|Y] = aY + b$.

Exercise 2.20 Let N be a bounded integer-valued random variable. Let $(X_n)_{n \geq 1}$ be a sequence of identically distributed integrable random variables. Assume that N, X_1, X_2, \dots are independent. Set

$$S = \sum_{n=1}^N X_n.$$

Prove that S is integrable and compute its expectation.

Exercise 2.21 Consider the probability space $([0, 1], \mathcal{B}_{[0,1]}, \lambda)$, where λ is the Lebesgue measure. Choose $n \geq 0$ and consider the σ -field

$$\mathcal{F}_n = \sigma \left(\left\{ \left[\frac{k}{2^n}, \frac{k+1}{2^n} \right) : k \in \{0 \dots 2^n - 1\} \right\} \right).$$

Define the random variable $X : [0, 1) \rightarrow \mathbb{R}$ by setting $X(t) = t$ for all $t \in [0, 1)$. Compute $\mathbb{E}[X|\mathcal{F}_n]$.

Exercise 2.22 Consider a random vector (X, Y, Z) whose distribution admits the density

$$f_{(X,Y,Z)}(x, y, z) = ce^{-z-2x} \mathbb{1}_{0 \leq x \leq y \leq z}$$

with respect to the Lebesgue measure on \mathbb{R}^3 , where c is a real constant. Determine the value of c , then compute $\mathbb{E}[X|Y, Z]$ and $\mathbb{E}[X|Y]$.

Exercise 2.23 Let X and Y be two square-integrable random variables such that $\mathbb{E}[X|Y] = Y$ and $\mathbb{E}[Y|X] = X$. Prove that $X = Y$.

Exercise 2.24 Consider a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Let \mathcal{B} be a sub- σ -field of \mathcal{A} . Propose a definition of the conditional expectation $\mathbb{E}[X|\mathcal{B}]$ for every random variable $X : (\Omega, \mathcal{A}) \rightarrow [0, +\infty]$, without any assumption of integrability. Investigate the existence and uniqueness of this conditional expectation, as well as its properties, in the spirit of Theorem 2.7.

Exercise 2.25 Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let \mathcal{B} be a sub- σ -field of \mathcal{A} . Check that $L^\infty(\Omega, \mathcal{B}, \mathbb{P})$ acts by multiplication on $L^1(\Omega, \mathcal{A}, \mathbb{P})$, in the sense that for all $X \in L^1(\Omega, \mathcal{A}, \mathbb{P})$ and all $Z \in L^\infty(\Omega, \mathcal{B}, \mathbb{P})$, the product ZX belongs to $L^1(\Omega, \mathcal{A}, \mathbb{P})$. Check that $L^\infty(\Omega, \mathcal{B}, \mathbb{P})$ acts also by multiplication on $L^1(\Omega, \mathcal{B}, \mathbb{P})$.

Find all linear maps

$$\mathcal{E} : L^1(\Omega, \mathcal{A}, \mathbb{P}) \rightarrow L^1(\Omega, \mathcal{B}, \mathbb{P})$$

which preserve the expectation and respect the action by multiplication of $L^\infty(\Omega, \mathcal{B}, \mathbb{P})$ in the sense that for all $X \in L^1(\Omega, \mathcal{A}, \mathbb{P})$ and all $Z \in L^\infty(\Omega, \mathcal{B}, \mathbb{P})$,

$$\mathcal{E}(ZX) = Z\mathcal{E}(X).$$

In particular, check that \mathcal{E} is automatically continuous.

Find all continuous linear maps $\tilde{\mathcal{E}} : L^1(\Omega, \mathcal{A}, \mathbb{P}) \rightarrow L^1(\Omega, \mathcal{B}, \mathbb{P})$ which respect the action by multiplication of $L^\infty(\Omega, \mathcal{B}, \mathbb{P})$.

3 Martingales

3.1 Martingales and conserved quantities

In the study of a deterministic dynamical system, a mechanical system for example, it is very important to identify quantities which stay constant through the evolution of the system. For example, let us think of the motion of a pendulum.

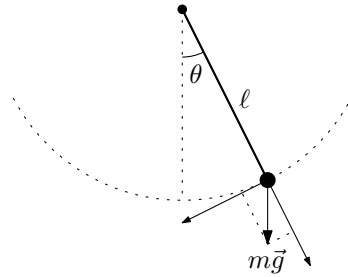


Figure 2: Newton's law writes $m\ell\ddot{\theta} = mg\sin\theta$, that is, $\ddot{\theta} = \frac{g}{\ell}\sin\theta$, where $g = 9.81 \text{ m s}^{-2}$ is the gravitation field on the surface of the Earth.

The differential equation $\ddot{\theta} = \frac{g}{\ell}\sin\theta$ is not easy to solve explicitly. Let us introduce the function $E = \frac{1}{2}m\ell\dot{\theta}^2 - mg\cos\theta$. It is immediately checked that $\dot{E} = 0$, and a very detailed study of the motion of the pendulum can be done from this single observation.

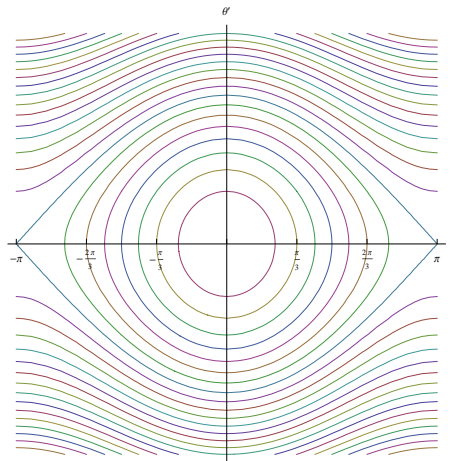


Figure 3: The curve $(\theta, \dot{\theta})$ corresponding to any possible motion of the pendulum is one of these curves (the right and left vertical sides of the picture, corresponding to $\theta = \pi$ and $\theta = -\pi$, should be identified). Which are the four radically different possible kinds of motion ?

In the evolution of a random dynamical system, it is unlikely that any quantity remains constant. However, certain quantities remain constant on average. For example, consider the simple random walk on \mathbb{Z} . That is, let $(X_n)_{n \geq 1}$ be a sequence of i.i.d. random variables such that $\mathbb{P}(X_1 = 1) = \mathbb{P}(X_1 = -1) = \frac{1}{2}$. Set $S_0 = 0$ and, for all $n \geq 1$, $S_n = X_1 + \dots + X_n$.

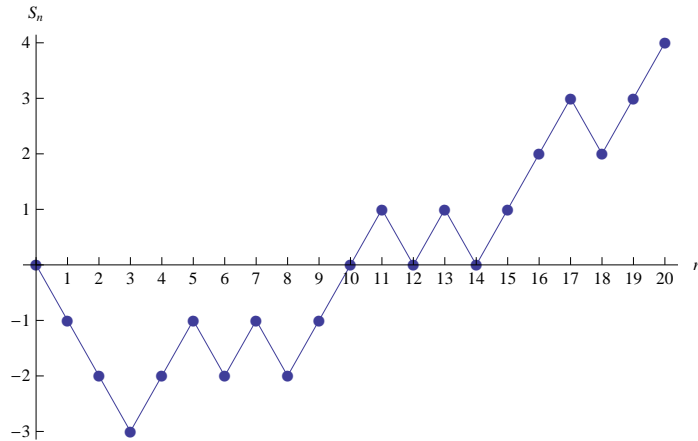


Figure 4: A sample path of the random dynamical system $(S_n)_{n \geq 0}$, also known as the simple random walk.

It is intuitively clear, and not difficult to prove, that the only functions $f : \mathbb{Z} \rightarrow \mathbb{R}$ such that $f(S_n)$ does not depend on n are the constant functions. There is nothing which, in this situation, plays exactly the same role as the mechanical energy in the case of the pendulum.

Exercise 3.1 Find a non-constant function $f : \mathbb{N} \times \mathbb{Z} \rightarrow \mathbb{R}$ such that $f(n, S_n)$ is constant.

However, you know that S_n itself is constant *on average*. More precisely, $\mathbb{E}[S_n] = 0$ does not depend on n . This observation is interesting but it is quite crude and, with some thought, it can be turned into a subtle and very fruitful one.

Indeed, look at the picture above. We know the value of S_n for $n \in \{0, \dots, 20\}$. What can we say about S_{21} ? Since $S_{20} = 4$, it has probability $\frac{1}{2}$ to be equal to 3, or to 5. In particular, its expected value is $\frac{1}{2}(3 + 5) = 4$. Here, the words “expected value” do not of course refer to the expectation, for we know that $\mathbb{E}[S_{21}] = 0$. They refer to the *conditional expectation* of S_{21} given the information to which we have access. This information consists in the values of $S_0(\omega), \dots, S_{20}(\omega)$ for the particular ω which corresponds to the particular realisation of the infinite experiment of which we see the first twenty steps.

Hence, the sentence “the expected value of S_{21} is 4” refers to the value at ω of the conditional expectation of S_{21} given the σ -field generated by S_0, \dots, S_{20} . It means that

$$\mathbb{E}[S_{21} | \sigma(S_0, \dots, S_{20})](\omega) = 4 = S_{20}(\omega).$$

How do we evaluate a random variable, which is only defined almost everywhere, at the point ω ? Well, the information to which we have access, contained in the picture above, describes exactly one event of the σ -field $\sigma(S_0, \dots, S_{20})$, namely the event

$$\{S_0 = 0, S_1 = -1, S_2 = -2, S_3 = -3, S_4 = -2, S_5 = -1, \dots, S_{19} = 3, S_{20} = 4\}.$$

We know that ω belongs to this event. We also know that $\mathbb{E}[S_{21} | \sigma(S_0, \dots, S_{20})]$ is constant on this event, equal to 4, which happens to be also the value of S_{20} .

Finally, the random dynamical system $(S_n)_{n \geq 0}$ has the following property: if for some $n \geq 0$ you observe S_0, \dots, S_n , then the values which you observe define an event, on which the

(conditional) average of S_{n+1} equals the value of S_n . This is beautifully expressed by the formula

$$\forall n \geq 0, \mathbb{E}[S_{n+1}|S_0, \dots, S_n] = S_n.$$

This is the main defining property of a *martingale*.

Martingales play for random dynamical systems the role played by conserved quantities for deterministic dynamical systems. For example, knowing that $(S_n)_{n \geq 0}$ is a martingale will allow us to answer the question : what is the probability that the first time at which S_n hits the value a is smaller than the first time at which S_n hits the value b ? You can try to find the answer for $a = 2$ and $b = -1$ and see that it is in general not easy to determine this probability.

Exercise 3.2 Check that $\mathbb{E}[S_n^2] = n$. What do you think about the process $(S_n^2 - n)_{n \geq 0}$?

In the classical case, it is sometimes not possible, or difficult, to identify a conserved quantity. This is for example the case for the pendulum if we take friction into account. But then, one can prove that $\dot{E} \leq 0$ and this is already a very interesting information. Similarly, martingales have variants called sub- and super-martingales, which correspond to Lyapunov functions of classical dynamical systems. Let us now turn to the systematic study of this class of processes.

3.2 Definition

Definition 3.1 Let (Ω, \mathcal{A}) be a measurable space. A filtration on (Ω, \mathcal{A}) is a non-decreasing sequence of sub- σ -fields of \mathcal{A} , that is, a sequence of sub- σ -fields of \mathcal{A} such that the inclusions

$$\mathcal{F}_0 \subset \mathcal{F}_1 \subset \dots \subset \mathcal{F}_n \subset \mathcal{F}_{n+1} \subset \dots \subset \mathcal{A}$$

hold. A probability space $(\Omega, \mathcal{A}, \mathbb{P})$ endowed with a filtration $(\mathcal{F}_n)_{n \geq 0}$ is called a filtered probability space.

Exercise 3.3 Recall the notation of Exercise 2.21. Prove that $(\mathcal{F}_n)_{n \geq 0}$ is a filtration. Is $\bigcup_{n \geq 0} \mathcal{F}_n$ a σ -field on $[0, 1)$?

The ambient σ -field \mathcal{A} can always be replaced by $\sigma\left(\bigcup_{n \geq 0} \mathcal{F}_n\right)$, and this is why we will sometimes not mention it in our definitions.

An important example of filtration is that generated by a sequence of random variables. Given a sequence $X = (X_n)_{n \geq 0}$ of random variables (and we will also call such a sequence a stochastic process), we can form a filtration $(\mathcal{F}_n^X)_{n \geq 0}$ by setting, for all $n \geq 0$,

$$\mathcal{F}_n^X = \sigma(X_0, \dots, X_n).$$

Definition 3.2 A stochastic process $X = (X_n)_{n \geq 0}$ defined on a filtered probability space $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ is said to be adapted if for all $n \geq 0$, the random variable X_n is measurable with respect to \mathcal{F}_n .

With the current notation, you can check that X is adapted if and only if $\mathcal{F}_n^X \subset \mathcal{F}_n$ for all $n \geq 0$.

Definition 3.3 (Martingales) Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let $X = (X_n)_{n \geq 0}$ be a stochastic process defined on this probability space. One calls X a martingale if the following conditions are satisfied.

1. X is adapted to $(\mathcal{F}_n)_{n \geq 0}$.
2. For all $n \geq 0$, X_n is integrable.
3. For all $n \geq 0$,

$$\mathbb{E}[X_{n+1} | \mathcal{F}_n] = X_n. \quad (\text{MG})$$

One calls X a supermartingale if it satisfies the same conditions, but with the equality (MG) replaced by the inequality

$$\mathbb{E}[X_{n+1} | \mathcal{F}_n] \leq X_n. \quad (\overline{\text{MG}})$$

One calls X a submartingale if it satisfies the same conditions with (MG) replaced by the other inequality

$$\mathbb{E}[X_{n+1} | \mathcal{F}_n] \geq X_n. \quad (\underline{\text{MG}})$$

The simple random walk is a fundamental example of martingale.

Exercise 3.4 Consider the simple random walk $S = (S_n)_{n \geq 0}$ as defined in the previous section. Recall in particular the i.i.d. sequence $(X_n)_{n \geq 1}$. Prove that the filtrations $(\mathcal{F}_n^X)_{n \geq 0}$ and $(\mathcal{F}_n^S)_{n \geq 0}$ are equal. Prove that S is a martingale with respect to $(\mathcal{F}_n^S)_{n \geq 0}$. What do you think about $(S_n^2)_{n \geq 1}$?

Exercise 3.5 Let $X = (X_n)_{n \geq 0}$ be a martingale with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$. Compute, for all $n, m \geq 0$, the conditional expectation $\mathbb{E}[X_n | \mathcal{F}_m]$. Prove that $\mathbb{E}[X_n] = \mathbb{E}[X_0]$ for all $n \geq 0$.

Exercise 3.6 Prove that if X is a submartingale (resp. a supermartingale), then the sequence $(\mathbb{E}[X_n])_{n \geq 0}$ is non-decreasing (resp. non-increasing). Mind the misleading vocabulary !

Exercise 3.7 Let $X = (X_n)_{n \geq 0}$ be a stochastic process. Prove that if X is a martingale with respect to a filtration $(\mathcal{F}_n)_{n \geq 0}$, then X is a martingale with respect to $(\mathcal{F}_n^X)_{n \geq 0}$. Is the converse true ? What about supermartingales and submartingales ?

Another important example is the following. Consider an integrable random variable Z and a filtration $(\mathcal{F}_n)_{n \geq 0}$ and define, for all $n \geq 0$, $X_n = \mathbb{E}[Z | \mathcal{F}_n]$. Then, for all $n \geq 0$, X_n is \mathcal{F}_n -measurable by definition of the conditional expectation, it is also integrable by definition of the conditional expectation. Thanks to the property labelled 5 in Theorem 2.7, we have

$$\mathbb{E}[X_{n+1} | \mathcal{F}_n] = \mathbb{E}[\mathbb{E}[Z | \mathcal{F}_{n+1}] | \mathcal{F}_n] = \mathbb{E}[Z | \mathcal{F}_n] = X_n.$$

Hence, $(X_n)_{n \geq 0}$ is a martingale with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$.

A martingale of the form $(\mathbb{E}[Z | \mathcal{F}_n])_{n \geq 0}$ for some integrable random variable Z is sometimes called a closed martingale.

Exercise 3.8 Is the simple random walk S a closed martingale ?

Exercise 3.9 Let again S be the simple random walk. Define $T = (T_n)_{n \geq 0}$ by setting $T_n = S_n$ if $n \leq 8888$ and $T_n = S_{8888}$ if $n > 8888$. Is T a martingale with respect to the filtration $(\mathcal{F}_n^S)_{n \geq 0}$? Is it a closed martingale ?

Proposition 3.4 Let $(X_n)_{n \geq 0}$ and $(Y_n)_{n \geq 0}$ be supermartingales. Let a be a real number.

1. If a is positive, then $(aX_n)_{n \geq 0}$ is a supermartingale.
2. If a is negative, then $(aX_n)_{n \geq 0}$ is a submartingale.
3. $(X_n + Y_n)_{n \geq 0}$ is a supermartingale.

The proposition obtained by interchanging the words supermartingale and submartingale is also true. In particular, any linear combination of martingales is a martingale.

The proof of this proposition is a verification. For the last sentence, observe that a stochastic process is a martingale if and only if it is both a supermartingale and a submartingale.

As a consequence of the conditional Jensen inequality (property 11 in Proposition 2.7), we have the following result.

Proposition 3.5 Let $(X_n)_{n \geq 0}$ be a martingale. Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be a convex function such that for all $n \geq 0$, $\phi(X_n)$ is integrable. Then $(\phi(X_n))_{n \geq 0}$ is a submartingale.

The proposition also holds if X is a submartingale and ϕ is non-decreasing.

Proof. Let us treat the case where X is a submartingale. The process $(\phi(X_n))_{n \geq 0}$ is adapted and integrable. For all $n \geq 0$, we have

$$\mathbb{E}[\phi(X_{n+1})|\mathcal{F}_n] \geq \phi(\mathbb{E}[X_{n+1}|\mathcal{F}_n]) \geq \phi(X_n),$$

where the first inequality is Jensen inequality and the second a consequence of the fact that ϕ is non-decreasing. This proves that $(\phi(X_n))_{n \geq 0}$ is a submartingale.

The case where X is a martingale is similar and easier. □

3.3 The case where the filtration is generated by partitions

Let us consider a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. On this probability space, let us consider a filtration $(\mathcal{F}_n)_{n \geq 0}$ and let us make the assumption that each σ -field \mathcal{F}_n is generated by a finite partition of Ω . In this case, the partition which generates \mathcal{F}_n is uniquely defined (as the set of elements of $\mathcal{F}_n \setminus \{\emptyset\}$ which are minimal for inclusion) and we will denote it by Π_n . The inclusion $\mathcal{F}_n \subset \mathcal{F}_{n+1}$ is equivalent to the fact that each element of Π_n is the union of some elements of Π_{n+1} (one says that Π_{n+1} is a finer partition of Ω than Π_n). For the sake of simplicity, let us also assume that $\mathcal{F}_0 = \{\emptyset, \Omega\}$, that is, $\Pi_0 = \{\Omega\}$.

Exercise 3.10 Check rigorously the two assertions made in the last paragraph.

There is a natural genealogical structure on $\bigcup_{n \geq 0} \Pi_n$. Indeed, any set $B \in \Pi_{n+1}$ is included in exactly one set $A \in \Pi_n$, and we could say that A is the father of B . This relation of fatherhood gives rise to a genealogical tree whose nodes are the elements of $\bigcup_{n \geq 0} \Pi_n$. In order to follow more easily this discussion, it may be useful to have a look at the picture on the next page.

Let us introduce a notation for the elements of $\bigcup_{n \geq 0} \Pi_n$ which takes its genealogical structure into account. The idea is that each element of Π_n will be labelled by a word of integers of length n which contains all its genealogy.

To start with, \mathcal{F}_0 is generated by $\Pi_0 = \{\Omega\}$. There is a unique word of length 0, it is the empty word, and we will accordingly set $A_\emptyset = \Omega$.

Now, \mathcal{F}_1 is generated by a partition $\Pi_1 = \{A_1, \dots, A_{n_\emptyset}\}$ of Ω . The number n_\emptyset happens to be the total number of blocks of Π_1 , but it is also more importantly the number of blocks of Π_1 contained in A_\emptyset . It is the number of children of A_\emptyset .

Then, the inductive rule is that the label of an element of Π_{n+1} is a word of $n+1$ integers whose n first letters are the label of the father of that element (the father belongs to Π_n) and whose last letter is the rank of this child in the progeny of the father. As you can guess and see on the picture, there is some freedom in the attribution of the last element, which amounts to a choice of a total order on the set of elements of Π_{n+1} which are contained in a particular element of Π_n .

Let us now consider a process $X = (X_n)_{n \geq 0}$ which is a martingale with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$. For each $n \geq 0$, X_n is measurable with respect to \mathcal{F}_n and this is equivalent to the fact that X_n is constant on each block $A_{i_1 \dots i_n}$ of Π_n . Let us denote by $x_{i_1 \dots i_n}$ the value which X_n takes on $A_{i_1 \dots i_n}$.

The martingale property can be easily formulated in terms of the numbers $x_{i_1 \dots i_n}$. Indeed, for each $n \geq 0$, the conditional expectation $\mathbb{E}[X_{n+1} | \mathcal{F}_n]$ is constant on each block $A_{i_1 \dots i_n}$ of Π_n , taking the value $\sum_{j=1}^k x_{i_1 \dots i_n j} \frac{\mathbb{P}(A_{i_1 \dots i_n j})}{\mathbb{P}(A_{i_1 \dots i_n})}$, where k is the number of children of $A_{i_1 \dots i_n}$.

The martingale property is thus equivalent to the equality

$$x_{i_1 \dots i_n} = \sum_{j=1}^k x_{i_1 \dots i_n j} \frac{\mathbb{P}(A_{i_1 \dots i_n j})}{\mathbb{P}(A_{i_1 \dots i_n})} = \frac{\sum_{j=1}^k x_{i_1 \dots i_n j} \mathbb{P}(A_{i_1 \dots i_n j})}{\sum_{j=1}^k \mathbb{P}(A_{i_1 \dots i_n j})}.$$

To summarise this discussion:

- we are considering the case of a filtration in which each σ -field is generated by a finite partition of Ω ,
- there is a genealogical structure on the set of all blocks of the partitions which generate the σ -fields of our filtration,
- a stochastic process adapted to this filtration is specified by its value on each of these blocks, that is, by a real number attached to each node of the genealogical tree,
- each node of the tree has a natural weight, which is the proportion of the probability of its father that it represents,
- a stochastic process adapted to this filtration is a martingale if and only if the value attached to each node is equal to the weighted average of the values attached to each of its children.

Exercise 3.11 Recall the notation of Exercise 2.21. Draw the first levels of the genealogical tree of the filtration $(\mathcal{F}_n)_{n \geq 0}$. Instead of labelling the sets by the integers $1, 2, 3, \dots$ use the integers starting from 0, so that $\Pi_1 = \{A_0, A_1, \dots\}$, $\Pi_2 = \{A_{00}, A_{01}, \dots, A_{10}, A_{11}, \dots\}$. Can you characterise the real numbers which belong to the set $A_{0100110}$? To the set $A_{i_1 \dots i_n}$?

Can you explain how martingales with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$ can be put in one-to-one correspondence with certain ways of writing real numbers in the circles of the following picture? What about submartingales and supermartingales?

Exercise 3.12 Let X be a non-negative martingale. By this we mean that for all $n \geq 0$, $X_n \geq 0$ almost surely. Prove that for all n and m such that $0 \leq n \leq m$, one has $X_m = 0$ on the event $\{X_n = 0\}$, that is, $X_m \mathbb{1}_{\{X_n = 0\}} = 0$.

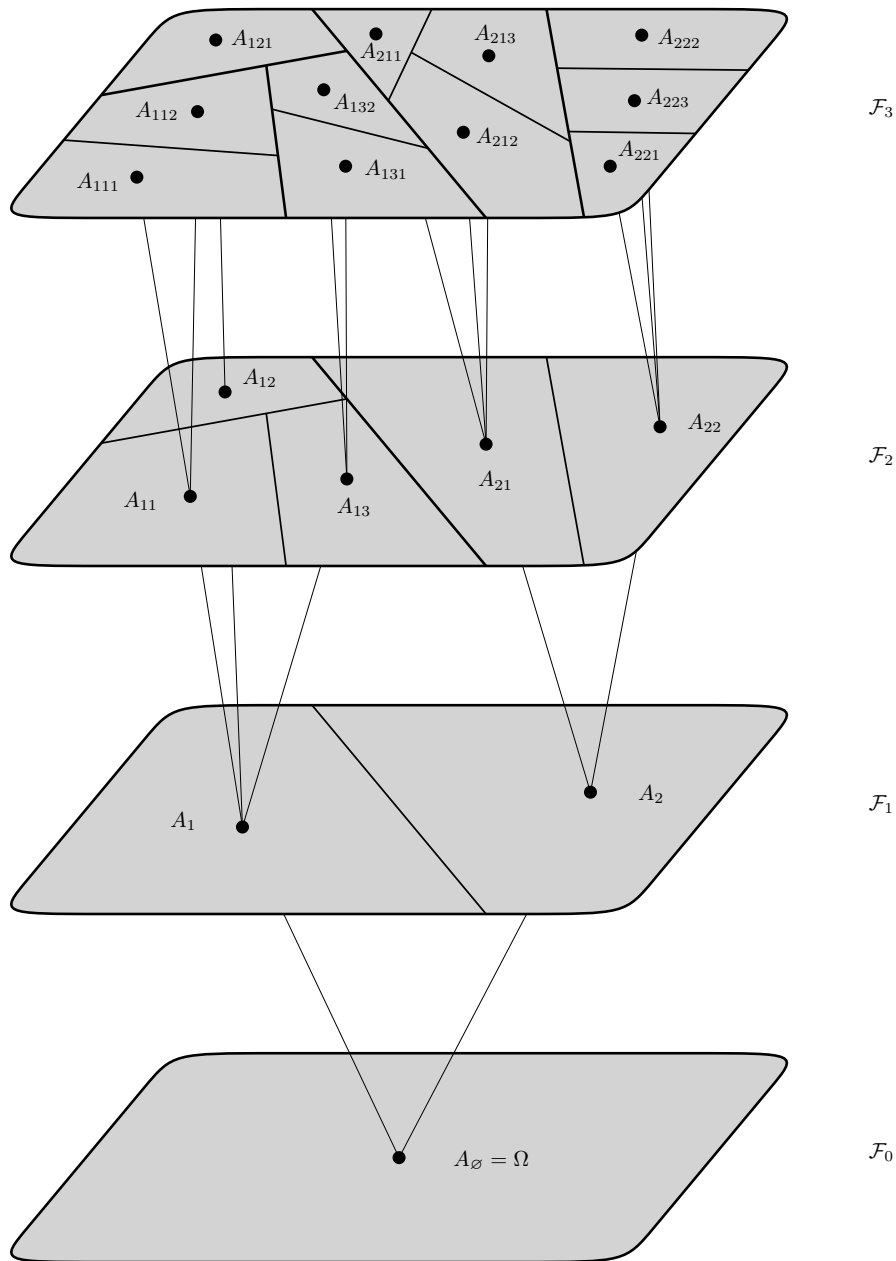


Figure 5: This illustrates the tree structure which underlies a filtration in which every σ -field is generated by a finite partition of Ω . A martingale is a process $X = (X_n)_{n \geq 0}$ such that on each block of the partition which generates \mathcal{F}_n , X_n is constant, and the average of X_{n+1} is equal to the value of X_n .

One usually phrases this property by saying that when a non-negative martingale hits zero, it stays equal to zero forever. This formulation implicitly considers the index n as a time variable, as is customary for stochastic processes.

Is this property also true of submartingales? Of supermartingales?

