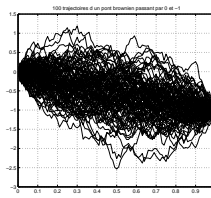


Randal DOUC

MONTE CARLO AND ADVANCED NUMERICAL SIMULATION METHODS



Contents

| | | |
|----------|---|----|
| 1 | Preliminaries | 3 |
| 2 | Refresh on Measure theory and Lebesgue integration | 5 |
| 2.1 | Sigma-fields, Measures and Probability | 5 |
| 2.2 | Integrals, random variables, expectation | 7 |
| 2.3 | Convergence | 9 |
| 2.4 | Some usual distributions | 10 |
| 3 | Introduction to the Monte Carlo methods | 11 |
| 3.1 | Principle of the method | 11 |
| 3.2 | On the use of the CLT for Monte Carlo methods | 14 |
| 3.3 | Take-home message | 15 |
| 3.4 | Highlights | 15 |
| 4 | Exact or approximate sampling | 17 |
| 4.1 | Exact Sampling | 17 |
| 4.2 | Approximate sampling | 23 |
| 4.3 | Take-home message | 25 |
| 5 | Metropolis-Hastings algorithms | 27 |
| 5.1 | Main notation | 27 |
| 5.2 | Definitions | 27 |
| 5.3 | Canonical space | 30 |
| 5.4 | At this point... .. | 32 |
| 5.5 | Metropolis-Hastings algorithms | 32 |
| 5.6 | Invariant probability measure: uniqueness | 35 |
| 5.7 | After studying this chapter... .. | 37 |
| 6 | Variance reduction techniques | 39 |
| 6.1 | Importance Sampling | 39 |
| 6.2 | Antithetic variates | 40 |
| 6.3 | Control Variates | 41 |
| 6.4 | Conditioning | 43 |
| 6.5 | Stratified sampling | 44 |
| 6.6 | Quasi Monte Carlo methods | 45 |
| 6.7 | Take-home message | 47 |
| 6.8 | Highlights | 47 |
| 7 | Exercises | 49 |

Plan of action

1. Day 1: 4H30
 - a. Lecture 1H30: Recap on Measures, Integration, Random Variables, independence.
 - b. Tutorials 1H30: Exercises.
 - c. Computer Session 1H30: Histograms and sampling of a random variable.
2. Day 2: 4H30
 - a. Lecture 1H30: LLN, Central Limit Theorem, Confidence Intervals.
 - b. Tutorials 1H30: Exact sampling, the quantile function. Rejection sampling: exercise.
 - c. Computer Session 1H30: LLN, CLT and confidence intervals. Exact sampling.
3. Day 3: 3H00
 - a. Lecture 1H30: Approximate Sampling, Importance Sampling. Monte Carlo by Markov chains.
 - b. Tutorials 1H30: Approximate Sampling, Markov chains, exercise.
4. Day 4: 4H30
 - a. Lecture 1H30: Variance reduction (I): antithetic, control variates.
 - b. Tutorials 1H30: Markov chains and antithetic variables.
 - c. Computer Session 1H30: Importance sampling and MCMC.
5. Day 5: 4H30
 - a. Lecture 1H30: Variance reduction (II): conditioning, stratified sampling
 - b. Tutorials 1H30: Exercises on variance reduction.
 - c. Computer Session 1H30: Variance reduction (II)

Chapter 1

Preliminaries

These lecture notes are built from various lecture notes, those of Bernard Lapeyre and Denis Talay, Eric Moulines and Gersende Fort, Emmanuel Temam, Jérôme Lelong, Bruno Bouchard, etc., all mixed up with various personal touches. It has been actively proofread by Cyrille Duguay.

Do not hesitate to point out the errors or typos that still remain and to propose any improvements on the content of this course.

Chapter 2

Refresh on Measure theory and Lebesgue integration

Contents

| | | |
|------------|---|-----------|
| 2.1 | Sigma-fields, Measures and Probability | 5 |
| 2.2 | Integrals, random variables, expectation | 7 |
| 2.3 | Convergence | 9 |
| 2.4 | Some usual distributions | 10 |

Keywords: *Sigma-fields, measures, probability measures, measurable functions, random variables, expectations. Weak convergence, almost sure convergence, Law of Large numbers, central limit theorems, confidence intervals.*

In this chapter, we provide a short recap on measure and integration theory with a special focus on their application in probability theory.

2.1 Sigma-fields, Measures and Probability

Let us start with the definition of sigma-fields...

Definition 2.1. Let Ω be a given set. We say that a family of sets $\mathcal{F} \subset \mathcal{P}(\Omega)$ is a sigma-field on Ω if and only if the three following properties are satisfied

- (i) $\Omega \in \mathcal{F}$
- (ii) if $A \in \mathcal{F}$ then $\bar{A} = \Omega \setminus A \in \mathcal{F}$.
- (iii) if for all $i \in \mathbb{N}$, $A_i \in \mathcal{F}$ then $\bigcap_{i \in \mathbb{N}} A_i \in \mathcal{F}$

We then say that (Ω, \mathcal{F}) is a measurable space.

A sigma-field is stable by complementary sets, countable intersection, countable union and also by taking the “set difference” \setminus in the sense that if $A, B \in \mathcal{F}$, then $A \setminus B \in \mathcal{F}$ (indeed, just write $A \setminus B = A \cap B^c$.)

► **Q 1:** Do those properties have special names?

The second property is often called *stability by complementary sets* and the last one *stability by countable intersection*. You may also find in the literature some other *equivalent* definitions:

- (i) $\emptyset \in \mathcal{F}$
- (ii) if $A \in \mathcal{F}$ then $\bar{A} = \Omega \setminus A \in \mathcal{F}$.
- (iii) if for all $i \in \mathbb{N}$, $A_i \in \mathcal{F}$ then $\bigcup_{i \in \mathbb{N}} A_i \in \mathcal{F}$

But I prefer the way it is expressed in Theorem 2.1.

► **Q 2:** Why do you need these properties?

In the theory of probability, a set A will correspond typically correspond to an event that may occur, it can be expressed as a constraint with respect to all the possibilities. The fact that we ask the stability by complementary sets or by countable intersections corresponds to considering either the complementary event or the fact that all the events A_i are satisfied.

► **Q 3:** Do you have any examples?

The smallest sigma-field is $\mathcal{F} = \{\Omega, \emptyset\}$ and the largest one is $\mathcal{P}(\Omega)$. Often, sigma-fields we are interested in are generated by some family of sets...

► **Q 4:** What do you mean ?

Say that you are interested in a family of sets $\mathcal{C} \subset \mathcal{P}(\Omega)$ but unfortunately, some of the 3 properties that define sigma-fields are not satisfied for \mathcal{C} . In that case, we can still include \mathcal{C} into a larger family so that it is a sigma-field. We can even consider the “smallest” one (in a sense to be defined), which contains all the sets in \mathcal{C} .

Definition 2.2. Let $\mathcal{C} \subset \mathcal{P}(\Omega)$. There exists a sigma-field, named $\sigma(\mathcal{C})$ which contains \mathcal{C} and which is minimal for the inclusion, that is, any other sigma-field that contains \mathcal{C} also contains $\sigma(\mathcal{C})$. We then say that $\sigma(\mathcal{C})$ is the sigma-field generated by \mathcal{C} .

A good exercise is to prove the property that appears in the above definition. This can be done by defining $\sigma(\mathcal{C})$ as the collection of all the sets which are common wrt to all the sigma fields that contain \mathcal{C} . In other words, define

$$\mathcal{A} = \cap \{ \mathcal{T} : \mathcal{C} \subset \mathcal{T} \text{ and } \mathcal{T} \text{ is a sigma-field} \}$$

Even though \mathcal{A} is defined with an uncountable intersection, you may check that \mathcal{A} is sigma-field, that this sigma-field contains all the sets in \mathcal{C} and that any other sigma-field that contains all the sets in \mathcal{C} necessarily contains all the sets in \mathcal{A} . All in all, \mathcal{A} is the smallest sigma-field satisfying this property and we can call it $\sigma(\mathcal{C})$.

► **Q 5:** What is your favorite example?

An important example is the case of open sets. If Ω is \mathbb{R}^k and \mathcal{C} is the family of open sets on Ω , then the sigma-field generated by open sets is called the Borel sigma-field and is noted $\mathcal{B}(\Omega)$.

► **Q 6:** For different family of sets, if you consider the sigma-fields generated by each of them, do you systematically find different sigma-fields?

Of course not... In practice, if we have two family of sets $\mathcal{C} \subset \mathcal{P}(\Omega)$ and $\mathcal{D} \subset \mathcal{P}(\Omega)$ and if we want to check if $\sigma(\mathcal{C}) = \sigma(\mathcal{D})$ then a necessary and sufficient condition for getting that is to check successively that $\mathcal{D} \subset \sigma(\mathcal{C})$ and $\mathcal{C} \subset \sigma(\mathcal{D})$. You can use this property for checking that the sigma-field generated by open sets (i.e. the Borel sigma-field) is also the sigma-field generated by closed sets.

► **Q 7:** That's nice. You can prove it easily?

Yes, please do so. It's a good way to check that everything is understandable.

Definition 2.3. Let (Ω, \mathcal{F}) be a measurable space. We say that a function $\mu : \mathcal{F} \rightarrow \bar{\mathbb{R}}^+ := \mathbb{R}^+ \cup \{\infty\}$ is a measure if it satisfies the *sigma-additivity property*, that is for any family of sets (A_i) such that $A_i \in \mathcal{F}$ for any $i \in \mathbb{N}$ and $A_i \cap A_j = \emptyset$ for all $i \neq j$, then

$$\mu \left(\bigcup_{i=0}^{\infty} A_i \right) = \sum_{i=0}^{\infty} \mu(A_i) \quad (2.1)$$

We then say that $(\Omega, \mathcal{F}, \mu)$ is a measured space. Moreover, if $\mu(\Omega) = 1$ then we say that μ is a probability measure.

Note that in the right hand side of (2.1), we sum quantities in $\bar{\mathbb{R}}^+$ that is, we use the convention that if $a \in \mathbb{R}^+$, $a + \infty = \infty$ and $\infty + \infty = \infty$.

► **Q 8:** The measure μ evaluated on a set A can be infinite?

Yes of course, $\mu(A)$ takes its values between 0 and $\mu(\Omega)$ actually... But if you consider a measure of probability (which is nothing but a particular measure), then the values of $\mu(A)$ are between 0 and 1.

► **Q 9:** Some useful properties?

- (i) If $A \subset B$, then $\mu(A) \leq \mu(B)$.
- (ii) $\mu(\cup_{i=1}^{\infty} A_i) \leq \sum_{i=1}^{\infty} \mu(A_i)$.
- (iii) $\mu(\cup_{i=1}^{\infty} A_i) = \lim_{n \rightarrow \infty} \mu(\cup_{i=1}^n A_i)$
- (iv) if $\mu(A_1) < \infty$, then, $\mu(\cap_{i=1}^{\infty} A_i) = \lim_{n \rightarrow \infty} \mu(\cap_{i=1}^n A_i)$

► **Q 10:** What are the typical measures I will deal with?

There are at least two fundamental examples of measures:

- (i) The Dirac measure on a which is defined by $\delta_a(A) = 0$ if $a \notin A$ and 1 otherwise.
- (ii) The Lebesgue measure λ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ which is defined with the following characterizing property: it is the only measure such that for any segment $A = [a, b]$, we have $\lambda(A) = b - a$. Then, you can show quite easily that the Lebesgue measure of any interval (the interval may be closed, open or none of them) is the length of the interval (which thus may be infinite, take $A = [1, \infty[$ for example).

From these measures, you can construct other measures, for example by multiplying them by some non negative measurable functions.

2.2 Integrals, random variables, expectation

► **Q 11:** Measurable functions?

Yes. The definition is below.

Definition 2.4. If (A, \mathcal{A}) and (B, \mathcal{B}) are measurable sets. We say that $h : A \rightarrow B$ is a $\mathcal{A} / \mathcal{B}$ measurable function if and only if for any $B \in \mathcal{B}$, $f^{-1}(B) \in \mathcal{A}$.

Of course, if $\mathcal{B} = \sigma(\mathcal{C})$, then instead of checking for any $B \in \mathcal{B}$, $f^{-1}(B) \in \mathcal{A}$, we may only check for any $C \in \mathcal{C}$, $f^{-1}(C) \in \mathcal{A}$. For a given measurable function f , we may define $\sigma\{f^{-1}(B) : B \in \mathcal{B}\}$, it turns out that it is a sigma field, called $\sigma(f)$. Measurable functions are linked with random variables...

► **Q 12:** Can you be more precise?

Here is the definition of a random variable.

Definition 2.5. Let (Ω, \mathcal{F}) be a measurable space and consider the measurable space $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. A random variable $X : \Omega \rightarrow \mathbb{R}$ is, by definition, a $\mathcal{F} / \mathcal{B}(\mathbb{R})$ measurable function.

► **Q 13:** You mean that a random variable is nothing more than a measurable function?

Yes, and it is always real valued. If X takes its values in \mathbb{R}^k , we don't speak about random variables but about random vectors and/or (in some books) about random elements. Now that we have defined general measures with the important particular case of the probability measures, we can define the *integral* associated to a measure μ and if μ is probability measure, we will defined the *expectation* of this random variable.

Recall that a measure μ associates to a set $A \in \mathcal{F}$, a real number $\mu(A)$. Now, we want μ to associate to any real-valued $\mathcal{F} / \mathcal{B}(\mathbb{R})$ -measurable function f a real number called $\mu(f)$ or $\int_{\Omega} f(w)\mu(dw)$. There is here an abuse of notation... Stricto sensu: $\mu(A)$ is well defined but $\mu(f)$ is an abuse of notation because between brackets, we write a function f and not a set A . To avoid confusions, most of the time we write $\int f(w)\mu(dw)$ or $\int f d\mu$ (to avoid w) instead of $\mu(f)$ but all these notation of course refer to the same object.

Actually, it will be not be possible to define $\mu(f)$ for any measurable function... Let us be more precise. The construction of the integral wrt μ will be done progressively. We start with the $\mu(\mathbb{1}_A)$. By definition, we set:

$$\mu(\mathbb{1}_A) = \mu(A)$$

Then, we will define

$$\mu\left(\sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}\right) = \sum_{i=1}^n \alpha_i \mu(A_i)$$

Then, we define for any measurable *non-negative* function f ,

$$\mu(f) = \sup \left\{ \mu\left(\sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}\right) : \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i} \leq f \right\}$$

Then, for any measurable function such that $\mu(|f|) < \infty$, we set

$$\mu(f) = \mu(f^+) - \mu(f^-)$$

► **Q 14:** You always integrate on the whole space?

Yes, but if you integrate a function f on a subset $\Omega_0 \in \mathcal{F}$ where $\Omega_0 \subset \Omega$, then by definition, it just means $\int f(w) \mathbb{1}_{\Omega_0}(w) \mu(dw)$. That is you integrate on the whole space but thanks to the indicator function $\mathbb{1}_{\Omega_0}$, only the values of f on Ω_0 are meaningful.

► **Q 15:** You told me that two examples of measures were important.
So what?

► **Q 16:** What are the integrals associated to those measures?

Our two important examples of integrals constructed from measures μ are

- (i) Integrals associated to Dirac measures... We can show that for any $a \in \Omega$ and any measurable function f , $\int f(w) \delta_a(dw) = f(a)$.
- (ii) Integrals associated to the Lebesgue measure... This is a common case, and instead of writing $\int f(w) \lambda(dw)$, we usually write $\int f(w) dw$.

► **Q 17:** OK for the construction of the integral but what are the essential properties?

I guess the very essential ones allow to interchange limit and integral. Two of them are essential:

- (i) The monotone convergence theorem: if $\{f_n, n \in \mathbb{N}\}$ is a family of measurable *non-negative* functions and $f_n \leq f_{n+1}$ for all large enough n , then $\int \lim_{n \rightarrow \infty} f_n(w) \mu(dw) = \lim_{n \rightarrow \infty} \int f_n(w) \mu(dw)$
- (ii) The Lebesgue dominated convergence theorem: if $\{f_n : n \in \mathbb{N}\}$ is a family of measurable functions such that $\lim_{n \rightarrow \infty} f_n(w)$ exists for μ -almost all $w \in \Omega$ and if $|f_n| \leq h$ where $\int h d\mu < \infty$, then $\int \lim_{n \rightarrow \infty} f_n(w) \mu(dw) = \lim_{n \rightarrow \infty} \int f_n(w) \mu(dw)$

But there are also essential properties: the linearity of the integral, or if $f \leq g$ then, $\mu(f) \leq \mu(g)$, or $|\mu(f)| \leq \mu(|f|)$. Or if $\mu(|f|) < \infty$ then, $|f(w)| < \infty$ for μ -almost all $w \in \Omega$.

► **Q 18:** You said that when you multiply a non-negative measurable function and a measure, it is a measure...
What do you mean exactly?

If f is a non-negative measurable function, then, the measure $f d\mu$ is defined by: $A \mapsto \int \mathbb{1}_A(w) f(w) \mu(dw)$. Therefore, with the two typical measures (δ_a and λ), you can define so many different measures $f d\delta_0$ or $f d\lambda$ by using a non-negative measurable function f .

► **Q 19:** You told me there is some link between the expectation operator associated to a probability measure and the integral associated to a measure.

The expectation is defined in the same way as integrals associated to some measure: it is constructed from a measured space $(\Omega, \mathcal{F}, \mathbb{P})$ and a random variable X (that is a $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable function). Then, by definition, $\mathbb{E}[X]$ is just the integral associated to the measure \mathbb{P} taken at the random variable X , that is, $\mathbb{E}[X] = \int_{\Omega} X(w) \mathbb{P}(dw)$.

Definition 2.6. The law of a random variable X on $(\Omega, \mathcal{F}, \mathbb{P})$ is the measure μ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ defined by: $\mu : A \mapsto \mathbb{P}(X \in A)$.

In this definition, we use $\mathbb{P}(X \in A)$ which means $\mathbb{P}(\Omega_0)$ where $\Omega_0 = \{w \in \Omega : X(w) \in A\} := \{X \in A\}$. We also call this measure, the push-forward measure of \mathbb{P} through the function X . It defines a measure on the arrival sigma field.

► **Q 20:** How do you check that two random variables are independent?