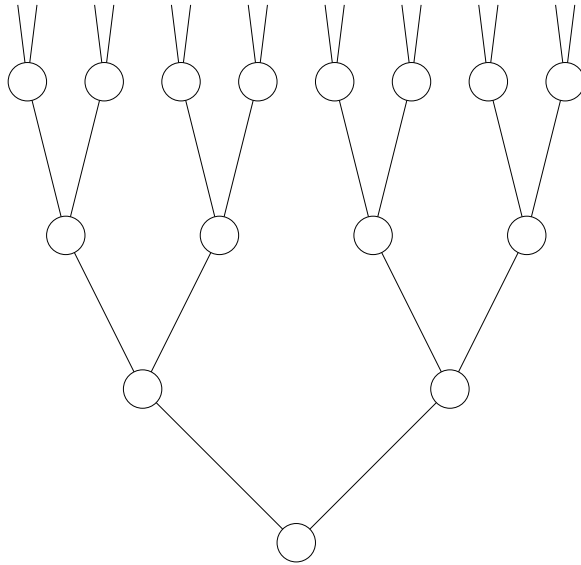


Figure 6: An infinite binary tree.

3.4 Stochastic integration

Recall the definition of the simple random walk $(S_n)_{n \geq 0}$. You can think of it as representing the (random) unfolding of a game. At each step of the game, a coin is tossed, and according to the result, the player wins or loses 1. There is a variant of the game where the player is allowed to bet any real number x before the coin is tossed. Then his or her gain or loss will be x , depending on the coin. Of course, if the game is to be fair, the player is only allowed to bet *before* the coin is tossed. This means that the amount which he or she bets for the n -th run of the game can be decided only on the basis of the results of the $n - 1$ first runs. This notion of fair betting is encoded by the following definition.

Definition 3.6 Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. A stochastic process $H = (H_n)_{n \geq 1}$ is said to be *previsible* if for each $n \geq 1$, the random variable H_n is measurable with respect to \mathcal{F}_{n-1} .

Exercise 3.13 Check that a previsible process is adapted. What can you say about a previsible martingale?

Let $(X_n)_{n \geq 0}$ be a martingale, which we think as representing the successive states of fortune of a player which bets 1 at each turn of the game. Let us consider another player which bets in a fair way, according to the values of a previsible process $(H_n)_{n \geq 1}$. What is his fortune at time n ?

Let us be more precise: X_0 is the initial fortune of the player and X_n his fortune after the n -th turn of the game. The gain of this first player during the n -th turn is thus $X_n - X_{n-1}$. The second player bets H_n just before this n -th turn, and gets $H_n(X_n - X_{n-1})$. Let us suppose that the second player starts with nothing. His fortune after the n -th turn is $H_1(X_1 - X_0) + \dots + H_n(X_n - X_{n-1})$.

Definition 3.7 (Discrete stochastic integral) Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let $X = (X_n)_{n \geq 0}$ be a martingale (or a sub- or super-martingale) and $H = (H_n)_{n \geq 1}$ a previsible process, both with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$. The stochastic integral of H with respect to X is the process $H \bullet X = ((H \bullet X)_n)_{n \geq 0}$ defined by $(H \bullet X)_0 = 0$ and, for all $n \geq 1$,

$$(H \bullet X)_n = \sum_{k=1}^n H_k (X_k - X_{k-1}).$$

The process $H \bullet X$ is also sometimes denoted by $\int H dX$. The strength of this construction lies in the following result, which says that if the correct assumptions of integrability are satisfied, then $H \bullet X$ is still a martingale when X is a martingale.

Theorem 3.8 Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let $X = (X_n)_{n \geq 0}$ be a martingale or a sub- or super-martingale and $H = (H_n)_{n \geq 1}$ be a previsible process.

1. If X is a martingale and each random variable H_n is bounded, then $H \bullet X$ is a martingale.
2. If X is a supermartingale and each random variable H_n is bounded and non-negative, then $H \bullet X$ is a supermartingale.
3. In the two previous assertions, the assumption “each random variable H_n is bounded” can be replaced by “all random variables X_n and H_n are square-integrable”.

Proof. For all $n \geq 1$, the random variable $(H \bullet X)_n$ is a function of the random variables $X_0, \dots, X_n, H_1, \dots, H_n$ which all are \mathcal{F}_n -measurable. Thus, $(H \bullet X)_n$ is \mathcal{F}_n -measurable and $H \bullet X$ is adapted.

The product of a bounded random variable with an integrable random variable is an integrable random variable. Also, by Hölder inequality, the product of two square-integrable random variables is an integrable random variable. Thus, in all cases considered in the statement, the random variable $(H \bullet X)_n$ is integrable for all $n \geq 0$.

Let us now check the main relation. Let us assume that X is a martingale. Then, for all $n \geq 0$, we must compute the conditional expectation

$$\begin{aligned} \mathbb{E}[(H \bullet X)_{n+1} | \mathcal{F}_n] &= \mathbb{E}\left[\underbrace{H_1(X_1 - X_0) + \dots + H_n(X_n - X_{n-1})}_{\mathcal{F}_n\text{-measurable}} + H_{n+1}(X_{n+1} - X_n) | \mathcal{F}_n \right] \\ &= (H \bullet X)_n + \mathbb{E}\left[\underbrace{H_{n+1}}_{\mathcal{F}_n\text{-meas.}} (X_{n+1} - X_n) | \mathcal{F}_n \right] \\ &= (H \bullet X)_n + H_{n+1} \underbrace{\mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n]}_{=0 \text{ because } X \text{ mart.}} \quad (\text{Thm. 2.7, 2}) \\ &= (H \bullet X)_n. \end{aligned}$$

If X is a supermartingale and H is non-negative, then only the last line changes. Indeed, in this case, $\mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n] \geq 0$, hence $H_{n+1} \mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n] \geq 0$ and finally $\mathbb{E}[(H \bullet X)_{n+1} | \mathcal{F}_n] \geq (H \bullet X)_n$. \square

Exercise 3.14 Choose an integer $N \geq 1$ and define H_n to be the constant random variable equal to 1 if $n \leq N$, and to 0 if $n > N$. Describe $H \bullet X$ in terms of X .

Exercise 3.15 (The clever gambler) *A clever gambler plays the fair coin tossing game described at the beginning of this section. He thinks: “The coin is random, hence it doesn’t like to repeat itself. Let me bet 1 on heads for the next turn whenever the coin gives tails, and conversely”. Define rigorously the previsible process H which corresponds to his strategy (you will have to make a choice for H_1). Describe as precisely as you can the distribution of the stochastic process $H \bullet S$. Is our player quite as clever as he thinks?*

Exercise 3.16 *Recall once again the notation of Exercise 2.21. For all $n \geq 0$, write $X_n = \mathbb{E}[X | \mathcal{F}_n]$. Prove that for every martingale M with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$ such that $M_0 = 0$, there exists a previsible process H such that $M = H \bullet X$. It may be useful to have in mind the point of view of Exercise 3.11.*

3.5 Almost sure convergence

As a beautiful application of the construction which we made in the previous section, let us prove one of the fundamental theorems on martingales. In the next statement, we use the notation $X_n^+ = \max(X_n, 0)$ for the positive part of a random variable X_n .

Theorem 3.9 (Almost sure convergence) *Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let $X = (X_n)_{n \geq 0}$ be a submartingale such that $\sup\{\mathbb{E}[X_n^+] : n \geq 0\} < \infty$. Then there exists an integrable random variable X_∞ such that the sequence $(X_n)_{n \geq 0}$ converges almost surely to X_∞ .*

This theorem is stated for submartingales. There is also a version for supermartingales, which one reads by replacing X by $-X$: the assumption must be replaced by $\sup\{\mathbb{E}[X_n^-] : n \geq 0\} < \infty$. For martingales, which are both submartingales and supermartingales, any of the two assumptions implies the conclusion.

The following corollary may seem slightly less general than Theorem 3.9, but it is in fact equivalent to it, and perhaps more easily remembered.

Corollary 3.10 *A supermartingale (resp. submartingale, resp. martingale) which is bounded in L^1 converges almost surely.*

The following exercise explains why the assumptions of Theorem 3.9 and Corollary 3.10 are equivalent. For the sake of practising the definitions, it is formulated in terms of supermartingales.

Exercise 3.17 *Let X be a supermartingale. Prove that the sequence $(\mathbb{E}[X_n^-])_{n \geq 0}$ is non-decreasing. Prove that $\sup\{\mathbb{E}[X_n^-] : n \geq 0\}$ is finite if and only if X is bounded in L^1 . Can we replace X_n^- by X_n^+ in this statement ?*

Exercise 3.18 *Let X be a (sub-, super-) martingale bounded in L^1 . Prove that the almost sure limit of X is an integrable random variable.*

A useful special case of Theorem 3.9 is the following. It is easily remembered by analogy with the fact that a non-increasing sequence of non-negative real numbers is convergent.

Corollary 3.11 *A non-negative supermartingale converges almost surely towards a limit which is an integrable random variable.*

Proof. For all $n \geq 0$, we have

$$\mathbb{E}[|X_n|] = \mathbb{E}[X_n] \leq \mathbb{E}[X_0],$$

so that a non-negative supermartingale is bounded in L^1 . The result now follows from Theorem 3.9. \square

The main tool for proving the almost sure convergence is the notion of upcrossing of an interval. Let us first define it for a deterministic sequence. Let $x = (x_n)_{n \geq 0}$ be a sequence of real numbers. Let $a < b$ be two real numbers. We call upcrossing of $[a, b]$ by the sequence x any couple (s, t) of integers such that

$$s < t \text{ and } x_s < a < b < x_t.$$

We say that two upcrossings (s_1, t_1) and (s_2, t_2) occur successively if $t_1 < s_2$ or $t_2 < s_1$.

Exercise 3.19 *Prove that the sequence x does not converge to an element of $[-\infty, +\infty]$ if and only if there exists two rational numbers a and b such that $a < b$ and such that there are infinitely many successive upcrossings of $[a, b]$ by x .*

For all $N \geq 1$, let us denote by $u_N(x; a, b)$ the number of successive upcrossings of $[a, b]$ by x before time N . More rigorously, $u_N(x; a, b)$ is the largest integer k such that there exists $2k$ integers $s_1, t_1, \dots, s_k, t_k$ such that

$$0 \leq s_1 < t_1 < \dots < s_k < t_k \leq N \text{ and for all } i \in \{1, \dots, k\}, x_{s_i} < a < b < x_{t_i}.$$

The sequence $(u_N(x; a, b))_{N \geq 1}$ is non-decreasing and we call $u_\infty(x; a, b)$ its limit, which belongs to $\mathbb{N} \cup \{+\infty\}$.

The statement of the last exercise is that x converges to an element of $[-\infty, +\infty]$ if and only if $u_\infty(x; a, b)$ is finite for all rational a and b .

Let us now turn to the probabilistic case. Let X be a stochastic process. We define for all $N \geq 1$ and for all real numbers $a < b$ the number of upcrossings $U_N(X; a, b)$ of $[a, b]$ by X before time N , which is now an integer-valued random variable. We define $U_\infty(X; a, b)$ as the almost sure limit of the non-decreasing sequence of random variables $(U_N(X; a, b))_{N \geq 1}$.

In order to prove Theorem 3.9, we are going to prove that its assumptions imply that $U_\infty(X; a, b) < +\infty$ for all rational $a < b$ almost surely. One needs to pay attention to the order of the quantifiers “for all rational $a < b$ ” and “almost surely”. Fortunately, there are countably many intervals with rational endpoints, and the two statements

$$\forall a < b \in \mathbb{Q}, \mathbb{P}(U_\infty(X; a, b) < \infty) = 1 \text{ and } \mathbb{P}(\forall a < b \in \mathbb{Q}, U_\infty(X; a, b) < \infty) = 1$$

are equivalent.

Exercise 3.20 *Check this equivalence.*

The main ingredient of the proof is an estimation of this number of upcrossings, which we will deduce from Theorem 3.8 by an appropriate choice of the previsible process H .

Proposition 3.12 (Doob's upcrossing lemma) *Let X be a supermartingale. Let a, b be two real numbers such that $a < b$. For all $n \geq 1$, one has*

$$\mathbb{E}[U_n(X; a, b)] \leq \frac{1}{b-a} \mathbb{E}[(X_n - a)^-].$$

Proof. Let us think in the same terms as before Definition 3.7. Here is the line of reasoning of a very clever gambler. “I know that the game is fair, so that whenever X_n has reached a value which is too low, it will have to compensate and increase until it takes again higher values. Let me choose two levels a and b , which I call respectively the low level and the high level. I will wait until the first time X reaches a level below a . Then I will start betting 1 at each turn, and stop as soon as X reaches a level above b . And then I will repeat this strategy. I am quite confident that I will make a lot of money that way.”

Well, we know by Theorem 3.8 that our gambler is, on average, not going to make money, but rather to lose money. However he is going to help us proving Doob's upcrossing lemma, which is certainly more important.

Let us define the previsible process H which corresponds to his strategy. We define it inductively. First set

$$H_1 = \mathbb{1}_{\{X_0 < a\}},$$

because we start betting at the first turn if and only if X_0 is below the low level a . Then suppose H_1, \dots, H_{n-1} have been defined, for some $n \geq 2$. If $H_{n-1} = 1$, then X has recently reached a level below a and we are currently betting until it exceeds b . Thus, $H_n = 1$ unless X has just reached such a level, that is, unless $X_{n-1} > b$. If on the contrary $H_{n-1} = 0$, then we are waiting until X passes below a . We thus have $H_n = 0$, unless $X_{n-1} < a$. Altogether, we set

$$H_n = \mathbb{1}_{\{H_{n-1}=1\}} \mathbb{1}_{\{X_{n-1} \leq b\}} + \mathbb{1}_{\{H_{n-1}=0\}} \mathbb{1}_{\{X_{n-1} < a\}}.$$

From this definition it follows that H_1 is $\sigma(X_0)$ -measurable, hence \mathcal{F}_0 -measurable, and for all $n \geq 2$, H_n is $\sigma(H_{n-1}, X_{n-1})$ -measurable. By induction, it follows that H_n is \mathcal{F}_{n-1} -measurable.

Finally, $H = (H_n)_{n \geq 0}$ is previsible, and since it is bounded and non-negative by definition, the second assertion of Theorem 3.8 applies, allowing us to claim that $H \bullet X$ is a supermartingale.

Incidentally, this is why our very clever gambler is losing money on average : his winnings are given by $H \bullet X$, but $\mathbb{E}[(H \bullet X)_n]$ is a non-increasing sequence.

But the main point is the following inequality : for all $n \geq 1$, we have

$$(H \bullet X)_n \geq (b-a)U_n(X; a, b) - (X_n - a)^-. \quad (2)$$

The first term is what motivates our gambler: each upcrossing of $[a, b]$ by X returns him at least $b - a$. The second term is what he forgot to take into account, and what restores the equity. Sometimes, X will reach a level below a , and it will reach much lower values before it rises again to a level higher than b . In the mean time, our gambler will have reached the limit of his own solvability.

For those of you who want to read a more rigorous proof of (2), let us introduce the sequence of times at which the gambler changes his bet. More precisely, let us set

$$S_1 = \inf\{n \geq 0 : H_{n+1} = 1\} = \inf\{n \geq 0 : X_n < a\}.$$

Then, set

$$T_1 = \inf\{n \geq S_1 : H_{n+1} = 0\} = \inf\{n \geq S_1 : X_n > b\}.$$

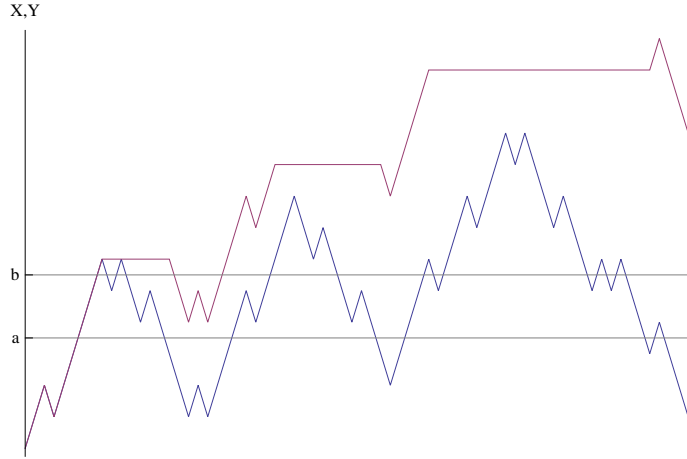


Figure 7: A sample path of X and the corresponding sample path of $H \bullet X$.

The first upcrossing of $[a, b]$ by X is thus the random interval $[S_1, T_1]$. During this interval, the gambler gains $X_{T_1} - X_{S_1}$. Then, suppose $S_1, T_1, \dots, S_k, T_k$ defined. We set

$$\begin{cases} S_{k+1} = \inf\{n \geq T_k : H_{n+1} = 1\} = \inf\{n \geq T_k : X_n < a\}, \\ T_{k+1} = \inf\{n \geq S_{k+1} : H_{n+1} = 0\} = \inf\{n \geq S_{k+1} : X_n > b\}. \end{cases}$$

Any of these random times can be infinite. If this occurs, the next times are not defined.

Now let us consider an integer n . There is, among the random times which we have defined, one which is largest among those which are smaller than n . There are two cases, depending on whether this last time is an S or a T .

1. The largest time is S_{k+1} for some k . In this case, $U_n(X; a, b) = k$. Moreover,

$$(H \bullet X)_n = (X_{T_1} - X_{S_1}) + \dots + (X_{T_k} - X_{S_k}) + (X_n - X_{S_{k+1}}).$$

The first k terms are larger than $b - a$ because $X_{T_l} > b$ and $X_{S_l} < a$ for all $l \in \{1, \dots, k\}$. The last term is larger than $X_n - a$, hence than $-(X_n - a)^-$. Indeed, recall that for all real number x , we use the notation $x^- = \max(-x, 0)$, so that the inequality $x \geq -x^-$ holds.

Thus, we have

$$(H \bullet X)_n \geq k(b - a) - (X_n - a)^+ = (b - a)U_n(X; a, b) - (X_n - a)^+$$

and (2) is proved.

2. The largest time is T_k for some k . In this case,

$$(H \bullet X)_n = (X_{T_1} - X_{S_1}) + \dots + (X_{T_k} - X_{S_k}) \geq (b - a)U_n(X; a, b)$$

and (2) is also proved.

3. There is actually a third case, where $S_1 > n$. In this case, both sides of (2) are equal to zero.

Now let us take the expectation on both sides of (2). We find

$$(b - a)\mathbb{E}[U_n(X; a, b)] - \mathbb{E}[(X_n - a)^-] \leq \mathbb{E}[(H \bullet X)_n] \leq \mathbb{E}[(H \bullet X)_0] = 0,$$

the last inequality being due to the fact that $H \bullet X$ is a supermartingale. This concludes the proof of the proposition. \square

Proof. (Theorem 3.9) Let X be a supermartingale such that $\sup\{\mathbb{E}[X_n^-] : n \geq 0\} < \infty$. Let us prove that it converges almost surely.

Consider $a, b \in \mathbb{Q}$ such that $a < b$. For all $n \geq 1$, Doob's lemma reads

$$\mathbb{E}[U_n(X; a, b)] \leq \frac{1}{b-a} \mathbb{E}[(X_n - a)^-].$$

On one hand, for all real number x , we have $(x - a)^- \leq x^- + |a|$. On the other hand, the non-decreasing sequence $(U_n(X; a, b))_{n \geq 1}$ converges to $U_\infty(X; a, b)$, so that the monotone convergence theorem entails

$$\mathbb{E}[U_\infty(X; a, b)] \leq \frac{1}{b-a} (\sup\{\mathbb{E}[X_n^-] : n \geq 0\} + |a|) < \infty$$

Hence,

$$\mathbb{P}(U_\infty(X; a, b) = \infty) = 0.$$

We have proved that for all rationals $a < b$, there are almost surely only finitely many upcrossings of $[a, b]$ by X . By the equivalence checked in Exercise 3.20, it follows that

$$\mathbb{P}(\forall a < b \in \mathbb{Q}, U_\infty(X; a, b) < \infty) = 1.$$

This in turn, by Exercise 3.19, implies that $(X_n)_{n \geq 0}$ converges almost surely. Let us denote by X_∞ its limit. There remains to prove that X_∞ is integrable.

Observe that for all $n \geq 0$, we have $\mathbb{E}[X_0] \geq \mathbb{E}[X_n] = \mathbb{E}[X_n^+] - \mathbb{E}[X_n^-]$, so that

$$\mathbb{E}[X_n^+] \leq \mathbb{E}[X_0] + \mathbb{E}[X_n^-] \leq \mathbb{E}[X_0] + \sup\{\mathbb{E}[X_n^-] : n \geq 0\}.$$

Hence, $\mathbb{E}[X_n^+]$ is also bounded, and we deduce that $\mathbb{E}[|X_n|]$ is bounded, by $\mathbb{E}[X_0] + 2 \sup\{\mathbb{E}[X_n^-] : n \geq 0\}$.

The sequence $(X_n)_{n \geq 0}$ is bounded in L^1 and converges almost surely to X_∞ . Hence, by Fatou's lemma,

$$\mathbb{E}[|X_\infty|] = \mathbb{E}[\liminf |X_n|] \leq \liminf \mathbb{E}[|X_n|] \leq \sup\{\mathbb{E}[|X_n|] : n \geq 0\} < \infty.$$

In other words, X_∞ is integrable and the proof of the theorem is finished. \square

Let us summarise the strategy of the proof.

1. We expressed the convergence of a sequence of real numbers in terms of upcrossings of intervals.
2. By integrating a well-chosen previsible process against our supermartingale, we derived a bound on the number of its upcrossings of a fixed interval.
3. Thanks to the countability of \mathbb{Q} , we deduced that a supermartingale which is bounded in L^1 has almost surely only finitely many upcrossings of every interval with rational endpoints.
4. We concluded that a supermartingale bounded in L^1 converges almost surely, and Fatou's lemma allowed us to prove that the limit is integrable.

Exercise 3.21 Let X be a supermartingale bounded in L^1 . For all real numbers a, b such that $a < b$, define

$$N_{a,b} = \{U_\infty(X; a, b) = \infty\}.$$

Compute $\mathbb{P}(\bigcup_{a < b} N_{a,b})$, where the union is taken over all real numbers $a < b$.

Exercise 3.22 Consider a sequence $(Y_n)_{n \geq 1}$ of independent random variables such that for all $n \geq 1$, one has

$$\mathbb{P}(Y_n = 0) = 1 - \frac{1}{n^2}, \quad \mathbb{P}(Y_n = e^n) = \frac{1}{2n^2} \quad \text{and} \quad \mathbb{P}(Y_n = -e^n) = \frac{1}{2n^2}.$$

Set $X_0 = 0$ and, for all $n \geq 1$, $X_n = Y_1 + \dots + Y_n$.

Prove that $(X_n)_{n \geq 0}$ is a martingale in its natural filtration. Prove that $(X_n)_{n \geq 0}$ converges almost surely to a random variable X_∞ . Prove that X_∞ is not integrable (Hint: consider the events $J_n = \{Y_1 = e, Y_n = e^n \text{ and } Y_k = 0 \text{ for all } k \notin \{1, n\}\}$).

Exercise 3.23 Let $S = (S_n)_{n \geq 0}$ be the simple random walk. Define $H = (H_n)_{n \geq 1}$ inductively by setting $H_1 = 1$ and for all $n \geq 2$,

$$H_n = \mathbb{1}_{\{H_{n-1}=1\}} \mathbb{1}_{\{S_{n-1} \neq -1\}}.$$

Explain with words what the process $H \bullet S$ is. Prove that $H \bullet S$ is a martingale. Using the following fact (which you can admit if you don't know it):

$$\mathbb{P}(\inf\{k > 0 : S_k = -1\} < \infty) = 1,$$

prove that $H \bullet S$ converges almost surely to a random variable C . What is this random variable? Is there convergence in L^1 of the sequence $(S_n)_{n \geq 1}$ towards C ?

3.6 Branching processes

Branching processes are a very important class of stochastic processes, and a rich source of examples of martingales. They constitute a simple model for the evolution of the size of a population, according to the following simple scheme.

The population starts with a certain number $\ell \geq 1$ of individuals, who constitute the 0-th generation. Then, at each step of the process, each individual of the existing population gives birth to a certain number of children, and disappears. The number of children of each individual is random, it has the same distribution for each individual, and the numbers of children of the individuals living at a certain generation are independent. Let us give a more formal definition of the process $(X_n)_{n \geq 0}$, where X_n is the size of the population at time n .

Let $(Z_{n,k})_{n,k \geq 0}$ be a family of independent identically distributed integer-valued random variables. Let us exclude the cases where these random variables are 0 almost surely, and 1 almost surely. Let us define inductively $X_0 = \ell$ and, for each $n \geq 0$,

$$X_{n+1} = \sum_{k=1}^{X_n} Z_{n,k}.$$

Let us define a filtration by setting $\mathcal{F}_0 = \{\emptyset, \Omega\}$ and, for all $n \geq 1$,

$$\mathcal{F}_n = \sigma(Z_{m,k} : m < n, k \geq 1).$$

Let m denote the common expectation of the random variables $Z_{n,k}$, that is, the average number of children of an individual of our population. It is a positive real number or $+\infty$, for the case $\mathbb{P}(Z_{n,k} = 0) = 1$ is excluded. Let us assume that $m < +\infty$.

The following result is extremely helpful in the study of branching processes, and it is an instance of the general idea that it is useful in the study of a random dynamical system to identify that certain quantities are martingales.

Proposition 3.13 *The process $(m^{-n}X_n)_{n \geq 0}$ is a martingale with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$.*

Proof. By definition, X_0 is \mathcal{F}_0 -measurable. Then, for all $n \geq 1$, X_n is a function of X_{n-1} and $\{Z_{n-1,k} : k \geq 1\}$. By induction, and by definition of \mathcal{F}_n , it follows that X_n is \mathcal{F}_n -measurable. Thus, $(X_n)_{n \geq 0}$ is adapted, and so is $(m^{-n}X_n)_{n \geq 0}$.

Let us now prove by induction on n that X_n is integrable. For $n = 0$, this is true. Let us assume that X_n is integrable. Then, since X_{n+1} is non-negative, we may consider its expectation, and we must prove that it is not equal to ∞ . We have

$$\mathbb{E}[X_{n+1}] = \sum_{x=0}^{\infty} \mathbb{E} \left[\sum_{k=1}^x Z_{n,k} \mathbb{1}_{X_n=x} \right].$$

Since X_n is \mathcal{F}_n -measurable and the random variables $\{Z_{n,k} : k \geq 1\}$ are independent of \mathcal{F}_n , we find

$$\begin{aligned} \mathbb{E}[X_{n+1}] &= \sum_{x=0}^{\infty} \sum_{k=1}^x \mathbb{E}[Z_{n,k}] \mathbb{P}(X_n = x) \\ &= \sum_{x=0}^{\infty} xm \mathbb{P}(X_n = x) \\ &= m \mathbb{E}[X_n]. \end{aligned}$$

Thus, $\mathbb{E}[X_n] = m^n \ell$ is finite.

Let us finally compute $\mathbb{E}[X_{n+1} | \mathcal{F}_n]$. We have

$$\begin{aligned} \mathbb{E}[X_{n+1} | \mathcal{F}_n] &= \mathbb{E} \left[\sum_{k=1}^{X_n} Z_{n,k} \middle| \mathcal{F}_n \right] \\ &= \sum_{k=1}^{X_n} \mathbb{E}[Z_{n,k}] \\ &= mX_n. \end{aligned}$$

It follows from this result that $(m^{-n}X_n)_{n \geq 0}$ is a martingale. □

Exercise 3.24 *Check the passage from the first to the second line in the last computation of this proof.*

We may apply Corollary 3.11 to the non-negative martingale $(m^{-n}X_n)_{n \geq 0}$, which thus converges almost surely towards an integrable random variable Y . There are three fairly different cases, depending on the value of m .

- The case $m < 1$. In this case, the almost sure convergence of $(m^{-n}X_n)_{n \geq 0}$ to Y implies that almost surely, $X_n = m^n(m^{-n}X_n)$ converges to $0 \times Y = 0$. Since $(X_n)_{n \geq 0}$ is an integer-valued process, this forces it to be almost surely stationary with limiting value 0.

Indeed, a sequence of integers is convergent if and only if it is stationary¹. Let us write our conclusion in symbols:

$$\mathbb{P}(\exists N \geq 1, \forall n \geq 1, X_n = 0) = 1.$$

In words, our conclusion is that when $m < 1$, the population gets extinct with probability 1. This case is called the *subcritical* case.