Actually, two random vectors X, Y defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$ but taking values in \mathbb{R}^p and \mathbb{R}^d are independent if and only one the equivalent properties are satisfied:

(i) For all $(A, B) \in \mathcal{B}(\mathbb{R}^p) \times \mathcal{B}(\mathbb{R}^d)$,

$$\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A)\mathbb{P}(Y \in B)$$

(ii) For all bounded or non-negative $\mathcal{B}(\mathbb{R}^p)/\mathcal{B}(\mathbb{R})$ -measurable functions $f : \mathbb{R}^p \rightarrow \mathbb{R}$ and for all $\mathcal{B}(\mathbb{R}^d)/\mathcal{B}(\mathbb{R})$ -measurable bounded or non-negative functions $g : \mathbb{R}^d \rightarrow \mathbb{R}$, we have

$$\mathbb{E}[f(X)g(Y)] = \mathbb{E}[f(X)]\mathbb{E}[g(Y)]$$

(iii) for all $(u, v) \in \mathbb{R}^p \times \mathbb{R}^d$,

$$\mathbb{E}\left[e^{iu^T X + iv^T Y}\right] = \mathbb{E}\left[e^{iu^T X}\right]\mathbb{E}\left[e^{iv^T Y}\right]$$

(iv) for all $(u, v) \in \mathbb{R}^p \times \mathbb{R}^d$ and all $(x, y) \in \mathbb{R}^2$,

$$\mathbb{P}(u^T X \leq x, v^T Y \leq y) = \mathbb{P}(u^T X \leq x)\mathbb{P}(v^T Y \leq y)$$

We have $X \stackrel{\mathcal{L}}{\equiv} Y$ (where X and Y are two random variables) if and only if any of the following conditions holds true:

- (i) $\mathbb{P}(X \in A) = \mathbb{P}(Y \in A)$ for any $A \in \mathcal{B}(\mathbb{R})$
- (ii) $\mathbb{P}(X \leq t) = \mathbb{P}(Y \leq t)$ for any $t \in \mathbb{R}$. Note that $t \mapsto \mathbb{P}(X \leq t)$ is the cumulative distribution function for the random variable X .
- (iii) $\mathbb{E}[e^{iuX}] = \mathbb{E}[e^{iuY}]$ for any $u \in \mathbb{R}$. Note that $u \mapsto \mathbb{E}[e^{iuX}]$ is the characteristic function for the random variable X .

2.3 Convergence

We start with some notation. In what follows,

- i.i.d means independent and identically distributed.
- r.v. means random variables.
- for $r, s \in \mathbb{N}$ such that $r \leq s$, we write $[r : s] = \{r, r+1, \dots, s\}$,
- $X \perp\!\!\!\perp Y$ means X and Y are independent random variables,
- $X \stackrel{\mathcal{L}}{\equiv} Y$ means X and Y have the same law.
- $\liminf_n a_n = \lim_{n \rightarrow \infty} (\inf_{k \geq n} a_k)$ and similarly, $\limsup_n a_n = \lim_{n \rightarrow \infty} (\sup_{k \geq n} a_k)$. Moreover, $\lim_n a_n$ exists if and only if $\liminf_n a_n = \limsup_n a_n$.
- for any $a \in \mathbb{R}$, $a^+ = \max(a, 0)$ and $a^- = \max(-a, 0) = -\min(a, 0)$ and we have $|a| = a^+ + a^-$ and $a = a^+ - a^-$.

Moreover, the following notions of convergence for random variables is used throughout these lecture notes.

► $X_n \xrightarrow{w} X$ means *convergence in distribution* (or "convergence en loi" in French). It is equivalent to any of the following statements.

- (a) for all bounded continuous functions h , we have $\lim_n \mathbb{E}[h(X_n)] = \mathbb{E}[h(X)]$.
- (b) for all $A \in \mathcal{B}(\mathbb{R})$ such that $\mathbb{P}(X \in \partial A) = 0$, we have $\lim_n \mathbb{P}(X_n \in A) = \mathbb{P}(X \in A)$.

- (c) for all $x \in \mathbb{R}$ such that $\mathbb{P}(X = x) = 0$, we have $\lim_n \mathbb{P}(X_n \leq x) = \mathbb{P}(X \leq x)$.
- (d) for all $u \in \mathbb{R}$, we have $\lim_n \mathbb{E}[e^{iuX_n}] = \mathbb{E}[e^{iuX}]$

By abuse of terminology, we may also say that X_n **weakly converges to X** instead of saying the distribution of X_n converges weakly to the distribution of X .

- $X_n \xrightarrow{\mathbb{P}\text{-prob}} X$ means *convergence in probability*: for all $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0.$$

- $X_n \xrightarrow{\mathbb{P}\text{-a.s.}} X$ means *almost sure convergence*:

$$\mathbb{P}(\lim_{n \rightarrow \infty} X_n = X) = 1.$$

Almost sure convergence implies convergence in probability. We also recall the following properties.

- (i) If $X_n \xrightarrow{w} X$ then for all continuous functions f , $f(X_n) \xrightarrow{w} f(X)$. Note that this property holds, when f is continuous (and not necessarily bounded), for example $f(u) = u^2$ so that $X_n^2 \xrightarrow{w} X^2$.
- (ii) **The Slutsky Lemma** If $X_n \xrightarrow{\mathbb{P}\text{-prob}} c$ where c is a constant and if $Y_n \xrightarrow{w} Y$, then $(X_n, Y_n) \xrightarrow{w} (c, Y)$ that is for all continuous functions f , $f(X_n, Y_n) \xrightarrow{w} f(c, Y)$.
- (iii) $X \sim N(0, 1)$ iff for all $u \geq 0$, $\mathbb{E}[e^{iuX}] = e^{-u^2/2}$. Moreover, $X \sim N(\mu, \sigma^2)$ iff for all $u \geq 0$, $\mathbb{E}[e^{iuX}] = e^{-u^2 \text{Var}(X)/2 + iu\mathbb{E}(X)}$ and in that case, $\sigma^2 = \text{Var}(X)$ and $\mu = \mathbb{E}(X)$.

2.4 Some usual distributions

| Name | Acronym | Parameter | density function: $f_X(x)$ | cdf: $F_X(x) = \int_{-\infty}^x f_X(u) du$ | Other properties |
|-------------|---------------------|--------------------------------------|--|---|--|
| Gaussian | $N(\mu, \sigma^2)$ | (μ, σ^2) | $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)}$ | No explicit expression | $\mathbb{E}[X] = \mu, \text{Var}X = \sigma^2$ |
| Exponential | $\exp(\lambda)$ | $\lambda > 0$ | $\lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}_+}(x)$ | $(1 - e^{-\lambda x}) \mathbb{1}_{\mathbb{R}_+}(x)$ | $\mathbb{E}[X] = 1/\lambda, \text{Var}X = 1/\lambda^2$ |
| Gamma | $\Gamma(k, \theta)$ | $(k, \theta) \in (\mathbb{R}_+^*)^2$ | $\frac{x^{k-1} e^{-x/\theta}}{\Gamma(k)\theta^k}$ | $\frac{\Gamma_x/\theta(k)}{\Gamma(k)}$ | |

In the above description,

- (i) if $X_i \sim \Gamma(k_i, \theta)$ and (X_i) are independent, then $\sum_{i=1}^n X_i \sim \Gamma(\sum_{i=1}^n k_i, \theta)$.
- (ii)

$$\Gamma(k) = \begin{cases} \int_0^\infty t^{k-1} e^{-t} dt & \text{if } k \in \mathbb{R}_+^* \\ k! & \text{if } k \in \mathbb{N}. \end{cases} \quad (\blacktriangleright \text{GAMMA FUNCTION})$$

$$\Gamma_x(k) = \int_0^x t^{k-1} e^{-t} dt \quad (\blacktriangleright \text{INCOMPLETE GAMMA FUNCTION})$$

Chapter 3

Introduction to the Monte Carlo methods

Contents

| | | |
|-----|---|----|
| 3.1 | Principle of the method | 11 |
| 3.2 | On the use of the CLT for Monte Carlo methods | 14 |
| 3.3 | Take-home message | 15 |
| 3.4 | Highlights | 15 |

Keywords: LLN, CLT, Slutski's lemma, unbiased estimator, Delta-method, confidence intervals.

Many probabilistic issues that arise in statistical applications boil down to the calculation of expectations. For example, in Bayesian inference for hidden models, we might be interested in the expectation of some function of the parameter wrt to a posteriori distribution. In that context, it's quite natural to focus on numerical methods, which allow to calculate these expectations, considering that in most cases closed-form formulae are not available.

3.1 Principle of the method

Monte Carlo methods are used to numerically calculate expectations. These methods are based on the celebrated Strong Law of the Large Numbers.

Theorem 3.1. (► STRONG LAW OF LARGE NUMBERS (LLN)). *Let $(X_n)_n$ be i.i.d random variables with the same law as X . If $\mathbb{E}[|X|] < \infty$, then*

$$\bar{X}_N = N^{-1} \sum_{i=1}^N X_i \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}[X]$$

The convergence also holds in L^1 : $\mathbb{E}[|\bar{X}_N - \mathbb{E}(X)|] \rightarrow 0$.

Remark 3.2 *The integrability assumption is mandatory, indeed we can show that if you consider an i.i.d. sequence of r.v. according to a Cauchy distribution, then the empirical mean does not converge. Please write an R program for illustrating this property.*

PROOF. We will prove only the a.s. convergence. We start with an elementary result:

Lemma 3.3 **A preliminary result:** *Let (Y_i) be iid random variables such that $\mathbb{E}[|Y_1|] < \infty$ and $\mathbb{E}[Y_1] > 0$, then a.s.,*

$$\liminf_n S_n/n \geq 0$$

where $S_n = \sum_{i=1}^n Y_i$.

► (Proof of the lemma) Set $L_n = \inf(S_k, k \in [1 : n])$, $L_\infty = \inf(S_k, k \in \mathbb{N}^*)$, $A = \{L_\infty = -\infty\}$. Let $\theta(y_1, y_2, \dots) = (y_2, y_3, \dots)$ be the shift operator. Then, a.s.,

$$\begin{aligned} L_n &= S_1 + \inf(0, S_2 - S_1, \dots, S_n - S_1) = Y_1 + \inf(0, L_{n-1} \circ \theta) \\ &\geq Y_1 + \inf(0, L_n \circ \theta) = Y_1 - L_n^- \circ \theta. \end{aligned}$$

where the inequality follows from the fact that $n \mapsto L_n$ is nonincreasing. This implies a.s. (since $L_n^- \circ \theta$ is a.s. finite)

$$\mathbf{1}_A Y_1 \leq \mathbf{1}_A L_n + \mathbf{1}_A L_n^- \circ \theta$$

Taking the expectation on both sides and then, using $\mathbb{P}(\mathbf{1}_A = \mathbf{1}_A \circ \theta) = 1$, and the strong stationarity of the sequence:

$$\mathbb{E}[\mathbf{1}_A Y_1] \leq \mathbb{E}[\mathbf{1}_A L_n] + \mathbb{E}[\mathbf{1}_A \circ \theta L_n^- \circ \theta] = \mathbb{E}[\mathbf{1}_A L_n] + \mathbb{E}[\mathbf{1}_A L_n^-] = \mathbb{E}[\mathbf{1}_A L_n^+] \rightarrow 0$$

where the right-hand side tends to 0 by the dominated convergence theorem since a.s. $\lim_n \mathbf{1}_A L_n^+ = \mathbf{1}_A L_\infty^+ = 0$ and $0 \leq \mathbf{1}_A L_n^+ \leq Y_1^+$. Finally $\mathbb{E}[\mathbf{1}_A Y_1] \leq 0$. Therefore, noting that $\mathbf{1}_A \circ \theta$ is independent from Y_1 ,

$$0 \geq \mathbb{E}[\mathbf{1}_A Y_1] = \mathbb{E}[\mathbf{1}_A \circ \theta Y_1] = \mathbb{E}[\mathbf{1}_A \circ \theta] \mathbb{E}[Y_1] = \underbrace{\mathbb{E}[\mathbf{1}_A]}_{\geq 0} \underbrace{\mathbb{E}[Y_1]}_{> 0}$$

This implies $\mathbb{P}(A) = 0$ and the lemma is proved. ◀

(Proof of the Theorem.) We now turn to the proof of the LLN. Without loss of generality, we assume that $\mathbb{E}[X_1] = 0$. Applying the lemma with $Y_i = X_i + \varepsilon$ (where $\varepsilon > 0$), we get $\liminf_n n^{-1} \sum_{i=1}^n X_i \geq -\varepsilon$ a.s. And applying again the lemma with $Y_i = -X_i + \varepsilon$, we get a.s., $\limsup_n n^{-1} \sum_{i=1}^n X_i \leq \varepsilon$ which finishes the proof since ε is arbitrary. ■

Since an estimator of $\mathbb{E}[X]$ is a random variable, it is important to provide confidence intervals. Indeed, two particular samplings may lead to very different estimates. The following theorem gives a result on the convergence speed of the estimator and then allows to provide confidence intervals. In practice, we should always balance the obtained estimator with the width of the confidence interval in order to justify the relevance of the announced result.

Theorem 3.4. (► CENTRAL LIMIT THEOREM (CLT)) *Let $(X_n)_n$ be a sequence of i.i.d random variables with the same law as X . If $\mathbb{E}[X^2] < \infty$, then*

$$\frac{\bar{X}_N - \mathbb{E}[X]}{\sqrt{\frac{\sigma^2}{N}}} \xrightarrow{w} \mathcal{N}(0, 1)$$

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

PROOF. Replacing if necessary X_i by $(X_i - \mathbb{E}[X])/\sigma$, we can assume that $\mathbb{E}[X] = 0$ and $\sigma^2 = \text{Var}(X) = \mathbb{E}[X^2] = 1$. In such a case, we only need to prove that

$$Z_N = \sum_{i=1}^N \frac{X_i}{\sqrt{N}} \xrightarrow{w} \mathcal{N}(0, 1)$$

To this aim, (as for all weak convergence results), it is sufficient to show that the characteristic function of Z_N , $u \mapsto \mathbb{E}(e^{iuZ_N})$ tends to the one of $\mathcal{N}(0, 1)$ that is $u \mapsto e^{-u^2/2}$. Write

$$\mathbb{E}(e^{iuZ_N}) = \mathbb{E}(e^{iu \sum_{i=1}^N X_i / \sqrt{N}}) = \left[\mathbb{E}(e^{iuX/\sqrt{N}}) \right]^N$$

$$\begin{aligned} \text{Taylor Exp.} \quad & \approx \left[\mathbb{E} \left(1 + \frac{iu}{\sqrt{N}} X - \frac{u^2}{2N} X^2 + \text{remainder terms} \right) \right]^N \approx \left[1 + \frac{iu}{\sqrt{N}} \underbrace{\mathbb{E}[X]}_0 - \frac{u^2}{2N} \underbrace{\mathbb{E}[X^2]}_1 \right]^N = \left(1 - \frac{u^2}{2N} \right)^N \rightarrow e^{-u^2/2} \end{aligned}$$

This completes the proof. Of course, the difficult part (the one we forgot to mention) consists in showing that the remainder terms do not come into play when N tends toward infinity (in the notation \approx). But the take-home message from this proof is that the CLT is obtained from the convergence of the characteristic function with a Taylor expansion of order 2. ■

Nevertheless, in many practical situations $\sigma^2 = \text{Var}(X)$ is an unknown quantity. We are then tempted to replace σ^2 by an estimator. A natural question is: does the previous theorem remain valid? The answer is yes, as shown by the following proposition.

Proposition 3.5 *Let $(X_n)_n$ be a series of r.v. i.i.d of the same law as X , such that $\mathbb{E}[X^2] < \infty$ and $\sigma^2 > 0$. Noting*

$$\tilde{\sigma}_N^2 = \frac{1}{N} \sum_{i=1}^N [X_i - \mathbb{E}[X]]^2, \quad \hat{\sigma}_N^2 = \frac{1}{N-1} \sum_{i=1}^N [X_i - \bar{X}_N]^2$$

then

$$\frac{\bar{X}_N - \mathbb{E}[X]}{\sqrt{\frac{\hat{\sigma}_N^2}{N}}} \xrightarrow{w} \mathcal{N}(0,1), \quad \frac{\bar{X}_N - \mathbb{E}[X]}{\sqrt{\frac{\tilde{\sigma}_N^2}{N}}} \xrightarrow{w} \mathcal{N}(0,1) \quad (3.1)$$

Lemma 3.6 *If $\mathbb{E}[X^2]$, the estimators $\tilde{\sigma}_N^2$ and $\hat{\sigma}_N^2$ are unbiased and strongly convergent estimators of σ^2 , i.e.*

$$\begin{aligned} \mathbb{E}[\tilde{\sigma}_N^2] &= \mathbb{E}[\hat{\sigma}_N^2] = \sigma^2, \\ \tilde{\sigma}_N^2 &\xrightarrow{\mathbb{P}\text{-a.s.}} \sigma^2, \quad \hat{\sigma}_N^2 \xrightarrow{\mathbb{P}\text{-a.s.}} \sigma^2 \end{aligned}$$

PROOF. Clearly, $\mathbb{E}[\tilde{\sigma}_N^2] = \mathbb{E}[X_1 - \mathbb{E}[X_1]]^2 = \sigma^2$. By the LGN, $\tilde{\sigma}_N^2 = \frac{1}{N} \sum_{i=1}^N \{X_i - \mathbb{E}[X]\}^2 \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}[X - \mathbb{E}[X]]^2 = \sigma^2$. Moreover, introducing $\mathbb{E}[X]$ in the square term $[X_i - \bar{X}_N]^2$, we get after straightforward algebra:

$$\hat{\sigma}_N^2 = \frac{1}{N-1} \left(\sum_{i=1}^N (X_i - \mathbb{E}[X])^2 - 2 \underbrace{\sum_{i=1}^N (X_i - \mathbb{E}[X])(\bar{X}_N - \mathbb{E}[X])}_{N(\bar{X}_N - \mathbb{E}[X])} + N(\bar{X}_N - \mathbb{E}[X])^2 \right) = \frac{\sum_{i=1}^N (X_i - \mathbb{E}[X])^2}{N-1} - \frac{N}{N-1} (\bar{X}_N - \mathbb{E}[X])^2 \quad (3.2)$$

The LGN shows that $\bar{X}_N \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}[X]$ and we deduce by applying again the LGN on the expression of $\hat{\sigma}_N^2$ obtained in (3.2), that $\hat{\sigma}_N^2 \xrightarrow{\mathbb{P}\text{-a.s.}} \sigma^2$. Finally, by (3.2),

$$\mathbb{E}(\hat{\sigma}_N^2) = \frac{1}{N-1} (N\sigma^2 - N\text{Var}\bar{X}_N) = \frac{1}{N-1} \left(N\sigma^2 - N\frac{\sigma^2}{N} \right) = \sigma^2 \quad \blacksquare$$

Remark 3.7 *If $\mathbb{E}[X]$ is known, it is better to use $\tilde{\sigma}_N^2$. Otherwise, we use $\hat{\sigma}_N^2$.*

The proof of the proposition 3.5 is a simple consequence of Slutski's Lemma (please, check it!), which we will now state without proof. Slutski's Lemma in many practical situations allows to obtain weak convergences of quantities of interest from other weak convergences provided that the quantities of interest only differ by variables that converges in probability to constants.

Theorem 3.8. (► SLUTSKI'S LEMMA) *If*

- i) $X_n \xrightarrow{w} X$,
- ii) $Y_n \xrightarrow{\mathbb{P}\text{-prob}} a$ where a is a constant

then for any continuous (and not necessarily bounded) function f

$$f(X_n, Y_n) \xrightarrow{w} f(X, a)$$

Another result is extremely useful for obtaining other weak convergence results associated with $g(X_n)$ from a weak convergence result associated to X_n .

Theorem 3.9. (► Δ -METHOD) *If*

- i) $\sqrt{n}(X_n - a) \xrightarrow{w} Z$,
- ii) $x \mapsto g(x)$ is differentiable at the point $x = a$

then

$$\sqrt{n}(g(X_n) - g(a)) \xrightarrow{w} g'(a)Z$$

PROOF. First we notice that i) implies that $X_n \xrightarrow{\mathbb{P}\text{-prob}} a$. This yields $\frac{g(X_n) - g(a)}{X_n - a} \xrightarrow{\mathbb{P}\text{-prob}} g'(a)$. We then apply Slutski's lemma, noting that

$$\sqrt{n}(g(X_n) - g(a)) = \underbrace{\sqrt{n}(X_n - a)}_{\xrightarrow{w} Z} \underbrace{\left(\frac{g(X_n) - g(a)}{X_n - a}\right)}_{\xrightarrow{\mathbb{P}\text{-prob}} g'(a)}$$

■

3.2 On the use of the CLT for Monte Carlo methods

Let X be a random variable and f be a measurable function such that $\mathbb{E}(f^2(X)) < \infty$. We wish to approximate $\mathbb{E}(f(X))$ by a Monte Carlo method. We can draw a N -sample (X_1, \dots, X_N) according to the law of X , and then set

$$S_N = \frac{1}{N} \sum_{i=1}^N f(X_i), \quad V_N = \frac{1}{N-1} \sum_{i=1}^N (f(X_i) - S_N)^2$$

Due to the LLN, S_n converges a.s. to $\mathbb{E}(f(X))$ and by Proposition 3.5,

$$\sqrt{N} \frac{S_N - \mathbb{E}(f(X))}{\sqrt{V_N}} \xrightarrow{w} \mathcal{N}(0, 1)$$

This last result then gives a confidence interval on the estimator S_n . Indeed, weak convergence implies that

$$\mathbb{P} \left(\sqrt{N} \frac{S_N - \mathbb{E}(f(X))}{\sqrt{V_N}} \in [-a, a] \right) \rightarrow \mathbb{P}(|G| \leq a)$$

where $G \sim \mathcal{N}(0, 1)$. So, $\left[S_N - a\sqrt{\frac{V_N}{N}}; S_N + a\sqrt{\frac{V_N}{N}} \right]$ is a confidence interval for $\mathbb{E}(f(X))$ of asymptotic level α if a is the quantile of order $1 - \alpha/2$ of the law $\mathcal{N}(0, 1)$, i.e. $\mathbb{P}(|G| \leq a) = 1 - \alpha/2$ with $G \sim \mathcal{N}(0, 1)$. In many examples, we take $\alpha = 0.05$ and in this case $a \approx 1.96$.

3.3 Take-home message

- a) The assumptions of the LLN and CLT should be perfectly known.
- b) The student should know exactly which estimators are biased or unbiased. The proof should be known.
- c) Use of the Slutski lemma and of the delta method to obtain other weak convergence from the CLT.
- d) Diverse methods to obtain a confidence interval.

3.4 Highlights

Monte Carlo methods. Source: Wikipedia

The term "Monte Carlo method" was coined in the 1940s by physicists working on nuclear weapon projects in the Los Alamos National Laboratory.

Enrico Fermi in the 1930s and Stanislaw Ulam in 1946 first had the idea. Ulam later contacted John Von Neumann to work on it.

Physicists at Los Alamos Scientific Laboratory were investigating radiation shielding and the distance that neutrons would likely travel through various materials. Despite having most of the necessary data, such as the average distance a neutron would travel in a substance before it collided with an atomic nucleus or how much energy the neutron was likely to give off following a collision, the problem could not be solved with analytical calculations. John von Neumann and Stanislaw Ulam suggested that the problem be solved by modeling the experiment on a computer using chance. Being secret, their work required a code name. Von Neumann chose the name "Monte Carlo". The name is a reference to the Monte Carlo Casino in Monaco where Ulam's uncle would borrow money to gamble.

Random methods of computation and experimentation (generally considered forms of stochastic simulation) can be arguably traced back to the earliest pioneers of probability theory (see, e.g., Buffon's needle, and the work on small samples by William Sealy Gosset), but are more specifically traced to the pre-electronic computing era. The general difference usually described about a Monte Carlo form of simulation is that it systematically "inverts" the typical mode of simulation, treating deterministic problems by first finding a probabilistic analog (see Simulated annealing). Previous methods of simulation and statistical sampling generally did the opposite: using simulation to test a previously understood deterministic problem. Though examples of an "inverted" approach do exist historically, they were not considered a general method until the popularity of the Monte Carlo method spread.

Monte Carlo methods were central to the simulations required for the Manhattan Project, though were severely limited by the computational tools at the time. Therefore, it was only after electronic computers were first built (from 1945 on) that Monte Carlo methods began to be studied in depth. In the 1950s they were used at Los Alamos for early work relating to the development of the hydrogen bomb, and became popularized in the fields of physics, physical chemistry, and operations research. The Rand Corporation and the U.S. Air Force were two of the major organizations responsible for funding and disseminating information on Monte Carlo methods during this time, and they began to find a wide application in many different fields.

Uses of Monte Carlo methods require large amounts of random numbers, and it was their use that spurred the development of pseudorandom number generators, which were far quicker to use than the tables of random numbers which had been previously used for statistical sampling.

Chapter 4

Exact or approximate sampling

Contents

| | | |
|------------|--|-----------|
| 4.1 | Exact Sampling | 17 |
| 4.1.1 | The inverse cumulative distribution function | 17 |
| 4.1.2 | The rejection sampling | 19 |
| 4.1.3 | Sampling from a conditional distribution | 20 |
| 4.1.4 | Other sampling methods | 22 |
| 4.2 | Approximate sampling | 23 |
| 4.2.1 | Importance Sampling | 23 |
| 4.2.2 | Other methods for approximate sampling | 25 |
| 4.3 | Take-home message | 25 |

Keywords: *Cumulative distribution function, generalized inverse, the rejection sampling, sampling by mapping, sampling from a conditional law, importance sampling.*

4.1 Exact Sampling

4.1.1 The inverse cumulative distribution function

Let Y be a real random variable with cumulative distribution function $F_Y : t \mapsto F_Y(t) = \mathbb{P}(Y \leq t)$. Immediate properties of the function F_Y are as follows: F_Y is non-decreasing, right-continuous, has a left-limit, tends toward 0 in $-\infty$ and toward 1 in ∞ . If F_Y is strictly increasing, it is easy to define the inverse function of F_Y . Unfortunately, it may be constant on some intervals (for example when Y is a discrete r.v.). As a result, we need to define a generalized inverse, which is also valid for non-decreasing functions (which can be constant over certain intervals):

Definition 4.1. (►GENERALIZED INVERSE) We define *the generalized inverse* of F_Y by: for any $x \in [0, 1]$,

$$F_Y^{-1}(x) = \inf\{y \in \mathbb{R} : F_Y(y) \geq x\}$$

Of course, if F_Y is invertible then the generalized inverse is the inverse in the classical sense. The generalized inverse F_Y^{-1} of the cumulative function F_Y is also called the **quantile function**.

Proposition 4.2 *If $U \sim \mathcal{U}[0,1]$, then $F_Y^{-1}(U) \stackrel{\mathcal{L}}{\equiv} Y$.*

PROOF. Using that the cumulative distribution functions are right-continuous, the following (intuitive) equivalence can be proved:

$$\forall u, v \in \mathbb{R}, \quad \{F_Y^{-1}(u) \leq v\} \iff \{u \leq F_Y(v)\}.$$

So, if U is a r.v. with uniform distribution on $[0,1]$, we have the equality

$$P(F_Y^{-1}(U) \leq y) = P(U \leq F_Y(y)) = F_Y(y),$$

which completes the proof. ■

As a byproduct of this Proposition, we can propose a sampling method by using the inverse cumulative distribution function (in the case where this function is explicitly available). We formalize it in the following Corollary.

Corollary 4.3 *Let f be a probability density. Let $F(t) = \int_{-\infty}^t f(u)du$ with generalized inverse F^{-1} . Suppose F^{-1} is explicitly available. Then we can draw a r.v. Y with density f by the following Algorithm 1:*

Algorithm 1 Sampling by the inverse cumulative distribution function

- 1: Draw $U \sim \mathcal{U}[0,1]$
 - 2: Set $Y = F^{-1}(U)$
-

Example 4.4 (► SIMULATION OF A DISCRETE LAW) *Let Y be a random variable with support on $(y_k)_{k \in \mathbb{N}}$, such that $P(Y = y_k) = p_k$. If $U \sim \mathcal{U}[0,1]$, then*

$$X = y_0 \mathbf{1}_{U \leq p_0} + \sum_{k \geq 1} y_k \mathbf{1}_{\sum_{i=0}^{k-1} p_i < U \leq \sum_{i=0}^k p_i} \stackrel{\mathcal{L}}{\equiv} Y$$

Indeed, we can easily check that for $k > 0$, we have

$$\mathbb{P}(X = y_k) = \mathbb{P}\left(\sum_{i=0}^{k-1} p_i < U \leq \sum_{i=0}^k p_i\right) = p_k$$

This intuitive sampling method is actually a sampling by inversion of the cumulative distribution function (please, check it).

A typical example is to draw Y according to the Bernoulli law with parameter p by drawing $U \sim \mathcal{U}[0,1]$ and setting $Y = \mathbf{1}_{U \leq p}$. Another important example is to sample a binomial law $\mathcal{B}(n, p)$ by setting $Y = \sum_{i=1}^n \mathbf{1}_{U_i \leq p}$ where U_i is $\mathcal{U}[0,1]$. Apart from the examples of discrete law, the generalized inverse is reduced in most other situations to the classical inverse. And if available in a closed-form, it is possible to use this sampling method. Let us look at an example of a "continuous" law.

Example 4.5 *The exponential distribution of parameter $\lambda > 0$ has a density of $f(x) = \lambda \exp(-\lambda x) \mathbf{1}_{\mathbb{R}^+}(x)$. The associated cumulative distribution function is $F(t) = -\exp(-\lambda t) + 1$ for $t \geq 0$ with inverse $F^{-1}(u) = -\frac{1}{\lambda} \ln(1-u)$. So, if $U \sim \mathcal{U}[0,1]$, $-\frac{1}{\lambda} \ln(1-U) \sim \exp(\lambda)$. Now $1-U$ also follows $\mathcal{U}[0,1]$. Since it is better to use as few calculations as possible, we set $F^{-1}(1-U) = -\frac{1}{\lambda} \ln(U) \sim \exp(\lambda)$.*

4.1.2 The rejection sampling

The inverse cumulative distribution function does not always admit an explicit form. We now present a method widely used in practice which consists in (1) proposing candidates according to another law and then (2) in defining a stopping rule, namely, a criterion which if met, allows us to select a variable which ("surprisingly") will be distributed according to the target distribution.

Proposition 4.6 *Let f and g be two probability densities satisfying for any $x \in \mathbb{R}$, $f(x) \leq Mg(x)$. We consider (U_i) and (X_i) two sequences of random variables such that*

- i) $U_i \sim \mathcal{U}[0,1]$ and $X_i \sim g$
- ii) the sequence of random variables $(U_s)_{s \geq 1}$ and $(X_t)_{t \geq 1}$ are independent.

Let $\nu = \inf \left\{ i \in \mathbb{N}, U_i \leq \frac{f(X_i)}{Mg(X_i)} \right\}$ and $Y = X_\nu$. Then,

- a) $\nu \sim \text{Geom}(M^{-1})$ and $Y \sim f$.
- b) ν and Y are independent variables.

We encourage the reader to read the following proof, rich in technical lessons.

PROOF. Note that the random vector $\begin{pmatrix} U \\ X \end{pmatrix}$ has the density $(u, x) \mapsto \mathbf{1}_{[0,1]}(u)g(x)$. This allows to write, knowing that the vectors $\begin{pmatrix} U_k \\ X_k \end{pmatrix}$ are i.i.d.,

$$\begin{aligned}
 \mathbb{P}(Y \in A, \nu = k) &= \mathbb{P}(X_k \in A, \nu = k) \\
 &= \mathbb{P}\left(U_1 > \frac{f(X_1)}{Mg(X_1)}, \dots, U_{k-1} > \frac{f(X_{k-1})}{Mg(X_{k-1})}, U_k \leq \frac{f(X_k)}{Mg(X_k)}, X_k \in A\right) \\
 &= \mathbb{P}\left(U_1 > \frac{f(X_1)}{Mg(X_1)}\right)^{k-1} \mathbb{P}\left(U_k \leq \frac{f(X_k)}{Mg(X_k)}, X_k \in A\right) \\
 &= \left(\iint_{\{u > \frac{f(x)}{Mg(x)}\}} \mathbf{1}_{[0,1]}(u)g(x) \, du dx \right)^{k-1} \iint_{\{u \leq \frac{f(x)}{Mg(x)}\} \cap \{x \in A\}} \mathbf{1}_{[0,1]}(u)g(x) \, du dx \\
 &= \left(\int \left(1 - \frac{f(x)}{Mg(x)}\right) g(x) \, dx \right)^{k-1} \int_{\{x \in A\}} \frac{f(x)}{Mg(x)} g(x) \, dx \\
 &= \left(1 - \frac{1}{M}\right)^{k-1} \frac{1}{M} \int_{\{x \in A\}} f(x) dx
 \end{aligned} \tag{4.1}$$

where the penultimate line comes from integration over u and the last line comes from the fact that f and g are probability densities. We then have by marginalization

$$\mathbb{P}(\nu = k) = \mathbb{P}(Y \in \mathbb{R}, \nu = k) = \left(1 - \frac{1}{M}\right)^{k-1} \frac{1}{M}$$

so $\nu \sim \text{Geom}(M^{-1})$, and therefore $\mathbb{P}(\nu < \infty) = \sum_{k=1}^{\infty} \mathbb{P}(\nu = k) = 1$. This implies:

$$\mathbb{P}(Y \in A) = \mathbb{P}(Y \in A, \nu < \infty) = \sum_{k=1}^{\infty} \mathbb{P}(Y \in A, \nu = k) = \sum_{k=1}^{\infty} \left(1 - \frac{1}{M}\right)^{k-1} \frac{1}{M} \int_{\{x \in A\}} f(x) dx = \int_{\{x \in A\}} f(x) dx$$

Therefore $Y \sim f$. Finally, (4.1) shows that ν and Y are independent random variables. ■

The immediate corollary of this proposition will be the rejection sampling method:

Corollary 4.7 (►REJECTION ALGORITHM) *Let f and g be two densities such that*

- i) we can sample according to the density g ,*
- ii) there is M such that for every $x \in \mathbb{R}$, $f(x) \leq Mg(x)$*
- iii) for any x such that $g(x) > 0$, the quantity $\frac{f(x)}{Mg(x)}$ is explicitly available*

Then, we can sample a random variable Y according to the density f by the algorithm 2.

Algorithm 2 The rejection sampling

- 1: draw $X \sim g$, $U \sim \mathcal{U}[0,1]$ independently
 - 2: **while** $U \geq \frac{f(X)}{Mg(X)}$ **do** draw $X \sim g$, $U \sim \mathcal{U}[0,1]$ independently
 - 3: **end while**
 - 4: Set $Y = X$
-

A fundamental difference with sampling by the inverse cumulative distribution function is the number of samplings required before producing a single r.v. distributed according to the target distribution: this waiting time before acceptance is here a random variable v of geometric distribution $\text{Geom}(M^{-1})$. The expectation of a geometric distribution is the inverse of its parameter, which gives $\mathbb{E}(v) = M$. The function f being given, we are therefore interested in a density g such that $f \leq Mg$ with M as small as possible (to minimize the waiting time before acceptance). Of course, we always have $M \geq 1$ (to see this, just integrate $f(x) \leq Mg(x)$ and note that f and g are probability densities) but we can't reach the limit $M = 1$. Otherwise, we would have $f(x) \leq g(x)$ and by integration $\int \underbrace{g(x) - f(x)}_{\geq 0} dx = 1 - 1 = 0$, hence $f = g$ and if we can draw according to g , that means we can draw according to f !