- The case $m = 1$. In this case, $(X_n)_{n \geq 1}$ itself converges almost surely to Y . Because X is integer-valued, this convergence implies that it is almost surely stationary. In symbols:

$$\mathbb{P}(\exists p \geq 0, \exists N \geq 1, \forall n \geq 1, X_n = p) = 1.$$

For each integer $p \geq 0$, we can consider the event

$$S_p = \{\exists N \geq 1, \forall n \geq 1, X_n = p\}$$

on which X is stationary at p . We claim that for all $p \geq 1$, the event S_p has probability 0. Indeed, choose $p \geq 1$ and rewrite the event S_p as

$$S_p = \{\exists N \geq 1, \forall n \geq N, Z_{n,1} + \dots + Z_{n,p} = p\} = \bigcup_{N \geq 1} \bigcap_{n \geq N} \{Z_{n,1} + \dots + Z_{n,p} = p\}.$$

The events $(\{Z_{n,1} + \dots + Z_{n,p} = p\})_{n \geq 1}$ are independent and they all have the same probability. Since the case $\mathbb{P}(Z_{n,k} = 1) = 1$ is excluded, and since $\mathbb{E}[Z_{n,k}] = 1$, we must have $\mathbb{P}(Z_{n,k} = 0) > 0$ (check this assertion!). So, in particular,

$$\mathbb{P}(Z_{n,1} + \dots + Z_{n,p} = p) \leq 1 - \mathbb{P}(Z_{n,1} + \dots + Z_{n,p} = 0) = 1 - \mathbb{P}(Z_{n,k} = 0)^p < 1.$$

The version of Borel-Cantelli's lemma for independent events now implies that $\mathbb{P}(S_p) = 0$.

Since for all $p \geq 1$, the event on which X is stationary with limiting value p has probability 0, it must be that X converges to 0 almost surely. In this case too, the population gets extinct with probability 1. This case is called the *critical* case.

- The case $m > 1$. In this *supercritical* case, which we will not treat in detail, one can show that there is a positive probability that the population survives forever.

Exercise 3.25 Prove that in the case where $m = 1$, the martingale $(X_n)_{n \geq 0}$ does not converge in L^1 to its almost sure limit. How do you intuitively understand this fact?

Exercise 3.26 Prove that in a subcritical branching process, the total population $\sum_{n \geq 0} X_n$ is not only finite almost surely, but has a finite expectation. Compute this expectation. What about the critical case?

¹Let me spell this out, since it is so important in the present problem. Let $a = (a_n)_{n \geq 0}$ be a sequence of elements of \mathbb{Z} . If a is stationary, it is of course convergent. Let us prove the converse. Assume that a is convergent. Then it is in particular Cauchy. Applying the definition of a Cauchy sequence with $\varepsilon = \frac{1}{2}$, we find that there exists an integer N such that for all $n \geq N$, we have $|a_n - a_N| < \frac{1}{2}$. Since a_n and a_N are both integers, they must be equal. Hence, a is stationary after rank N .

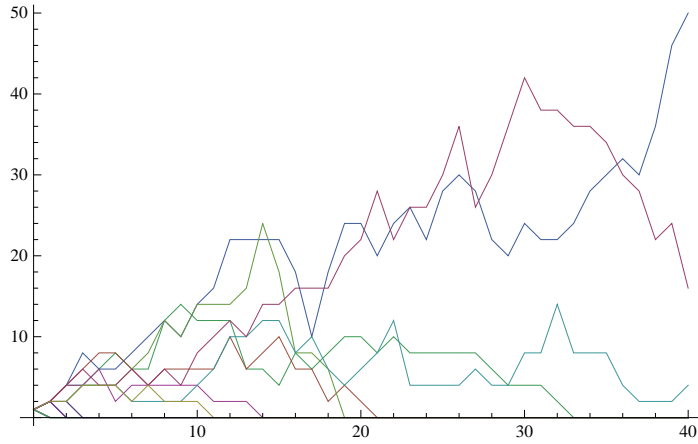


Figure 8: Twenty sample paths of X when $\ell = 1$ and $\mathbb{P}(Z_{n,k} = 0) = \mathbb{P}(Z_{n,k} = 2) = \frac{1}{2}$.

Exercise 3.27 *This exercise is a preparation for the next. Prove that every random variable Y with values in \mathbb{N} satisfies the inequality*

$$\mathbb{P}(Y \geq 1) \geq \frac{\mathbb{E}[Y]^2}{\mathbb{E}[Y^2]}.$$

Exercise 3.28 *Prove that for a subcritical branching process, there exists a positive real $\alpha > 0$ such that for all $n \geq 0$,*

$$\mathbb{P}(X_n \geq 1) \leq e^{-\alpha n}.$$

In words, the probability of survival of the population up to time n decays exponentially fast.

On the other hand, consider a critical branching process such that the random variables $Z_{n,k}$ are square-integrable. Set $\sigma^2 = \text{Var}(Z_{n,k})$. Prove that

$$\mathbb{E}[X_n^2] = n\ell\sigma^2 + \ell(\ell - \sigma^2)$$

and deduce that there exists a positive real $\beta > 0$ such that

$$\mathbb{P}(X_n \geq 1) \geq \frac{\beta}{n}.$$

3.7 Convergence in L^1

Critical branching processes are examples of non-negative martingales which converge almost surely (as they must according to Corollary 3.11), but not in L^1 (see Exercise 3.25). In this section, we clarify what it means for a martingale to be convergent in L^1 .

Theorem 3.14 *Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let $X = (X_n)_{n \geq 0}$ be a martingale. The following two conditions are equivalent:*

1. *The martingale X converges to a random variable X_∞ almost surely and in L^1 .*

2. *There exists an integrable random variable Z such that $X_n = \mathbb{E}[Z | \mathcal{F}_n]$ for all $n \geq 0$.*

Moreover, if these conditions are satisfied, then one can take $Z = X_\infty$ in 2.

Proof. 1 \Rightarrow 2. Choose $n \geq 0$. For all $m \geq n$, we have $X_n = \mathbb{E}[X_m | \mathcal{F}_n]$. We have

$$\begin{aligned} \mathbb{E}[|\mathbb{E}[X_m | \mathcal{F}_n] - \mathbb{E}[X_\infty | \mathcal{F}_n]|] &= \mathbb{E}[|\mathbb{E}[X_m - X_\infty | \mathcal{F}_n]|] \\ &\leq \mathbb{E}[\mathbb{E}[|X_m - X_\infty| | \mathcal{F}_n]] \\ &= \mathbb{E}[|X_m - X_\infty|], \end{aligned}$$

so that $\mathbb{E}[X_m | \mathcal{F}_n]$ converges to $\mathbb{E}[X_\infty | \mathcal{F}_n]$ in L^1 as n tends to infinity. Since $\mathbb{E}[X_m | \mathcal{F}_n]$ does in fact not depend on $m \geq n$, we even have $\mathbb{E}[X_m | \mathcal{F}_n] = \mathbb{E}[X_\infty | \mathcal{F}_n]$. Finally,

$$X_n = \mathbb{E}[X_\infty | \mathcal{F}_n].$$

2 \Rightarrow 1. For all $n \geq 0$, one has

$$\mathbb{E}[|X_n|] = \mathbb{E}[|\mathbb{E}[Z | \mathcal{F}_n]|] \leq \mathbb{E}[\mathbb{E}[|Z| | \mathcal{F}_n]] = \mathbb{E}[|Z|].$$

Hence, $(X_n)_{n \geq 0}$ is bounded in L^1 and, thanks to Theorem 3.9, converges almost surely to an integrable random variable X_∞ . We need to prove that $(X_n)_{n \geq 0}$ converges in L^1 to X_∞ .

Let us do this first in the case where Z is bounded by a constant M , that is, under the assumption that $\mathbb{P}(|Z| \leq M) = 1$. Then, for all $n \geq 1$, the positivity of the conditional expectation implies $\mathbb{P}(|X_n| \leq M) = 1$. We can conclude, by the dominated convergence theorem, that $(X_n)_{n \geq 1}$ converges in L^1 to X_∞ .

Let us now treat the general case. Let us fix $\varepsilon > 0$. There exists a constant $M > 0$ such that

$$\mathbb{E}[|Z - Z\mathbb{1}_{|Z| \leq M}|] < \varepsilon.$$

Let us choose such an M and set $\tilde{Z} = Z\mathbb{1}_{|Z| \leq M}$. For all $n \geq 0$, let us define $\tilde{X}_n = \mathbb{E}[\tilde{Z} | \mathcal{F}_n]$. Then for all $n \geq 0$,

$$\mathbb{E}[|X_n - \tilde{X}_n|] = \mathbb{E}[|\mathbb{E}[Z - \tilde{Z} | \mathcal{F}_n]|] \leq \mathbb{E}[|Z - \tilde{Z}|] < \varepsilon.$$

Since \tilde{Z} is a bounded random variable, we know from our study of the bounded case that the martingale $(\tilde{X}_n)_{n \geq 0}$ converges in L^1 . Thus, it is a Cauchy sequence in L^1 , which is to say that there exists $n_0 \geq 1$ such that for all $n, m \geq n_0$, one has

$$\mathbb{E}[|\tilde{X}_n - \tilde{X}_m|] < \varepsilon.$$

Then, for all $n, m \geq n_0$, we have

$$\mathbb{E}[|X_n - X_m|] \leq \mathbb{E}[|X_n - \tilde{X}_n|] + \mathbb{E}[|\tilde{X}_n - \tilde{X}_m|] + \mathbb{E}[|\tilde{X}_m - X_m|] < 3\varepsilon.$$

Hence, the sequence $(X_n)_{n \geq 0}$ is Cauchy in L^1 . Thus, it converges in L^1 , and it must be towards X_∞ . \square

In this proof, we used several facts which it may be useful to review.

Exercise 3.29 Prove that $\|\mathbb{E}[X | \mathcal{F}] - \mathbb{E}[Y | \mathcal{F}]\|_{L^1} \leq \|X - Y\|_{L^1}$. In other words, as a mapping of L^1 into itself, conditional expectation is 1-Lipschitz. Prove in fact that, as a linear mapping from L^1 into itself, it has norm exactly 1. What about the norm of $\mathbb{E}[\cdot | \mathcal{F}]$ as a linear operator on L^p ?

Exercise 3.30 Let Z be an integrable random variable. Let $\varepsilon > 0$ be fixed. Prove that there exists $M > 0$ such that $\|Z - Z\mathbb{1}_{|Z| \leq M}\|_{L^1} < \varepsilon$.

Exercise 3.31 Prove that if a sequence of random variables converges almost surely and in L^1 , it must be to the same limit.

The next exercise provides us with a refinement of the preceding theorem.

Exercise 3.32 Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let Z be an integrable random variable. Let X_∞ the almost sure and L^1 limit of the martingale $(\mathbb{E}[Z|\mathcal{F}_n])_{n \geq 0}$. Does the equality $X_\infty = Z$ necessarily hold ?

Define the σ -field $\mathcal{F}_\infty = \sigma\left(\bigcup_{n \geq 0} \mathcal{F}_n\right)$. Prove that X_∞ is \mathcal{F}_∞ -measurable. Prove that for all $A \in \bigcup_{n \geq 0} \mathcal{F}_n$, one has

$$\int_A Z \, dP = \int_A X_\infty \, dP.$$

Prove, using the monotone class theorem, that the same equality holds for all $A \in \mathcal{F}_\infty$. What is the conclusion of this argument ?

So far, we have studied the almost sure convergence of martingales, observed that it is not always a convergence in L^1 , and understood that it is equivalent for a martingale to be a closed martingale or to be convergent in L^1 . If time allows, we shall come back to the question of convergence in L^1 . Before that, we want to study convergence in L^p for $p > 1$. As often, the case $p = 2$ is particularly pleasant, but there is a very beautiful theory for general p which relies on Doob's maximal inequality. In order to study it, we will need a fundamental tool in the study of martingales, which is the theory of stopping. This is the subject of the next section.

3.8 Stopping times

Let us remember our understanding of a martingale in the picture involving a gambler. A gambler does not usually play forever, he stops playing at a certain point. The reason why he stops may be a mixture of various ingredients : lack of time, lack of money, bad luck, sense that he has won all that he could, time for dinner. In any case, the time at which the gambler stops playing can depend on how the game went : it is normally a random variable. However, if things are to be fair, the decision of stopping cannot depend on any information about the future of the game. The random time at which the player stops must thus have the following property : the decision of stopping or not at time n must be taken in a deterministic way from the information available at time n . The following definition encodes this property.

Definition 3.15 Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. A random variable $T : \Omega \rightarrow \mathbb{N} \cup \{\infty\}$ is a stopping time if for all $n \geq 0$, the event $\{T = n\}$ belongs to \mathcal{F}_n .

For example, a random time which is almost surely constant is a stopping time.

Exercise 3.33 Prove that if T is a stopping time, then $\{T = \infty\}$ belongs to \mathcal{A} , and in fact to $\sigma\left(\bigcup_{n \geq 0} \mathcal{F}_n\right)$

Exercise 3.34 Prove that T is a stopping time if and only if for all $n \geq 0$, the event $\{T \leq n\}$ belongs to \mathcal{F}_n .

Observe that the definition of a stopping time does not involve the probability \mathbb{P} . It only depends on the filtration. Here are a few simple but useful properties of stopping times.

Lemma 3.16 Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let S, T be stopping times. Then $S \wedge T = \min(S, T)$, $S \vee T = \max(S, T)$, $S + T$ are stopping times.

The proof is left as an exercise. A fundamental example of stopping time is the hitting time of a Borel subset by an adapted process.

Exercise 3.35 Let $X = (X_n)_{n \geq 0}$ be an adapted process. Let B be a Borel subset of \mathbb{R} . Prove that $H = \inf\{n \geq 0 : X_n \in B\}$, with the convention $\inf \emptyset = \infty$, is a stopping time.

As their name indicates, stopping times are meant to stop processes, in particular martingales. In fact, any integer-valued random variable can be used to stop any stochastic process, according to the following definition.

Definition 3.17 Let $X = (X_n)_{n \geq 0}$ be a stochastic process. Let T be a random variable with values in $\mathbb{N} \cup \{\infty\}$. The stopped process $X^T = (X_{T \wedge n})_{n \geq 0}$ is defined by setting, for all $n \geq 0$ and all $\omega \in \Omega$,

$$X_{T \wedge n}(\omega) = X_{T(\omega) \wedge n}(\omega) = \begin{cases} X_n(\omega) & \text{if } T(\omega) > n, \\ X_m(\omega) & \text{if } T(\omega) = m \leq n. \end{cases}$$

If T is finite almost surely, then the random variable X_T is well defined by the formula

$$X_T(\omega) = X_{T(\omega)}(\omega).$$

Observe that in general, for X_T to be defined, one need to make sure that T is finite almost surely. If one knows that the sequence $(X_n)_{n \geq 0}$ converges almost surely to X_∞ , one might however remove this assumption and set $X_T = X_\infty$ on the event $\{T = \infty\}$.

The first important result about stopping of martingales by stopping times is the following and we will deduce it from Theorem 3.8.

Proposition 3.18 Let X be a supermartingale. Let T be a stopping time. The process X^T is a supermartingale.

Proof. In the picture of the gambler, stopping a process at time T amounts to playing 1 until time T and then playing 0 forever. Let us define accordingly, for all $n \geq 1$,

$$H_n = \mathbb{1}_{\{T \geq n\}} = \mathbb{1}_{\{T \leq n-1\}^c}.$$

We play 1 during the n -th turn if and only if we did not decide to stop just after the end of the $n-1$ -th turn. By construction, H is previsible. It is bounded and non-negative. Thus, $H \bullet X$ is

a supermartingale. Moreover, for almost all ω , we have

$$\begin{aligned}
(H \bullet X)_n(\omega) &= \sum_{k=1}^n H_k(\omega)(X_k(\omega) - X_{k-1}(\omega)) \\
&= \sum_{k=1}^n \mathbb{1}_{T(\omega) \geq k} (X_k(\omega) - X_{k-1}(\omega)) \\
&= \sum_{k=1}^{T(\omega) \wedge n} X_k(\omega) - X_{k-1}(\omega) \\
&= X_{T(\omega) \wedge n}(\omega) - X_0(\omega).
\end{aligned}$$

Thus, X^T itself is a supermartingale. □

In particular, $\mathbb{E}[X_{T \wedge n}] \leq \mathbb{E}[X_0]$ and if X is a martingale, the equality holds : one has $\mathbb{E}[X_{T \wedge n}] = \mathbb{E}[X_0]$. It is tempting to let n tend to infinity in this equation, but one must be very cautious in doing so. In order to be able to say that $\mathbb{E}[X_T] = \mathbb{E}[X_0]$, one must usually combine properties of X and T , in one of many possible ways. Let us see what can go wrong.

Firstly, as we already mentioned, the sequence $(X_{T \wedge n})_{n \geq 0}$ may not converge almost surely. Indeed, the event $\{T = +\infty\}$ may have positive probability, and X may not converge on it. However, if $T < \infty$ almost surely, it is true that $X_{T \wedge n}$ converges almost surely to X_T . Nevertheless, this convergence may not happen in L^1 . The classical counter-example is that of the simple random walk (see Exercise 3.23). Indeed, let T be the hitting time of -1 for S , that is, $T = \inf\{n : S_n = -1\}$. Then T is finite almost surely, and S^T converges almost surely to -1 , although $0 = \mathbb{E}[S_{T \wedge n}] \neq \mathbb{E}[S_T] = -1$.

Let us give three simple sets of conditions under which the expected result holds.

Theorem 3.19 *Let X be a supermartingale. Let T be a stopping time. If one of the following conditions is satisfied :*

1. T is bounded,
2. T is integrable and there exists $M > 0$ such that for all $n \geq 0$, $|X_{n+1} - X_n| \leq M$,
3. T is almost surely finite and X^T is bounded,

then X_T is integrable and the inequality $\mathbb{E}[X_T] \leq \mathbb{E}[X_0]$ holds.

Observe that the three successive sets of assumptions are in decreasing order of strength on T , and of increasing order of strength on X .

Proof. 1. If T is bounded by an integer N , then $\mathbb{E}[X_0] \geq \mathbb{E}[X_{T \wedge N}] = \mathbb{E}[X_T]$.

2. Since T is integrable, it is almost surely finite, so that $X_{T \wedge n}$ converges almost surely to X_T . Moreover, for all $n \geq 1$,

$$|X_{T \wedge n} - X_0| \leq \sum_{k=0}^{(T \wedge n) - 1} |X_{k+1} - X_k| \leq MT.$$

Hence, the sequence $(X_{T \wedge n})_{n \geq 0}$ is dominated by the integrable random variable $MT + |X_0|$ and the dominated convergence theorem allows us to conclude that $X_{T \wedge n}$ converges in L^1 to X_T . In

particular, X_T is integrable and $\mathbb{E}[X_T] = \lim_{n \rightarrow \infty} \mathbb{E}[X_{T \wedge n}] \leq \mathbb{E}[X_0]$.

3. If T is finite almost surely, then $X_{T \wedge n}$ converges almost surely to X_T . If moreover X^T is bounded, then the dominated convergence theorem ensures that the convergence holds in L^1 , and we conclude as in the previous case. \square

Exercise 3.36 Let S be the simple random walk. Consider two integers a and b such that $a < 0 < b$. Prove that S will almost surely visit the set $\{a, b\}$ (you can prove this without using any sophisticated results on the simple random walk). Compute the probability that S hits a before hitting b . Application : a gambler starts with a fortune of 1. He has decided to stop playing as soon as his fortune reaches 100, or when he has no money anymore. What is the probability that he returns home with empty pockets ?

Exercise 3.37 Prove that a random variable T with values in $\mathbb{N} \cup \{+\infty\}$ is a stopping time if and only if the process $H = (H_n)_{n \geq 1}$ defined by $H_n = \mathbb{1}_{\{T \geq n\}}$ is previsible.

Exercise 3.38 Let $H = (H_n)_{n \geq 1}$ be a stochastic process. Prove that the following two assertions are equivalent.

1. H is adapted.
2. For every adapted process X , the process $H \bullet X$ is adapted.

Exercise 3.39 Let T be a random variable with values in $\mathbb{N} \cup \{+\infty\}$. Prove that the following assertions are equivalent.

1. For all $n \geq 0$, the event $\{T = n\}$ belongs to \mathcal{F}_{n+1} .
2. For every adapted process X , the stopped process X^T is adapted.
3. For every martingale X , the stopped process X^T is adapted.

3.9 Convergence in L^p for $p > 1$

In this section, we shall prove that a supermartingale which is bounded in L^p for some real $p > 1$ converges almost surely and in L^p . We already know that it converges almost surely, for boundedness in L^p implies boundedness in L^1 . The point is to establish the convergence in L^p . For this, we shall use Doob's maximal inequality, for which we need the following lemma.

Lemma 3.20 Let X be a submartingale. Let S and T be two bounded stopping times such that $S \leq T$ almost surely. Then $\mathbb{E}[X_S] \leq \mathbb{E}[X_T]$.

Proof. Let our gambler start playing at time S and stop at time T . This amounts to defining, for all $n \geq 1$,

$$H_n = \mathbb{1}_{\{S < n \leq T\}} = \mathbb{1}_{\{S \leq n-1\}} \mathbb{1}_{\{T \leq n-1\}}^c.$$

The way we wrote it makes it clear that $H = (H_n)_{n \geq 0}$ is previsible. The same computation as in the proof of Proposition 3.18 yields

$$(H \bullet X)_n = X_{T \wedge n} - X_{S \wedge n}.$$

Let N be an integer such that $S \leq T \leq N$ almost surely. Then $(H \bullet X)_N = X_T - X_S$. Since H is bounded and non-negative, $H \bullet X$ is a submartingale, so that

$$\mathbb{E}[X_T - X_S] = \mathbb{E}[(H \bullet X)_N] \geq \mathbb{E}[(H \bullet X)_0] = 0,$$

as expected. □

Let us use this lemma to prove Doob's maximal inequality.

Proposition 3.21 (Doob's maximal inequality) *Let X be a submartingale. For all integer $n \geq 0$ and all real number a , the following inequalities hold :*

$$a\mathbb{P}\left(\sup_{0 \leq k \leq n} X_k \geq a\right) \leq \mathbb{E}[X_n \mathbb{1}_{\{\sup_{0 \leq k \leq n} X_k \geq a\}}] \leq \mathbb{E}[X_n^+].$$

Proof. Choose a real number a . Define T as the first hitting time of $[a, +\infty)$ for X :

$$T = \inf\{n \geq 0 : X_n \geq a\}.$$

Consider now an integer $n \geq 0$. The random variable $X_{T \wedge n}$ is equal to X_T , hence greater or equal to a , on the event $\{T \leq n\}$. On the complement of this event, $X_{T \wedge n}$ is equal to X_n . Observing that the events $\{T \leq n\}$ and $\{\sup_{0 \leq k \leq n} X_k \geq a\}$ are the same, we thus find

$$X_{T \wedge n} \geq a \mathbb{1}_{\{\sup_{0 \leq k \leq n} X_k \geq a\}} + X_n \mathbb{1}_{\{\sup_{0 \leq k \leq n} X_k < a\}}.$$

Taking the expectation of both sides of this inequality, we find

$$a\mathbb{P}\left(\sup_{0 \leq k \leq n} X_k \geq a\right) \leq \mathbb{E}[X_{T \wedge n}] - \mathbb{E}[X_n \mathbb{1}_{\{\sup_{0 \leq k \leq n} X_k < a\}}].$$

The previous lemma applied to the bounded stopping times $T \wedge n$ and n yields the inequality $\mathbb{E}[X_{T \wedge n}] \leq \mathbb{E}[X_n]$, which combined with the previous one implies that

$$a\mathbb{P}\left(\sup_{0 \leq k \leq n} X_k \geq a\right) \leq \mathbb{E}[X_n(1 - \mathbb{1}_{\{\sup_{0 \leq k \leq n} X_k < a\}})] = \mathbb{E}[X_n \mathbb{1}_{\{\sup_{0 \leq k \leq n} X_k \geq a\}}].$$

This proves the first inequality. The second follows from an instance of the inequality $Z \mathbb{1}_A \leq Z^+$ which is true for any random variable Z and any event A . □

Although one would quite naturally apply this result with a positive value for a , it seems that the result and the proof do not depend on any assumption on the sign of a . It is interesting to take a moment to think about what the theorem says when $a = 0$ (and even in this case the theorem is not trivial), and when $a < 0$.

Let us now take one further step towards the L^p convergence. For this, let us introduce a notation. If $X = (X_n)_{n \geq 0}$ is a stochastic process, let us define its maximal process $X^* = (X_n^*)_{n \geq 0}$ by setting, for all $n \geq 0$,

$$X_n^* = \sup_{0 \leq k \leq n} |X_k|.$$

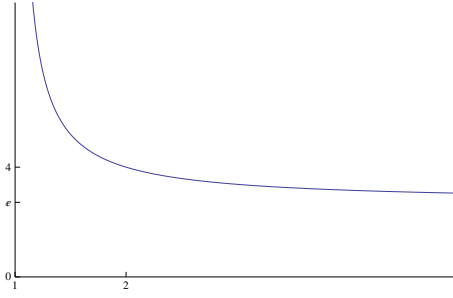


Figure 9: The graph of the function $p \mapsto \left(\frac{p}{p-1}\right)^p$.

Proposition 3.22 *Let X be a martingale. Let $p > 1$ be a real number. For all integer $n \geq 1$, one has*

$$\mathbb{E}[(X_n^*)^p] \leq \left(\frac{p}{p-1}\right)^p \mathbb{E}[|X_n|^p].$$

Observe that we are only considering non-negative quantities in this statement. Hence, we do not need to assume that X belongs to L^p .

Proof. By the previous proposition applied to the submartingale $(|X_n|)_{n \geq 0}$, we have, for all $a > 0$,

$$a\mathbb{P}(X_n^* \geq a) \leq \mathbb{E}[|X_n| \mathbb{1}_{\{X_n^* \geq a\}}].$$

Let us multiply both sides of this inequality by a^{p-2} and integrate with respect to a from 0 to $+\infty$. On the left, we find

$$\int_0^{+\infty} a^{p-1} \mathbb{P}(X_n^* \geq a) da = \mathbb{E} \left[\int_0^{X_n^*} a^{p-1} da \right] = \frac{1}{p} \mathbb{E}[(X_n^*)^p].$$

On the right, we find

$$\begin{aligned} \int_0^{+\infty} a^{p-2} \mathbb{E}[|X_n| \mathbb{1}_{\{X_n^* > a\}}] da &= \mathbb{E} \left[|X_n| \int_0^{X_n^*} a^{p-2} da \right] \\ &= \frac{1}{p-1} \mathbb{E}[|X_n| (X_n^*)^{p-1}] \\ &\leq \frac{1}{p-1} \mathbb{E}[|X_n|^p]^{\frac{1}{p}} \mathbb{E}[(X_n^*)^p]^{\frac{p-1}{p}}, \end{aligned}$$

where the last inequality is Hölder's inequality. Combining the two expressions, we find the expected inequality. \square

As n tends to infinity, X_n^* increases and converges almost surely to

$$X_\infty^* = \sup\{|X_n| : n \geq 0\}.$$

Let us now prove the convergence result in L^p .

Theorem 3.23 *Let X be a martingale. Let $p > 1$ be a real number. Assume that X is bounded in L^p . Then X converges almost surely and in L^p towards a random variable X_∞ which satisfies*

$$E[|X_\infty|^p] = \sup\{\mathbb{E}[|X_n|^p] : n \geq 0\}.$$

Moreover, one has

$$\mathbb{E}[(X_\infty^*)^p] \leq \left(\frac{p}{p-1}\right)^p \mathbb{E}[|X_\infty|^p].$$

Proof. Let us assume that X is bounded in L^p . In particular, X is bounded in L^1 , hence it converges almost surely to a random variable X_∞ .