4.1.3 Sampling from a conditional distribution

4.1.3.1 Link with the rejection sampling

Let X be a r.v. with density g , which can be sampled. If you want to sample a r.v. according to the distribution of X conditionally on the event $\{X \in [a, b]\}$, an intuitive idea would be to sample following g and to keep only the candidates that fall into $[a, b]$. We will show that this method actually corresponds to a particular case of the rejection sampling. The distribution of X conditionally on the event $\{X \in [a, b]\}$ has the density

$$f(x) = \mathbf{1}_{[a,b]}(x) \frac{g(x)}{\int_a^b g(u) du}.$$

Indeed, write

$$\mathbb{P}(X \in A | X \in [a, b]) = \frac{\mathbb{P}(X \in A \cap [a, b])}{\mathbb{P}(X \in [a, b])} = \frac{\int_A \mathbf{1}_{[a,b]}(x) g(x) dx}{\int_a^b g(u) du} = \int_A f(x) dx.$$

Now, note that

$$f(x) \leq Mg(x) \quad \text{avec} \quad M = \frac{1}{\int_a^b g(u) du}$$

Therefore, applying the rejection sampling, draw $U \sim \mathcal{U}[0, 1]$ and $X \sim g$ independently and set $Y = X$ only if

$$U \leq \frac{f(X)}{Mg(X)} = \mathbf{1}_{[a,b]}(X) = \begin{cases} 1 & \text{if } X \in [a,b] \\ 0 & \text{otherwise} \end{cases}.$$

Finally if $X \in [a, b]$, the candidate is accepted with probability 1, and if $X \notin [a, b]$, it is refused with probability 1. This is exactly the intuitive idea of the algorithm. In this example, $M = 1 / \int_a^b g(u) du$ is not available in closed-form in general; still, as seen before, the ratio $f(X)/(Mg(X))$ can be calculated and this allows to apply the rejection sampling. Indeed if $\mathbb{P}(X \in [a, b])$ is small then, we will wait long before accepting a candidate. To apply the rejection sampling with a target density with support inside $[a, b]$, it is often better to propose according to a density with support inside $[a, b]$.

4.1.3.2 Link with sampling by inversion of the cumulative distribution function

If F_X et F_X^{-1} are available in closed-form, it is more efficient to sample a random variable according to the conditional law of X with respect to the event $\{X \in [a, b]\}$ in the following way

- i) draw $U \sim \mathcal{U}[0, 1]$
- ii) set $Y = F_X^{-1}(F_X(a) + U[F_X(b) - F_X(a)])$.

Algorithm 3 Sample from a conditional distribution

- 1: draw $U \sim \mathcal{U}[0, 1]$
 - 2: Set $Y = F_X^{-1}(F_X(a) + U[F_X(b) - F_X(a)])$.
-

Indeed, write for all $t \in [a, b]$,

$$\begin{aligned} \mathbb{P}(Y \leq t) &= \mathbb{P}(F_X^{-1}(F_X(a) + U[F_X(b) - F_X(a)]) \leq t) = \mathbb{P}(F_X(a) + U[F_X(b) - F_X(a)] \leq F_X(t)) \\ &= \mathbb{P}\left(U \leq \frac{F_X(t) - F_X(a)}{F_X(b) - F_X(a)}\right) = \frac{F_X(t) - F_X(a)}{F_X(b) - F_X(a)} = \frac{\int_a^t g(x) dx}{\int_a^b g(u) du} = \int_{-\infty, t[} f(x) dx \end{aligned}$$

An alternative way to see this is the following: the target density is

$$g(x) = \frac{f(x) \mathbb{1}_{x \in [a,b]}}{\int_a^b f(t) dt} = \frac{f(x) \mathbb{1}_{x \in [a,b]}}{F_X(b) - F_X(a)}$$

This target density has the cumulative distribution function G defined by for $x \in [a, b]$,

$$G(x) = \frac{\int_a^x f(t) dt}{F_X(b) - F_X(a)} = \frac{F(x) - F(a)}{F_X(b) - F_X(a)}$$

Since $\frac{F_X(x) - F_X(a)}{F_X(b) - F_X(a)} = u$ is equivalent to $F_X(x) = F_X(a) + u[F_X(b) - F_X(a)]$, we deduce that G has a generalized inverse which is:

$$G^{-1}(u) = F_X^{-1}(F_X(a) + u[F_X(b) - F_X(a)]).$$

This completes the proof.

4.1.4 Other sampling methods

4.1.4.1 Sampling by mapping

We also can try to show that the target distribution is the one of random variables, which are obtained by a mapping (typically C^1 -mappings but not always) of r.v. which can be sampled easily (typically r.v. with uniform distribution).

Lemma 4.8 CHANGE OF VARIABLES- *Let ϕ be C^1 -mapping from an open set O of \mathbb{R}^d to an open set O' of \mathbb{R}^d . Assume that $V \sim g$ where g has a support in O , then $U = \phi(V)$ has the density $u \mapsto f(u) = g \circ \phi^{-1}(u) |\det(\nabla \phi^{-1}(u))|$*

It is useless to learn this lemma by heart since we can recover it quite intuitively by the change of variable formula for multiple integrals:

$$\mathbb{E}(h(U)) = \mathbb{E}(h \circ \phi(V)) = \int \underbrace{h \circ \phi(v)}_u g(v) dv = \int h(u) \times g(\phi^{-1}(u)) \underbrace{\left| \frac{\partial v}{\partial u} \right|}_{|\det \nabla(\phi^{-1}(u))|} du = \int h(u) f(u) du$$

Thanks to this lemma, we can state the following proposition which proposes a sampling method for obtaining a couple of independent gaussian variables.

Proposition 4.9 BOX MULLER- *Let U and V be two i.i.d.r.v. with distribution $\mathcal{U}[0, 1]$. Set*

$$X = \sqrt{-2 \ln U} \cos(2\pi V), \quad Y = \sqrt{-2 \ln U} \sin(2\pi V)$$

Then X and Y are i.i.d. with distribution $\mathcal{N}(0, 1)$.

PROOF. We can easily show by the change of variable formula that if X and Y are i.i.d. of distribution $\mathcal{N}(0, 1)$ then, denoting by (R, θ) the polar coordinates of (X, Y) , the couple of r.v. (R, θ) is independent, $R^2 \sim \exp(1/2)$ and $\theta \sim \mathcal{U}[0, 2\pi]$. ■

To draw gaussian r.v., we can also use the inversion of the cumulative distribution function of $\mathcal{N}(0, 1)$. Set $\mathcal{N}(t) = \int_{-\infty}^t (e^{-u^2/2} / \sqrt{2\pi}) du$ which is not explicit (nor its inverse) but it may happen that the numerical values of $\mathcal{N}^{-1}(t)$ can be approximated at any order in some libraries and in that case, we can draw a gaussian distribution by using the inverse cumulative distribution function.

4.1.4.2 Sampling a gaussian vector

The Box and Muller algorithm allows to sample easily a vector of distribution $\mathcal{N}(0, I)$ where I is the identity matrix. If we then intend to sample a general gaussian vector of distribution $\mathcal{N}(\mu, \Sigma)$ then note that Σ being symmetric, and nonnegative, there exists a real valued triangular matrix A such that $AA^T = \Sigma$ (see the **Cholesky** decomposition), then the random vector $\mu + AG$ with $G \sim \mathcal{N}(0, I)$ is distributed according to $\mathcal{N}(\mu, \underbrace{\Sigma}_{AA^T})$.

4.1.4.3 Sampling by marginalisation

Other methods for sampling from a given distributions exist: for example, it happens that the target density f can be written as $f(u) = \int g(u, v) dv$ where g is itself a density. In that case, if we can sample according to

g then we sample $(U, V) \sim g$ and U will have the density f . Often, $g(u, v) = h(v)p(u|v)$ and h and p can be sampled. We draw $V \sim h$. Then, setting $V = v$, we draw $U \sim p(\cdot|v)$. When V is a discrete-valued random variable, we then say that f is a *mixing distribution*.

4.2 Approximate sampling

4.2.1 Importance Sampling

In most situations, we do not really intend to sample according to a distribution but we want an approximation of the expectation of a functional of a r.v. with a given law, i.e. $\mathbb{E}_f(h(Y)) = \int h(y)f(y)dy$ (for the notation \mathbb{E}_f , see the footnote ¹). In that case, Importance Sampling (IS) boils down to sample (X_1, \dots, X_N) according to a proposal density $g > 0$, $X_i \sim g$, and then, to use the approximation

$$\mathbb{E}_f(h(Y)) = \int h(y)f(y)dy \approx \bar{S}_N = \frac{\sum_{i=1}^N \frac{f(X_i)}{g(X_i)} h(X_i)}{N}$$

Conditions for applying importance sample are then

- i) we have $\int_{g(x)=0} f(x)dx = 0$.
- ii) we can draw from g .
- iii) for all x , $f(x)/g(x)$ has a closed-form expression.

The first assumption implies

$$\mathbb{E}_g(f(X)h(X)/g(X)) = \int \frac{f(x)}{g(x)} h(x)g(x)dx = \int_{g(x) \neq 0} h(x)f(x)dx = \mathbb{E}_f(h(Y)).$$

The importance sampling estimator \bar{S}_N is clearly unbiased $\mathbb{E}(\bar{S}_N) = \mathbb{E}_f(h(Y))$, strongly convergent under the assumption $\mathbb{E}_g \left[\frac{f(X)}{g(X)} |h(X)| \right] = \mathbb{E}_f[|h(X)|] < \infty$. The LLN indeed show that

$$\frac{\sum_{i=1}^N \frac{f(X_i)}{g(X_i)} h(X_i)}{N} \xrightarrow{\mathbb{P}\text{-a.s.}} \int \frac{f(x)}{g(x)} h(x)g(x)dx = \int h(y)f(y)dy$$

We can see Importance Sampling as a method where we draw X_i according to a proposal distribution g and the error is then corrected, (because we have not drawn according to f) by associating to each X_i a weight $f(X_i)/g(X_i)$.

We can have the impression that IS can be applied more often than the rejection sampling since in IS, we do not need to assume that $\sup_x \frac{f(x)}{g(x)} < \infty$ (this condition can be a bit restrictive sometimes)... It is true but the advantage of the rejection sampling is that it produces an exact sample according to f !

4.2.1.1 Optimisation of the proposal density

Under the assumption that $\mathbb{E}_g \left[\frac{f^2(X)}{g^2(X)} h^2(X) \right] = \mathbb{E}_f \left[\frac{f(X)}{g(X)} h^2(X) \right] < \infty$, the CLT gives the quality of the approximation as follows

$$\sqrt{N} \left(\frac{\sum_{i=1}^N \frac{f(X_i)}{g(X_i)} h(X_i)}{N} - \int h(y)f(y)dy \right) \xrightarrow{w} \mathcal{N} \left(0, \text{Var}_g \left(\frac{f(x)}{g(x)} h(x) \right) \right)$$

¹ When we put an f under the expectation, this means that we take the expectation with respect to a random variable of density f

To target a density g , that gives rise to the most precise approximation, we should then minimize the quantity $\text{Var}_g\left(\frac{f(x)}{g(x)}h(x)\right)$.

Proposition 4.10

$$\inf \left\{ \text{Var}_g \left(\frac{f(x)}{g(x)} h(x) \right); g \text{ density} \right\} = \left(\int f(x) |h(x)| dx \right)^2 - \left(\int f(x) h(x) dx \right)^2$$

The infimum is attained with a density g^* defined by $g^*(x) = f(x)|h(x)| / \int f(u)|h(u)| du$.

PROOF. Write

$$\text{Var}_g \left(\frac{f(x)}{g(x)} h(x) \right) = \int \frac{f(x)}{g(x)} h^2(x) f(x) dx - \left(\int h(y) f(y) dy \right)^2$$

The second term does not depend on g , we can thus minimize the first term of the right-hand side. By the Cauchy-Schwarz inequality;

$$\left(\int f(x) |h(x)| dx \right)^2 = \left(\int \frac{f(x)}{\sqrt{g(x)}} |h(x)| \sqrt{g(x)} dx \right)^2 \leq \left(\int \frac{f^2(x)}{g(x)} |h(x)|^2 dx \right) \underbrace{\left(\int (\sqrt{g(x)})^2 dx \right)}_1 = \int \frac{f(x)}{g(x)} h^2(x) f(x) dx \quad (4.2)$$

Moreover, the equality in the Cauchy Schwarz inequality holds for $\sqrt{g^*(x)} \propto \frac{f(x)}{\sqrt{g^*(x)}} |h(x)|$, which can also be written as, knowing that g^* is a density: $g^*(x) = f(x)|h(x)| / \int f(u)|h(u)| du$.

An alternative (and more direct) proof is as follows: Note that

$$\text{Var}_g \left(\frac{f(x)}{g(x)} |h(x)| \right) = \int \frac{f(x)}{g(x)} h^2(x) f(x) dx - \left(\int |h(y)| f(y) dy \right)^2$$

And since the variance is non negative, we get

$$\int \frac{f(x)}{g(x)} h^2(x) f(x) dx \geq \left(\int |h(y)| f(y) dy \right)^2$$

This shows (4.2). This inequality becomes an equality for g^* such that $\text{Var}_{g^*} \left(\frac{f(x)}{g^*(x)} |h(x)| \right) = 0$, that is, if $x \mapsto \frac{f(x)}{g^*(x)} |h(x)|$ is g^* -a.s. a constant, that is if $g^* \propto |h|f$ which corresponds to $g^*(x) = f(x)|h(x)| / \int f(u)|h(u)| du$. Finally, with this definition of g^* , we have

$$\int \frac{f(x)}{g(x)} h^2(x) f(x) dx \geq \int \frac{f(x)}{g^*(x)} h^2(x) f(x) dx$$

■

Unfortunately, this lemma has no immediate practical consequences since it is not obvious to know how to sample according to the density distribution $g^*(x) = f(x)|h(x)| / \int f(u)|h(u)| du$ and even if it were the case, to be able to use sampling importance, it would have been necessary to know how to calculate explicitly

$$\frac{f(x)}{g^*(x)} = \frac{\int f(u)|h(u)| du}{|h(x)|}$$

which we don't usually know how to do, since we're trying to give a numerical value to $\int f(u)h(u) du$. However, the message of this lemma is that, when considering the optimal approximation of $\mathbb{E}_f(h(Y))$, the density that corresponds to the minimal variance is not necessarily f .

4.2.2 Other methods for approximate sampling

There are many other approximate sampling algorithms, in particular the MCMC algorithms (Monte Carlo by Markov Chains) whose general principle consists in building a Markov chain whose stationary distribution is of density f . The approximation of $\mathbb{E}_f(h(X))$ by $N^{-1} \sum_{i=1}^N h(X_i)$ then comes from a law of large numbers for a Markov chain (we can no longer use the usual LLN because the r.v. are not i.i.d.).

Other hybrid methods judiciously combine Importance Sampling and MCMC algorithms. That being said, the question remains open for choosing the most efficient estimation method that approximates $\mathbb{E}_f(h(X))$; it is still the subject of intensive research, particularly when the variable X evolves in a high dimensional space.

4.3 Take-home message

- a) Sampling by the inverse cumulative distribution function: the student should know how to show it when the cumulative distribution function is invertible. He should know how to distinguish discrete and continuous variables.
- b) Know how to re-prove the proposition that justifies the rejection sampling. The rejection algorithm requires a number of random samplings before producing a sample from the target distribution.
- c) Know how to make change of variables to do the mapping sampling.
- d) On the use of importance sampling: the student should know the justification for the convergence and the CLT. Optimization of asymptotic variance.

Chapter 5

Metropolis-Hastings algorithms

Keywords: *Markov chains, Markov property, Metropolis Hastings, canonical space.*

This chapter introduces very basic tools in the Markov Chain Monte Carlo (MCMC) theory. We only focus on the most fundamental properties that will be useful for a first understanding of Metropolis-Hastings (MH) algorithms. Let us start with a gentle and smooth introduction to Markov chains.

5.1 Main notation

Let (X, \mathcal{X}) be a measurable space, i.e. \mathcal{X} is a σ -algebra on X , and consider the following notations.

- $M_+(X)$ is the set of (non-negative) measures on (X, \mathcal{X}) .
- $M_1(X)$ is the set of probability measures on (X, \mathcal{X}) .
- $F(X)$ is the set of real-valued measurable functions f on X and $F_+(X)$ the set of non-negative measurable functions on X .
- If $k \leq \ell$, $u_{k:\ell}$ means (u_k, \dots, u_ℓ) and $u_{k:\infty}$ means $(u_{k+\ell})_{\ell \in \mathbb{N}}$.

Other notation will be introduced progressively.

5.2 Definitions

We first describe a Markov kernel, which will then be fundamental for the definition of a Markov chain.

Definition 5.1. We say that $P : X \times \mathcal{X} \rightarrow \mathbb{R}^+$ is a Markov kernel, if for all $(x, A) \in X \times \mathcal{X}$,

- $X \ni y \mapsto P(y, A)$ is $\mathcal{X} / \mathcal{B}(\mathbb{R}^+)$ measurable,
- $\mathcal{X} \ni B \mapsto P(x, B)$ is a probability measure on (X, \mathcal{X}) .

In words, for all $(x, A) \in X \times \mathcal{X}$, as a function of the first component only, $P(\cdot, A)$ is measurable and as a function of the second component only, $P(x, \cdot)$ is a probability measure. In particular, $P(x, X) = 1$ for all $x \in X$. Since $P(x, \cdot)$ is a measure, we also use the infinitesimal notation: $P(x, dy)$. For example,

$$P(x, A) = \int_X \mathbb{1}_A(y) P(x, dy) = \int_A P(x, dy).$$

In almost all the course, a Markov kernel P allows to move a point x from a measurable space (X, \mathcal{X}) to another point on the same measurable space, that is, P is defined on $X \times \mathcal{X}$ but we can more generally define a Markov kernel from a measurable space (X, \mathcal{X}) to another measurable space (Y, \mathcal{Y}) . In such case, P will be a Markov kernel on $X \times \mathcal{Y}$. We can now move on to the definition of a Markov chain.

Definition 5.2. Let $\{X_k : k \in \mathbb{N}\}$ be a sequence of random variables on the same probability space $(\Omega, \mathcal{G}, \mathbb{P})$ and taking values on X , we say that $\{X_k : k \in \mathbb{N}\}$ is a Markov chain with Markov kernel P and initial distribution $\nu \in M_1(X)$ if and only if

- (i) for all $(k, A) \in \mathbb{N} \times \mathcal{X}$, $\mathbb{P}(X_{k+1} \in A | X_{0:k}) = P(X_k, A)$, \mathbb{P} -a.s.
- (ii) $\mathbb{P}(X_0 \in A) = \nu(A)$.

Note that in the definition we consider $\mathbb{P}(X_{k+1} \in A | X_{0:k})$, that is, the conditional probability is with respect to the sigma-field $\sigma(X_{0:k})$. We can actually replace $\sigma(X_{0:k})$ by \mathcal{F}_k as soon as we know that $(X_k)_{k \geq 0}$ is $(\mathcal{F}_k)_{k \geq 0}$ -adapted.