The previous proposition, the fact that $(X_n^*)_{n \geq 0}$ is a non-decreasing sequence which converges towards X_∞^* and the monotone convergence theorem imply that

$$\mathbb{E}[(X_\infty^*)^p] \leq \left(\frac{p}{p-1}\right)^p \sup\{\mathbb{E}[|X_n|^p] : n \geq 0\} < \infty.$$

Hence, the almost sure convergence of X_n to X_∞ is dominated in L^p by X_∞^* . By the dominated convergence theorem, the convergence holds in L^p . In particular, $\mathbb{E}[|X_\infty|^p]$ is the limit of the sequence $(\mathbb{E}[|X_n|^p])_{n \geq 0}$ which, by Jensen's inequality, is non-decreasing. Hence,

$$E[|X_\infty|^p] = \sup\{\mathbb{E}[|X_n|^p] : n \geq 0\}.$$

The last inequality follows immediately. □

The results of this section, as well as those of Section 3.5, fall in a category that one could call “rigidity results” for martingales. Doob’s upcrossing lemma, which was the main result needed to prove the almost sure convergence of martingales bounded in L^1 , can be rephrased by saying that if a martingale oscillates a lot, then its L^1 norm becomes large. Doob’s maximal inequality says that if the largest value taken by a martingale is large in L^p , then the martingale itself is large in L^p . Neither of these results hold, even in very weak forms, for arbitrary stochastic processes.

Exercise 3.40 *The two equalities*

$$\sum_{n=1}^{\infty} \frac{1}{n} = \infty \quad \text{and} \quad \sum_{n=1}^{\infty} \frac{(-1)^n}{n} = -\log 2$$

are well known. Consider now a sequence $(\varepsilon_n)_{n \geq 1}$ of i.i.d. random variables such that $\mathbb{P}(\varepsilon_1 = 1) = \mathbb{P}(\varepsilon_1 = -1) = \frac{1}{2}$. Study the random series

$$\sum_{n=1}^{\infty} \frac{\varepsilon_n}{n} \quad \text{and more generally} \quad \sum_{n=1}^{\infty} \frac{\varepsilon_n}{n^s},$$

where s is an arbitrary complex number.

3.10 Square-integrable martingales

As often in analysis when something can be done for all L^p spaces, the case $p = 2$ enjoys special properties. Martingales are no exception, the more so since the conditional expectation, which lies at the principle of the notion of martingale, is closely related to the geometry of L^2 . Let us illustrate this by the following elementary but useful result.

Proposition 3.24 *The increments of a square-integrable martingale are orthogonal in L^2 . More precisely, if $X = (X_n)_{n \geq 0}$ is a martingale on $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ such that $\mathbb{E}[X_n^2] < \infty$ for all $n \geq 0$, then for all integers m, n, p such that $m \leq n \leq p$, one has*

$$\mathbb{E}[(X_n - X_m)(X_p - X_n)] = 0.$$

Proof. Consider three integers $m \leq n \leq p$. We have

$$\begin{aligned} \mathbb{E}[(X_n - X_m)(X_p - X_n)] &= \mathbb{E}[\mathbb{E}[(X_n - X_m)(X_p - X_n) | \mathcal{F}_n]] \\ &= \mathbb{E}[(X_n - X_m)\mathbb{E}[X_p - X_n | \mathcal{F}_n]] \\ &= \mathbb{E}[(X_n - X_m) \cdot 0] \\ &= 0, \end{aligned}$$

as expected. □

Exercise 3.41 *Is the converse true? Is any square-integrable process $X = (X_n)_{n \geq 0}$ such that for all $m \leq n \leq p$ one has $\mathbb{E}[(X_n - X_m)(X_p - X_n)] = 0$ necessarily a martingale?*

Exercise 3.42 *Let H be a Hilbert space. A curve $x : \mathbb{R} \rightarrow H, t \mapsto x_t$ is called a helix if for all reals $s \leq t \leq u$, the vectors $x_t - x_s$ and $x_u - x_t$ are perpendicular. Let $x = (x_t)_{t \geq 0}$ be a helix in a Hilbert space H . Prove that for all reals $r \leq s \leq t \leq u$, the vectors $x_s - x_r$ and $x_u - x_t$ are perpendicular. Prove that H is infinite dimensional. Find an example of a helix in your favorite separable infinite-dimensional Hilbert space, and prove that there exists a helix in any infinite-dimensional Hilbert space.*

A collection of orthogonal vectors in a Hilbert space begs us for an application of the Pythagorean theorem.

Proposition 3.25 *Let X be a square-integrable martingale. For all $n \geq 0$, one has*

$$\mathbb{E}[X_n^2] = \mathbb{E}[X_0^2] + \sum_{k=0}^{n-1} \mathbb{E}[(X_{k+1} - X_k)^2].$$

In particular, X is bounded in L^2 if and only if

$$\sum_{k=0}^{\infty} \mathbb{E}[(X_{k+1} - X_k)^2] < +\infty.$$

Proof. This follows immediately from the previous proposition and the Pythagorean theorem. □

Exercise 3.43 Prove that if $(Z_n)_{n \geq 0}$ is a sequence of independent random variables such that $\mathbb{E}[Z_n] = 0$ for all $n \geq 0$ and the series $\sum_{n \geq 0} \text{Var}(Z_n)$ converges, then the series $\sum_{n \geq 0} Z_n$ converges almost surely and in L^2 .

The following result is fundamental. It is not specific of the square-integrable case, but we will see that it has very important consequences in this case.

Proposition 3.26 (Doob's decomposition) Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let $X = (X_n)_{n \geq 0}$ be an adapted integrable stochastic process.

1. There exists a martingale $M = (M_n)_{n \geq 0}$ with $M_0 = 0$ and a previsible process $A = (A_n)_{n \geq 0}$ such that for all $n \geq 0$, one has

$$X_n = X_0 + M_n + A_n.$$

Moreover, this decomposition is unique.

2. The process X is a submartingale if and only if the process A is non-decreasing, that is, for all $n \geq 0$, $\mathbb{P}(A_n \leq A_{n+1}) = 1$.

Proof. If such a decomposition exists, then we must have $A_0 = 0$ and, for all $n \geq 0$,

$$\mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n] = A_{n+1} - A_n.$$

This formula defines inductively a previsible process $(A_n)_{n \geq 0}$. Moreover, we must have $M_n = X_n - X_0 - A_n$ and the equality

$$\mathbb{E}[(X_{n+1} - A_{n+1}) - (X_n - A_n) | \mathcal{F}_n] = 0$$

shows that $(M_n)_{n \geq 0}$ is a martingale.

The second assertion follows immediately from the first line of computation above. \square

For square-integrable martingales, this decomposition leads to the definition of a very important object.

Definition 3.27 Let X be a square-integrable martingale. Let $X_n^2 = X_0^2 + M_n + A_n$ be the Doob decomposition of the integrable submartingale $(X_n^2)_{n \geq 0}$. The process $(A_n)_{n \geq 0}$ is called the increasing process associated to X and is often denoted by $\langle X \rangle = (\langle X \rangle_n)_{n \geq 0}$.

The process $\langle X \rangle$ plays a crucial role in the study of continuous-time martingales and in the construction of the stochastic integral, one of the main results being the continuous-time analogue of the following property.

Exercise 3.44 Let X be a martingale. Let H be a bounded previsible process. Prove that for all $n \geq 0$,

$$\langle H \bullet X \rangle_n = \sum_{k=1}^n H_k^2 (\langle X \rangle_k - \langle X \rangle_{k-1}).$$

This can be written informally as

$$\left\langle \int_0^\cdot H dX \right\rangle_n = \int_0^n H^2 d\langle X \rangle.$$

Using the increasing process associated with a square-integrable martingale, we can refine the theorem of convergence. For this, let us introduce

$$\langle X \rangle_\infty = \lim_{n \rightarrow \infty} \langle X \rangle_n,$$

which exists almost surely as the limit of a non-decreasing sequence.

Exercise 3.45 Check that $\mathbb{E}[X_n^2] = \mathbb{E}[\langle X \rangle_n]$ for all $n \geq 0$. Prove that X is bounded in L^2 if and only if $\langle X \rangle_\infty$ is integrable.

Theorem 3.28 Let X be a square-integrable martingale. The sequence $(X_n)_{n \geq 0}$ converges almost surely on the event $\{\langle X \rangle_\infty < \infty\}$.

Proof. Choose an integer $k \geq 0$ and define

$$T_k = \inf\{n \geq 0 : \langle X \rangle_{n+1} > k\}.$$

It is a stopping time, because it is the hitting time of the Borel subset $(k, +\infty)$ of \mathbb{R} by the adapted process $(\langle X \rangle_{n+1})_{n \geq 0}$. We claim that the process $\langle X \rangle^{T_k}$ is the increasing process associated to X^{T_k} .

Firstly, since $X^2 - \langle X \rangle$ is a martingale and T_k a stopping time, the process

$$(X^2 - \langle X \rangle)^{T_k} = (X^{T_k})^2 - \langle X \rangle^{T_k}$$

is a martingale. It remains to prove that $\langle X \rangle^{T_k}$ is previsible.

Choose $n \geq 1$ and B a Borel subset of \mathbb{R} . Then

$$\begin{aligned} \{\langle X \rangle_n^{T_k} \in B\} &= \{\langle X \rangle_{T_k \wedge n} \in B\} \\ &= (\{T_k \geq n\} \cap \{\langle X \rangle_n \in B\}) \cup (\{T_k < n\} \cap \{\langle X \rangle_{T_k} \in B\}) \\ &= (\{T_k \leq n-1\}^c \cap \{\langle X \rangle_n \in B\}) \cup \bigcup_{m=0}^{n-1} (\{T_k = m\} \cap \{\langle X \rangle_m \in B\}) \end{aligned}$$

and this way of writing this event makes it apparent that it belongs to \mathcal{F}_{n-1} .

The increasing process of the martingale X^{T_k} is bounded by k by construction, so that this martingale is bounded in L^2 , from which it follows that it converges almost surely.

On the event $\langle X \rangle_\infty \leq k$, the stopping time T_k is equal to ∞ , and the processes X and X^{T_k} are equal. Hence, X itself converges almost surely on $\{\langle X \rangle_\infty \leq k\}$.

Finally, the equality $\{\langle X \rangle_\infty < \infty\} = \bigcup_{k \geq 1} \{\langle X \rangle_\infty \leq k\}$ implies that X converges almost surely on the whole event $\{\langle X \rangle_\infty < \infty\}$. \square

Exercise 3.46 Give an example of a sequence $(z_n)_{n \geq 0}$ of real numbers such that $\sum_{n \geq 0} z_n^2$ converges, but $\sum_{n \geq 0} z_n$ does not.

Let $(Z_n)_{n \geq 0}$ be a sequence of independent random variables such that $\mathbb{E}[Z_n] = 0$ for all $n \geq 0$. Set $X_0 = 0$ and, for all $n \geq 1$, $X_n = Z_1 + \dots + Z_n$. Compute $\langle X \rangle$. What does the last theorem say in this situation?

3.11 Uniform integrability

In this section, we will give a much more detailed answer to the question of knowing when a martingale which is bounded in L^1 converges in L^1 . The crucial tool is the notion of uniform integrability.

Definition 3.29 Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. A family $(X_i)_{i \in I}$ of integrable random variables is said to be uniformly integrable if for all $\varepsilon > 0$ there exists $M > 0$ such that

$$\forall i \in I, \mathbb{E} [|X_i| \mathbb{1}_{\{|X_i| > M\}}] < \varepsilon.$$

An equivalent formulation of this definition is

$$\lim_{M \rightarrow \infty} \left(\sup_{i \in I} \mathbb{E} [|X_i| \mathbb{1}_{\{|X_i| > M\}}] \right) = 0.$$

Exercise 3.47 Prove that a uniformly integrable family is bounded in L^1 .

Exercise 3.48 Let $(X_i)_{i \in I}$ be a family of integrable random variables. Prove that it is uniformly integrable as soon as one of the following assumptions is satisfied.

- The index set I is finite.
- There exists an integrable random variable Z such that $|X_i| \leq Z$ for all $i \in I$.
- The random variable $\sup_{i \in I} |X_i|$ is integrable.
- There exists $p > 1$ such that the family $(X_i)_{i \in I}$ is bounded in L^p .

Consider the probability space $([0, 1], \mathcal{B}_{[0,1]}, \lambda)$. For each $n \geq 1$, set $m = \lfloor \log_2 n \rfloor$, so that $n = 2^m + k$ with $k \in \{0, \dots, 2^m - 1\}$, and define the random variable

$$X_n = \frac{2^m}{\log(m+1)} \mathbb{1}_{\left[\frac{k}{2^m}, \frac{k+1}{2^m}\right]}.$$

Prove that the family $(X_n)_{n \geq 1}$ is uniformly integrable but does not satisfy any of the properties above. Incidentally, does the sequence $(X_n)_{n \geq 1}$ converge, and in which sense?

The name *uniform integrability* is justified by the following proposition.

Proposition 3.30 Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let $(X_i)_{i \in I}$ be a family of random variables which is bounded in L^1 . The following two assertions are equivalent.

1. The family $(X_i)_{i \in I}$ is uniformly integrable.
2. For all $\varepsilon > 0$ there exists $\delta > 0$ such that for all $A \in \mathcal{A}$ such that $\mathbb{P}(A) < \delta$, one has

$$\forall i \in I, \int_A |X_i| d\mathbb{P} < \varepsilon.$$

Proof. $1 \Rightarrow 2$. Choose $\varepsilon > 0$. Since the family $(X_i)_{i \in I}$ is uniformly integrable, there exists a real $M > 0$ such that for all $i \in I$, $\mathbb{E}[|X_i| \mathbb{1}_{\{|X_i| > M\}}] < \frac{\varepsilon}{2}$. Let A be any event in \mathcal{A} such that $\mathbb{P}(A) < \frac{\varepsilon}{2M}$. Then

$$\begin{aligned} \int_A |X_i| d\mathbb{P} &= \int_{A \cap \{|X_i| \leq M\}} |X_i| d\mathbb{P} + \int_{A \cap \{|X_i| > M\}} |X_i| d\mathbb{P} \\ &\leq \int_A M d\mathbb{P} + \int_{\{|X_i| > M\}} |X_i| d\mathbb{P} \\ &< M\mathbb{P}(A) + \frac{\varepsilon}{2} \\ &< \varepsilon. \end{aligned}$$

2 \Rightarrow 1. Let $K > 0$ be such that $\mathbb{E}[|X_i|] \leq K$ for all $i \in I$. Choose $\varepsilon > 0$. Let $\delta > 0$ be such that for all $A \in \mathcal{A}$, $\mathbb{P}(A) < \delta$ implies $\int_A |X_i| d\mathbb{P} < \varepsilon$ for all $i \in I$. Set $M = \frac{2K}{\delta}$. Then for all $i \in I$,

$$\mathbb{P}(|X_i| > M) \leq \frac{1}{M} \int_{\{|X_i| > M\}} |X_i| d\mathbb{P} \leq \frac{K}{M} \leq \frac{\delta}{2} < \delta,$$

so that

$$\int_{\{|X_i| > M\}} |X_i| d\mathbb{P} < \varepsilon,$$

and the proof is finished. \square

Exercise 3.49 *When did we use the assumption that the family $(X_i)_{i \in I}$ is bounded in L^1 ? Prove that if the probability space $(\Omega, \mathcal{A}, \mathbb{P})$ is such that Ω is a finite set, then the second assertion is true for any family of random variables. Under which assumption on the probability space $(\Omega, \mathcal{A}, \mathbb{P})$ can one deduce from the second assertion that the family $(X_i)_{i \in I}$ is bounded in L^1 ?*

Corollary 3.31 *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let Z be an integrable variable. The family $\{\mathbb{E}[Z|\mathcal{B}] : \mathcal{B} \text{ sub-}\sigma\text{-field of } \mathcal{A}\}$ is uniformly integrable.*

Proof. For each sub- σ -field \mathcal{B} of \mathcal{A} , the random variable $\mathbb{E}[Z|\mathcal{B}]$ is integrable. Let us now choose ε . Thanks to the uniform integrability of the family which consists in the single random variable Z , let us consider $\delta > 0$ such that $\mathbb{P}(A) < \delta$ implies $\int_A |Z| d\mathbb{P} < \varepsilon$. Set $M = \frac{2}{\delta} \mathbb{E}[|Z|]$. Let \mathcal{B} be a sub- σ -field of \mathcal{A} . We claim that

$$\mathbb{E} [|\mathbb{E}[Z|\mathcal{B}]| \mathbb{1}_{\{|\mathbb{E}[Z|\mathcal{B}]| > M\}}] < \varepsilon.$$

Indeed, we have

$$\mathbb{P}(|\mathbb{E}[Z|\mathcal{B}]| > M) \leq \frac{1}{M} \mathbb{E}[|\mathbb{E}[Z|\mathcal{B}]|] \leq \frac{1}{M} \mathbb{E}[\mathbb{E}[|Z||\mathcal{B}]] = \frac{1}{M} \mathbb{E}[|Z|] < \delta.$$

Hence,

$$\int_{\{|\mathbb{E}[Z|\mathcal{B}]| > M\}} |\mathbb{E}[Z|\mathcal{B}]| d\mathbb{P} \leq \int_{\{|\mathbb{E}[Z|\mathcal{B}]| > M\}} |Z| d\mathbb{P} < \varepsilon$$

and the result is proved. \square

The reason why uniform integrability is so useful in the context of convergence of martingales is the following result, which is a stronger form of the dominated convergence theorem (but it is only valid on finite measure spaces).

Theorem 3.32 *Let $(X_n)_{n \geq 0}$ be a sequence of integrable random variables. Let X_∞ be a random variable. The following assertions are equivalent.*

1. *The sequence $(X_n)_{n \geq 0}$ converges in L^1 to X_∞ .*
2. *The sequence $(X_n)_{n \geq 0}$ is uniformly integrable and converges in probability to X_∞ .*

Proof. $1 \Rightarrow 2$. We know that convergence in L^1 implies convergence in probability. Let us prove that $(X_n)_{n \geq 0}$ is uniformly integrable.

Choose $\varepsilon > 0$. Let $\delta_1 > 0$ be such that $\mathbb{P}(A) < \delta_1$ implies $\int_A |X_\infty| d\mathbb{P} < \frac{\varepsilon}{2}$. Let n_0 be such that for all $n \geq n_0$, $\mathbb{E}[|X_n - X_\infty|] < \frac{\varepsilon}{2}$. Let $\delta_2 > 0$ be such that for all $n \leq n_0$ and all A with $\mathbb{P}(A) < \delta_2$, we have $\int_A |X_n| d\mathbb{P} < \varepsilon$.

Now set $\delta = \min(\delta_1, \delta_2)$. Choose $A \in \mathcal{A}$ such that $\mathbb{P}(A) < \delta$. Choose $n \geq 0$. If $n \leq n_0$, then since $\mathbb{P}(A) < \delta_2$, we have $\int_A |X_n| d\mathbb{P} < \varepsilon$. If $n \geq n_0$, then

$$\int_A |X_n| d\mathbb{P} \leq \int_A |X_n - X_\infty| d\mathbb{P} + \int_A |X_\infty| d\mathbb{P} < \mathbb{E}[|X_n - X_\infty|] + \frac{\varepsilon}{2} < \varepsilon.$$

$2 \Rightarrow 1$. From Proposition 3.30 it follows that the family $(X_n - X_m)_{n,m \geq 0}$ is uniformly integrable. Choose $\varepsilon > 0$. Let $M > 0$ be such that $\mathbb{E}[|X_n - X_m| \mathbb{1}_{\{|X_n - X_m| > M\}}] < \varepsilon$ for all $n, m \geq 0$. Thus,

$$\begin{aligned} \mathbb{E}[|X_n - X_m|] &= \mathbb{E}[|X_n - X_m| \mathbb{1}_{\{|X_n - X_m| < \varepsilon\}}] + \mathbb{E}[|X_n - X_m| \mathbb{1}_{\{\varepsilon \leq |X_n - X_m| \leq M\}}] \\ &\quad + \mathbb{E}[|X_n - X_m| \mathbb{1}_{\{|X_n - X_m| > M\}}] \\ &\leq 2\varepsilon + \mathbb{E}[|X_n - X_m| \mathbb{1}_{\{\varepsilon \leq |X_n - X_m| \leq M\}}] \\ &\leq 2\varepsilon + M\mathbb{P}(|X_n - X_m| > \varepsilon). \end{aligned}$$

Let n_0 be such that for all $n \geq n_0$, $\mathbb{P}(|X_n - X_\infty| > \frac{\varepsilon}{2}) < \frac{\varepsilon}{2M}$. Then for all $n, m \geq n_0$, we have

$$\mathbb{P}(|X_n - X_m| > \varepsilon) \leq \mathbb{P}\left(|X_n - X_\infty| > \frac{\varepsilon}{2}\right) + \mathbb{P}\left(|X_m - X_\infty| > \frac{\varepsilon}{2}\right) < \frac{\varepsilon}{M}.$$

Thus, for all $n, m \geq n_0$ we have

$$\mathbb{E}[|X_n - X_m|] \leq 3\varepsilon.$$

The sequence $(X_n)_{n \geq 0}$ is a Cauchy sequence in L^1 . Hence, it converges in L^1 , and its limit must be X_∞ . \square

The proof of the following result is left as an exercise.

Theorem 3.33 *Let $X = (X_n)_{n \geq 0}$ be a martingale. The following assertions are equivalent.*

1. *The martingale X converges in L^1 .*
2. *The family $(X_n)_{n \geq 0}$ is uniformly integrable.*
3. *There exists an integrable random variable such that $X_n = \mathbb{E}[Z | \mathcal{F}_n]$ for all $n \geq 0$.*

Moreover, if any of these assertions holds, then the martingale converges almost surely.

We would like to complete this picture by extending the stopping theorem to the case of uniformly integrable martingales. For this, we need to introduce the notion of σ -field associated with a stopping time.

Definition 3.34 *Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, \mathbb{P})$ be a filtered probability space. Let T be a stopping time. The collection of events*

$$\mathcal{F}_T = \{A \in \mathcal{A} : \forall n \geq 0, A \cap \{T = n\} \in \mathcal{F}_n\}$$

is a sub- σ -field of \mathcal{A} called the σ -field of the past up to time T .

Exercise 3.50 Prove that \mathcal{F}_T is indeed a σ -field. Compute \mathcal{F}_T when $T = n$ almost surely. Compute \mathcal{F}_T when $T = \infty$ almost surely.

Exercise 3.51 Prove that a subset A of Ω belongs to \mathcal{F}_T if and only if there exists a sequence $(A_n)_{n \geq 0}$ of events, and an event A_∞ , such that $A_\infty \in \mathcal{A}$ and $A_n \in \mathcal{F}_n$ for every $n \geq 0$, such that

$$A = (A_\infty \cap \{T = \infty\}) \cup \bigcup_{n=0}^{\infty} (A_n \cap \{T = n\}).$$

Lemma 3.35 Let S and T two stopping times. If $S \leq T$, then $\mathcal{F}_S \subset \mathcal{F}_T$.

Proof. Consider $A \in \mathcal{F}_S$. Choose $n \geq 0$. Then the event

$$A \cap \{T = n\} = \bigcup_{k=0}^{\infty} (A \cap \{T = n\} \cap \{S = k\}) = \bigcup_{k=0}^n (A \cap \{S = k\} \cap \{T = n\})$$

belongs to \mathcal{F}_n . Hence, A belongs to \mathcal{F}_T . \square

Lemma 3.36 Let $X = (X_n)_{n \geq 0}$ be an adapted stochastic process. Let T be a stopping time. If one of the following assumptions is satisfied :

1. T is finite almost surely,
2. X_n converges almost surely to X_∞ ,

then X_T is well defined and \mathcal{F}_T -measurable.

Proof. We know already that X_T is well defined under any of the two assumptions. Let us prove that X_T is \mathcal{F}_T measurable. For this, let us consider a Borel subset B of \mathbb{R} and an integer $n \geq 0$. We have

$$\{X_T \in B\} \cap \{T = n\} = \{X_n \in B\} \cap \{T = n\} \in \mathcal{F}_n.$$

Since this holds for all $n \geq 0$, $\{X_T \in B\}$ belongs to \mathcal{F}_T . \square

Let us now state and prove the stopping theorem for uniformly integrable martingales.

Theorem 3.37 Let X be a uniformly integrable martingale. Let T be a stopping time. Then

$$X_T = \mathbb{E}[X_\infty | \mathcal{F}_T].$$

In particular, $\mathbb{E}[X_T] = \mathbb{E}[X_\infty] = \mathbb{E}[X_n]$ for all $n \geq 0$. Moreover, if S and T are two stopping times such that $S \leq T$, then

$$X_S = \mathbb{E}[X_T | \mathcal{F}_S].$$

Proof. Let us check that X_T is integrable. We have

$$\begin{aligned} \mathbb{E}[|X_T|] &= \sum_{n=0}^{\infty} \mathbb{E}[|X_n| \mathbb{1}_{\{T=n\}}] + \mathbb{E}[|X_\infty| \mathbb{1}_{\{T=\infty\}}] \\ &\leq \sum_{n=0}^{\infty} \mathbb{E}[\mathbb{E}[|X_\infty| | \mathcal{F}_n] \mathbb{1}_{\{T=n\}}] + \mathbb{E}[|X_\infty| \mathbb{1}_{\{T=\infty\}}] \\ &= \sum_{n=0}^{\infty} \mathbb{E}[|X_\infty| \mathbb{1}_{\{T=n\}}] + \mathbb{E}[|X_\infty| \mathbb{1}_{\{T=\infty\}}] \\ &= \mathbb{E}[|X_\infty|] < \infty. \end{aligned}$$

Let us now choose $A \in \mathcal{F}_T$. We have

$$\begin{aligned}
\mathbb{E}[X_\infty \mathbb{1}_A] &= \sum_{n=0}^{\infty} \mathbb{E}[X_\infty \mathbb{1}_{A \cap \{T=n\}}] + \mathbb{E}[X_\infty \mathbb{1}_{A \cap \{T=\infty\}}] \\
&= \sum_{n=0}^{\infty} \mathbb{E}[\mathbb{E}[X_\infty | \mathcal{F}_n] \mathbb{1}_{A \cap \{T=n\}}] + \mathbb{E}[X_\infty \mathbb{1}_{A \cap \{T=\infty\}}] \\
&= \sum_{n=0}^{\infty} \mathbb{E}[X_n \mathbb{1}_{A \cap \{T=n\}}] + \mathbb{E}[X_\infty \mathbb{1}_{A \cap \{T=\infty\}}] \\
&= \mathbb{E}[X_T \mathbb{1}_A].
\end{aligned}$$

This proves that $X_T = \mathbb{E}[X_\infty | \mathcal{F}_T]$.

It follows that $\mathbb{E}[X_T] = \mathbb{E}[X_\infty]$. We already know that $\mathbb{E}[X_n] = \mathbb{E}[\mathbb{E}[X_\infty | \mathcal{F}_n]] = \mathbb{E}[X_\infty]$. Finally, if $S \leq T$, then since $\mathcal{F}_S \subset \mathcal{F}_T$, we have

$$X_S = \mathbb{E}[X_\infty | \mathcal{F}_S] = \mathbb{E}[\mathbb{E}[X_\infty | \mathcal{F}_T] | \mathcal{F}_S] = \mathbb{E}[X_T | \mathcal{F}_S]$$

and the proof is finished. \square

Exercise 3.52 *Let X be a uniformly integrable martingale. Let T be a stopping time. Prove that $X_{T \wedge n}$ converges almost surely and in L^1 to X_T .*

3.12 Backward martingales

The name *backward martingales* may be misleading. We are not going to run time backwards but rather allow time to be unbounded below.

Definition 3.38 *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. A backward filtration is a non-decreasing sequence $(\mathcal{F}_n)_{n \leq 0}$ of sub- σ -fields of \mathcal{A} :*

$$\dots \subset \mathcal{F}_{-n-1} \subset \mathcal{F}_{-n} \subset \dots \subset \mathcal{F}_{-1} \subset \mathcal{F}_0.$$

A backward martingale on the backward filtered probability space $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \leq 0}, \mathbb{P})$ is a sequence $X = (X_n)_{n \leq 0}$ such that for all $n \leq 0$, X_n is integrable and \mathcal{F}_n -measurable, and

$$\mathbb{E}[X_n | \mathcal{F}_{n-1}] = X_{n-1}.$$

Having time unbounded below instead of unbounded above makes a huge difference.