What does it mean exactly? Well, recall that if $\{\mathcal{F}_k : k \in \mathbb{N}\}$ is a sequence of embedded sigma-fields on X (that is, $\mathcal{F}_k \subset \mathcal{F}_{k+1}$ for all k), then $\{\mathcal{F}_k : k \in \mathbb{N}\}$ is called a **filtration** on X and we say that $(X_k)_{k \geq 0}$ is $(\mathcal{F}_k)_{k \geq 0}$ -adapted if X_k is $\mathcal{G}/\mathcal{F}_k$ -measurable for all $k \in \mathbb{N}$. Of course, the most natural filtration for $\{X_k : k \in \mathbb{N}\}$ is indeed $\mathcal{F}_k = \sigma(X_{0:k})$ and unsurprisingly, we call it the **natural filtration**. But other possibilities exist, where \mathcal{F}_k is enlarged to include some other variables alongside with $X_{0:k}$. For example, let $\{Y_k : k \in \mathbb{N}\}$ be any other sequence of random variables on $(\Omega, \mathcal{G}, \mathbb{P})$ and taking values in X (we do not assume anything on the relation between $\{X_k : k \in \mathbb{N}\}$ and $\{Y_k : k \in \mathbb{N}\}$). A typical example corresponds to $X_{k+1} = g(X_{0:k}, Y_{0:k})$, but we do not even need to assume that for the moment. Set $\mathcal{F}_k = \sigma(X_{0:k}, Y_{0:k})$ and assume that

$$\mathbb{P}(X_{k+1} \in A | \mathcal{F}_k) = P(X_k, A), \quad \mathbb{P} - \text{a.s.} \quad (5.1)$$

Since X_ℓ is \mathcal{F}_ℓ -measurable and $\mathcal{F}_\ell \subset \mathcal{F}_k$ for $\ell \leq k$, we deduce that $\sigma(X_{0:k}) \subset \mathcal{F}_k$. This allows to apply the tower property, which yields

$$\begin{aligned} \mathbb{P}(X_{k+1} \in A | X_{0:k}) &= \mathbb{E}[\mathbb{P}(X_{k+1} \in A | \mathcal{F}_k) | X_{0:k}] \\ &= \mathbb{E}[P(X_k, A) | X_{0:k}] = P(X_k, A), \quad \mathbb{P} - \text{a.s.} \end{aligned}$$

and therefore if we assume (5.1), then as soon as $(X_k)_{k \geq 0}$ is $(\mathcal{F}_k)_{k \geq 0}$ -adapted, we can conclude that $\{X_k : k \in \mathbb{N}\}$ is a Markov chain with Markov kernel P . Why is it useful? Well, sometimes you define iteratively X_{k+1} using other variables rather than $X_{0:k}$ only and therefore, considering $\mathbb{P}(X_{k+1} \in A | \mathcal{F}_k)$ is easier to deal with. Let us see it in action through a very simple example.

Example 5.3 Let $\{\varepsilon_k : k \geq 1\}$ be i.i.d. random variables on \mathbb{R}^p with density f with respect to the Lebesgue measure on \mathbb{R}^p , and let $X_0 \sim \mu$. We assume that X_0 is independent of $\{\varepsilon_k : k \in \mathbb{N}\}$. Define

$$X_{k+1} = aX_k + b\varepsilon_{k+1}, \quad k \in \mathbb{N}.$$

Set $\mathcal{F}_0 = \sigma(X_0)$ and for $k \geq 1$, $\mathcal{F}_k = \sigma(X_0, \varepsilon_{1:k})$. Since X_ℓ is a deterministic function of X_0 and $\varepsilon_{1:\ell}$, we deduce that $(X_k)_{k \geq 0}$ is $(\mathcal{F}_k)_{k \geq 0}$ -adapted. Therefore, we only need to check (5.1). Now, for any non-negative or bounded measurable function h on X ,

$$\mathbb{E}[h(X_{k+1}) | \mathcal{F}_k] = \int_{-\infty}^{\infty} \underbrace{h(aX_k + b\varepsilon)}_y f(\varepsilon) d\varepsilon = \int_{-\infty}^{\infty} \underbrace{h(y) f\left(\frac{y - aX_k}{b}\right)}_{P(X_k, dy)} \frac{1}{b^p} dy,$$

where the last equality follows from an adequate change of variable. Therefore, $\{X_k : k \in \mathbb{N}\}$ is a Markov chain with Markov kernel

$$(x, A) \mapsto P(x, A) = \int_A f\left(\frac{y - ax}{b}\right) \frac{1}{b^p} dy.$$

In this example, we can check that \mathcal{F}_k is actually the natural filtration of $\{X_k : k \in \mathbb{N}\}$ but we even do not need to check this property for getting that $\{X_k : k \in \mathbb{N}\}$ is a Markov chain.

What have we learned so far? At this stage, we are able to solve typical exercises where some random variables are given and the question is to determine whether these random variables form a Markov chain or not and if yes, what is the expression of the associated Markov kernel.

5.2.1 Additional notation

We are now ready (and eager) to absorb frantically other notation... For all $\mu \in \mathcal{M}_+(\mathbb{X})$, all Markov kernels P, Q on $\mathbb{X} \times \mathcal{X}$, and all measurable non-negative or bounded functions h on \mathbb{X} , we use the following convention and notation.

- μP is the (positive) measure: $\mathcal{X} \ni A \mapsto \mu P(A) = \int \mu(dx)P(x, A)$,
- PQ is the Markov kernel: $(x, A) \mapsto \int_{\mathbb{X}} P(x, dy)Q(y, A)$,
- Ph is the measurable function $x \mapsto \int_{\mathbb{X}} P(x, dy)h(y)$.

It is easy to check that if μ is a probability measure, then μP is also a probability measure (since $\mu P(\mathbb{X}) = \int_{\mathbb{X}} \mu(dx)P(x, \mathbb{X}) = \int_{\mathbb{X}} \mu(dx) = 1$). With this notation, using Fubini's theorem,

$$\begin{aligned} \mu(P(Qh)) &= (\mu P)(Qh) = (\mu(PQ))h \\ &= \mu((PQ)h) = \int \cdots \int_{\mathbb{X}^3} \mu(dx)P(x, dy)Q(y, dz)h(z). \end{aligned}$$

Therefore, all these parenthesis can be discarded and we can write μPQh without any ambiguity. That is excellent news because it is simpler to deal with expressions without all these parenthesis. To sum up, measures act on the left side of a Markov kernel whereas functions acts on the right side. To make sure you have mastered all the notation, check your understanding with the following equalities $\delta_x P(A) = P(x, A) = P\mathbb{1}_A(x)$.

To finish up with notation, we now define the iterates of a Markov kernel P , which will come in very handy thereafter: for a given Markov kernel P on $\mathbb{X} \times \mathcal{X}$, define $P^0 = I$ where I is the identity kernel: $(x, A) \mapsto \mathbb{1}_A(x)$, and set for $k \geq 0$, $P^{k+1} = P^k P$.

Lemma 5.4 Let $\{X_k : k \in \mathbb{N}\}$ be a Markov chain on the same probability space $(\Omega, \mathcal{G}, \mathbb{P})$ and taking values on \mathbb{X} , with Markov kernel P and with initial distribution $\nu \in \mathcal{M}_1(\mathbb{X})$. Then, for any $n \in \mathbb{N}$, the law of $X_{0:n}$ is $\nu(dx_0) \prod_{i=0}^{n-1} P(x_i, dx_{i+1})$ (with the convention that $\prod_{i=0}^{-1} = 1$).

PROOF. Recall that for all $(n, A) \in \mathbb{N} \times \mathcal{X}$, $\mathbb{P}(X_{n+1} \in A | X_{0:n}) = P(X_n, A)$, \mathbb{P} -a.s. or equivalently for all non-negative measurable functions h_{k+1} on \mathbb{X} , $\mathbb{E}[h_{n+1}(X_{n+1}) | X_{0:n}] = Ph_{n+1}(X_n)$ \mathbb{P} -a.s. We now show by induction that for all $n \in \mathbb{N}$,

$$(H_n) \text{ the law of } X_{0:n} \text{ is } \nu(dx_0) \prod_{i=0}^{n-1} P(x_i, dx_{i+1}).$$

We first note that (H_0) is true since by assumption, $X_0 \sim \nu$. Assume now that (H_n) holds for some $n \in \mathbb{N}$. Then, for all non-negative measurable functions h_0, \dots, h_{n+1} on \mathbb{X} , the tower property yields

$$\mathbb{E} \left[\prod_{i=0}^{n+1} h_i(X_i) \right] = \mathbb{E} \left[\left(\prod_{i=0}^n h_i(X_i) \right) \mathbb{E}[h_{n+1}(X_{n+1}) | X_{0:n}] \right] = \mathbb{E} \left[\left(\prod_{i=0}^n h_i(X_i) \right) Ph_{n+1}(X_n) \right],$$

and since the inner term in the rhs only depends on $X_{0:n}$, we can apply (H_n) and thus,

$$\begin{aligned} \mathbb{E} \left[\prod_{i=0}^{n+1} h_i(X_i) \right] &= \int \cdots \int_{\mathbb{X}^{n+1}} \left[\nu(dx_0) \prod_{i=0}^{n-1} P(x_i, dx_{i+1}) \right] \left(\prod_{i=0}^n h_i(x_i) \right) Ph_{n+1}(x_n) \\ &= \int \cdots \int_{\mathbb{X}^{n+1}} \left[\nu(dx_0) \prod_{i=0}^n P(x_i, dx_{i+1}) \right] \left(\prod_{i=0}^{n+1} h_i(x_i) \right), \end{aligned}$$

showing that the law of $X_{0:n+1}$ is $\nu(dx_0) \prod_{i=0}^n P(x_i, dx_{i+1})$ and (H_{n+1}) is thus proved. \blacksquare

As a consequence, the marginal law of X_n is given by integrating $\nu(dx_0) \prod_{i=0}^{n-1} P(x_i, dx_{i+1})$ over $x_{0:n-1}$ and thus, for all $A \in \mathcal{X}$,

$$\mathbb{P}(X_n \in A) = \int \cdots \int_{\mathcal{X}^{n+1}} \mathbb{1}_A(X_n) \nu(dx_0) \prod_{i=0}^{n-1} P(x_i, dx_{i+1}) = \nu P^n(A),$$

that is, νP^n is the distribution of X_n or in a compact notation, $\boxed{X_n \sim \nu P^n}$.

5.3 Canonical space

So far, the situation is the following: the $\{X_k : k \in \mathbb{N}\}$ are already given and we check the two items in the definition of a Markov chain (Theorem 5.2) with Markov kernel P and initial distribution ν . We now turn to the reverse situation where a couple (ν, P) of initial distribution and Markov kernel are given beforehand and we intend to construct the random variables $\{X_k : k \in \mathbb{N}\}$ on some convenient (common) probability space $(\Omega, \mathcal{F}, \mathbb{P})$ such that $\{X_k : k \in \mathbb{N}\}$ is a Markov chain with Markov kernel P and initial distribution ν .

5.3.1 A simpler problem

We start with an easier problem where we only want to construct $\{X_k : k \in [0 : n-1]\}$ where n is some given positive integer. That is, we only consider a **finite range of integers** such that the first item in Theorem 5.2 is satisfied. For a given $\nu \in \mathbb{M}_1(\mathcal{X})$, define the triplet $(\Omega_n, \mathcal{G}_n, \mathbb{P}_{\nu, n})$ as follows:

- $\Omega_n = \mathcal{X}^{n+1}$, $\mathcal{G}_n = \mathcal{X}^{\otimes(n+1)}$ and $\mathbb{P}_{\nu, n}$ is the probability measure defined on $(\Omega_n, \mathcal{G}_n)$ by

$$\mathcal{G}_n \ni A \mapsto \mathbb{P}_{\nu, n}(A) = \int \cdots \int_{\mathcal{X}^{n+1}} \mathbb{1}_A(\omega_{0:n}) \nu(d\omega_0) \prod_{i=1}^n P(\omega_{i-1}, d\omega_i),$$

and for $\omega \in \Omega_n$, set $X_k(\omega) = \omega_k$, (that $X_k(\omega)$ is the projection of the k -th component of ω).

We aim to show that for all $A \in \mathcal{X}$ and all $k \in [0 : n-1]$, we have $\mathbb{P}_{\nu, n}(X_{k+1} \in A | X_{0:k-1}) = P(X_k, A)$, $\mathbb{P}_{\nu, n}$ -a.s. To do so, write for any $k \in [0 : n-1]$, any non-negative measurable function h on \mathcal{X}^{k+1} and any $A \in \mathcal{X}$

$$\begin{aligned} \mathbb{E}_{\nu, n} [h(X_{0:k}) \mathbb{1}_A(X_{k+1})] &= \int \cdots \int_{\mathcal{X}^{k+2}} h(\omega_{0:k}) \mathbb{1}_A(\omega_{k+1}) \nu(d\omega_0) \prod_{i=1}^{k+1} P(\omega_{i-1}, d\omega_i) \\ &= \int \cdots \int_{\mathcal{X}^{k+2}} h(\omega_{0:k}) P(\omega_k, A) \nu(d\omega_0) \prod_{i=1}^k P(\omega_{i-1}, d\omega_i) \\ &= \mathbb{E}_{\nu, n} [h(X_{0:k}) P(X_k, A)]. \end{aligned}$$

Since h is arbitrary, this is equivalent to saying that $\mathbb{P}_{\nu, n}(X_{k+1} \in A | X_{0:k}) = P(X_k, A)$, \mathbb{P} -a.s.

5.3.2 The general case

We now consider the general case where $k \in \mathbb{N}$ instead of $k \in [0 : n-1]$. Define the coordinate process (X_n) by $X_n(\omega) = \omega_n$ for all $\omega \in \mathcal{X}^{\mathbb{N}}$. We will sometimes use $X_{k:\ell} : \omega \mapsto (\omega_k, \dots, \omega_\ell)$ for $k \leq \ell$ and by extension $X_{k:\infty} : \omega \mapsto (\omega_k, \dots, \omega_\ell, \omega_{\ell+1}, \dots)$. In particular, $X_{0:\infty}(\omega) = I(\omega) = \omega$ where I is the identity function.

Theorem 5.5. (The canonical space) Let (X, \mathcal{X}) be a measurable space and let P be a Markov kernel on $X \times \mathcal{X}$. For every probability measure $\nu \in M_1(X)$, there exists a unique probability measure \mathbb{P}_ν on the canonical space $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}})$ such that, under \mathbb{P}_ν , the coordinate process $\{X_n : n \in \mathbb{N}\}$ is a Markov chain with Markov kernel P and initial distribution ν .

This result is often referred to as the *Ionescu-Tulcea theorem*. Its proof goes far beyond the scope of this course and we will admit it here. Some other, much simpler proofs exist and are based on the Kolmogorov extension theorem, but they hold at the price of additive assumptions on the space (in contrast, we only assume here that (X, \mathcal{X}) is a measurable space, which is quite minimal). In the canonical representation, we therefore set $\Omega = X^\mathbb{N}$, $\mathcal{G} = \mathcal{X}^{\otimes \mathbb{N}}$ and $\mathbb{P} = \mathbb{P}_\nu$. In the particular case where the initial distribution ν is a Dirac mass, we use the compact notation $\boxed{\mathbb{P}_x = \mathbb{P}_{\delta_x}}$. Thus, the theorem allows to define not only one probability measure but a family of probability measures $(\mathbb{P}_\nu)_{\nu \in M_1(X)}$ on the space of trajectories.

What are the relations between the probability measures $(\mathbb{P}_\nu)_{\nu \in M_1(X)}$? A consequence of this theorem is that for all $A \in \mathcal{X}^{\otimes (n+1)}$, $\mathbb{P}_\nu(X_{0:n} \in A) = \int_A \nu(d\omega_0) \prod_{i=1}^n P(\omega_{i-1}, d\omega_i)$. Replacing ν by δ_{x_0} and comparing the two obtained expressions, we get

$$\mathbb{P}_\nu(X_{0:n} \in A) = \int_X \nu(dx_0) \mathbb{P}_{x_0}(X_{0:n} \in A).$$

We can actually extend this result to any $A \in \mathcal{X}^{\otimes \mathbb{N}}$ by replacing the $n+1$ -tuple $X_{0:n}$ by the (infinite) trajectory $X_{0:\infty}$. First note the following equalities: $A = \{\omega \in A\} = \{\omega \in \Omega : X_{0:\infty}(\omega) \in A\} = \{X_{0:\infty} \in A\}$.

We then obtain the following identity: for all $A \in \mathcal{X}^{\otimes \mathbb{N}}$,

$$\mathbb{P}_\nu(A) = \mathbb{P}_\nu(X_{0:\infty} \in A) = \int_X \nu(dx_0) \mathbb{P}_{x_0}(X_{0:\infty} \in A) = \int_X \nu(dx_0) \mathbb{P}_{x_0}(A) \quad (5.2)$$

5.3.3 The Markov property.

Define the shift operator S by $S : X^\mathbb{N} \ni \omega \mapsto \omega' \in X^\mathbb{N}$ where $\omega = (\omega_i)_{i \in \mathbb{N}}$ and $\omega' = (\omega_{i+1})_{i \in \mathbb{N}}$

Theorem 5.6. (The Markov property) For any $\nu \in M_1(X)$, any non-negative or bounded function h on $X^\mathbb{N}$ and any $n \in \mathbb{N}$,

$$\mathbb{E}_\nu \left[h \circ S^k | \mathcal{F}_k \right] = \mathbb{E}_{X_k} [h], \quad \mathbb{P}_\nu - a.s. \quad (5.3)$$

where $\mathcal{F}_k = \sigma(X_{0:k})$.

The expression of the Markov property, as it stands, may seem a bit cryptic. Recalling the definition of $X_{k:\infty}$, we have $\mathbb{P}_\nu - a.s.$,

$$\begin{aligned} \mathbb{E}_\nu [h(X_{k:\infty}) | \mathcal{F}_k] &= \mathbb{E}_\nu \left[h \circ S^k \circ X_{0:\infty} | \mathcal{F}_k \right] = \mathbb{E}_\nu \left[h \circ S^k \circ I | \mathcal{F}_k \right] \\ &= \mathbb{E}_\nu \left[h \circ S^k | \mathcal{F}_k \right] = \mathbb{E}_{X_k} [h] = \mathbb{E}_{X_k} [h \circ I] \\ &= \mathbb{E}_{X_k} [h(X_{0:\infty})]. \end{aligned}$$

Therefore, for any $\nu \in M_1(X)$,

$$\mathbb{E}_\nu [h(X_{k:\infty}) | \mathcal{F}_k] = \mathbb{E}_{X_k} [h(X_{0:\infty})], \quad \mathbb{P}_\nu - a.s.$$

is another **equivalent expression of the Markov property**. This expression may seem easier to deal with but the reader has to be at ease with both formulations. A stronger version of the Markov property exists and is called (as expected) the strong Markov property: its statement is (5.3) with the exception that k is replaced by a stopping time τ and that the identity only holds on the event $\{\tau < \infty\}$. The strong Markov property is extremely important in Markov Chain theory but, quite surprisingly, it is not needed in this basic course.

5.4 At this point...

- a) I can write perfectly expression of Markov kernels if it is asked.
- b) I understand the different notation $\mu PQf...$
- c) I perfectly understand the canonical space.
- d) I can understand (5.2).
- e) I perfectly understand the Markov property (5.3).

5.5 Metropolis-Hastings algorithms

5.5.1 Invariant probability measures: existence

Definition 5.7. We say that $\pi \in M_1(X)$ is an invariant probability measure for the Markov kernel P on $X \times \mathcal{X}$ if $\pi P = \pi$.

In words, if (X_k) is a Markov chain with Markov kernel P and assuming that $X_0 \sim \pi$, then for all $k \geq 1$, we have $X_k \sim \pi$. (This is due to the fact that applying P^k on both sides of $\pi P = \pi$ shows that $\pi P^{k+1} = \pi P^k$ and therefore, for all $k \in \mathbb{N}$, $\pi P^k = \pi$). This result on the (marginal) distribution of X_k may be extended to n -tuples.

More precisely, it can be readily checked that if π is an *invariant probability measure* for P , then the sequence of random variables $\{X_k : k \in \mathbb{N}\}$ is a *strongly stationary sequence* under \mathbb{P}_π (in the sense that for all $n, p \in \mathbb{N}^*$, and all n -tuple $k_{1:n}$, the random vector $(X_{k_1}, \dots, X_{k_n})$ follows the same distribution as $(X_{k_1+p}, \dots, X_{k_n+p})$).

We now introduce the notion of reversibility for a Markov kernel. This will be of crucial importance for designing Markov kernels with a given invariant probability measure.

Definition 5.8. Let $\pi \in M_1(X)$ and P be a Markov kernel on $X \times \mathcal{X}$. We say that P is π -reversible if and only if (with infinitesimal notation)

$$\pi(dx)P(x, dy) = \pi(dy)P(y, dx), \quad (5.4)$$

that is, for all measurable bounded or non-negative functions h on $(X^2, \mathcal{X}^{\otimes 2})$,

$$\iint_{\mathcal{X}^2} h(x, y) \pi(\mathrm{d}x) P(x, \mathrm{d}y) = \iint_{\mathcal{X}^2} h(x, y) \pi(\mathrm{d}y) P(y, \mathrm{d}x). \quad (5.5)$$

In words, a Markov kernel P is π -reversible if and only if the probability measure $\pi(\mathrm{d}x)P(x, \mathrm{d}y)$ is symmetric with respect to (x, y) .

Proposition 5.9 *Let P be a Markov kernel on $\mathcal{X} \times \mathcal{X}$. Let $\pi \in \mathcal{M}_1(\mathcal{X})$ such that P is π -reversible, then the Markov kernel P is π -invariant.*

PROOF. For any $A \in \mathcal{X}$, we have by the reversibility relation

$$\pi P(A) = \iint_{\mathcal{X}^2} \mathbb{1}_A(y) \pi(\mathrm{d}x) P(x, \mathrm{d}y) = \iint_{\mathcal{X}^2} \mathbb{1}_A(y) \pi(\mathrm{d}y) P(y, \mathrm{d}x) = \int_A \pi(\mathrm{d}y) \underbrace{P(y, \mathcal{X})}_1 = \pi(A),$$

which finishes the proof. ■

What are the consequences? If you want to check easily that a kernel P is π -invariant, it is sufficient to check that it is π -reversible.

5.5.2 Metropolis-Hastings (MH) algorithms

In this section, we are given a probability measure $\pi \in \mathcal{M}_1(\mathcal{X})$ and the idea now is to construct a Markov chain $\{X_k : k \in \mathbb{N}\}$ admitting π as invariant probability measure, in which case we say that π is a target distribution. In other words, we try to find a Markov kernel P on $\mathcal{X} \times \mathcal{X}$ such that P is π -invariant. The reason for that is that an invariant probability measure will be a good candidate for the “limiting” distribution of $\{X_k : k \in \mathbb{N}\}$ (in some sense to be defined) and this in turn, will allow us to provide an approximation $\pi(h)$:

$$\pi(h) = \int_{\mathcal{X}} h(x) \pi(\mathrm{d}x) \approx n^{-1} \sum_{k=0}^{n-1} h(X_k).$$

5.5.2.1 Construction of the kernel

For simplicity we now assume that π has a density with respect to some dominating σ -finite measure λ and by abuse of notation, we also denote by π this density, that is we write $\pi(\mathrm{d}x) = \pi(x)\lambda(\mathrm{d}x)$ and we assume that this density π is **positive**.

Moreover, let Q be Markov kernel on $\mathcal{X} \times \mathcal{X}$ such that $Q(x, \mathrm{d}y) = q(x, y)\lambda(\mathrm{d}y)$, that is, for any $x \in \mathcal{X}$, $Q(x, \cdot)$ is also dominated by λ and denoting by $q(x, \cdot)$ this density, we assume for simplicity that $q(x, y)$ is **positive** for all $x, y \in \mathcal{X}$. At this stage, there is almost no link between Q and the target distribution π .

For a given function $\alpha : \mathcal{X}^2 \rightarrow [0, 1]$, consider the following Algorithm 4.

In words, Q allows to propose a candidate for the next value of the Markov chain (X_k) and this candidate will be accepted or refused according to a probability that depends on the function α .

We will now choose conveniently α in such a way that (X_k) is a Markov chain with invariant probability measure π . The concept of π -reversibility will help us. To do so, let us assume that for all $x, y \in \mathcal{X}$,

$$\pi(x)\alpha(x, y)q(x, y) = \pi(y)\alpha(y, x)q(y, x), \quad (5.6)$$

and let us show that it implies that the Markov kernel P associated to (X_k) is π -reversible.

Algorithm 4 The Metropolis-Hastings algorithm

```

1: Input:  $n$ 
2: Output:  $X_0, \dots, X_n$ 
3: At  $t = 0$ , draw  $X_0$  according to some arbitrary distribution
4: for  $t \leftarrow 0, \dots, n-1$  do
5:   Draw independently  $Y_{t+1} \sim Q(X_t, \cdot)$  and  $U_{t+1} \sim \text{Unif}(0, 1)$ 
6:   Set  $X_{t+1} = \begin{cases} Y_{t+1} & \text{if } U_{t+1} \leq \alpha(X_t, Y_{t+1}) \\ X_t & \text{otherwise} \end{cases}$ 
7: end for

```

First, we write down the Markov kernel associated to (X_k) : in passing, this is an excellent opportunity to check if we are able to express explicitly a Markov kernel by analyzing conveniently the update transition. Denote $\mathcal{F}_t = \sigma(X_0, U_{1:t}, Y_{1:t})$ and note that (X_t) is adapted to the filtration (\mathcal{F}_t) (which is equivalent to $\sigma(X_{0:t}) \subset \mathcal{F}_t$). Then, setting $\bar{\alpha}(x) = 1 - \int_{\mathbb{X}} Q(x, dy) \alpha(x, y)$, we have for any bounded or non-negative measurable function h on \mathbb{X} and any $t \in \mathbb{N}$,

$$\begin{aligned}
\mathbb{E}[h(X_{t+1}) | \mathcal{F}_t] &= \mathbb{E}[\mathbb{1}_{\{U_{t+1} < \alpha(X_t, Y_{t+1})\}} h(Y_{t+1}) | \mathcal{F}_t] + \mathbb{E}[\mathbb{1}_{\{U_{t+1} \geq \alpha(X_t, Y_{t+1})\}} h(X_t) | \mathcal{F}_t] \\
&= \int_{\mathbb{X}} Q(X_t, dy) \alpha(X_t, y) h(y) + \bar{\alpha}(X_t) h(X_t) \\
&= \int_{\mathbb{X}} \underbrace{Q(X_t, dy) \alpha(X_t, y) + \bar{\alpha}(X_t) \delta_{X_t}(dy)}_{P_{(\pi, Q)}^{MH}(X_t, dy)} h(y) = P_{(\pi, Q)}^{MH} h(X_t).
\end{aligned}$$

Therefore, $\{X_t : t \in \mathbb{N}\}$ is a Markov chain with Markov kernel

$$P_{(\pi, Q)}^{MH}(x, dy) = Q(x, dy) \alpha(x, y) + \bar{\alpha}(x) \delta_x(dy) \quad (5.7)$$

Lemma 5.10 *The Markov kernel $P_{(\pi, Q)}^{MH}$ is π -reversible if and only if*

$$\pi(dx) Q(x, dy) \alpha(x, y) = \pi(dy) Q(y, dx) \alpha(y, x). \quad (5.8)$$

In this literature, (5.8) is often called the **detailed balance condition**. **PROOF.** First, note that

$$\pi(dx) \bar{\alpha}(x) \delta_x(dy) = \pi(dy) \bar{\alpha}(y) \delta_y(dx) \quad (5.9)$$

Indeed, for any measurable function h on \mathbb{X}^2 , we have

$$\begin{aligned}
\iint_{\mathbb{X}^2} h(x, y) \pi(dx) \bar{\alpha}(x) \delta_x(dy) &= \int_{\mathbb{X}} h(x, x) \pi(dx) \bar{\alpha}(x) \\
&= \int_{\mathbb{X}} h(y, y) \pi(dy) \bar{\alpha}(y) = \iint_{\mathbb{X}^2} h(x, y) \pi(dy) \bar{\alpha}(y) \delta_y(dx).
\end{aligned}$$

Combining (5.7) with (5.9), we obtain that $P_{(\pi, Q)}^{MH}$ is π -reversible if and only if the detailed balance condition (5.8) is satisfied. This completes the proof. \blacksquare

5.5.2.2 Acceptance probability

We now make use of Lemma 5.10 in order to find an explicit expression of the acceptance probability α . We have the following lemma.

Lemma 5.11 Denote $\alpha^{MH}(x, y) = \min\left(\frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}, 1\right)$ and $\alpha^b(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y) + \pi(y)q(y, x)}$. Then, any $\alpha \in \{\alpha^{MH}, \alpha^b\}$ satisfies the detailed balance condition (5.8). Moreover, any other $\alpha \in [0, 1]$ that satisfies the detailed balance condition is dominated by α^{MH} in the sense that: for $\lambda^{\otimes 2}$ -almost all $x, y \in \mathcal{X}$,

$$\alpha(x, y) \leq \alpha^{MH}(x, y) \quad (5.10)$$

PROOF. The fact that any $\alpha \in \{\alpha^{MH}, \alpha^b\}$ satisfies $\pi(x)q(x, y)\alpha(x, y) = \pi(y)q(y, x)\alpha(y, x)$ for $\lambda^{\otimes 2}$ -almost all $x, y \in \mathcal{X}$, is immediate, by replacing α by its expression. It remains to check (5.10). Assume now that for $\lambda^{\otimes 2}$ -almost all $x, y \in \mathcal{X}$,

$$\pi(x)q(x, y)\alpha(x, y) = \pi(y)q(y, x)\alpha(y, x).$$

then, using that $\alpha(y, x) \leq 1$ shows that $\alpha(x, y) \leq \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}$. Moreover, $\alpha(x, y) \leq 1$ and this finally implies

$$\alpha(x, y) \leq \min\left(\frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}, 1\right) = \alpha^{MH}(x, y),$$

which completes the proof. ■

According to (5.10), α^{MH} is actually the highest acceptance probability among the acceptance probabilities such that $P_{(\pi, Q)}^{MH}$ is π -reversible and therefore, this acceptance is widely used in practice (in the sense that we expect that a Markov kernel that accepts often, explores the space more rapidly and therefore is preferable to another one with less acceptance probability). In what follows, unless otherwise stated, we *implicitly assume that the Markov kernel $P_{(\pi, Q)}^{MH}$ is associated to the acceptance probability α^{MH} .*

Example 5.12 (The independence sampler) If the proposition kernel is $Q(x, dy) = q(y)\lambda(dy)$ where q is a density wrt λ on \mathcal{X} , then at each time step, the proposed candidate is drawn irrespective of the current value of the Markov chain (this is because, $Q(x, dy)$ does not depend on x), that is, in the step 2(a) of Algorithm 4, we draw $Y_{t+1} \sim q(\cdot)$. In such case, the acceptance probability is $\alpha(x, y) = \min\left(\frac{\pi(y)q(x)}{\pi(x)q(y)}, 1\right)$ and the Metropolis-Hastings algorithm is called the Independence Sampler.

Another important example is the following.

Example 5.13 (The random walk MH sampler) If $\mathcal{X} = \mathbb{R}^p$ and if the proposition kernel is $Q(x, dy) = q(y - x)\lambda(dy)$ where q is a symmetric density wrt λ on \mathcal{X} , (by symmetric, we mean that $q(u) = q(-u)$ for all $u \in \mathcal{X}$) then at each time step in Algorithm 4, we draw a candidate $Y_{t+1} \sim q(y - X_t)\lambda(dy)$. In such case, the acceptance probability is $\alpha(x, y) = \min\left(\frac{\pi(y)}{\pi(x)}, 1\right)$ and the associated algorithm is called the (symmetric) Random Walk Metropolis-Hasting. Another way of writing the proposition update is $Y_{t+1} = X_t + \eta_k$ where $\eta_k \sim q(\cdot)$.

5.6 Invariant probability measure: uniqueness

We start with a very simple lemma that will be useful for finding sufficient conditions for uniqueness.

Lemma 5.14 If P admits two distinct invariant probability measures, it also admits distinct invariant probability measures π_0 and π_1 that are mutually singular, i.e., such that there exists $A \in \mathcal{X}$ such that $\pi_0(A) = \pi_1(A^c) = 0$.

PROOF. Let ζ_0, ζ_1 be two distinct invariant probability measures for P . Both have densities with respect to some common dominating measure (for example, taking $\zeta = \zeta_0 + \zeta_1$, we have that ζ dominates both ζ_0 and ζ_1 , which can be seen from

the implication $\zeta(A) = 0 \Rightarrow (\zeta_1(A) = 0 \text{ and } \zeta_2(A) = 0)$ for any $A \in \mathcal{X}$ and according to the Radon Nikodym theorem, if a measure dominates another one, the latter has a density with respect to the former). Write then $\zeta_0(dx) = f_0(x)\zeta(dx)$ and $\zeta_1(dx) = f_1(x)\zeta(dx)$ where f_0, f_1 are non-negative measurable functions on \mathbf{X} . Define the positive part $(\zeta_1 - \zeta_0)^+$ and the negative part $(\zeta_1 - \zeta_0)^-$ of the signed measure $\zeta_1 - \zeta_0$ by $(\zeta_1 - \zeta_0)^+(dx) = [f_1(x) - f_0(x)]^+\zeta(dx)$ and $(\zeta_1 - \zeta_0)^-(dx) = [f_1(x) - f_0(x)]^-\zeta(dx)$. Then,

$$\begin{aligned} (\zeta_1 - \zeta_0)^+ P \mathbb{1}_A &= \int_{\mathbf{X}} \zeta(dx) [f_1(x) - f_0(x)]^+ P(x, A) \\ &\geq \int_{\mathbf{X}} \zeta(dx) [f_1(x) - f_0(x)] P(x, A) \\ &\geq \zeta_1 P(A) - \zeta_0 P(A) = \zeta_1(A) - \zeta_0(A). \end{aligned}$$

Therefore, $(\zeta_1 - \zeta_0)^+ P$ is a (non-negative) measure that is greater than the signed measure $\zeta_1 - \zeta_0$. Since the positive part $(\zeta_1 - \zeta_0)^+$ is also the smallest (non-negative) measure that is greater than $\zeta_1 - \zeta_0$, we conclude that $(\zeta_1 - \zeta_0)^+ \leq (\zeta_1 - \zeta_0)^+ P$. The measure $(\zeta_1 - \zeta_0)^+ P - (\zeta_1 - \zeta_0)^+$ is therefore non-negative and we have

$$[(\zeta_1 - \zeta_0)^+ P - (\zeta_1 - \zeta_0)^+](\mathbf{X}) = \int_{\mathbf{X}} (\zeta_1 - \zeta_0)^+(dx) \underbrace{P(x, \mathbf{X})}_1 - (\zeta_1 - \zeta_0)^+(\mathbf{X}) = 0.$$

Finally, $(\zeta_1 - \zeta_0)^+ = (\zeta_1 - \zeta_0)^+ P$. The probability measure $\pi_0 = \frac{(\zeta_1 - \zeta_0)^+}{(\zeta_1 - \zeta_0)^+(\mathbf{X})}$ is thus an invariant probability measure for P . Replacing $(\zeta_1 - \zeta_0)^+$ by $(\zeta_1 - \zeta_0)^-$, we obtain in the same way that $\pi_1 = \frac{(\zeta_1 - \zeta_0)^-}{(\zeta_1 - \zeta_0)^-(\mathbf{X})}$ is an invariant probability measure. We can easily check that taking $A = \{f_0 \geq f_1\}$, we have $\pi_0(A) = \pi_1(A^c) = 0$, showing that these probability measures are mutually singular. ■

To be rigorous, in the course of the proof, we actually need some results on the positive and negative part of a signed-measure. The interested reader may work on the following exercise to fully understand the previous proof:

Exercise 5.15. Define $M_s(\mathbf{X})$ the set of signed-measures. Let $\mu \in M_s(\mathbf{X})$ and assume that $\mu \preceq \zeta$ where $\zeta \in M_+(\mathbf{X})$ (in the sense that we have the implication: if for some $A \in \mathcal{X}$, $\zeta(A) = 0$, then $\mu(A) = 0$). According to the Radon-Nikodym theorem, there exists a measurable function h such that $\mu(dx) = h(x)\zeta(dx)$. Define $\mu^+(dx) = |h(x)|\zeta(dx)$.