Lemma 3.39 *A backward martingale is uniformly integrable.*

Proof. Indeed, we have for all $n \leq 0$ the equality $X_n = \mathbb{E}[X_0 | \mathcal{F}_n]$, and the result is a consequence of Corollary 3.31. \square

It is thus not too surprising that the problem of convergence of backward martingales is simpler than that of martingales.

Theorem 3.40 *Let X be a backward martingale. Set $\mathcal{F}_{-\infty} = \bigcap_{n \leq 0} \mathcal{F}_n$. Then, as n tends to $-\infty$, X_n converges almost surely and in L^1 to $\mathbb{E}[X_0 | \mathcal{F}_{-\infty}]$.*

Proof. The proof that X converges almost surely is the same as that for usual martingales. Indeed, for all $N \leq 0$, the sequence X_N, X_{N+1}, \dots, X_0 is a usual martingale and for all $a < b$, the number of upcrossings of this martingale can be estimated by Doob's upcrossing lemma. As N tends to $-\infty$, this number of upcrossings converges to the number of upcrossings of the full backward martingale X , and turns out to be finite almost surely for all $a < b$. Thus, X_n converges almost surely as n tends to $-\infty$, towards a random variable which we denote by $X_{-\infty}$.

Since, by the preceding lemma, X is uniformly integrable, it also converges to $X_{-\infty}$ in L^1 .

The random variable $X_{-\infty}$ is \mathcal{F}_n -measurable for each $n \leq 0$, because it is the limit of the sequence $X_n, X_{n-1}, X_{n-2}, \dots$. Thus, $X_{-\infty}$ is $\mathcal{F}_{-\infty}$ measurable. Now if A belongs to $\mathcal{F}_{-\infty}$, then the L^1 convergence of X_n to $X_{-\infty}$ implies

$$\int_A X_{-\infty} d\mathbb{P} = \lim_{n \rightarrow -\infty} \int_A X_n d\mathbb{P} = \lim_{n \rightarrow -\infty} \int_A \mathbb{E}[X_0 | \mathcal{F}_n] d\mathbb{P} = \lim_{n \rightarrow -\infty} \int_A X_0 d\mathbb{P} = \int_A X_0 d\mathbb{P},$$

because A belongs to \mathcal{F}_n for all $n \leq 0$. Hence, $X_{-\infty} = \mathbb{E}[X_0 | \mathcal{F}_{-\infty}]$ and the proof is finished. \square

This result of convergence has several spectacular consequences, including a proof of the strong law of large numbers. This proof however requires some preparatory results.

Lemma 3.41 *Let X, Y_1, \dots, Y_n and X', Y'_1, \dots, Y'_n be random variables. Assume that the random vectors (X, Y_1, \dots, Y_n) and (X', Y'_1, \dots, Y'_n) have the same distribution. Assume that X and X' are integrable. Assume finally that $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a measurable function such that*

$$\mathbb{E}[X | \sigma(Y_1, \dots, Y_n)] = h(Y_1, \dots, Y_n).$$

Then

$$\mathbb{E}[X' | \sigma(Y'_1, \dots, Y'_n)] = h(Y'_1, \dots, Y'_n),$$

with the same function h .

Proof. Let B be a Borel subset of \mathbb{R}^n . The assumptions imply that

$$\int_{\Omega} X \mathbb{1}_B(Y_1, \dots, Y_n) dP = \int_{\Omega} h(Y_1, \dots, Y_n) \mathbb{1}_B(Y_1, \dots, Y_n) dP.$$

Denoting by μ the distribution of (X, Y_1, \dots, Y_n) (which is a probability measure on \mathbb{R}^{n+1}) and by ν the distribution of (Y_1, \dots, Y_n) (which is a probability measure on \mathbb{R}^n), this rewrites as

$$\int_{\mathbb{R}^{n+1}} x \mathbb{1}_B(y_1, \dots, y_n) \mu(dx, dy_1, \dots, dy_n) = \int_{\mathbb{R}^n} h(y_1, \dots, y_n) \mathbb{1}_B(y_1, \dots, y_n) \nu(dy_1, \dots, dy_n).$$

Since (X, Y_1, \dots, Y_n) and (X', Y'_1, \dots, Y'_n) have the same distribution, the last equality also says that

$$\int_{\Omega} X' \mathbb{1}_B(Y'_1, \dots, Y'_n) dP = \int_{\Omega} h(Y'_1, \dots, Y'_n) \mathbb{1}_B(Y'_1, \dots, Y'_n) dP,$$

that is,

$$\mathbb{E}[X' | \sigma(Y'_1, \dots, Y'_n)] = h(Y'_1, \dots, Y'_n)$$

as expected. \square

This lemma has the following consequence.

Lemma 3.42 *Let $(X_n)_{n \geq 1}$ be an i.i.d. sequence of integrable random variables. For all $n \geq 1$, set $S_n = X_1 + \dots + X_n$. Then for all $n \geq 1$, and all $k \in \{1, \dots, n\}$,*

$$\mathbb{E}[X_k | S_n] = \frac{S_n}{n} = \mathbb{E}[X_k | \sigma(S_n, S_{n+1}, S_{n+2}, \dots)].$$

Proof. For all $k, l \in \{1, \dots, n\}$, the vectors (X_k, S_n) and (X_l, S_n) have the same distribution. Indeed, assuming $k < l$, the vectors (X_1, \dots, X_n) and $(X_1, \dots, X_l, \dots, X_k, \dots, X_n)$ with X_k and X_l exchanged have the same distribution, so that (X_k, S_n) and (X_l, S_n) , which are respectively obtained from these two vectors by applying the function

$$(x_1, \dots, x_n) \mapsto (x_k, x_1 + \dots + x_n),$$

also have the same distribution.

Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a measurable function such that $\mathbb{E}[X_1 | S_n] = h(S_n)$. By the previous lemma, we have $\mathbb{E}[X_k | S_n] = h(S_n)$ for all $k \in \{1, \dots, n\}$. Hence,

$$S_n = \mathbb{E}[S_n | S_n] = \sum_{k=1}^n \mathbb{E}[X_k | S_n] = nh(S_n),$$

so that $h(S_n) = \frac{S_n}{n}$. This proves the first equality.

In order to prove the second, observe that $\sigma(S_n, S_{n+1}, \dots) = \sigma(S_n, X_{n+1}, X_{n+2}, \dots)$. We need to prove that for all event $A \in \sigma(S_n, X_{n+1}, X_{n+2}, \dots)$, we have

$$\mathbb{E}[X_k \mathbb{1}_A] = \frac{1}{n} \mathbb{E}[S_n \mathbb{1}_A].$$

Let us denote by \mathcal{C} the class of events $A \in \mathcal{A}$ for which this equality holds. It is a λ -system. Let us consider an event $B \in \sigma(S_n)$ and an event $C \in \sigma(X_{n+1}, X_{n+2}, \dots)$. Since the σ -fields $\sigma(X_k, S_n)$ and $\sigma(X_{n+1}, X_{n+2}, \dots)$ are independent, we have

$$\mathbb{E}[X_k \mathbb{1}_B \mathbb{1}_C] = \mathbb{E}[X_k \mathbb{1}_B] \mathbb{E}[\mathbb{1}_C] = \frac{1}{n} \mathbb{E}[S_n \mathbb{1}_B] \mathbb{E}[\mathbb{1}_C] = \frac{1}{n} \mathbb{E}[S_n \mathbb{1}_B \mathbb{1}_C].$$

Hence, the λ -system \mathcal{C} contains the π -system of all events of the form $B \cap C$ with $B \in \sigma(S_n)$ and $C \in \sigma(X_{n+1}, X_{n+2}, \dots)$. Thus, by the monotone class theorem, it contains $\sigma(S_n, X_{n+1}, X_{n+2}, \dots)$ and the second equality is proved. \square

Theorem 3.43 (Strong law of large numbers) *Let $(X_n)_{n \geq 1}$ be an i.i.d. sequence of integrable random variables. Then as n tends to infinity,*

$$\frac{X_1 + \dots + X_n}{n} \longrightarrow \mathbb{E}[X_1]$$

almost surely and in L^1 .

Proof. Set $S_0 = 0$ and, for all $n \geq 1$, $S_n = X_1 + \dots + X_n$. For all $n \geq 0$, set $\mathcal{F}_{-n} = \sigma(S_n, S_{n+1}, \dots) = \sigma(S_k : k \geq n)$. Then $(\mathcal{F}_n)_{n \leq 0}$ is a backward filtration. For each $n \geq 1$, we have

$$\mathbb{E} \left[\frac{S_n}{n} \middle| \mathcal{F}_{-n-1} \right] = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[X_k | \mathcal{F}_{-n-1}] = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[X_k | S_{n+1}, S_{n+2}, \dots] = \frac{S_{n+1}}{n+1}.$$

In other words, $(\frac{S_{-n}}{-n})_{n \leq 1}$ is a backward martingale with respect to its natural filtration. This implies that the sequence $(\frac{S_n}{n})_{n \geq 0}$ converges almost surely and in L^1 . Its limit is

$$Y = \lim_{n \rightarrow \infty} \frac{X_1 + \dots + X_n}{n},$$

but for each $n_0 \geq 0$, it is also

$$Y = \lim_{n \rightarrow \infty} \frac{X_{n_0} + \dots + X_n}{n}$$

almost surely, so that Y is measurable with respect to $\sigma(X_n : n \geq n_0)$. Since this holds for all $n_0 \geq 1$, the random variable Y is measurable with respect to the tail σ -field

$$\bigcap_{n \geq 1} \sigma(X_k : k \geq n)$$

and Kolmogorov's 0-1 law asserts that this σ -field is trivial. Hence, the limit Y is a constant. Since $\mathbb{E}[Y] = \mathbb{E}[X_1]$, this constant must be $\mathbb{E}[X_1]$. \square

4 Markov chains

4.1 Introduction

Just as martingales, Markov chains are a class of stochastic processes whose definition involves conditional expectation in a crucial way. However, in contrast with martingales, which are a technical tool for the study of dynamical systems, Markov chains are a model for actual random phenomena. The class of Markov chains, or more generally Markov processes, has the main features of a good and successful model, namely: simplicity and versatility. It is mathematically simple enough to be studied in great detail, and allows one to give interesting models of a whole range of real situations.

Both points should however be nuanced. On one hand, the general theory of Markov processes is a very sophisticated theory, of which we are going to study a few fundamental ideas in the simplest interesting situation. On the other hand, Markovian models have their limit when applied to reality, and must often be enhanced to fit any particular actual random phenomenon.

The general idea of Markov processes is that of a random motion in a so-called *state space*, that we will denote by E . This random motion can either be thought of as the random motion of a particle in the space E , or as the random evolution in time of the state of a system, each point of E describing such a state. The second point of view is more general than the first, since the position of a particle is a special case of the state of a system. In any case, the main defining property of the random motion, or random evolution, is that it is *memoryless*, in the sense that the way it moves, or evolves, immediately after a given instant of time depends on the history of the movement, or evolution, only through its present location, or state. Another way of stating this property is to say that at every instant of time, the future evolution is independent of the past evolution conditional on the present state.

Markov processes come in several variants, depending on the time being discrete or continuous, and on the space E being discrete or continuous. We are going to study the simplest case, where E is discrete and time is discrete. Although technically simple, this case will allow us to meet most of the fundamental ideas of the theory of Markov processes.

4.2 First definition and first properties

In this chapter, we will call *state space* and denote by E an arbitrary non-empty finite or countable set. Whenever it is necessary, this set will be endowed with the σ -field $\mathcal{P}(E)$ of all subsets of E .

A Markov chain is a sequence of random variables with values in E which satisfies a certain property which we described as absence of memory. In dealing with this property, we will make use of conditional expectations in a way that differs slightly from the way we discussed it in the first chapter of these notes.

Given a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, a sub- σ -field \mathcal{B} of \mathcal{A} and an event $A \in \mathcal{A}$, we will use the notation

$$\mathbb{P}(A|\mathcal{B}) = \mathbb{E}[\mathbb{1}_A|\mathcal{B}]$$

for the conditional probability of A given \mathcal{B} . Let us emphasize that this is a random variable, and not a number. In the simple case where \mathcal{B} is the σ -field $\{\emptyset, B, B^c, \Omega\}$ generated by a single event B such that $0 < \mathbb{P}(B) < 1$, we have (see Section 2.1)

$$\mathbb{P}(A|\mathcal{B}) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \mathbb{1}_B + \frac{\mathbb{P}(A \cap B^c)}{\mathbb{P}(B^c)} \mathbb{1}_{B^c}$$

almost surely. The number $\frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ is traditionally denoted by $\mathbb{P}(A|B)$, and called the conditional probability of A given B , with the important precaution that it is defined only when $\mathbb{P}(B) > 0$.

A situation that we will meet often is summarised in the following exercise.

Exercise 4.1 *Let A be an event and Y a random variable with values in a finite or countable set D endowed with the σ -field $\mathcal{P}(D)$. Suppose that we want to understand the conditional probability $\mathbb{P}(A|\sigma(Y))$, which incidentally is often denoted by $\mathbb{P}(A|Y)$. What is usually easy to compute is, for an element $y \in D$, the conditional probability $\mathbb{P}(A|\{Y = y\})$, which is usually denoted by $\mathbb{P}(A|Y = y)$. However, this conditional probability is defined only if $\mathbb{P}(Y = y) > 0$. Prove that*

$$\mathbb{P}(A|Y) = \sum_{\substack{y \in D \\ \mathbb{P}(Y=y) > 0}} \mathbb{P}(A|Y = y) \mathbb{1}_{\{Y=y\}}.$$

The point here is that we consider in the sum only those values y for which $\mathbb{P}(Y = y) > 0$.

Apart from the conditional expectation, the main ingredient in the definition of a Markov chain is the *transition kernel*.

Definition 4.1 *A transition kernel on E is a function*

$$\begin{aligned} P : E \times E &\longrightarrow [0, 1] \\ (x, y) &\longmapsto p(x, y) \end{aligned}$$

such that

$$\forall x \in E, \sum_{y \in E} p(x, y) = 1.$$

The value of P at a couple (x, y) is usually denoted by $p(x, y)$, but also sometimes by $P(x, y)$.

Informally, conditional on the system currently being in the state x , the probably of this state becoming y at the next instant of time is $p(x, y)$.

Exercise 4.2 Check that it is equivalent to consider a transition kernel on E or to consider a map $E \rightarrow \text{Prob}(E)$ from E into the set of all probability measures on $(E, \mathcal{P}(E))$.

A transition kernel can be thought of as, and indeed, when E is finite, identified with a matrix whose rows and columns are indexed by elements of E . The conditions on this matrix are that it should have non-negative entries and that the sum of the entries in each row should be 1.

Just as matrices, transition kernels can be multiplied: if P and P' are transition kernels, then $PP' : E \rightarrow [0, +\infty]$ defined by

$$\forall x, y \in E, (PP')(x, y) = \sum_{z \in E} P(x, z)P'(z, y)$$

is a transition kernel.

Exercise 4.3 Check this statement, and check that the product thus defined is associative on the set of transition kernels, in the sense that if P, P', P'' are transition kernels, then $P(P'P'') = (PP')P''$.

The transition kernel $I : E \times E \rightarrow [0, 1]$ defined by

$$I(x, y) = \delta_{x,y} = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{otherwise} \end{cases}$$

is the unit element of the associative product on the set of all transition kernels. It is of course the identity matrix if E is finite, and the natural analogue of it when E is infinite.

Given a transition kernel P , we define for every non-negative integer n a new transition kernel P^n by setting

$$P^0 = I \text{ and for all } n \geq 1, P^n = \underbrace{P \dots P}_{n \text{ times}}.$$

We can now give a first definition of a Markov chain.

Definition 4.2 (Markov chains, first definition) Let E be a non-empty, finite or countable set endowed with the σ -field of all its subsets. Let P be a transition kernel on E . Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let $X = (X_n)_{n \geq 0}$ be a sequence of random variables defined on this probability space, with values in E .

The sequence X is a Markov chain on E with transition kernel P if the following condition holds:

$$\forall n \geq 0, \forall y \in E, \mathbb{P}(X_{n+1} = y | X_0, \dots, X_n) = p(X_n, y). \quad (3)$$

According to the result of Exercise 4.1, this condition can be written in a slightly longer but more elementary way as follows : for all $n \geq 0$, all $y \in E$ and all $x_0, \dots, x_n \in E$ such that $\mathbb{P}(X_0 = x_0, \dots, X_n = x_n) > 0$,

$$\mathbb{P}(X_{n+1} = y | X_0 = x_0, \dots, X_n = x_n) = p(x_n, y).$$

Let us immediately give an equivalent characterisation of Markov chains.

Proposition 4.3 *With the notation of Definition 4.2, X is a Markov chain on E with transition kernel P if and only if the following condition holds:*

$$\forall n \geq 0, \forall x_0, \dots, x_n \in E, \mathbb{P}(X_0 = x_0, \dots, X_n = x_n) = \mathbb{P}(X_0 = x_0)p(x_0, x_1) \dots p(x_{n-1}, x_n). \quad (4)$$

Proof. Let us start by the ‘only if’ part, that is, the implication \Rightarrow . We prove (4) by induction on n . For $n = 0$, it reduces to $\mathbb{P}(X_0 = x_0) = \mathbb{P}(X_0 = x_0)$, which is true. Let us now assume that (4) has been proved up to rank $n - 1$ for some $n \geq 1$ and let us consider x_0, \dots, x_n in E . If $\mathbb{P}(X_0 = x_0, \dots, X_{n-1} = x_{n-1}) = 0$, then both sides of (4) vanish, the left-hand side because it is the probability of an event included in a negligible event, and the right-hand side because the product of all factors except the last, being equal by induction to $\mathbb{P}(X_0 = x_0, \dots, X_{n-1} = x_{n-1})$, is equal to 0.

If, on the other hand, $\mathbb{P}(X_0 = x_0, \dots, X_{n-1} = x_{n-1}) > 0$, then

$$\begin{aligned} \mathbb{P}(X_0 = x_0, \dots, X_n = x_n) &= \mathbb{P}(X_0 = x_0, \dots, X_{n-1} = x_{n-1}) \\ &\quad \mathbb{P}(X_n = x_n | X_0 = x_0, \dots, X_{n-1} = x_{n-1}) \\ &= \mathbb{P}(X_0 = x_0)p(x_0, x_1) \dots p(x_{n-2}, x_{n-1})p(x_{n-1}, x_n), \end{aligned}$$

as expected.

Let us now prove the ‘if’ part, that is, the implication \Leftarrow . According to the remark made after the definition of Markov chains, it suffices to consider x_0, \dots, x_n, y in E such that $\mathbb{P}(X_0 = x_0, \dots, X_n = x_n) > 0$ and compute $\mathbb{P}(X_{n+1} = y | X_0 = x_0, \dots, X_n = x_n)$. Thanks to (4), this conditional expectation is equal to

$$\frac{\mathbb{P}(X_0 = x_0)p(x_0, x_1) \dots p(x_{n-1}, x_n)p(x_n, y)}{\mathbb{P}(X_0 = x_0)p(x_0, x_1) \dots p(x_{n-1}, x_n)} = p(x_n, y),$$

and this finishes the proof. □

It follows from this proposition that if X is a Markov chain on E with transition kernel P , then

$$\forall n \geq 0, \forall y \in E, \mathbb{P}(X_{n+1} = y | X_n) = p(X_n, y). \quad (5)$$

However, it is important to realise that this property is much weaker than the property which defines Markov chains. Exercise 4.5 below gives an example of a process that is not a Markov chain but satisfies (5).

Exercise 4.4 *Check the fact that (5) holds for a Markov chain.*

Exercise 4.5 *Consider the state space $E = \{a, b, c\}$. Define a random variable Z with values in E such that $\mathbb{P}(Z = b) = \mathbb{P}(Z = c) = \frac{1}{2}$. Define a sequence $X = (X_n)_{n \geq 0}$ of random variables with values in E such that*

$$\forall n \geq 0, X_n = \begin{cases} a & \text{if } n \text{ is even,} \\ Z & \text{if } n \text{ is odd.} \end{cases}$$

Prove that there exists a unique transition kernel P such that X satisfies the property (5), and prove that X is not a Markov chain with transition kernel P .

Let us give two simple consequences of this characterisation of Markov chains.

Proposition 4.4 *Let X be a Markov chain on E with transition kernel P . For all $n \geq 0$ and all $y \in E$, one has*

$$\mathbb{P}(X_n = y|X_0) = P^n(X_0, y).$$

Proof. We need to prove that for every element x_0 of E such that $\mathbb{P}(X_0 = x_0) > 0$, the equality $\mathbb{P}(X_n = y|X_0 = x_0) = P^n(x_0, y)$ holds. But for such an x_0 , we have

$$\begin{aligned} \mathbb{P}(X_0 = x_0)\mathbb{P}(X_n = y|X_0 = x_0) &= \mathbb{P}(X_0 = x_0, X_n = y) \\ &= \sum_{x_1, \dots, x_{n-1} \in E} \mathbb{P}(X_0 = x_0, \dots, X_{n-1} = x_{n-1}, X_n = y) \\ &= \mathbb{P}(X_0 = x_0) \sum_{x_1, \dots, x_{n-1} \in E} p(x_0, x_1) \dots p(x_{n-1}, y) \\ &= \mathbb{P}(X_0 = x_0)P^n(x_0, y), \end{aligned}$$

as expected. □

Exercise 4.6 *Let X be a Markov chain on E with transition kernel P . Check that for every integer $\ell \geq 0$, the sequence $(X_{n\ell})_{n \geq 0}$ is a Markov chain on E with transition kernel P^ℓ .*

Proposition 4.5 *Let X be a Markov chain on E with transition kernel P . Let $N \geq 0$ be an integer. For every $n \geq 0$, set $Y_n = X_{N+n}$. Then $Y = (Y_n)_{n \geq 0}$ is a Markov chain on E with transition kernel P .*

Proof. Let us consider $n \geq 0$ and $y_0, \dots, y_n \in E$. Firstly, we know from the previous proposition that

$$\begin{aligned} \mathbb{P}(Y_0 = y_0) &= \sum_{x_0 \in E} \mathbb{P}(X_0 = x_0, Y_0 = y_0) \\ &= \sum_{x_0 \in E} \mathbb{P}(X_0 = x_0, X_N = y_0) \\ &= \sum_{x_0 \in E} \mathbb{P}(X_0 = x_0)P^N(x_0, y_0). \end{aligned}$$

Secondly,

$$\begin{aligned} \mathbb{P}(Y_0 = y_0, \dots, Y_n = y_n) &= \mathbb{P}(X_N = y_0, \dots, X_{N+n} = y_n) \\ &= \sum_{x_0, \dots, x_{N-1} \in E} \mathbb{P}(X_0 = x_0, \dots, X_{N-1} = x_{N-1}, X_N = y_0, \dots, X_{N+n} = y_n) \\ &= \sum_{x_0, \dots, x_{N-1} \in E} \mathbb{P}(X_0 = x_0)p(x_0, x_1) \dots p(x_{N-1}, y_0)p(y_0, y_1) \dots p(y_{n-1}, y_n) \\ &= \sum_{x_0 \in E} \mathbb{P}(X_0 = x_0)P^N(x_0, y_0)p(y_0, y_1) \dots p(y_{n-1}, y_n) \\ &= \mathbb{P}(Y_0 = y_0)p(y_0, y_1) \dots p(y_{n-1}, y_n), \end{aligned}$$

from which we see that Y is a Markov chain with transition kernel P . □

Exercise 4.7 Check that for a Markov chain X with transition kernel P , one has, for all integers $0 \leq n \leq m$ and all $y \in E$

$$\mathbb{P}(X_m = y | X_n) = P^{m-n}(X_n, y).$$

Let us conclude this first section by giving a few examples of Markov chains.

- **Independent random variables.** Let $X = (X_n)_{n \geq 0}$ be a sequence of i.i.d. random variables with values in E , with common distribution μ . Then X is a Markov chain with transition kernel P given by

$$\forall x, y \in E, P(x, y) = \mu(y),$$

where $\mu(y)$ is a notation for $\mu(\{y\})$.

- **Random walks on \mathbb{Z}^d .** Let $d \geq 1$ be an integer. Let μ be a probability measure on \mathbb{Z}^d . Let $(\xi_i)_{i \geq 1}$ be an i.i.d. sequence of random variables with values in \mathbb{Z}^d and common distribution μ . Let X_0 be a random variable with values in \mathbb{Z}^d independent of $(\xi_i : i \geq 1)$. For every $n \geq 1$, set

$$X_n = X_0 + \xi_1 + \dots + \xi_n.$$

Then $X = (X_n)_{n \geq 0}$ is a Markov chain on \mathbb{Z}^d with transition kernel P given by

$$\forall x, y \in E, P(x, y) = \mu(y - x).$$

This Markov chain is called the *random walk on \mathbb{Z}^d* with jump distribution μ .

Let (e_1, \dots, e_d) denote the canonical basis of \mathbb{Z}^d . In the special case where

$$\mu = \frac{1}{2d} \sum_{i=1}^d (\delta_{e_i} + \delta_{-e_i}),$$

the random walk is called the *simple random walk on \mathbb{Z}^d* . The simple random walk is one of the most classical objects of the theory of probability.

- **Random walk on a graph.** Let A be a subset of the set $\mathcal{P}_2(E)$ of pairs of elements of E . We think of E as the set of vertices of a graph, and of each element $\{x, y\} \in A$ as an unoriented edge joining the vertices x and y .

For each $x \in E$, we define the set N_x of neighbours of x in the graph (E, A) by

$$N_x = \{y \in E : \{x, y\} \in A\}.$$

We make the assumption that for all $x \in E$, the set N_x is non-empty and finite:

$$\forall x \in E, 0 < |N_x| < \infty.$$

The *random walk on the graph (E, A)* is the Markov chain with transition kernel P given by

$$\forall x, y \in E, P(x, y) = \begin{cases} \frac{1}{|N_x|} & \text{if } y \in N_x \\ 0 & \text{otherwise.} \end{cases}$$

We did not prove yet that such a Markov chain exists, and we will do it soon.

- **Branching processes.** Recall from Section 3.6 the definition of a branching process. With the notation used there, the sequence $(X_n)_{n \geq 0}$ is a Markov chain on \mathbb{N} with transition kernel

$$\forall n, m \in \mathbb{N}, P(n, m) = \mu^{*n}(m),$$

where μ is the reproduction law of the branching process, that is, the common distribution of all the variables $(Z_{n,k})_{n,k \geq 0}$, and for every integer $n \geq 0$, μ^{*n} is the n -th convolution power of μ , that is, the distribution of the sum of n independent random variables with distribution μ , for example $Z_{0,1} + \dots + Z_{0,n}$. In particular, $\mu^{*0} = \delta_0$.

Exercise 4.8 Check that a sequence of i.i.d. random variables, a random walk on \mathbb{Z}^d , a branching process, are indeed Markov chains with the claimed kernels.

4.3 Construction of Markov chains

In this section, we prove that on any finite or countable state space, any transition kernel is the transition kernel of a Markov chain.

Proposition 4.6 Let E be a finite or countable set. Let P be a transition kernel on E . Let x_0 be an element of E . There exists a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and a sequence $X = (X_n)_{n \geq 0}$ of E -valued random variables defined on this probability space such that $X_0 = x_0$ a.s. and X is a Markov chain with transition kernel P .

In order to construct such a Markov chain, we will need a source of randomness, which will be provided by the next lemma.

Proposition 4.7 There exists a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and a sequence $U = (U_n)_{n \geq 0}$ of i.i.d. random variables with common distribution equal to the uniform distribution on the interval $[0, 1]$.

It may be that you have always taken this existence result for granted: in this case, you should take a moment to think about the fact that it actually needs a proof.

Proof. Take $(\Omega, \mathcal{A}, \mathbb{P}) = ([0, 1], \mathcal{B}_{[0,1]}, \lambda)$, where λ is the Lebesgue measure. For all $n \geq 1$, define a random variable B_n on this probability space by setting

$$\forall t \in [0, 1], B_n(t) = \lfloor 2^n t \rfloor - 2 \lfloor 2^{n-1} t \rfloor = \sum_{k=0}^{2^{n-1}-1} \mathbb{1}_{\left[\frac{2k+1}{2^n}, \frac{2k+2}{2^n}\right)}.$$

Then $(B_n)_{n \geq 1}$ is an i.i.d. sequence with Bernoulli distribution of parameter $\frac{1}{2}$:

$$\forall n \geq 1, \mathbb{P}(B_n = 0) = \mathbb{P}(B_n = 1) = \frac{1}{2}.$$

Now, for all $n \geq 1$, let p_n denote the n -th prime number, so that $p_1 = 2, p_2 = 3, p_3 = 5$ and so on. For all $n \geq 0$, define

$$U_n = \sum_{k=1}^{\infty} 2^{-k} B_{p_{n+1}}^k.$$

Then, by computing for instance its characteristic function, one checks that U_n is uniformly distributed on $[0, 1]$ for all $n \geq 0$. Moreover, each variable U_n is built from a subset of the variables $(B_r)_{r \geq 0}$ that is disjoint from the subset used to build all the other U_m 's. Thus, $(U_n)_{n \geq 0}$ is an independent sequence of random variables. \square

We now turn to the proof of the proposition.

Proof of Proposition 4.6. Take $(\Omega, \mathcal{A}, \mathbb{P}) = ([0, 1], \mathcal{B}_{[0,1]}, \lambda)$ as before and consider a sequence $(U_n)_{n \geq 0}$ of i.i.d. random variables with uniform distribution on $[0, 1]$.

We claim that there exists a function $F : E \times [0, 1] \rightarrow E$ with the property that

$$\forall x, y \in E, \lambda(\{t \in [0, 1] : F(x, t) = y\}) = p(x, y).$$

In order to build such a function, we order E and write its elements as $E = \{y_1, y_2, \dots\}$ in an arbitrary fashion. Now choose $x \in E$ and $t \in [0, 1)$. There exists a unique integer $k \geq 1$ such that

$$\sum_{l=1}^{k-1} p(x, y_l) \leq t < \sum_{l=1}^k p(x, y_l)$$

and we define

$$F(x, t) = y_k.$$

With this definition,

$$\{t \in [0, 1] : F(x, t) = y\} = \left(\sum_{l=1}^{k-1} p(x, y_l), \sum_{l=1}^k p(x, y_l) \right]$$

is indeed a subset of $[0, 1]$ with Lebesgue measure $p(x, y_k)$.

Let us now define $X_0 = y_{k_0}$ and inductively, for all $n \geq 0$,

$$X_{n+1} = F(X_n, U_n).$$

Let k_0 be the integer such that $x_0 = y_{k_0}$. Then, for all $n \geq 1$ and all $k_1, \dots, k_n \geq 1$,

$$\begin{aligned} \mathbb{P}(X_0 = y_{k_0}, X_1 = y_{k_1}, \dots, X_{k_n} = y_{k_n}) &= \prod_{m=1}^n \lambda \left(\left(\sum_{l=1}^{k_m-1} p(y_{k_{m-1}}, y_l), \sum_{l=1}^{k_m} p(y_{k_{m-1}}, y_l) \right) \right) \\ &= p(y_{k_0}, y_{k_1}) \cdots p(y_{k_{n-1}}, y_{k_n}) \end{aligned}$$

and X is a Markov chain issued from $y_{k_0} = x_0$ with transition kernel P . \square

We would like now to take a slightly different point of view on Markov chains. So far, we thought of a Markov chain as a sequence of random variables with values in E , and we would like to become familiar with the idea that it is, or at least can be thought of as a single random variable with values in the space $E^{\mathbb{N}}$ of sequences of elements of E . Informally, this amounts to realising that $X = (X_n)_{n \geq 0}$ is really a function of two variables, namely $(n, \omega) \mapsto X_n(\omega)$, and changing the priority between the two variables n and ω .

In order to make this idea more precise, we need to define a measurable space with underlying set $E^{\mathbb{N}}$. Concretely, we need to endow the set

$$E^{\mathbb{N}} = \{\omega = (\omega_i)_{i \geq 0} : \forall i \geq 0, \omega_i \in E\}$$

with a σ -field. For this, we start by listing natural functions on this set that we would want to be measurable.

Definition 4.8 *Let E be a finite or countable set. On the set $E^{\mathbb{N}}$, which in the present context is called the canonical space, we define for each $n \geq 0$ the function*

$$\begin{aligned} \widehat{X}_n : E^{\mathbb{N}} &\longrightarrow E \\ \omega = (\omega_i)_{i \geq 0} &\longmapsto \widehat{X}_n(\omega) = \omega_n. \end{aligned}$$

The collection $\widehat{X} = (\widehat{X}_n)_{n \geq 0}$ of E -valued functions on $E^{\mathbb{N}}$ is called the canonical process.

We think of an element ω of $E^{\mathbb{N}}$ as the complete record of all the successive positions of our particle (or states of our system), from the origin to the end of times. Then, $\widehat{X}_n(\omega)$ is the position at time n of our particle.

In the next definition, recall that the set E is endowed with the σ -field $\mathcal{P}(E)$.

Definition 4.9 *For every $n \geq 0$, we define on $E^{\mathbb{N}}$ the σ -field*

$$\mathcal{C}_n = \sigma(\widehat{X}_0, \dots, \widehat{X}_n).$$

The elements of the π -system

$$\bigcup_{n=0}^{\infty} \mathcal{C}_n$$

are called cylinders on $E^{\mathbb{N}}$, and the σ -field

$$\mathcal{C} = \sigma\left(\bigcup_{n=0}^{\infty} \mathcal{C}_n\right) = \sigma(\widehat{X}_n : n \geq 0)$$

is called the cylinder σ -field on $E^{\mathbb{N}}$.

In the present context, \mathcal{C} turns out to be the natural σ -field on $E^{\mathbb{N}}$. This means, in English, and without being too precise, that the functions of a trajectory that we want to consider (and to call measurable) are those which are expressible in terms of finitely many successive positions of the trajectory, or which can be approximated in an appropriate sense by such functions.

It is a useful observation that for all $n \geq 0$, the σ -field \mathcal{C}_n is generated by the partition

$$\begin{aligned} E^{\mathbb{N}} &= \bigsqcup_{x_0, \dots, x_n \in E} \{\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n\} \\ &= \bigsqcup_{x_0, \dots, x_n \in E} \{\omega \in E^{\mathbb{N}} : \omega_0 = x_0, \dots, \omega_n = x_n\} \\ &= \bigsqcup_{x_0, \dots, x_n \in E} \{x_0\} \times \dots \times \{x_n\} \times E^{\mathbb{N} \setminus \{0, \dots, n\}}. \end{aligned}$$

Lemma 4.10 *A map f from a measurable space (Ω, \mathcal{A}) into $(E^{\mathbb{N}}, \mathcal{C})$ is measurable if and only if for all $n \geq 0$, the map $\widehat{X}_n \circ f$ is measurable from (Ω, \mathcal{A}) to $(E, \mathcal{P}(E))$.*

Proof. The ‘only if’ part of the statement is a consequence of the fact that a composition of measurable maps is measurable. Let us prove the ‘if’ part: let us assume that for all $n \geq 0$, the map $\widehat{X}_n \circ f$ is measurable. Let us define the class

$$\mathcal{I} = \{C \in \mathcal{C} : f^{-1}(X) \in \mathcal{A}\}$$

of subsets of $E^{\mathbb{N}}$. This is sometimes called the *image σ -field*² of \mathcal{A} by f . It is indeed a σ -field (check it!) and it is the largest sub- σ -field of \mathcal{C} that makes f measurable. Of course, we want to show that $\mathcal{I} = \mathcal{C}$. For this, it suffices to show that $\mathcal{C}_n \subset \mathcal{I}$ for all $n \geq 0$. According to the observation made just before stating the present lemma, it is enough to prove that for all $n \geq 0$ and all $x_0, \dots, x_n \in E$, the set $\{\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n\}$ belongs to \mathcal{I} . This is indeed the case, because

$$\begin{aligned} f^{-1}(\{\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n\}) &= f^{-1}\left(\bigcap_{i=0}^n \{\widehat{X}_i = x_i\}\right) \\ &= \bigcap_{i=0}^n f^{-1}(\{\widehat{X}_i = x_i\}) \\ &= \bigcap_{i=0}^n (\widehat{X}_i \circ f)^{-1}(\{x_i\}) \end{aligned}$$

belongs to \mathcal{A} by assumption. □

We are now in possession of a filtered measurable space with an adapted process on it, namely $(E^{\mathbb{N}}, \mathcal{C}, (\mathcal{C}_n)_{n \geq 0}, \widehat{X} = (\widehat{X}_n)_{n \geq 0})$. The last lemma will allow us to connect this nice space with our more familiar notion of sequence of E -valued random variables.

Definition 4.11 *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space on which we are given a sequence $X = (X_n)_{n \geq 0}$ of E -valued random variables. The trajectory map of this stochastic process is the map*

$$\begin{aligned} X : \Omega &\longrightarrow E^{\mathbb{N}} \\ \omega &\longmapsto (X_n(\omega))_{n \geq 0}. \end{aligned}$$

A technical remark is in order at this point. Each random variable X_n is really an equivalence class of functions, any two of which coincide outside a \mathbb{P} -negligible subset of Ω . Because we are considering a countable collection of random variables, the map X is indeed well defined \mathbb{P} -almost surely, in the sense that two distinct choices of a sequence of representatives of the sequence of random variables $(X_n)_{n \geq 0}$ yields two functions X that may be distinct, but coincide outside a \mathbb{P} -negligible subset of Ω . This is one respect in which the theory of continuous time random processes is technically more demanding than its discrete time analogue.

It follows immediately from Lemma 4.10 and from the identity

$$\forall n \geq 0, X_n = \widehat{X}_n \circ X$$

that the trajectory map is measurable with respect to the σ -fields \mathcal{A} and \mathcal{C} .

²More precisely, \mathcal{I} is the intersection of \mathcal{C} and the image σ -field of \mathcal{A} by f .

Exercise 4.9 What is the trajectory map of the canonical process $(\widehat{X}_n)_{n \geq 0}$ defined on the canonical space $(E^{\mathbb{N}}, \mathcal{C})$?

We can now give an upgraded and more canonical version of Proposition 4.6.

Theorem 4.12 Let E be a finite or countable set. Let P be a transition kernel on E . Let x be an element of E . There exists on the measurable space $(E^{\mathbb{N}}, \mathcal{C})$ a unique probability measure $\widehat{\mathbb{P}}_x$ such that $\widehat{X}_0 = x$ $\widehat{\mathbb{P}}_x$ -a.s. and under $\widehat{\mathbb{P}}_x$, the canonical process $\widehat{X} = (\widehat{X}_n)_{n \geq 0}$ is a Markov chain on E with transition kernel P .

Proof. To prove the existence of $\widehat{\mathbb{P}}_x$, let us apply Proposition 4.6 to find a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and a Markov chain $X = (X_n)_{n \geq 0}$ defined on this space, issued from x and with transition kernel P . Let us consider the trajectory map X of this process and define

$$\widehat{\mathbb{P}}_x = \mathbb{P} \circ X^{-1}.$$

In English, $\widehat{\mathbb{P}}_x$ is the law of the trajectory of the process $X = (X_n)_{n \geq 0}$. Let us check that $(\widehat{X}_n)_{n \geq 0}$ is under $\widehat{\mathbb{P}}_x$ a Markov chain on E issued from x with transition kernel P . For this, let us choose $n \geq 0$ and x_0, \dots, x_n in E . According to Proposition 4.3, we must prove that

$$\widehat{\mathbb{P}}_x(\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n) = \delta_{x, x_0} p(x_0, x_1) \dots p(x_{n-1}, x_n).$$

The following computation is elementary in that it consists exclusively in unfolding definitions:

$$\begin{aligned} \widehat{\mathbb{P}}_x(\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n) &= \mathbb{P}(X^{-1}(\{\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n\})) \\ &= \mathbb{P}(X^{-1}(\{\widehat{X}_0 = x_0\}) \cap \dots \cap X^{-1}(\{\widehat{X}_n = x_n\})) \\ &= \mathbb{P}((\widehat{X}_0 \circ X)^{-1}(\{x_0\}) \cap \dots \cap (\widehat{X}_n \circ X)^{-1}(\{x_n\})) \\ &= \mathbb{P}(X_0^{-1}(\{x_0\}) \cap \dots \cap X_n^{-1}(\{x_n\})) \\ &= \mathbb{P}(X_0 = x_0, \dots, X_n = x_n) \\ &= \delta_{x, x_0} p(x_0, x_1) \dots p(x_{n-1}, x_n), \end{aligned}$$

as expected³.

Let us turn to the uniqueness of $\widehat{\mathbb{P}}_x$. Suppose that $\widehat{\mathbb{Q}}_x$ is a probability measure on $(E^{\mathbb{N}}, \mathcal{C})$ under which the canonical process is a Markov chain on E issued from x and with transition kernel P . The for all $n \geq 0$ and all $x_0, \dots, x_n \in E$, we have

$$\widehat{\mathbb{P}}_x(\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n) = \widehat{\mathbb{Q}}_x(\widehat{X}_0 = x_0, \dots, \widehat{X}_n = x_n).$$

According to the remark made just before the statement of Lemma 4.10, this implies that for all $n \geq 0$, the probability measures $\widehat{\mathbb{P}}_x$ and $\widehat{\mathbb{Q}}_x$ agree on \mathcal{C}_n . Hence, they agree on the class $\bigcup_{n \geq 0} \mathcal{C}_n$ of all cylinder sets. We conclude the proof by applying the monotone class theorem, as follows.

On any measurable space, the class of measurable sets on which two probability measures agree is a monotone class (also called a λ -system). Applying this general fact in our particular situation, we find that the class of all elements of \mathcal{C} on which $\widehat{\mathbb{P}}_x$ and $\widehat{\mathbb{Q}}_x$ agree is a λ -system.

³The last computation is a proof of the fact, which we could have used directly, that, because $\widehat{\mathbb{P}}_x = \mathbb{P} \circ X^{-1}$ and $(X_0, \dots, X_n) = (\widehat{X}_0, \dots, \widehat{X}_n) \circ X$, the random variable $(\widehat{X}_0, \dots, \widehat{X}_n)$ has under $\widehat{\mathbb{P}}_x$ the same distribution as the random variable (X_0, \dots, X_n) under \mathbb{P} .

We just proved that this λ -system contains $\bigcup_{n \geq 0} \mathcal{C}_n$, which is a π -system. The monotone class theorem states that a λ -system which contains a π -system also contains the σ -field generated by this π -system. Hence, the class of sets on which $\hat{\mathbb{P}}_x$ and $\hat{\mathbb{Q}}_x$ agree contains the σ -field generated by all cylinder sets, which, by definition, is \mathcal{C} . Hence, $\hat{\mathbb{P}}_x$ and $\hat{\mathbb{Q}}_x$ agree on \mathcal{C} : they are equal. \square

We are now in position of giving a second, more sophisticated definition of a Markov chain.

Definition 4.13 (Markov chains, second definition) *Let E be a non-empty, finite or countable set. Let P be a transition kernel on E . A Markov chain with transition kernel P on E is a quintuple $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, (\mathbb{P}_x)_{x \in E}, X = (X_n)_{n \geq 0})$ in which $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0})$ is a filtered measurable space on which $X = (X_n)_{n \geq 0}$ is an adapted E -valued process which, for all x , is under \mathbb{P}_x , and in the sense of Definition 4.2, a Markov chain on E issued from x and with transition kernel P .*

Our discussion of the canonical space shows that every transition kernel on E is the transition kernel of a Markov chain on E in this more sophisticated sense.

Given a Markov chain in the sense of the definition above, we define, for every probability measure μ on E , the measure \mathbb{P}_μ by

$$\mathbb{P}_\mu = \int_E \mathbb{P}_x d\mu(x) = \sum_{x \in E} \mu(\{x\}) \mathbb{P}_x.$$

Under \mathbb{P}_μ , the process X is a Markov chain with initial distribution μ and transition kernel P .

There is a last improvement that we should make to our definition. There is a slight inconsistency in Definition 4.13 because we are given the filtration $(\mathcal{F}_n)_{n \geq 0}$ on our measurable space, and in Definition 4.2, and more precisely in Equation (3), we use the natural filtration of the process $(X_n)_{n \geq 0}$. We are assuming that the process $(X_n)_{n \geq 0}$ is adapted with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$, which means by definition that its natural filtration is included in the filtration $(\mathcal{F}_n)_{n \geq 0}$, but there is no reason why these two filtrations should be equal. This leads us to the following, final definition of a Markov chain.

Definition 4.14 (Markov chains, third definition) *Let E be a non-empty, finite or countable set. Let P be a transition kernel on E . A Markov chain with transition kernel P on E is a quintuple $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, (\mathbb{P}_x)_{x \in E}, X = (X_n)_{n \geq 0})$ in which $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0})$ is a filtered measurable space on which $X = (X_n)_{n \geq 0}$ is an adapted E -valued process such that for all $x, y \in E$ and all $n \geq 0$, one has*

$$\mathbb{P}_x(X_{n+1} = y | \mathcal{F}_n) = p(X_n, y).$$

It is important to see that a Markov chain in the third sense is a Markov chain in the second sense. Indeed, for a Markov chain in the second sense, we have, for all $x, y \in E$ and $n \geq 0$

$$\begin{aligned} \mathbb{P}_x(X_{n+1} = y | X_0, \dots, X_n) &= \mathbb{E}_x[\mathbb{E}_x[X_{n+1} = y | \mathcal{F}_n] | X_0, \dots, X_n] = \mathbb{E}_x[p(X_n, y) | X_0, \dots, X_n] \\ &= p(X_n, y). \end{aligned}$$

On the other hand, the 'if' part of Proposition 4.3 is not true any more: it is still true for a Markov chain in the third sense on E with transition matrix P that for all $n \geq 0$ and all $x_0, \dots, x_n \in E$ one has

$$\mathbb{P}_{x_0}(X_0 = x_0, \dots, X_n = x_n) = p(x_0, x_1) \dots p(x_{n-1}, x_n) \tag{6}$$

but knowing that this equality holds does not imply that $(X_n)_{n \geq 0}$ is a Markov in the third sense: this depends indeed on the filtration $(\mathcal{F}_n)_{n \geq 0}$.

Exercise 4.10 Consider a Markov chain in the third sense in the situation where for all $n \geq 0$, one has $\mathcal{F}_n = \sigma(X_0, \dots, X_{n+1})$. Prove that the transition matrix of this Markov chain can only take the values 0 and 1.

Exercise 4.11 Consider a Markov chain in the second sense with a transition matrix which takes at least one value in the interval $(0, 1)$. Assume that for all $n \geq 0$, the random variable X_{n+1} is measurable with respect to \mathcal{F}_n . Prove that the process $(X_n)_{n \geq 0}$ satisfies (6) but is not a Markov chain in the third sense with respect to the filtration $(\mathcal{F}_n)_{n \geq 0}$.

Exercise 4.12 Consider a Markov in the third sense. Prove that for all $x \in E$, all $n, m \geq 0$ and all $x_n, \dots, x_{n+m} \in E$, we have \mathbb{P}_x -almost surely

$$\mathbb{P}_x(X_n = x_n, \dots, X_{n+m} = x_{n+m} | \mathcal{F}_n) = \mathbb{1}_{\{X_n = x_n\}} p(x_n, x_{n+1}) \dots p(x_{n+m-1}, x_{n+m}).$$

The distinction between the second and third definition of Markov chains is a bit subtle, and can be ignored in a first time. In what follows, we stick to the second definition and make some comments to indicate how certain proofs should be modified to fit the third definition.

4.4 The Markov property

The essence of a Markov chain is the fact that it is memoryless, a property that can be rephrased by saying that it starts afresh at every instant. The Markov property expresses this absence of memory in a very effective way, and extends it to random times: we will prove that, in a certain sense, a Markov chain starts afresh at every stopping time.

In order to articulate the Markov property, we need to introduce one last piece of structure on the canonical space.

Definition 4.15 The shift operator on the canonical space is the map

$$\begin{aligned} \theta : E^{\mathbb{N}} &\longrightarrow E^{\mathbb{N}} \\ \omega = (\omega_i)_{i \geq 0} &\longmapsto \theta(\omega) = (\omega_{i+1})_{i \geq 0}. \end{aligned}$$

We also define $\theta_0 = \text{id}_{E^{\mathbb{N}}}$ and, for all $n \geq 2$,

$$\theta_n = \theta^n = \underbrace{\theta \circ \dots \circ \theta}_{n \text{ times}}.$$

For all $n \geq 0$, we have simply $\theta_n(\omega) = (\omega_{n+i})_{i \geq 0}$.

Exercise 4.13 Check that for all $n \geq 0$, the map θ_n is measurable with respect to the σ -field \mathcal{C} .

Finally, let us introduce a classical piece of notation. Let us consider a Markov chain in the sophisticated sense of Definition 4.13. For all $x \in E$, we naturally denote by \mathbb{E}_x the expectation with respect to \mathbb{P}_x . Let us extend this notation as follows. Assume that Z is an E -valued

random variable on (Ω, \mathcal{A}) . Then, for all non-negative random variable Y on (Ω, \mathcal{A}) , we define $\mathbb{E}_Z[Y]$ by the formula

$$\mathbb{E}_Z[Y] = \sum_{x \in E} \mathbb{E}_x[Y] \mathbb{1}_{\{Z=x\}}.$$

This can also be written

$$\mathbb{E}_Z[Y] = h(Z), \text{ where } h(x) = \mathbb{E}_x[Y].$$

Let us emphasise that $\mathbb{E}_Z[Y]$ is not a number in general, but a random variable.

We can now state the Markov property. The following statement is called the *weak* Markov property because it involves ‘only’ deterministic times. Immediately after proving it, we shall strengthen it into the so-called strong Markov property, which covers the case of random times.

Theorem 4.16 (Weak Markov property) *Let P be a transition kernel on the state space E . Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, (\mathbb{P}_x)_{x \in E}, X = (X_n)_{n \geq 0})$ be a Markov chain on E with transition kernel P .*

Let $F : (E^{\mathbb{N}}, \mathcal{C}) \rightarrow (\mathbb{R}_+, \mathcal{B}_{\mathbb{R}_+})$ be a measurable non-negative function. For all integer $n \geq 0$ and all $x \in E$, we have the equality

$$\mathbb{E}_x[F(\theta_n(X)) | \mathcal{F}_n] = \mathbb{E}_{X_n}[F(X)] \quad \mathbb{P}_x\text{-a.s.} \quad (7)$$

of non-negative random variables on $(\Omega, \mathcal{A}, \mathbb{P}_x)$.

Proof. The right-hand side of (7) is a function of X_n , so that it is \mathcal{F}_n -measurable. What remains to prove is that for every $A \in \mathcal{F}_n$, both sides of (7) have the same \mathbb{P}_x -integral over A . Since both sides of (7) depend linearly on F , and since any measurable positive function can be written as the pointwise limit of an increasing sequence of simple functions, it is enough to prove (7) when F is an indicator function, that is, $F = \mathbb{1}_C$ for some $C \in \mathcal{C}$. Thus, we need to prove that

$$\forall A \in \mathcal{F}_n, \forall C \in \mathcal{C}, \quad \mathbb{E}_x[\mathbb{1}_C(\theta_n(X)) \mathbb{1}_A] = \mathbb{E}_x[\mathbb{E}_{X_n}[\mathbb{1}_C(X)] \mathbb{1}_A]. \quad (8)$$

Here we use a monotone class argument to reduce the problem further. The class of all sets $C \in \mathcal{C}$ such that the last equality holds is a λ -system. In order to prove that it contains \mathcal{C} , it suffices to prove that it contains the π -system $\bigcup_{m \geq 0} \mathcal{C}_m$, which generates the σ -field \mathcal{C} . Thus, it is enough to prove that the equality holds when $C \in \mathcal{C}_m$ for an arbitrary m . Finally, according to the remark made before Lemma 4.10, it suffices to prove that for all choices of $m \geq 0$, $x_0, \dots, x_m \in E$, $y_0, \dots, y_m \in E$, the equality holds for

$$A = \{X_0 = x_0, \dots, X_n = x_n\} \text{ and } C = \{\widehat{X}_0 = y_0, \dots, \widehat{X}_m = y_m\}.$$

In this case, we can compute both sides of (8). The left-hand side is equal to

$$\begin{aligned} \mathbb{E}_x[\mathbb{1}_C(\theta_n(X)) \mathbb{1}_A] &= \mathbb{P}_x(X_0 = x_0, \dots, X_n = x_n, X_n = y_0, \dots, X_{n+m} = y_m) \\ &= \delta_{x, x_0} p(x_0, x_1) \dots p(x_{n-1}, x_n) \delta_{x_n, y_0} p(y_0, y_1) \dots p(y_{m-1}, y_m). \end{aligned}$$

The right-hand side is

$$\begin{aligned} \mathbb{E}_x[\mathbb{E}_{X_n}[\mathbb{1}_C(X)] \mathbb{1}_A] &= \mathbb{E}_x[\mathbb{E}_{x_n}[\mathbb{1}_C(X)] \mathbb{1}_A] \\ &= \mathbb{E}_x[\mathbb{1}_A] \mathbb{E}_{x_n}[\mathbb{1}_C(X)] \\ &= \mathbb{P}_x(X_0 = x_0, \dots, X_n = x_n) \mathbb{P}_{x_n}(X_0 = y_0, \dots, X_m = y_m) \\ &= \delta_{x, x_0} p(x_0, x_1) \dots p(x_{n-1}, x_n) \delta_{x_n, y_0} p(y_0, y_1) \dots p(y_{m-1}, y_m), \end{aligned}$$

and the equality is proved. \square

Exercise 4.14 Prove the weak Markov property for Markov chains in the third sense (Definition 4.14). Hint: use the result of Exercise 4.12. You will see that the proof is slightly easier than the proof above.

Let us extend the Markov property to random times. For this, we need to remember the definition of a stopping time (see Definition 3.15) and the definition of the σ -field of events prior to a stopping time (see Definition 3.34). We need also to define the shift by a stopping time.

For this, let us consider again a Markov chain in the sense of Definition 4.13. Let T be a stopping time on the measurable space $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0})$. On the event $\{T < \infty\}$, we denote by $\theta_T(X)$ the map from (Ω, \mathcal{A}) to $(E^{\mathbb{N}}, \mathcal{C})$ defined by

$$(\theta_T(X))(\omega) = \theta_{T(\omega)}(X(\omega)) = (X_{T+i}(\omega))_{i \geq 0}.$$

Finally, let us agree on the following convention: if $h : \mathbb{N} \rightarrow \mathbb{R}$ is a function and T is a stopping time, we use the notation

$$h(T) \mathbb{1}_{\{T < \infty\}},$$

although T might take the value ∞ and $h(\infty)$ is not defined, to indicate

$$h(T) \mathbb{1}_{\{T < \infty\}} = \sum_{n \geq 0} h(n) \mathbb{1}_{\{T=n\}}.$$

To be clear, this random variable vanishes on the event $\{T = \infty\}$. With all this preparation, we can state the strong Markov property.

Theorem 4.17 (Strong Markov property) Let P be a transition kernel on the state space E . Let $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, (\mathbb{P}_x)_{x \in E}, X = (X_n)_{n \geq 0})$ be a Markov chain on E with transition kernel P .

Let $F : (E^{\mathbb{N}}, \mathcal{C}) \rightarrow (\mathbb{R}_+, \mathcal{B}_{\mathbb{R}_+})$ be a measurable non-negative function. Let T be a stopping time on $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0})$. For all $x \in E$, we have the equality

$$\mathbb{E}_x[F(\theta_T(X)) \mathbb{1}_{\{T < \infty\}} | \mathcal{F}_T] = \mathbb{E}_{X_T}[F(X)] \mathbb{1}_{\{T < \infty\}} \quad \mathbb{P}_x\text{-a.s.} \quad (9)$$

of non-negative random variables on $(\Omega, \mathcal{A}, \mathbb{P}_x)$.

Proof. Just as in the proof of the weak Markov property, the right-hand side of (9), which is a function of X_T , is \mathcal{F}_T -measurable, and it suffices to prove that for every $A \in \mathcal{F}_T$, both sides of (9) have the same \mathbb{P}_x -integral over A .