1. Show that the measure μ^+ is well-defined (in the sense that the measure $|h(x)|\zeta(dx)$ does not depend on the measure ζ , provided that the ζ dominates μ .)
2. Show that for ζ -almost all $x \in \mathbf{X}$, $|h(x)| \leq 1$.
3. Assume that there exists $\nu \in M_+(\mathbf{X})$ such that for all $A \in \mathcal{X}$, we have $\mu(A) \leq \nu(A)$. Show that $\mu^+(A) \leq \nu(A)$ for all $A \in \mathcal{X}$.

We now make use of Lemma 5.14 in order to give a sufficient condition for uniqueness.

Proposition 5.16 Assume that there exists a non-null measure $\mu \in M_+(\mathbf{X})$ satisfying the following property:

- For all $A \in \mathcal{X}$ such that $\mu(A) > 0$ and for all $x \in \mathbf{X}$, there exists $n \in \mathbb{N}$ such that $P^n(x, A) > 0$.

Then, P admits at most one invariant probability measure.

If the assumption of Proposition 5.16 holds, we say that P is μ -irreducible and in such case, μ is called an **irreducibility measure** for P . **PROOF.** The proof is by contradiction. Assume that there exists two distinct invariant probability measures. According to Lemma 5.14, we can consider two invariant probability measures π_1 and π_2 that are mutually singular. Under the assumptions of the Proposition, let $A \in \mathcal{X}$ such that $\mu(A) > 0$. Then, for any $i \in \{1, 2\}$, we have

$$0 < \int_{\mathbf{X}} \pi_i(dx) \underbrace{\sum_{n=0}^{\infty} P^n(x, A)}_{>0} = \sum_{n=0}^{\infty} \pi_i P^n(A) = \sum_{n=0}^{\infty} \pi_i(A),$$

which in turn implies that $\pi_i(A) > 0$. The contraposed implication gives that if for some $i \in \{1, 2\}$, $\pi_i(A) = 0$, then $\mu(A) = 0$. Now, since $\{\pi_i : i \in \{1, 2\}\}$ are mutually singular, there exists $A \in \mathcal{X}$ such that $\pi_1(A) = \pi_2(A^c) = 0$ and this shows that $\mu(A) = \mu(A^c) = 0$ which is impossible. ■

5.6.1 Application to Metropolis-Hastings algorithms.

We have already seen that $P_{(\pi, Q)}^{MH}$ is π -invariant and we have assumed that $Q(x, dy) = q(x, y)\lambda(dy)$ and $\pi(dy) = \pi(y)\lambda(dy)$ and for simplicity, we said that for all $x, y \in X$, $q(x, y) > 0$ and $\pi(y) > 0$. This in turn implies that $\alpha(x, y) = \alpha^{MH}(x, y) > 0$ and therefore if $\lambda(A) > 0$, then for all $x \in X$,

$$P(x, A) \geq \int_A \underbrace{q(x, y)\alpha(x, y)}_{>0} \lambda(dy) > 0.$$

This shows that P is λ -irreducible and therefore π is the unique invariant probability measure for P .

5.7 After studying this chapter...

- a) I can understand the link between reversibility and invariant measure.
- b) If asked, I can check that MH chains are reversible.
- c) I understand the form of the acceptance probability.
- d) Independence sampling and random walk MH have no secrets for me and I am able to implement them.
- e) If asked, I can show that an invariant measure is unique by using Proposition 5.16
- f) I find that Metropolis-Hasting algorithms are magical.

Chapter 6

Variance reduction techniques

Contents

| | | |
|------------|--|-----------|
| 6.1 | Importance Sampling | 39 |
| 6.2 | Antithetic variates | 40 |
| 6.3 | Control Variates | 41 |
| 6.4 | Conditioning | 43 |
| 6.5 | Stratified sampling | 44 |
| 6.6 | Quasi Monte Carlo methods | 45 |
| | 6.6.1 Weak discrepancy sequences | 45 |
| 6.7 | Take-home message | 47 |
| 6.8 | Highlights | 47 |

Keywords: *Importance sampling, antithetic variates, stratification, control variates, conditioning, weak discrepancy, Van der Corput sequences, Halton sequences.*

We have seen in the previous chapters some tools for sampling exactly or approximately according to a target distribution defined by a density f or by the solution of a SDE (in particular, we have seen discretization schemes for some SDE), the final goal remaining to give the most precise numerical value of a quantity of the type $\mathbb{E}[h(Y)] = \int h(y)f(y)dy$ where Y has the density f .

This chapter is devoted to various methods that allow to produce estimators \bar{S}_n of $\mathbb{E}[h(Y)]$ such that the quadratic deviation for a fixed n is the smallest possible. All the estimators we will see here will be unbiased and strongly convergent: $\mathbb{E}[\bar{S}_n] = \mathbb{E}[h(Y)]$ et $\bar{S}_n \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}[h(Y)]$, so that the standard deviation writes

$$\mathbb{E} \left[\bar{S}_n - \underbrace{\mathbb{E}[h(Y)]}_{\mathbb{E}[\bar{S}_n]} \right]^2$$

which is also the variance $\text{Var}(\bar{S}_n)$. Let us now look at some sampling methods for reducing the variance compared to a standard Monte Carlo estimator.

6.1 Importance Sampling

We refer the reader to Section (4.2.1) where the importance sampling is introduced and where it is shown that the asymptotic variance of the importance sampling estimator did not necessarily reach its optimal variance when the instrumental distribution is the target distribution. Let us recall the method. Consider a r.v. $Y \sim f$ where f is a density. Suppose we want to approximate $\mathbb{E}[h(Y)]$. An (instrumental) density g is then selected. And we draw an n -sample (X_1, \dots, X_n) according to the distribution of density g and we set

$$\bar{S}_n = \frac{1}{n} \sum_{i=1}^n \frac{f(X_i)}{g(X_i)} h(X_i)$$

On the contrary to the chapter 4, we will focus on the non-asymptotic variance (instead of the asymptotic variance) of \bar{S}_n .

$$\text{Var}(\bar{S}_n) = \frac{\text{Var}\left(\frac{f(X)}{g(X)}h(X)\right)}{n} = \frac{\int f(x)\frac{f(x)}{g(x)}h^2(x)dx - (\int h(x)f(x)dx)^2}{n}$$

We thus obtain a reduction of the variance in comparison with a standard Monte Carlo estimator iff

$$\int f(x)\frac{f(x)}{g(x)}h^2(x)dx \leq \int f(x)h^2(x)$$

6.2 Antithetic variates

Suppose we wish to evaluate $\mathbb{E}(h(U))$ where $U \sim \mathcal{U}[0, 1]$ by a Monte Carlo method. If we have an n -sample (U_1, \dots, U_n) of i.i.d. r.v. according to the distribution of U , a classical method would be to consider

$$S_n = \frac{1}{n} \sum_{i=1}^n h(U_i).$$

Nevertheless, one can note that if $U \sim \mathcal{U}[0, 1]$ then $1 - U$ also follows the distribution $\mathcal{U}[0, 1]$. It is then tempting to consider

$$S'_n = \frac{1}{n} \sum_{i=1}^n h(1 - U_i).$$

or

$$\bar{S}_{2n} = \frac{1}{2n} \sum_{i=1}^n (h(U_i) + h(1 - U_i)) \quad (6.1)$$

Clearly, \bar{S}_{2n} is an *unbiased* and *strongly convergent* estimator of $\mathbb{E}(h(U))$. Let us calculate now the variance of \bar{S}_{2n} :

$$\text{Var}(\bar{S}_{2n}) = \frac{1}{4n} (2\text{Var}(h(U)) + 2\text{Cov}(h(U), h(1 - U))) = \text{Var}(S_{2n}) + \frac{1}{2n} \text{Cov}(h(U), h(1 - U))$$

If $\text{Cov}(h(U), h(1 - U))$ is negative, then the speed of convergence of \bar{S}_{2n} is better than the one of S_{2n} . The following lemma gives a sufficient condition that ensures the negativity of the covariance of two functions of the same random variable.

Lemma 6.1 CORRELATION INEQUALITY- *Let X be a r.v and h, g be two monotone functions such that one is decreasing and the other one increasing, then $\text{Cov}(h(X), g(X)) \leq 0$.*

PROOF. Assume first that $\mathbb{E}(g(X)) = 0$ and $\mathbb{E}(h(X)) = 0$. Let x and y be two real numbers. We have

$$(h(x) - h(y))(g(x) - g(y)) \leq 0.$$

If X and Y are independent, with the same distribution, then

$$\mathbb{E}[(h(X) - h(Y))(g(X) - g(Y))] \leq 0.$$

Expanding this quantity, we obtain

$$\mathbb{E}[h(X)g(X)] + \mathbb{E}[h(Y)g(Y)] - \mathbb{E}[h(X)]\mathbb{E}[g(Y)] - \mathbb{E}[h(Y)]\mathbb{E}[g(X)] \leq 0$$

which is the desired result by noting that $X \stackrel{\mathcal{L}}{\equiv} Y$. The proof is completed for $\mathbb{E}(g(X)) = \mathbb{E}(h(X)) = 0$. For the general case, we set $\bar{g}(x) = g(x) - \mathbb{E}(g(X))$ and $\bar{h}(x) = h(x) - \mathbb{E}(h(X))$ and we apply the previous result. ■

Corollary 6.2 *If $U \sim \mathcal{U}[0, 1]$ and if h is nondecreasing then, $\text{Cov}(h(U), h(1-U)) \leq 0$.*

PROOF. Take $g(x) = h(1-x)$ in the previous lemma. ■

Example 6.3 *Let $\psi(x) = (\lambda e^{\sigma x} - K)^+$. We want to approximate $\mathbb{E}[\psi(G)]$ or $G \sim \mathcal{N}(0, 1)$ by a Monte Carlo method (of course, we know that Black and Scholes provides an exact closed-form formula, without any approximation). In this case, noting that G and $-G$ have the same law, one can choose an estimator of the form*

$$\bar{S}_n = \frac{1}{2n} \sum_{i=1}^n (\psi(G_i) + \psi(-G_i))$$

where $G_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. Clearly the proof of $\text{Cov}(\psi(G), \psi(-G)) \leq 0$ is a direct consequence of the Lemma 6.1 with $f(x) = \psi(x)$ and $g(x) = \psi(-x)$.

More generally, we can use a Monte Carlo method with antithetic variate to approximate $\mathbb{E}[h(Y)]$ if $h: \mathbb{R} \rightarrow \mathbb{R}$ is monotone and if there is a decreasing function ϕ such as $h(Y) \stackrel{\mathcal{L}}{\equiv} h \circ \phi(Y)$.

Exercise 6.4. GENERALISATION- Let Y be a r.v. with values in \mathbb{R}^d . Let $h: \mathbb{R}^d \rightarrow \mathbb{R}$ and $\tilde{h}: \mathbb{R}^d \rightarrow \mathbb{R}$ such that h is increasing wrt to any of its coordinates and \tilde{h} is decreasing wrt to any of its coordinates. Assume furthermore that

$$h(Y) \stackrel{\mathcal{L}}{\equiv} \tilde{h}(Y).$$

Show that under some integrability assumptions (that should be explicit), the quantity $\frac{1}{2n} \sum_{i=1}^n [h(Y_i) + \tilde{h}(Y_i)]$ is an unbiased and strongly convergent estimator, with a lower variance than the one of the standard Monte Carlo estimator $\frac{1}{n} \sum_{i=1}^n h(Y_i)$.

6.3 Control Variates

In order to reduce the variance of the Monte Carlo method for the calculation of $\mathbb{E}(Y)$, another approach consists in finding a r.v X , such that its expectation **is available in closed-form** and such that $\boxed{\text{Var}(Y - \alpha X) \leq \text{Var}(Y)}$ for some $\alpha \in \mathbb{R}$.

We then consider (Y_1, \dots, Y_n) an n -sample according to the distribution of Y and (X_1, \dots, X_n) an n -sample according to the distribution of X . We then set

$$\bar{S}_n = \frac{1}{n} \sum_{i=1}^n [Y_i - \alpha(X_i - \mathbb{E}[X])] = \frac{1}{n} \sum_{i=1}^n Y_i - \alpha \left(\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X] \right)$$

Of course, \bar{S}_n is an *unbiased* and *strongly convergent* estimator of $\mathbb{E}(Y)$. Indeed by the LLN, we have $\bar{S}_n \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}(Y)$. Now, setting $S_n = \frac{1}{n} \sum_{i=1}^n Y_i$

$$\text{Var}(\bar{S}_n) = \frac{\text{Var}[Y - \alpha(X - \mathbb{E}[X])]}{n} = \frac{\text{Var}[Y - \alpha X]}{n} \leq \frac{\text{Var}[Y]}{n} = \text{Var}(S_n)$$

which shows that the variance is smaller than the one obtained by a standard Monte Carlo method. We then say that X is a *control variate* for Y .

Unfortunately, there is no general method for creating a control variate X starting from Y . This is on a case-by-case basis. However, to obtain a "good" control variate X , we must keep in mind that $Y - \alpha X$ should "vary" less than Y alone and so, X must "follow" more or less the evolution of Y - replacing if necessary X by $-X$, (if Y is large, also should be X and if Y small, also should be X) while X still has the advantage over Y to have an expectation which can be explicitly calculated.

6.3.0.1 Optimisation of α

Once the control variate X chosen, the problem is then to select the most efficient coefficient α . We have the following lemma which turns out to be an orthogonal projection result:

Lemma 6.5 *Let Y and X be two square integrable r.v integrable, then $\text{Var}(Y - \alpha^*X) = \inf_{\alpha \in \mathbb{R}} \text{Var}(Y - \alpha X)$ where we have set*

$$\alpha^* = \frac{\text{Cov}(Y, X)}{\text{Var}(X)}$$

PROOF. Note that, using the notation $\|\cdot\| = \sqrt{\langle \cdot, \cdot \rangle}$ where $\langle U, V \rangle = \mathbb{E}(UV)$ is the usual scalar product on L^2 random variables,

$$\text{Var}(Y - \alpha X) = \text{Var}[Y - \mathbb{E}(Y) - \alpha(X - \mathbb{E}[X])] = \|Y - \mathbb{E}(Y) - \alpha(X - \mathbb{E}[X])\|^2$$

Therefore, the optimum coefficient $\alpha = \alpha^*$ is obtained by searching in the space $\mathcal{H} = \{\alpha[X - \mathbb{E}[X]]; \alpha \in \mathbb{R}\}$ the closest element from $Y - \mathbb{E}(Y)$ for the norm associated to a scalar product: this is exactly the orthogonal projection of the vector $Y - \mathbb{E}(Y)$ on the subspace \mathcal{H} . We then obtain (by the classical formula for the projection on an orthogonal basis, knowing that here, the basis consists in only one vector!)

$$\mathcal{P}_{\perp}^{Y - \mathbb{E}(Y)}[\mathcal{H}] = \langle Y - \mathbb{E}(Y), \frac{X - \mathbb{E}[X]}{\|X - \mathbb{E}[X]\|} \rangle \frac{X - \mathbb{E}[X]}{\|X - \mathbb{E}[X]\|} = \frac{\langle Y - \mathbb{E}(Y), X - \mathbb{E}[X] \rangle}{\langle X - \mathbb{E}[X], X - \mathbb{E}[X] \rangle} (X - \mathbb{E}[X]) = \frac{\text{Cov}(Y, X)}{\text{Var}(X)} (X - \mathbb{E}[X])$$

hence the value of α^* . ■

From this lemma, the estimator issued from the control variate technique writes

$$\bar{s}_n = \frac{1}{n} \sum_{i=1}^n Y_i + \frac{\text{Cov}(Y, X)}{\text{Var}(X)} \left(\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X] \right)$$

For this estimator to be really used in practise, it is therefore necessary to know explicitly not only $\mathbb{E}[X]$ but also $\text{Var}(X)$ and $\text{Cov}(Y, X)$; and the latter is not known in general since $\mathbb{E}(Y)$ is not known itself (this is precisely what we have been looking for since the beginning of this course). We are therefore forced to estimate this covariance, taking care that this additional estimation error does not harm the decrease in variance.

Remark 6.6 *We can see the Monte Carlo methods with antithetic variates as some particular cases of control variates technique. To see this, let us take again the example we used for antithetic variates (6.1):*

$$\frac{1}{n} \sum_{i=1}^n \frac{h(U_i) + h(1 - U_i)}{2} = \frac{1}{n} \sum_{i=1}^n h(U_i) - \underbrace{\frac{h(U_i) - h(1 - U_i)}{2}}_{X_i}$$

where $\mathbb{E}(X_i) = 0$, which is thus known exactly and this situation finally falls down into the control variates technique.

Exercise 6.7. THE CALL-PUT PARITY- In a simplified manner, the payoff of a selling option (a Call) can be written as $(\lambda e^{\sigma G} - K)_+$ where $G \sim \mathcal{N}(0, 1)$. Noting that $(\lambda e^{\sigma G} - K)_+ - (K - \lambda e^{\sigma G})_+ = \lambda e^{\sigma G} - K$, propose a control variate technique for the calculation of $\mathbb{E}(\lambda e^{\sigma G} - K)_+$.

Exercise 6.8. KEMNA AND VORST TECHNIQUE FOR ASIAN OPTIONS.- We want to approximate the price of an asian put in a Black and Scholes model where the price of the risky asset is given by $S_t = xe^{(r-\frac{\sigma^2}{2})t+\sigma W_t}$. The payoff of the asian put under the risk neutral probability is equal to

$$P = \mathbb{E} \left[e^{-rT} \left(K - \frac{1}{T} \int_0^T S_u du \right)_+ \right]$$

Estimate this expectation using a Monte Carlo method with control variates: $X = e^{-rT} (K - \exp(Z))_+$ where $Z = \frac{1}{T} \int_0^T \ln S_u du$. Explain the method and calculate $\mathbb{E}[X]$.

Exercise 6.9. BASKET OPTIONS- Consider d assets such that the price at time T writes in the following way (for $i = 1, \dots, d$),

$$S_T^i = x_i \exp \left[\left(\underbrace{r - \frac{1}{2} \sum_{j=1}^p \sigma_{ij}^2}_{\beta_i} \right) T + \sum_{j=1}^p \sigma_{ij} W_T^j \right]$$

We want approximate the price of a basket put: $\mathbb{E}(K - Y)_+$ where $Y = \sum_{i=1}^d a_i S_T^i$ (with $a_i \geq 0$ and $\sum_{i=1}^d a_i = 1$). For this, we set $m = \sum_{i=1}^d a_i x_i$ and we approximate Y/m with $X = \exp \left[\sum_{i=1}^d \frac{a_i x_i}{m} \left(\beta_i T + \sum_{j=1}^p \sigma_{ij} W_T^j \right) \right]$. Explain the method, give the control variate associated to X and calculate the expectation.