Let us choose $A \in \mathcal{F}_T$. For all $n \geq 0$, we have

$$\mathbb{E}_x[F(\theta_T(X)) \mathbb{1}_{\{T=n\}} \mathbb{1}_A] = \mathbb{E}_x[F(\theta_T(X)) \mathbb{1}_{\{T=n\} \cap A}].$$

Using the weak Markov property and the fact that $\{T = n\} \cap A$ belongs to \mathcal{F}_n , we find

$$\mathbb{E}_x[F(\theta_T(X)) \mathbb{1}_{\{T=n\}} \mathbb{1}_A] = \mathbb{E}_x[\mathbb{E}_{X_n}[F(X)] \mathbb{1}_{\{T=n\} \cap A}] = \mathbb{E}_x[\mathbb{E}_{X_T}[F(X)] \mathbb{1}_{\{T=n\}} \mathbb{1}_A].$$

Summing this equality over n , we find

$$\mathbb{E}_x[F(\theta_T(X)) \mathbb{1}_{\{T < \infty\}} \mathbb{1}_A] = \mathbb{E}_x[\mathbb{E}_{X_T}[F(X)] \mathbb{1}_{\{T < \infty\}} \mathbb{1}_A],$$

and the proof is finished. \square

We will see over time that the theorem that we just proved expresses the Markov property in a very convenient and powerful way. For the moment, let us give a simple consequence.

Corollary 4.18 *We use the notation of Theorem 4.17. Let x, y be elements of E and let T be a stopping time such that $T < \infty$ and $X_T = y$ \mathbb{P}_x -a.s. Then under \mathbb{P}_x , $\theta_T(X)$ is independent of \mathcal{F}_T and has the same distribution as X under \mathbb{P}_y .*

Proof. The statement is equivalent to saying that for all non-negative measurable function F on $E^{\mathbb{N}}$ and all $B \in \mathcal{F}_T$,

$$\mathbb{E}_x[F(\theta_T(X))\mathbb{1}_B] = \mathbb{P}_x(B)\mathbb{E}_y[F(X)].$$

Let us compute the expectation on the left-hand side by conditioning with respect to \mathcal{F}_T and using the strong Markov property. We find

$$\mathbb{E}_x[F(\theta_T(X))\mathbb{1}_B] = \mathbb{E}_x[\mathbb{E}_{X_T}[F(X)]\mathbb{1}_B] = \mathbb{E}_x[\mathbb{E}_y[F(X)]\mathbb{1}_B] = \mathbb{P}_x(B)\mathbb{E}_y[F(X)],$$

as expected. □

4.5 Recurrent and transient states

Equipped with the powerful tools developed in the last section, we can embark on the study of the recurrence properties of the states of a Markov chain. From this point on, and until further notice, we fix a Markov chain $(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, (\mathbb{P}_x)_{x \in E}, X = (X_n)_{n \geq 0})$ on the state space E with transition kernel P .

For every state $x \in E$, let us introduce a random variable N_x with values in $\mathbb{N} \cup \{\infty\}$ defined by

$$N_x = \sum_{n=0}^{\infty} \mathbb{1}_{\{X_n=x\}}$$

and a stopping time

$$T_x = \inf\{n \geq 1 : X_n = x\}.$$

The random variable N_x is the total number of visits of the chain to the state x and the stopping time T_x is the first return time of the chain to x . Note that $T_x \geq 1$ by construction, and the usual convention $\inf \emptyset = \infty$ applies.

It is a simple but important observation that the following equality of events holds:

$$\{N_x \geq \mathbb{1}_{\{X_0=x\}} + 1\} = \{T_x < \infty\}.$$

These two events can be described as 'the chain visits the state x at least once strictly after time to 0'.

We will prove that, starting from x , either the Markov chain visits x infinitely often almost surely, or it visits x a finite number of times which moreover has finite expectation.

Before that, let us introduce the canonical versions of N_x and T_x : let us define on the canonical space

$$\widehat{N}_x = \sum_{n=0}^{\infty} \mathbb{1}_{\{\widehat{X}_n=x\}} \quad \text{and} \quad \widehat{T}_x = \inf\{n \geq 1 : \widehat{X}_n = x\}.$$

We have $N_x = \widehat{N}_x(X)$ and $T_x = \widehat{T}_x(X)$.

Proposition 4.19 *Let x be an element of E . Exactly one of the following two situations occurs.*

1. $\mathbb{P}_x(T_x < \infty) = 1$. *In this case, $N_x = \infty$ \mathbb{P}_x -a.s. and one says that x is recurrent.*

2. $\mathbb{P}_x(T_x < \infty) < 1$. *In this case, $N_x < \infty$ \mathbb{P}_x -a.s. and one says that x is transient.*

Moreover, in this case,

$$\mathbb{E}_x[N_x] = \frac{1}{\mathbb{P}_x(T_x = \infty)}.$$

Proof. Let $k \geq 1$ be an integer and let us compute $\mathbb{P}_x(N_x \geq k + 1)$.

$$\begin{aligned} \mathbb{P}_x(N_x \geq k + 1) &= \mathbb{P}_x(\widehat{N}_x(X) \geq k + 1) \\ &= \mathbb{P}_x(T_x < \infty, \widehat{N}_x(\theta_{T_x}(X)) \geq k) \\ &= \mathbb{E}_x[\mathbb{1}_{\{\widehat{N}_x \geq k\}}(\theta_{T_x}(X)) \mathbb{1}_{\{T_x < \infty\}}] \\ &= \mathbb{E}_x[\mathbb{E}_{T_x}[\mathbb{1}_{\{\widehat{N}_x \geq k\}}(X)] \mathbb{1}_{\{T_x < \infty\}}] \\ &= \mathbb{P}_x(N_x \geq k) \mathbb{P}_x(T_x < \infty). \end{aligned}$$

Since $\mathbb{P}_x(N_x \geq 1) = 1$, we get by induction

$$\forall k \geq 1, \mathbb{P}_x(N_x \geq k) = \mathbb{P}_x(T_x < \infty)^{k-1}.$$

If $\mathbb{P}_x(T_x < \infty) = 1$, we deduce that $\mathbb{P}_x(N_x = \infty) = 1$. This is the recurrent case.

On the other hand, if $\mathbb{P}_x(T_x < \infty) < 1$, then

$$\begin{aligned} \mathbb{E}_x[N_x] &= \sum_{k \geq 1} \mathbb{P}_x(N_x \geq k) \\ &= \sum_{k \geq 0} \mathbb{P}_x(T_x < \infty)^k \\ &= \frac{1}{1 - \mathbb{P}_x(T_x < \infty)}, \end{aligned}$$

so that $\mathbb{E}_x[N_x] < \infty$ and N_x is finite \mathbb{P}_x -almost surely. □

The dichotomy between recurrent and transient states suggests the definition of the following function.

Definition 4.20 *The Green function of the chain is the function $G : E \times E \rightarrow [0, \infty]$ defined by*

$$\forall x, y \in E, \quad G(x, y) = \mathbb{E}_x[N_y].$$

Proposition 4.21 1. *For all $x, y \in E$, we have*

$$G(x, y) = \sum_{n=0}^{\infty} P^n(x, y).$$

2. *The state $x \in E$ is recurrent if and only if $G(x, x) = \infty$.*

3. *If $x \neq y$, then*

$$G(x, y) = \mathbb{P}_x(T_y < \infty) G(y, y).$$

In particular, $G(x, y) \leq G(y, y)$.

4. *For all $x \in E$,*

$$G(x, x) = 1 + \mathbb{P}_x(T_x < \infty) G(x, x).$$

Proof. 1. By definition of G and the monotone convergence theorem,

$$G(x, y) = \sum_{n=0}^{\infty} \mathbb{E}_x[\mathbb{1}_{\{X_n=y\}}] = \sum_{n=0}^{\infty} \mathbb{P}_x[X_n = y] = \sum_{n=0}^{\infty} P^n(x, y).$$

2. By the previous proposition, x is recurrent if and only if $\mathbb{E}_x[N_x] < \infty$.

3. Under \mathbb{P}_x , we have almost surely $N_y = \widehat{N}_y(\theta_{T_y}(X)) \mathbb{1}_{\{T_y < \infty\}}$, so that

$$G(x, y) = \mathbb{E}_x[N_y] = \mathbb{E}_x[\mathbb{E}_y[\widehat{N}_y(X)] \mathbb{1}_{\{T_y < \infty\}}] = \mathbb{P}_x(T_y < \infty) \mathbb{E}_y[N_y] = \mathbb{P}_x(T_y < \infty) G(y, y).$$

4. Under \mathbb{P}_x , we have almost surely $N_x = 1 + \widehat{N}_x(\theta_{T_x}(X)) \mathbb{1}_{\{T_x < \infty\}}$, and the relation on the Green functions follows. \square

The Green function allows us to define a binary relation on E : given two states x and y , we say that x leads to y , and we write $x \rightarrow y$, if $G(x, y) > 0$. Thus, x leads to y if and only if starting from x , there is a positive probability⁴ of visiting y . If $x \neq y$, then the following four assertions are equivalent:

$$x \rightarrow y \iff G(x, y) > 0 \iff \mathbb{P}_x(N_y \geq 1) > 0 \iff \mathbb{P}_x(T_y < \infty) > 0.$$

Lemma 4.22 *On E , the binary relation \rightarrow is reflexive and transitive.*

Proof. For every $x \in E$, one has $G(x, x) \geq 1$, so that $x \rightarrow x$. This means that \rightarrow is a reflexive relation.

To prove that it is transitive, consider three states x, y and z such that $x \rightarrow y$ and $y \rightarrow z$. If any two of these states are equal, then it is true that $x \rightarrow z$. On the other hand, if they are pairwise distinct, then

$$\begin{aligned} G(x, z) &\geq \mathbb{P}_x(N_z \geq 1) \\ &= \mathbb{P}_x(T_z < \infty) \\ &\geq \mathbb{P}_x(T_y < \infty, \widehat{T}_z(\theta_{T_y}(X)) < \infty) \\ &= \mathbb{P}_x(T_y < \infty) \mathbb{P}_y(T_z < \infty) \\ &= \mathbb{P}_x(N_y \geq 1) \mathbb{P}_y(N_z \geq 1) \\ &> 0 \end{aligned}$$

and $x \rightarrow z$. \square

It is not true in general that the relation \rightarrow is symmetric. We will now prove that its restriction to the subset of all recurrent states is symmetric.

Proposition 4.23 *Let x and y be two states. Assume that x is recurrent and x leads to y . Then y is recurrent, y leads to x , and*

$$\mathbb{P}_x(T_y < \infty) = \mathbb{P}_y(T_x < \infty) = 1.$$

⁴The notation may be misleading: $x \rightarrow y$ means that starting from x , a visit to y occurs with positive probability. It does not mean that starting from x , a visit to y is certain. This is the difference between $\mathbb{P}_x(N_y \geq 1) > 0$ and $\mathbb{P}_x(N_y \geq 1) = 1$.

Proof. If $y = x$, then there is nothing to prove. Let us assume that $y \neq x$. Then

$$\begin{aligned} 0 = \mathbb{P}_x(T_x = \infty) &\geq \mathbb{P}_x(T_y < \infty, \widehat{T}_x(\theta_{T_y}(X)) = \infty) \\ &= \mathbb{P}_x(T_y < \infty)\mathbb{P}_y(T_x = \infty). \end{aligned}$$

Since, by assumption, $\mathbb{P}_x(T_y < \infty) > 0$, we find $\mathbb{P}_y(T_x = \infty) = 0$, that is, $\mathbb{P}_y(T_x < \infty) = 1$. In particular, $y \rightarrow x$.

Since $G(x, y) > 0$ and $G(y, x) > 0$, there exists two integers $a, b \geq 1$ such that

$$P^a(x, y) > 0 \text{ and } P^b(y, x) > 0.$$

Hence,

$$\begin{aligned} G(y, y) &= \sum_{n \geq 0} P^n(y, y) \\ &\geq \sum_{n \geq a+b} P^n(y, y) \\ &\geq \sum_{c \geq 0} P^b(y, x)P^c(x, x)P^a(x, y) \\ &= P^a(x, y)G(x, x)P^b(y, x) \\ &= \infty \end{aligned}$$

and y is recurrent.

Since y is recurrent and $y \rightarrow x$, we proved already that $\mathbb{P}_x(T_y < \infty) = 1$. □

Let us write $x \sim y$ if $x \rightarrow y$ and $y \rightarrow x$. The relation \sim is an equivalence relation on E . It follows from the last proposition that the relations \rightarrow and \sim coincide on the set of all recurrent states. Note also that this proposition implies that if x is recurrent and y is transient, then $G(x, y) = 0$.

Proposition 4.24 *Under the assumptions of the previous proposition, we have*

$$\mathbb{P}_x(N_y = \infty) = \mathbb{P}_y(N_x = \infty) = 1.$$

In particular,

$$G(x, y) = G(y, x) = \infty.$$

Proof. Indeed, y is recurrent and for all $k \geq 1$,

$$\mathbb{P}_x(N_y \geq k) = \mathbb{P}_x(T_y < \infty)\mathbb{P}_y(N_y \geq k) = 1.$$

Hence, $\mathbb{P}_x(N_y = \infty) = 1$. □

We can now state the following result which summarises our study.

Theorem 4.25 (Classification of states) *Let R be the subset of E consisting of all recurrent states. Let*

$$R = \bigsqcup_{i \in I} R_i$$

be the partition of R in equivalence classes for the relation \sim . Consider a state $x \in E$.

1. If x is recurrent, let $i \in I$ be such that $x \in R_i$. Then \mathbb{P}_x -a.s.,

$$N_y = \infty \text{ for all } y \in R_i \text{ and } N_z = 0 \text{ for all } z \in E \setminus R_i.$$

In English, starting from a recurrent state, the chain stays forever in the class of its initial state and visits infinitely often every state of this class.

2. If x is transient, define $T = \inf\{n \geq 0 : X_n \in R\}$. Then

$$\mathbb{P}_x(\{T = \infty \text{ and } N_y < \infty \text{ for all } y \in E\} \cup \{T < \infty \text{ and } \exists j \in I, \forall n \geq T, X_n \in R_j\}) = 1.$$

In English, starting from a transient state, either the chain never visits a recurrent state, in which case it visits every state a finite number of times, or it eventually visits a recurrent state, in which case it gets stuck forever in the class of the first recurrent state which it visits.

The equivalence classes of R under the relation \sim are called the *recurrence classes* of the chain.

Proof. 1. Consider $y \in R_i$. Then, according to Proposition 4.24, we have $\mathbb{P}_x(N_y = \infty) = 1$. Consider now $z \in E \setminus R_i$. If z is transient, then $G(x, z) = 0$ by Proposition 4.23. If z is recurrent, then $G(x, z) = 0$ by definition of the equivalence class \sim .

2. Consider a transient state y . The third assertion of Proposition 4.21 asserts that $G(x, y) \leq G(y, y)$, so that $G(x, y) < \infty$. In particular, N_y is finite \mathbb{P}_x -almost surely.

On the event $\{T < \infty\}$, let J be the random element of I such that $X_T \in R_J$. Then

$$\mathbb{P}_x(T < \infty \text{ and } \forall n \geq T, X_n \in R_J) = \mathbb{E}_x[\mathbb{1}_{\{T < \infty\}} \mathbb{P}_{X_T}(\forall n \geq T, X_n \in R_J)].$$

By the first part of the theorem, and since $X_T \in R_J$ by definition of J , we have $\mathbb{P}_{X_T}(\forall n \geq T, X_n \in R_J) = 1$ on the event $T < \infty$. Thus,

$$\mathbb{P}_x(T < \infty \text{ and } \forall n \geq T, X_n \in R_J) = \mathbb{P}(T < \infty),$$

which is what we wanted to prove. □

Definition 4.26 *The Markov chain is said to be irreducible if $x \rightarrow y$ for all $x, y \in E$.*

Corollary 4.27 *Let us assume that the Markov chain is irreducible. Then we are in exactly one of the following two situations.*

- *All states are recurrent, there is only one recurrence class and*

$$\mathbb{P}_x(\forall y \in E, N_y = \infty) = 1.$$

- *All states are transient and*

$$\mathbb{P}_x(\forall y \in E, N_y < \infty) = 1.$$

If E is a finite set, then we are in the first situation.

Proof. If there is one recurrent state, then by Proposition 4.23, all states are recurrent and there is only one class. Moreover, Proposition 4.24 implies that every state is visited infinitely often \mathbb{P}_x -almost surely for every $x \in E$.

If there is no recurrent state, then all states are transient and, for all states x, y , we have $G(x, y) \leq G(y, y) < \infty$, so that N_y is finite \mathbb{P}_x -almost surely.

Finally, if E is finite, there exists at least one state that is visited infinitely often, and we must be in the first situation. \square

In the first situation, one says that the Markov chain is *irreducible and recurrent*. Let us study an important example.

Theorem 4.28 (Random walks on \mathbb{Z}) *Let μ be a probability measure on \mathbb{Z} . Let P be the transition kernel of the random walk on \mathbb{Z} with jump distribution μ :*

$$\forall x, y \in \mathbb{Z}, \quad p(x, y) = \mu(y - x).$$

Let ξ be a random variable with distribution μ . Let us assume that $\mathbb{E}[|\xi|] < \infty$.

1. *If $\mathbb{E}[\xi] \neq 0$, then every state is transient.*
2. *If $\mathbb{E}[\xi] = 0$, then every state is recurrent. Moreover, the chain is irreducible if and only if the subgroup of \mathbb{Z} generated by $\{x \in \mathbb{Z} : \mu(x) > 0\}$ is \mathbb{Z} itself.*

Proof. 1. For all $x \in \mathbb{Z}$, the strong law of large numbers implies that $|X_n|$ tends to $+\infty$ as n tends to infinity, \mathbb{P}_x -almost surely. Hence, N_x is finite \mathbb{P}_x -almost surely and x is transient.

2. There is an invariance by translation of the problem which implies that all states have the same nature. It is thus sufficient to prove that 0 is recurrent.

Let us consider an integer $p \geq 0$ and a positive real $\varepsilon > 0$, both to be specified later, and let us estimate

$$\sum_{|x| \leq \varepsilon p} G(0, x)$$

in two ways. Firstly, we want to say that this sum is not too big. For this, we argue that for every $x \in \mathbb{Z}$, we have $G(x, x) = G(0, 0)$ by invariance by translation, and

$$G(0, x) \leq G(x, x) = G(0, 0).$$

Thus,

$$\sum_{|x| \leq \varepsilon p} G(0, x) \leq (2\varepsilon p + 1)G(0, 0).$$

Secondly, we want to say that the same sum is not too small. For this, we use the weak law of large numbers, according to which $\frac{X_n}{n}$ converges to 0 in probability. This implies that there exists n_0 , which depends on ε , such that for all $n \geq n_0$,

$$\mathbb{P}_0(|X_n| \leq \varepsilon n) \geq \frac{1}{2}.$$

For all $p \geq n_0$, we have then

$$\begin{aligned}
\sum_{|x| \leq \varepsilon p} G(0, x) &= \sum_{|x| \leq \varepsilon p} \sum_{n=0}^{\infty} P^n(0, x) \\
&\geq \sum_{|x| \leq \varepsilon p} \sum_{n=n_0}^p P^n(0, x) \\
&\geq \sum_{n=n_0}^p \mathbb{P}_0(|X_n| \leq \varepsilon p) \\
&\geq \frac{p - n_0 + 1}{2}.
\end{aligned}$$

Thus, we proved that for all $\varepsilon > 0$, there exists an n_0 such that for all $p \geq n_0$,

$$\frac{p - n_0 + 1}{2} \leq \sum_{|x| \leq \varepsilon p} G(0, x) \leq (2\varepsilon p + 1)G(0, 0).$$

Thus, for all $\varepsilon > 0$, we have

$$G(0, 0) \geq \lim_{p \rightarrow \infty} \frac{p - n_0 + 1}{4\varepsilon p + 2} = \frac{1}{4\varepsilon}.$$

This implies that $G(0, 0) = \infty$, and 0 is recurrent.

There remains to study the irreducibility of the chain. Let us denote by S the set $\{x \in \mathbb{Z} : \mu(x) > 0\}$.

Let us first assume that the chain is irreducible. Then in particular 0 leads to 1, which means that there exists $n \geq 0$ such that $P^n(0, 1) > 0$. This in turns means that there exists integers x_1, \dots, x_{n-1} such that $p(0, x_1)p(x_1, x_2) \dots p(x_{n-1}, 1) > 0$. Thus, S contains the integers $x_1, x_2 - x_1, \dots, x_{n-1} - x_{n-2}, 1 - x_{n-1}$, which add up to 1. The subgroup of \mathbb{Z} generated by S is thus \mathbb{Z} .

Conversely, let us assume that S generates \mathbb{Z} . Then, let us choose $z \in \mathbb{Z}$ and let us prove that $0 \rightarrow z$. Then, there exists x_1, \dots, x_k and y_1, \dots, y_l in S such that

$$z = x_1 + \dots + x_k - y_1 - \dots - y_l.$$

Let us write $x = x_1 + \dots + x_k$ and $y = y_1 + \dots + y_l$. Now, we have on one hand

$$P^k(0, x) \geq \mu(x_1) \dots \mu(x_k) > 0,$$

so that $0 \rightarrow x$, and on the other hand

$$P^l(z, x) \geq p(z, z + y_1) \dots p(z + y_1 + \dots + y_{l-1}, x) = \mu(y_1) \dots \mu(y_l) > 0,$$

so that $z \rightarrow x$. Since z is recurrent, it follows that $x \rightarrow z$, and finally $0 \rightarrow z$. \square

Exercise 4.15 *With the notation of the theorem, give an example of a probability measure μ such that the random walk is transient, but such that the support of μ generates \mathbb{Z} as an additive group.*

4.6 Invariant measures

As in the previous section, we fix once and for all a Markov chain

$$(\Omega, \mathcal{A}, (\mathcal{F}_n)_{n \geq 0}, (\mathbb{P}_x)_{x \in E}, X = (X_n)_{n \geq 0})$$

on the state space E with transition kernel P .

Definition 4.29 *A measure μ on E is an invariant measure of the transition kernel P if μ is not the zero measure, μ gives finite mass to every singleton, and*

$$\forall y \in E, \sum_{x \in E} \mu(x)p(x, y) = \mu(y). \quad (10)$$

The conditions that μ is not zero and gives finite mass to every singleton can be written symbolically as

$$\exists x \in E, \mu(x) > 0 \text{ and } \forall x \in E, \mu(x) < \infty.$$

The relation (10) can be written matrixially, as follows. Let us remember that P can be thought of as a possibly infinite square matrix with as many rows and columns as E has elements. Let us think of a measure on E as a row, the width (or length) of which is $|E|$, the number of elements of E . In a dual fashion, it would be appropriate to think of a function on E as a column with height EI . Then, if μ is a measure and f a function on E , the matrix product μf is a 1×1 matrix, that is, a number, which is none but the integral of f with respect to μ :

$$\mu f = \int_E f d\mu.$$

Exercise 4.16 *Let f be a function on E , seen as a column. The matrix product Pf is well defined and it is a column. Thus, it represents a function on E . What is this function?*

Equation (10) can be written

$$\mu P = \mu.$$

This shows that there is a purely linear algebraic approach to the study of invariant measures. Indeed, if E is finite, μ is an invariant measures if and only if μ^* , the column obtained by transposing μ , is an eigenvector with non-negative coefficients of the transposed kernel P^* , associated with the eigenvalue 1. The piece of linear algebra which deals with stochastic matrices is called the Perron-Frobenius theory, and we shall discuss it later.

Let us give an example of an invariant measure. Let us consider a random walk on \mathbb{Z}^d , with arbitrary jump distribution. Thus, η is a probability measure on \mathbb{Z}^d and for all $x, y \in \mathbb{Z}^d$, we have $p(x, y) = \eta(y - x)$. Then the equality

$$\sum_{x \in \mathbb{Z}^d} p(x, y) = \sum_{x \in \mathbb{Z}^d} \eta(y - x) = \sum_{x \in \mathbb{Z}^d} \eta(x) = 1,$$

valid for all $y \in \mathbb{Z}^d$, shows that the counting measure μ , given by

$$\forall x \in \mathbb{Z}^d, \mu(x) = 1$$

is an invariant measure of this Markov chain.

It is an elementary observation that any positive multiple of an invariant measure is still an invariant measure. If there exists an invariant measure μ with the property that $\mu(E) < \infty$, it is natural to normalise it and to define

$$\pi = \frac{1}{\mu(E)},$$

which is an invariant probability measure on E . Then, for every bounded function f on E , the matricial relation

$$\mu P f = \mu f$$

can be written in terms of the Markov chain as

$$\mathbb{E}_\pi[f(X_1)] = \int_E f d\pi,$$

where \mathbb{E}_π denotes the expectation with respect to the probability measure

$$\mathbb{P}_\pi = \sum_{x \in E} \pi(x) \mathbb{P}_x.$$

In other words, if the initial distribution of the Markov chain is an invariant probability measure, then the Markov chain is stationary, in the sense that the distribution of X_n does not depend on n .

There is another instructive way of thinking of (10). Suppose that we consider the evolution of a very large assembly of particles (or sand grains, or people). At the initial time, the quantity of particles (or sand, or people) that is present at each state x is measured by the number $\mu(x)$. For example, $\mu(x)$ is the number of litres of water present at the state x at time 0. Then, between time 0 and time 1, from every state x , and for every state y , a proportion $p(x, y)$ of the water present at x moves to y . Thus, on one hand, the water received by the state y in this process is

$$\sum_{x \in E} \mu(x) p(x, y).$$

On the other hand, the water lost by the same state y is

$$\sum_{x \in E} \mu(y) p(y, x),$$

which incidentally is equal to $\mu(y)$. To say that μ is invariant is to say that the water received by y exactly compensates the water lost by y : this is the equality expressed by (10).

One could demand that more is true, and demand that in the process just described, for all states x and y , the water received by y from x exactly compensates the water sent from y to x .

Definition 4.30 *The measure μ on E is said to be reversible with respect to P if μ is not the zero measure, μ gives finite mass to every singleton, and*

$$\forall x, y \in E, \mu(x) p(x, y) = \mu(y) p(y, x). \quad (11)$$

Our discussion should make it clear that the following statement holds.

Proposition 4.31 *A reversible measure is an invariant measure.*

Exercise 4.17 *Prove directly this proposition.*

Let us give an example of a reversible measure. Let us consider on \mathbb{Z} the random walk with jump distribution η given by

$$\eta(1) = p \text{ and } \eta(-1) = q = 1 - p$$

for some $p \in (0, 1)$. This is called the p -biased random walk on \mathbb{Z} .

Then the measure μ given by

$$\mu(i) = \left(\frac{p}{q}\right)^i \tag{12}$$

is reversible. If $p \neq q$, that is, if $p \neq \frac{1}{2}$, this measure is not proportional to the counting measure on \mathbb{Z} , of which we know already that it is invariant.

Exercise 4.18 *Describe, for every $p \in (0, 1)$, the set of all invariant measures of the p -biased random walk on \mathbb{Z} .*

Exercise 4.19 *Recall the random walk on a graph described page 63. Prove that the formula*

$$\mu(x) = |A_x|$$

defines an invariant measure on E .

When a reversible measure exists, it is in general *much* easier to find by solving (11) than by solving (10). In other words, when looking for an invariant measure, one should always start by looking for a reversible measure.

Exercise 4.20 *Draw a small connected graph. Find an invariant probability measure for the random walk on this graph by solving (10) (for this to be tractable by hand, your graph should not have more than a few vertices, a dozen would typically be too much, unless your graph has a lot of symmetry and you find a way to use it). Compare with the effort needed to find a reversible measure for the same random walk.*

Another instance of the fact that reversible measures are nice to work with is given by the following exercise.

Exercise 4.21 *Under the assumption that the Markov chain is irreducible, prove that any two reversible measures are proportional.*

Is it true, under the same assumption of irreducibility, that any two invariant measures are proportional?

We will soon prove that if the chain is irreducible *and recurrent*, then any two invariant measures are proportional, but this will require more than a two-line proof.

Let us leave reversible measures aside and consider general invariant measures again. The fundamental result is the following.

Theorem 4.32 *Let $x \in E$ be a recurrent state. The formula*

$$\mu(y) = \mathbb{E}_x \left[\sum_{i=0}^{T_x-1} \mathbb{1}_{\{X_i=y\}} \right]$$

defines an invariant measure on E . Moreover, the support of this measure is

$$\{y \in E : \mu(y) > 0\} = \{y \in E : x \rightarrow y\} = \{y \in E : x \sim y\},$$

the communication class of x .

Proof. The first observation is that $\mu(x) = 1$ by definition of T_x . Hence, μ is not identically zero. Let us prove that μ satisfies (10). To start with, whether $y = x$ or $y \neq x$, we have

$$\mu(y) = \mathbb{E}_x \left[\sum_{i=1}^{T_x} \mathbb{1}_{\{X_i=y\}} \right].$$

Now, let us compute by considering not only the i -th state of the walk, but also the state immediately before.

$$\begin{aligned} \mu(y) &= \mathbb{E}_x \left[\sum_{i=1}^{T_x} \sum_{z \in E} \mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{X_i=y\}} \right] \\ &= \sum_{i=1}^{\infty} \sum_{z \in E} \mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{X_i=y\}} \mathbb{1}_{\{i \leq T_x\}} \right]. \end{aligned}$$

Now, for all $i \geq 1$, $y, z \in E$, we have

$$\mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{X_i=y\}} \mathbb{1}_{\{i \leq T_x\}} \right] = \mathbb{E}_x \left[\mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{X_i=y\}} \mathbb{1}_{\{i \leq T_x\}} \mid \mathcal{F}_{i-1} \right] \right].$$

Using the fact that $\{i \leq T_x\} = \{T_x \leq i-1\}^c$ belongs to \mathcal{F}_{i-1} , we find

$$\begin{aligned} \mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{X_i=y\}} \mathbb{1}_{\{i \leq T_x\}} \right] &= \mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{E}_x \left[\mathbb{1}_{\{X_i=y\}} \mid \mathcal{F}_{i-1} \right] \mathbb{1}_{\{i \leq T_x\}} \right] \\ &= \mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} p(X_{i-1}, y) \mathbb{1}_{\{i \leq T_x\}} \right] \\ &= p(z, y) \mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{i \leq T_x\}} \right]. \end{aligned}$$

Thus,

$$\begin{aligned} \mu(y) &= \sum_{z \in E} p(z, y) \sum_{i=1}^{\infty} \mathbb{E}_x \left[\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{i \leq T_x\}} \right] \\ &= \sum_{z \in E} p(z, y) \mathbb{E}_x \left[\sum_{i=1}^{T_x} \mathbb{1}_{\{X_{i-1}=z\}} \right] \\ &= \sum_{z \in E} p(z, y) \mathbb{E}_x \left[\sum_{i=0}^{T_x-1} \mathbb{1}_{\{X_{i-1}=z\}} \right] \\ &= \sum_{z \in E} \mu(z) p(z, y). \end{aligned}$$

Note that, by construction, for all $y \in E$,

$$\mu(y) \leq \mathbb{E}_x \left[\sum_{i=0}^{\infty} \mathbb{1}_{\{X_i=y\}} \right] = \mathbb{E}_x[N_y] = G(x, y).$$

In particular, if $x \not\rightarrow y$, then $\mu(y) = 0$.

Let now y be a state such that $x \rightarrow y$ and $y \rightarrow x$. Thus, there exist integers n and m such that $P^n(x, y) > 0$ and $P^m(y, x) > 0$. Since $\mu P^n = \mu P^m = \mu$, we have in particular

$$\mu(y) = \sum_{z \in E} \mu(z) P^n(z, y) \geq \mu(x) P^n(x, y) > 0,$$

so that $\mu(y) > 0$, and

$$1 = \mu(x) = \sum_{z \in E} \mu(z) P^m(z, x) \geq \mu(y) P^m(y, x),$$

so that $\mu(y) < \infty$.

This concludes the proof that μ is an invariant measure with support equal to the communication class of x . \square

Exercise 4.22 *Does every Markov chain admit an invariant measure ?*

Exercise 4.23 *Prove that any invariant measure of an irreducible Markov chain gives a positive mass to every state. In other words, such a measure μ satisfies $\forall x \in E, \mu(x) > 0$.*

The second fundamental result is the following.

Theorem 4.33 *Assume that the Markov chain is irreducible and recurrent. Then any two invariant measures are proportional.*

Proof. Let μ be an invariant measure. Let $x \in E$ be such that $\mu(x) > 0$. Dividing μ by $\mu(x)$, we may and will assume that $\mu(x) = 1$. This normalisation being made, we want to prove that μ is equal to the measure ν defined by

$$\forall y \in E, \nu(y) = \mathbb{E}_x \left[\sum_{i=0}^{T_x-1} \mathbb{1}_{\{X_i=y\}} \right].$$

For this, we will prove that $\mu \geq \nu$, in the sense that

$$\forall y \in E, \mu(y) \geq \nu(y). \tag{13}$$

Let us take this inequality for granted for one moment and explain how it implies $\mu = \nu$. The point is that, up to the fact that it may be identically zero, $\mu - \nu$ is an invariant measure of our Markov chain, which gives mass 0 to the state x , and which therefore must vanish identically⁵. More precisely, consider $y \in E$ and an integer m such that $P^m(y, x) > 0$. We have

$$0 = (\mu - \nu)(x) = \sum_{z \in E} (\mu - \nu)(z) P^m(z, x) \geq (\mu - \nu)(y) P^m(y, x),$$

⁵This may be a good time to search Exercise 4.23 if you did not yet do so.

from which it follows that $\mu(y) = \nu(y)$.

Let us now prove (13). For this, we prove that for all $n \geq 0$ and all $y \in E$,

$$\mu(y) \geq \mathbb{E}_x \left[\sum_{i=0}^{(T_x-1) \wedge n} \mathbb{1}_{\{X_i=y\}} \right]. \quad (14)$$

If $y = x$, then both sides are equal to 1, for all $n \geq 0$, and the inequality holds.

Let us assume that $y \neq x$ and prove the result by induction on n . If $n = 0$, the right-hand side of (14) is equal to 0 and the inequality holds. Let us assume that the result is proved at rank n . Then, proceeding in a way that is very similar to the proof of Theorem 4.32, we have

$$\begin{aligned} \mathbb{E}_x \left[\sum_{i=0}^{(T_x-1) \wedge (n+1)} \mathbb{1}_{\{X_i=y\}} \right] &= \mathbb{E}_x \left[\sum_{i=1}^{T_x \wedge (n+1)} \mathbb{1}_{\{X_i=y\}} \right] \\ &= \sum_{z \in E} \sum_{i=1}^{n+1} \mathbb{E}_x [\mathbb{1}_{\{X_{i-1}=z\}} \mathbb{1}_{\{X_i=y\}} \mathbb{1}_{\{i \leq T_x\}}] \\ &= \sum_{z \in E} \mathbb{E}_x \left[\sum_{i=1}^{T_x \wedge (n+1)} \mathbb{1}_{\{X_{i-1}=z\}} \right] p(z, y) \\ &= \sum_{z \in E} \mathbb{E}_x \left[\sum_{i=0}^{(T_x-1) \wedge n} \mathbb{1}_{\{X_{i-1}=z\}} \right] p(z, y) \\ &\leq \sum_{z \in E} \mu(z) p(z, y) \\ &= \mu(y). \end{aligned}$$

Using the monotone convergence theorem to let n tend to infinity in (14), we obtain (13), and the proof is finished. \square

Exercise 4.24 Compare your understanding of the invariant measures of the p -biased random walk on \mathbb{Z} with the previous theorem.

We proved that an irreducible recurrent Markov chain admits, up to a multiplicative constant, a unique invariant measure. We will distinguish between the cases where these invariant measures are finite, and the case where they are infinite.

Proposition 4.34 Assume that the chain is irreducible and recurrent. Then exactly one of the following two situations occurs.

1. All invariant measures have infinite total mass and for all $x \in E$, we have $\mathbb{E}_x[T_x] = \infty$. In this case, the Markov chain is called null recurrent.
2. All invariant measures have finite total mass. There exists a unique invariant probability measure π . For all $x \in E$, we have

$$\pi(x) > 0, \quad \mathbb{E}_x[T_x] < \infty \quad \text{and} \quad \pi(x) = \frac{1}{\mathbb{E}_x[T_x]}.$$

In this case, the chain is called positive recurrent.

Proof. Let x be a state and let μ be the unique invariant measure on E such that $\mu(x) = 1$, which is given by Theorem 4.32. We have

$$\mu(E) = \sum_{y \in E} \mathbb{E}_x \left[\sum_{i=0}^{T_x-1} \mathbb{1}_{\{X_i=y\}} \right] = \mathbb{E}_x \left[\sum_{i=0}^{T_x-1} \sum_{y \in E} \mathbb{1}_{\{X_i=y\}} \right] = \mathbb{E}_x[T_x].$$

Since all invariant measures are proportional, they are either all finite, or all infinite. If they are all infinite, the previous computation shows that $\mathbb{E}_x[T_x] = \infty$ for all $x \in E$.

If they are all finite, the same computation shows that $\mathbb{E}_x[T_x]$ is finite for all $x \in E$. Moreover, the unique invariant probability measure is

$$\pi = \frac{1}{\mu(E)} \mu = \frac{1}{\mathbb{E}_x[T_x]} \mu,$$

from which it follows that

$$\pi(x) = \frac{1}{\mathbb{E}_x[T_x]},$$

and the proof is finished. \square

Proposition 4.35 *Assume that the chain is irreducible. If there exists an invariant probability measure, then the chain is recurrent.*

Proof. Let y be a state such that $\pi(y) > 0$. Let us compute $G(y, y)$. Remember that for every state x , we have $G(x, y) \leq G(y, y)$. Thus,

$$G(y, y) = \sum_{x \in E} \pi(x) G(y, y) \geq \sum_{x \in E} \pi(x) G(x, y).$$

Now,

$$\sum_{x \in E} \pi(x) G(x, y) = \sum_{x \in E} \sum_{n=0}^{\infty} \pi(x) P^n(x, y) = \sum_{n=0}^{\infty} (\pi P^n)(y) = \sum_{n=0}^{\infty} \pi(y) = \infty.$$

Thus, $G(y, y) = \infty$ and y is recurrent. Thus, since the chain is irreducible, all states are recurrent, and the chain is positive recurrent. \square

4.7 Brief summary

Let us summarise the results of the last two sections. Firstly, about recurrence, transience, and communication between states.

- The state space E of a Markov chain is partitioned into two disjoint subsets: the set of recurrent states and the set of transient states. Recall that x is recurrent, by definition, if and only if

$$\mathbb{P}_x(T_x < \infty) = 1.$$

- On E , there is the relation \rightarrow defined by

$$x \rightarrow y \Leftrightarrow G(x, y) > 0.$$

It is reflexive and transitive. There is also the relation \sim defined by

$$x \sim y \Leftrightarrow (x \rightarrow y \text{ and } y \rightarrow x).$$

It is an equivalence relation, the classes of which are called the communication classes. If E is itself a communication class, the chain is said to be irreducible.

- A communication class consists either exclusively of recurrent states, or exclusively of transient states. One speaks of recurrent and transient classes.
- Given two communication classes C and D , the existence of $x \in C$ and $y \in D$ such that $x \rightarrow y$ implies that for every $x \in C$ and every $y \in D$, one has $x \rightarrow y$. In this case, one writes $C \succ D$. This defines a binary relation on the set E/\sim of communication classes.
- The relation \succ on E/\sim is a (partial) order⁶ of which the recurrent classes are minimal elements. Let us emphasize that there can exist transient minimal classes. It is also possible that there be no minimal class at all.

Now about invariant measures.

- The support of any invariant measure is a union of communication classes. Moreover, if the support of an invariant measure contains a class C , then it also contains every class D such that $C \succ D$.

Let us now consider irreducible chains. There are three cases: transient, null recurrent and positive recurrent.

- Transient irreducible chains can have no invariant measure at all, a unique (up to multiplication) invariant measure, or several non-proportional invariant measures. In any case, any invariant measure of a transient irreducible chain has infinite total mass.
- Recurrent irreducible chains admit, up to multiplication, exactly one invariant measure.
- Null recurrent chains are the recurrent irreducible chains for which all invariant measures are infinite. The expected return time at every state is infinite⁷.
- Positive recurrent chains are the recurrent irreducible chains for which all invariant measures are finite, or equivalently for which there exists an invariant probability measure. If π is this invariant probability measure, then the relation

$$\pi(x)\mathbb{E}_x[T_x] = 1$$

holds for every state x .

⁶As the notation suggests, we think of C being ‘larger’ than D if $C \succ D$.

⁷It is however not true that the hitting time of every state starting from every other is infinite.

4.8 The ergodic theorem

In this section and the next, we will study the asymptotic behaviour of our Markov chain when time tends to infinity.

The main question is to determine, for all x and y , the behaviour as n tends to infinity of $P^n(x, y) = \mathbb{P}_x(X_n = y)$.

If y is transient, then $G(x, y) \leq G(y, y) < \infty$, so that

$$\lim_{n \rightarrow \infty} P^n(x, y) = 0.$$

With recurrent states, the situation is more interesting.

Theorem 4.36 *Consider a Markov chain that is irreducible and recurrent. Let μ be an invariant measure of this chain. Let $f : E \rightarrow \mathbb{R}^+$ and $g : E \rightarrow \mathbb{R}^+$ be two non-negative real-valued functions on E . Assume that $\int_E g \, d\mu > 0$ and that at least one of the functions f and g has finite integral with respect to μ . Then for all $x \in E$,*

$$\frac{\sum_{i=0}^n f(X_i)}{\sum_{i=0}^n g(X_i)} \xrightarrow{n \rightarrow \infty} \frac{\int_E f \, d\mu}{\int_E g \, d\mu} \quad \mathbb{P}_x - a.s.$$

Proof. Let $x \in E$ be a state. Let us introduce the sequence of stopping times

$$T_x^{(0)} = 0 \text{ and, for all } k \geq 1, T_x^{(k)} = \inf\{n > T_x^{(k-1)} : X_n = x\}.$$

Note that for all $k \geq 1$,

$$T_x^{(k)} = \widehat{T}_x(\theta_{T_x^{(k-1)}}(X)).$$

In particular, $T_x^{(1)} = T_x$ and the sequence $(T_x^{(k)})_{k \geq 1}$ is a sequence of independent and identically distributed random variables.

Since x is recurrent, we have

$$\mathbb{P}_x(\forall k \geq 0, T_x^{(k)} < \infty) = 1.$$

For all $k \geq 1$, let us define

$$\xi_k = \sum_{n=T_x^{(k-1)}}^{T_x^{(k)}-1} f(X_n),$$

the sum of the values of f on the states visited by the chain during its k -th excursion from x .

The Markov property implies that the sequence $(\xi_k)_{k \geq 1}$ of non-negative random variables is i.i.d. The strong law of large numbers asserts that

$$\frac{\xi_1 + \dots + \xi_k}{k} \xrightarrow[k \rightarrow \infty]{\mathbb{P}_x - a.s.} \mathbb{E}[\xi_1].$$

On the other hand,

$$\begin{aligned} \mathbb{E}[\xi_1] &= \mathbb{E} \left[\sum_{n=0}^{T_x-1} \sum_{y \in E} \mathbb{1}_{\{X_n=y\}} f(y) \right] \\ &= \sum_{y \in E} \mathbb{E} \left[\sum_{n=0}^{T_x-1} \mathbb{1}_{\{X_n=y\}} \right] f(y) \\ &= \int_E f \, d\nu, \end{aligned}$$

where ν is the unique invariant measure on E such that $\nu(x) = 1$. This measure is none other than $\frac{1}{\mu(x)}\mu$, so that

$$\mathbb{E}[\xi_1] = \frac{1}{\mu(x)} \int_E f d\mu.$$

For all $n \geq 0$, let us define

$$N_{x,n} = \sum_{i=0}^n \mathbb{1}_{\{X_i=x\}},$$

the number of visits to x before n , of which we think as the number of the excursion at x which is going on at time n . Since x is recurrent, we have

$$N_{x,n} \xrightarrow[n \rightarrow \infty]{\mathbb{P}_x - a.s.} +\infty.$$

By construction, we have, for all $n \geq 0$,

$$T_x^{(N_{x,n}-1)} \leq n < T_x^{(N_{x,n})} \quad \mathbb{P}_x - a.s.$$

Thus, for all $n \geq 0$, we have

$$\xi_1 + \dots + \xi_{N_{x,n}-1} \leq \sum_{i=0}^n f(X_i) \leq \xi_1 + \dots + \xi_{N_{x,n}}.$$

Dividing by $N_{x,n}$ and letting n tend to infinity, we find

$$\frac{1}{N_{x,n}} \sum_{i=0}^n f(X_i) \xrightarrow[n \rightarrow \infty]{} \frac{1}{\mu(x)} \int_E f d\mu \quad \mathbb{P}_x - a.s.$$

The same argument can be applied to g , and by dividing the two a.s. convergences, one obtains the expected result. \square

Corollary 4.37 *Suppose that the Markov chain is irreducible and recurrent.*

1. *Assume that the chain is positive recurrent. Let π denote the unique invariant probability measure. Then for all probability measure ν on E and all $y \in E$, we have*

$$\frac{1}{n} \sum_{i=0}^{n-1} \mathbb{1}_{\{X_i=y\}} \xrightarrow[n \rightarrow \infty]{} \pi(y) \quad \mathbb{P}_\nu - a.s.$$

In particular, for every non-negative function f on E ,

$$\frac{1}{n} \sum_{i=0}^{n-1} f(X_i) \xrightarrow[n \rightarrow \infty]{} \int_E f d\pi \quad \mathbb{P}_\nu - a.s.$$

2. *Assume that the chain is null recurrent. Then for all probability measure ν on E and all $y \in E$, we have*

$$\frac{1}{n} \sum_{i=0}^{n-1} \mathbb{1}_{\{X_i=y\}} \xrightarrow[n \rightarrow \infty]{} 0. \quad \mathbb{P}_\nu - a.s.$$

Proof. It suffices to apply the ergodic theorem to the function $g = 1$. □

This corollary explains the names *positive recurrent* and *null recurrent*: in the positive recurrent case, the chain spends a positive proportion of the time in each state, whereas in the null recurrent case, the asymptotic proportion of time spent in any given state is 0.

This corollary shows also, in the positive recurrent case, that for all $x, y \in E$, the sequence $(P^n(x, y))_{n \geq 0}$ converges in the sense of Cesàro to $\pi(y)$. We will now study when this convergence holds in the usual sense.

Consider, as an example, the transition kernel

$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

on the set $E = \{1, 2\}$. The invariant probability measure of this Markov chain is the uniform measure $\frac{1}{2}\delta_1 + \frac{1}{2}\delta_2$, but under \mathbb{P}_1 for instance, the distribution of the chain does not converge to the invariant distribution. Instead, it is alternatively equal to δ_1 and δ_2 . This example shows that a phenomenon of cyclicity, or periodicity, can prevent the convergence of a Markov chain to its invariant measure. We will study this phenomenon and prove that it is the only possible obstruction to the convergence in distribution of a Markov chain.

We will need a small amount of arithmetic. Firstly, we say that a subset I of \mathbb{N} is a *semigroup* if it contains 0 and is stable under addition:

$$\forall x, y \in I, x + y \in I.$$

We will need two properties of semigroups.

Lemma 4.38 *Let $I \subset \mathbb{N}$ be a semigroup. Let d be the g.c.d. of the elements of I .*

1. *The subgroup of \mathbb{Z} generated by I is the set*

$$I - I = \{x - y : x, y \in \mathbb{Z}\}.$$

2. *The equality $I - I = d\mathbb{Z}$ holds.*

3. *If $d = 1$, then there exists an integer n_0 such that I contains every integer larger than n_0 .*

Proof. 1. The subgroup of \mathbb{Z} generated by I certainly contains $I - I$. Our claim is thus equivalent to the fact that $I - I$ is a subgroup of \mathbb{Z} . Since $I - I$ contains 0, it suffices to check that it is stable by subtraction. But for all $n, m, n', m' \in I$, we have

$$(n - m) - (n' - m') = \underbrace{(n + m')}_{\in I} - \underbrace{(n' + m)}_{\in I} \in I - I.$$

2. There exists an integer e such that $I - I = e\mathbb{Z}$. On one hand, since 0 belongs to I , we have $I \subset I - I \subset e\mathbb{Z}$, so that e is a common divisor of I . On the other hand, any common divisor of I is also a common divisor of $I - I$, hence of e . Thus, $e = d$.

3. Let us assume that $d = 1$. Then by the first assertion, I contains two consecutive integers, say i and $i + 1$. Hence, I , being stable under addition, contains the set

$$\{ai + b(i + 1) : a, b \geq 0\}.$$

We claim that this set contains all integers larger than i^2 . Indeed, by reducing n modulo i^2 and then the remainder modulo i , write $n = qi^2 + ri + s$ with $q \geq 1$ and $r, s \in \{0, \dots, i-1\}$. Then $r - s > -i$ and the equality

$$n = (qi + (r - s))i + s(i + 1)$$

shows that n belongs to I . □

Definition 4.39 *Let $x \in E$ be a recurrent state. Define the set*

$$I_x = \{n \geq 0 : P^n(x, x) > 0\}.$$

The g.c.d. of this set is called the period of x and it is denoted by d_x .

The assumption that x is recurrent implies that I_x is an infinite set, because

$$\infty = G(x, x) = \sum_{n \geq 0} P^n(x, x) = \sum_{n \in I_x} P^n(x, x) \leq |I_x|.$$

In particular, $I_x \neq \{0\}$ and its g.c.d. is well defined.

Let us observe that I_x contains 0 and is stable by addition. Indeed, for all $n, m \geq 0$, we have

$$P^{n,m}(x, x) \geq P^n(x, x)P^m(x, x),$$

so that $n + m$ belongs to I_x as soon as n and m do. In particular, according to Lemma 4.38,

$$I_x - I_x = d_x \mathbb{Z}.$$

Proposition 4.40 *If the chain is irreducible and recurrent, then all states have the same period.*

Proof. Consider two states x and y . By irreducibility, there exists integers k and l such that $P^k(x, y) > 0$ and $P^l(y, x) > 0$. It follows that

$$k + I_y + l \subset I_x.$$

Thus, for all $n, m \in I_y$, we have

$$n - m = (k + n + l) - (k + m + l) \in I_x - I_x = d_x \mathbb{Z}.$$

This shows that d_x is a common divisor of all the elements of I_y , so that d_x divides d_y . By symmetry, d_y divides d_x , and $d_x = d_y$. □

Definition 4.41 *An irreducible recurrent chain is said to be aperiodic if the common period of all states is 1.*

Proposition 4.42 *Assume that the Markov chain is irreducible, recurrent and aperiodic. For all $x, y \in E$, there exists an integer n_0 such that for all $n \geq n_0$, $P^n(x, y) > 0$.*

Proof. Consider $x, y \in E$. Let $n_2 \geq 0$ be such that $P^{n_2}(x, y) > 0$. Such an integer n_2 exists because the chain is irreducible. Then, since $d_x = 1$, there exists n_1 such that I_x contains every integer larger than n_1 . Set $n_0 = n_1 + n_2$. For all $n \geq n_0$, write $n = m + n_2$ with $m \geq n_1$. Then

$$P^n(x, y) \geq P^m(x, x)P^{n_2}(x, y) > 0,$$

and the result is proved. \square

Theorem 4.43 *Assume that the chain is irreducible, positive recurrent, and aperiodic. Then, for all $x \in E$, we have*

$$\lim_{n \rightarrow \infty} \sum_{y \in E} |P^n(x, y) - \pi(y)| = 0,$$

where π is the unique invariant probability.

Proof. Let us define the Markovian kernel \bar{P} on $E \times E$ by

$$\bar{p}((x_1, x_2), (y_1, y_2)) = p(x_1, y_1)p(x_2, y_2).$$

Let $(\bar{\Omega}, \bar{\mathcal{F}}, (\bar{\mathcal{F}}_n)_{n \geq 0}, (\bar{\mathbb{P}}_{(x_1, x_2)})_{(x_1, x_2) \in E^2}, \bar{X} = (X_n^1, X_n^2)_{n \geq 0})$ be a Markov chain with transition kernel \bar{P} .

This chain \bar{X} is irreducible. Indeed, consider (x_1, x_2) and (y_1, y_2) in $E \times E$. By aperiodicity, there exists n_1 and n_2 such that for all $n \geq n_1$ (resp. $n \geq n_2$), we have $P^n(x_1, y_1) > 0$ (resp. $P^n(x_2, y_2) > 0$). For $n = \max(n_1, n_2)$, we have $\bar{P}^n((x_1, x_2), (y_1, y_2)) > 0$.

Moreover, the probability measure $\pi \otimes \pi$ is invariant for this chain. Indeed, for all $(y_1, y_2) \in E^2$,

$$\begin{aligned} \sum_{(x_1, x_2) \in E \times E} (\pi \otimes \pi)(x_1, x_2) \bar{P}((x_1, x_2), (y_1, y_2)) &= \sum_{(x_1, x_2) \in E \times E} \pi(x_1)p(x_1, y_1)\pi(x_2)p(x_2, y_2) \\ &= \pi(y_1)\pi(y_2) \\ &= (\pi \otimes \pi)(y_1, y_2). \end{aligned}$$

It follows from the fact that it admits an invariant probability measure that the chain is positive recurrent.

For all $x, y \in E$,

$$P^n(x, y) - \pi(y) = \bar{\mathbb{P}}_{\pi \otimes \delta_x}(X_n^2 = y) - \bar{\mathbb{P}}_{\pi \otimes \delta_x}(X_n^1 = y) = \bar{\mathbb{E}}_{\pi \otimes \delta_x}[\mathbb{1}_{\{X_n^2=y\}} - \mathbb{1}_{\{X_n^1=y\}}].$$

Let us consider the stopping time

$$T = \inf\{n \geq 0 : X_n^1 = X_n^2\}.$$

We have

$$\begin{aligned} P^n(x, y) - \pi(y) &= \bar{\mathbb{E}}_{\pi \otimes \delta_x}[\mathbb{1}_{\{T > n\}}(\mathbb{1}_{\{X_n^2=y\}} - \mathbb{1}_{\{X_n^1=y\}})] \\ &\quad + \sum_{k=0}^n \sum_{z \in E} \bar{\mathbb{E}}_{\pi \otimes \delta_x}[\mathbb{1}_{\{T=k, X_k^1=X_k^2=z\}}(\mathbb{1}_{\{X_n^2=y\}} - \mathbb{1}_{\{X_n^1=y\}})]. \end{aligned}$$

For $k \in \{0, \dots, n-1\}$ and $z \in E$, we have

$$\begin{aligned} \bar{\mathbb{E}}_{\pi \otimes \delta_x} [\mathbb{1}_{\{T=k, X_k^1=X_k^2=z\}} \mathbb{1}_{\{X_n^2=y\}}] &= \bar{\mathbb{E}}_{\pi \otimes \delta_x} [\mathbb{1}_{\{T=k, X_k^1=X_k^2=z\}}] P^{n-k}(z, y) \\ &= \bar{\mathbb{E}}_{\pi \otimes \delta_x} [\mathbb{1}_{\{T=k, X_k^1=X_k^2=z\}} \mathbb{1}_{\{X_n^1=y\}}], \end{aligned}$$

so that the double sum of the last expression vanishes. Thus,

$$\sum_{y \in E} |P^n(x, y) - \pi(y)| \leq 2 \bar{\mathbb{P}}_{\pi \otimes \delta_x}(T > n).$$

Since the chain \bar{X} is recurrent, T is finite $\bar{\mathbb{P}}_{\pi \otimes \delta_x}$ -almost surely, and the last quantity tends to 0 as n tends to infinity. \square

References

- [1] Patrick Billingsley. *Probability and Measure*. Wiley, 3rd edition, 1995.
- [2] Joe Diestel. Uniform integrability: an introduction. *Rend. Istit. Mat. Univ. Trieste*, 23(1):41–80 (1993), 1991. School on Measure Theory and Real Analysis (Grado, 1991).
- [3] Rick Durrett. *Probability: theory and examples*, volume 31 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, fourth edition, 2010.
- [4] James R. Norris. *Markov chains*, volume 2 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998. Reprint of 1997 original.
- [5] David Williams. *Probability with martingales*. Cambridge University Press, 1991.