6.4 Conditioning

We wish to approximate $\mathbb{E}(g(X, Y))$ where the random vector (X, Y) has the density $(x, y) \mapsto f(x, y)$. Suppose that the function h defined by $h(X) = \mathbb{E}[g(X, Y)|X]$ is **available in a closed-form** and $\mathbb{E}|g(X, Y)| < \infty$; assume in addition that the distribution of X can be *sampled*. Then, the quantity $\mathbb{E}(g(X, Y))$ can be approximated by the estimator

$$\bar{S}_n = \frac{1}{N} \sum_{i=1}^N h(X_i)$$

where (X_i) are i.i.d. with density $x \mapsto \int f(x, y) dy$ (obtained by marginalization). \bar{S}_n is a *strongly convergent* and *unbiased* estimator of $\mathbb{E}(g(X, Y))$ with a lower variance than

$$S_n = \frac{1}{N} \sum_{i=1}^N g(X_i, Y_i)$$

where the (X_i, Y_i) are i.i.d. with density f . Indeed, the consistency of the estimator is given by the LLN combined with $\mathbb{E}(h(X)) = \mathbb{E}[\mathbb{E}(g(X, Y)|X)] = \mathbb{E}(g(X, Y))$. Moreover, by the tower property and the formula for conditional variances:

$$\mathbb{V}\text{ar}(g(X, Y)) = \underbrace{\mathbb{V}\text{ar}[\mathbb{E}(g(X, Y)|X)]}_{h(X)} + \underbrace{\mathbb{E}[\mathbb{V}\text{ar}(g(X, Y)|X)]}_{\geq 0} \geq \mathbb{V}\text{ar}(h(X))$$

Remark 6.10 According to this last inequality, the variance of conditional expectation is always less than the variance (without conditioning). This property is also used in statistics under the name of Rao-Blackwell's method or "Rao-Blackwellisation", it is not the subject of this course but we advise the reader to make the link between these methods and the Rao-Blackwell estimator whose definition and properties can be found in the literature...

Exercise 6.11. A STOCHASTIC VOLATILITY MODEL- CONDITIONING- We consider a financial asset such that the price satisfies the SDE:

$$dS_t = S_t(rdt + \sigma_t dW_t)$$

where the process $(\sigma_t)_{t \geq 0}$ is assumed continuous, random and independent from $(W_t)_{t \geq 0}$. Then, S_t writes

$$S_t = x \exp \left(rt - \int_0^t \frac{\sigma_s^2}{2} ds + \int_0^t \sigma_s dW_s \right)$$

We want to approximate $\mathbb{E}(e^{-rT}(S_T - K)_+)$, the price of an european call with strike K .

1. Show that conditionally on all the trajectory (σ_t) , $\int_0^t \sigma_s dW_s$ follows a gaussian distribution with a variance to be defined.
2. Deduce that $\mathbb{E} [e^{-rT}(S_T - K)_+ | (\sigma_t)_{0 \leq t \leq T}]$ can be explicitly written as a Black and Scholes formula.
3. Explain the estimation of a call by a Monte Carlo method with conditioning.

6.5 Stratified sampling

Let $Y \sim f$ be a r.v. with values in \mathbb{R}^d . We wish to approximate $\mathbb{E}[h(Y)]$. We note that

$$\mathbb{E}[h(Y)] = \int h(x)f(x)dx = \sum_{i=1}^p \underbrace{\left(\int_{S_i} f(u)du \right)}_{\alpha_i} \int h(x) \underbrace{\left[\frac{f(x)\mathbf{1}_{S_i}(x)}{\int_{S_i} f(u)du} \right]}_{\text{densité de } Y|_{Y \in S_i}} dx$$

where $(S_i)_{1 \leq i \leq p}$ are "regions" or "strata" of \mathbb{R}^d ; in mathematical terms, the $(S_i)_{1 \leq i \leq p}$ form a partition of \mathbb{R}^d . If the values of $\boxed{(\alpha_i)_{i=1, \dots, p}}$ are known and if we can draw from the distribution of Y conditionally on $\{Y \in S_i\}$ then, we can set:

$$\bar{S}_n = \sum_{i=1}^p \alpha_i \left[\frac{\sum_{j=1}^{n_i} h(Y_{i,j})}{n_i} \right],$$

with the conditions

- i) $\sum_{i=1}^p n_i = n$
- ii) $Y_{i,j} \stackrel{\mathcal{L}}{\equiv} Y|_{Y \in S_i}$ and $(Y_{i,j})$ are independent

Moreover, set $S_n = \frac{\sum_{i=1}^p h(X_i)}{n}$ with $X_i \stackrel{i.i.d.}{\sim} f$. Using the notation

$$\sigma_i^2 = \text{Var}(h(Y)|Y \in S_i), \quad \mu_i = \mathbb{E}[h(Y)|Y \in S_i]$$

we get

$$\text{Var}(\bar{S}_n) = \sum_{i=1}^p \alpha_i^2 \frac{\sigma_i^2}{n_i}, \quad \text{Var}(S_n) = \frac{1}{n} \left[\underbrace{\mathbb{E}[h^2(Y)]}_{\sum_{i=1}^p \alpha_i (\sigma_i^2 + \mu_i^2)} - \left(\underbrace{\mathbb{E}[h(Y)]}_{\sum_{i=1}^p \alpha_i \mu_i} \right)^2 \right] \quad (6.2)$$

We now consider two subproblems:

- i) For a given number of simulations n , what value should we take for n_i , i.e. how many samples n_i should we use in each region S_i ?
- ii) For given choice of the allocation numbers (n_i) , is the resulting variance really lower than a usual Monte Carlo method?

6.5.0.1 Proportional allocation

Choose $\frac{n_i}{n} = \alpha_i$. Then, (6.2) writes

$$\begin{aligned}\mathbb{V}\text{ar}(\bar{S}_n) &= \frac{1}{n} \sum_{i=1}^p \alpha_i \sigma_i^2 \\ \mathbb{V}\text{ar}(S_n) &= \frac{1}{n} \sum_{i=1}^p \alpha_i (\sigma_i^2 + \mu_i^2) - \left(\sum_{i=1}^p \alpha_i \mu_i \right)^2 = \mathbb{V}\text{ar}(\bar{S}_n) + \frac{1}{n} \left(\sum_{i=1}^p \alpha_i \mu_i^2 - \left(\sum_{i=1}^p \alpha_i \mu_i \right)^2 \right) \geq \mathbb{V}\text{ar}(\bar{S}_n)\end{aligned}$$

where the last equality follows from $\sum_{i=1}^p \alpha_i \mu_i^2 - \left(\sum_{i=1}^p \alpha_i \mu_i \right)^2 \geq 0$ since (α_i) defines a probability distribution on $\{1, \dots, p\}$. The stratified sampling thus induces a lower variance in the case of proportional allocation.

6.5.0.2 Optimal allocation

We now aim at finding the optimal (n_i^*) such that (see (6.2)),

$$\sum_{i=1}^p \alpha_i^2 \frac{\sigma_i^2}{n_i^*} = \inf \left\{ \sum_{i=1}^p \alpha_i^2 \frac{\sigma_i^2}{n_i}; \sum_{i=1}^p n_i = n \right\}$$

Set $\frac{n_i^*}{n} = \frac{\alpha_i \sigma_i}{\sum_{j=1}^p \alpha_j \sigma_j}$ and let us check that this choice is optimal. Indeed by Cauchy-Schwarz's inequality,

$$\sum_{i=1}^p \alpha_i^2 \frac{\sigma_i^2}{n_i^*} = \frac{1}{n} \left(\sum_{i=1}^p \alpha_i \sigma_i \right)^2 = \frac{1}{n} \left(\sum_{i=1}^p \sqrt{n_i} \frac{\alpha_i \sigma_i}{\sqrt{n_i}} \right)^2 \leq \frac{1}{n} \left(\sum_{i=1}^p n_i \right) \left(\sum_{i=1}^p \left(\frac{\alpha_i \sigma_i}{\sqrt{n_i}} \right)^2 \right) = \sum_{i=1}^p \alpha_i^2 \frac{\sigma_i^2}{n_i}$$

which indeed shows the optimality of the n_i^* . That being said, the proportional allocation has the advantage (when compared to the optimal allocation) to be independent of the function h whereas the n_i^* are defined in terms of σ_i , which in turn, depend on h .

Remark 6.12 *The meticulous reader can note that the proof of the optimal allocation shares some strong similarities with the one of the optimal variance in the Importance Sampling techniques (see Proposition 4.10). Let us gather the two optimisation problem in the following table (by setting $\beta_i = n_i/n$):*

| | |
|------------------------------|--|
| ALLOCATION OPTIMALE- | $\inf \left\{ \sum_{i=1}^p \frac{\alpha_i^2 \sigma_i^2}{\beta_i}; \beta_i \geq 0, \sum_{i=1}^p \beta_i = 1 \right\}$ |
| VARIANCE OPTIMALE POUR L'IS- | $\inf \left\{ \int \frac{f^2(x)h^2(x)}{g(x)} dx; g(x) \geq 0, \int g(x)dx = 1 \right\}$ |

And we can see that the two optimisation problems are the same, one being the continuous version of the other one.

6.6 Quasi Monte Carlo methods

6.6.1 Weak discrepancy sequences

We are looking for (x_i) with values in $[0, 1]$ such that

$$\frac{1}{n} \sum_{i=1}^n h(x_i) \rightarrow_{n \rightarrow \infty} \mathbb{E}(h(U))$$

where $U \sim \mathcal{U}[0, 1]$ and such that the convergence is quicker than the one given by the CLT. The sequence (x_i) being deterministic, we say that the approximation by $\frac{1}{n} \sum_{i=1}^n h(x_i)$ is of quasi-Monte Carlo type.

Definition 6.13. "SUITES EQUIRÉPARTIES"- A sequence $(x_i)_i$ with values in $[0, 1]^d$ is "équirépartie" on $[0, 1]^d$ if it satisfies one of the three equivalent properties :

- i) For all bounded continuous functions h on $[0, 1]^d$, $n^{-1} \sum_{i=1}^n h(x_i) \rightarrow \int_{[0,1]^d} h(u) du$
- ii) $\forall y = (y^1, \dots, y^d) \in [0, 1]^d$, $n^{-1} \sum_{i=1}^n \mathbf{1}_{[0,y^1] \times \dots \times [0,y^d]}(x_i) \rightarrow \text{Volume}([0, y^1] \times \dots \times [0, y^d]) = \prod_{j=1}^d y^j$.
- iii) Defining the *discrepancy at the origin*, \mathcal{D}_n^* , of the sequence (x_i) by

$$\mathcal{D}_n^* = \sup_{y=(y^1, \dots, y^d) \in [0,1]^d} \left| n^{-1} \sum_{i=1}^n \mathbf{1}_{[0,y^1] \times \dots \times [0,y^d]}(x_i) - \text{Volume}([0, y^1] \times \dots \times [0, y^d]) \right|.$$

we have $\mathcal{D}_n^* \rightarrow 0$.

Remark 6.14 The sequences are obtained as the output of r.v. with uniform distribution on $[0, 1]^d$ are "équirépartie". They indeed satisfy the first condition due to the LLN.

There exist also other "équirépartie" sequences (x_n) such that for a certain class of functions h , we have the following result

$$\left| \frac{1}{n} \sum_{i=1}^n h(x_i) - \mathbb{E}(h(U)) \right| \leq c(h) \frac{(\ln n)^d}{n} \quad (6.3)$$

Such sequences are called *sequences with weak discrepancy*. The speed of convergence is therefore of order $(\ln(n))^d/n$ which is much better than the ones obtained by standard Monte Carlo methods (where the speed, given by the CLT, is of order $1/\sqrt{n}$).

6.6.1.1 The Van der Corput sequence

We first give an example in dimension 1. Let p be a prime number. Each number n can be written in basis p as a sequence $(a_{i,n})_{i \geq 0}$ of integers in $\{0, \dots, p-1\}$, eventually equal to 0 such that

$$n = \sum_{i=0}^{\infty} a_{i,n} p^i \quad \text{où} \quad a_{i,n} \in \{0, \dots, p-1\}$$

Set $x_n^p = \sum_{i=0}^{\infty} \frac{a_{i,n}}{p^{i+1}}$. This sequence $(x_n^p)_{n \geq 0}$ is known as a Van Der Corput sequence, it is "équirépartie" on $[0, 1]$ with weak discrepancy.

The support of the Van Der Corput sequence associated to the number p is the set of the extreme values of intervals obtained by splitting the interval $[0, 1]$ into p intervals of the same length and by repeating the procedure on each subinterval.

6.6.1.2 Halton Sequence

It is a generalization of the Van der Corput sequence in dimension d . We consider \mathbb{R}^d and let (p_1, \dots, p_d) be the d first prime numbers, then the sequence of vectors $(x_n^{p_1}, \dots, x_n^{p_d})$ (where x_n^p is the n -th term in the Van der Corput sequence associated to the prime number p) is an "équirépartie" sequence on $[0, 1]^d$ with weak discrepancy.

There also exists other weak discrepancy sequences as the Faure sequence or the Sobol sequences, whose description and properties can be easily found in the literature.

Algorithm 5 Calcul de x_n^p

```

1: Input:  $n, p$ .
2:  $power = p$ .
3:  $vdc = (n \bmod p) / power$ .
4:  $n = \text{floor}(n/p)$ .
5: while  $n > 0$  do  $power = power * p$ 
6:    $vdc = vdc + (n \bmod p) / power$ 
7:    $n = \text{floor}(n/p)$ .
8: end while
9: Output:  $vdc$ .

```

Remark 6.15 We stress that these sequences are not random, it is therefore not possible to consider confidence intervals using the CLT, which does not mean anything for a deterministic sequence. Nevertheless, the inequality (6.3) gives an upper-bound for the error on the condition that we know exactly $c(h)$, which is not the case in general.

6.7 Take-home message

- a) Be able to distinguish all the means for reducing the variance: importance sampling, anti-athetic variates, control variates, conditioning, stratified sampling.
- b) For each method, be able to show the convergence properties and the asymptotic normality.
- c) Optimality for the importance sampling, for the control variates and the stratified sampling.
- d) Terminology of the quasi Monte Carlo sequences. Definition of the Van der Corput sequences and of the Halton sequences.

6.8 Highlights**6.8.0.1 Van der Corput.** source: Wikipedia

Johannes Gualtherus van der Corput (Rotterdam, September 4, 1890 - Amsterdam, September 16, 1975) was a Dutch mathematician, working in the field of analytic number theory.

He was appointed professor at the University of Groningen in 1923, and at the University of Amsterdam in 1946. He was one of the founders of the Mathematisch Centrum in Amsterdam, of which he also was the first director. From 1953 on he worked in the United States at the University of California, Berkeley and the University of Wisconsin-Madison. Among his students were J. F. Koksma and J. Popken.



Chapter 7

Exercises

Exercise 7.1. Let $p \in \mathbb{R}_*^+$.

1. Let $\Omega = [0, 1]$, $\mathcal{F} = \mathcal{B}(\Omega)$ and $\mathbb{P} = \lambda|_{\Omega}$ the Lebesgue measure restricted to $[0, 1]$. Define $X : \omega \mapsto X(\omega) = \mathbb{1}_{[0,p]}(\omega)$. What is the law of X ?
2. Let $\Omega = \{0, 1\}$, $\mathcal{F} = \mathcal{P}(\Omega)$ and \mathbb{P} such that $\mathbb{P}(\{1\}) = 1 - \mathbb{P}(\{0\}) = p$. Let $Y : \omega \mapsto Y(\omega) = \omega$. Show that the law of Y is the same as the one of X .

Exercise 7.2. Show that the Borel sigma field $\mathcal{B}(\mathbb{R})$ is generated by

- (i) the closed sets
- (ii) $\{[a, b] : a, b \in \mathbb{R}\}$
- (iii) $\{]-\infty, b] : a, b \in \mathbb{R}\}$

Exercise 7.3. Show that any continuous function $f : \mathbb{R} \rightarrow \mathbb{R}$ is $\mathcal{B}(\mathbb{R})/\mathcal{B}(\mathbb{R})$ -measurable.

Exercise 7.4. Let $\tau = \{A \in \mathcal{B}(\mathbb{R}) : A = -A\}$ where $-A = \{-x : x \in A\}$.

1. Show that $A \in \tau$ if and only if we have $(x \in A) \Rightarrow (-x \in A)$.
2. Show that τ is a sigma field on \mathbb{R} .
3. Consider the mappings $f : x \mapsto e^x$, $g : x \mapsto x^3$ and $h : x \mapsto \cos(x)$.
 - (a) Are they $\mathcal{B}(\mathbb{R})/\tau$ -measurable?
 - (b) Are they $\tau/\mathcal{B}(\mathbb{R})$ -measurable?
 - (c) Are they τ/τ -measurables?
4. Find all the functions that are $\tau/\mathcal{B}(\mathbb{R})$ -measurable.
5. Find all the functions that are τ/τ -measurables.

Exercise 7.5. Set for any $x \geq 1$: $f(x) = \sum_{n=1}^{\infty} n e^{-nx}$, what is the value of $\int_1^{\infty} f(x) dx$?

Exercise 7.6. Calculate the limit as n goes to the infinity of the quantity $\int_0^{\infty} \arctan(nx) e^{-x^n} dx$.

Index

- Δ -method, 14
- FAURE, sequence, 46
- HALTON, sequence, 46
- RAO BLACKWELL, 43
- SLUTSKI'S Lemma, 13
- SOBOL, sequence, 46
- VAN DER CORPUT, suite de , 46

- Acceptance probability, 34
- Antithetic, variates, 40

- Box Muller sampling, 22

- canonical space, 30, 31
- Central Limit theorem, 12
- Change of variables, 22
- Conditional law sampling, 20
- Conditioning, 43
- Confidence Interval, 14
- Control Variates, 41

- Discrepancy at the origin, 46
- Distribution
 - Exponential, 10
 - Gamma, 10
 - Gaussian, 10

- Exponential distribution, 18

- filtration, 28
 - natural, 28

- Generalized Inverse, 17

- Importance Sampling, 23
- independence sampler, 35
- invariant
 - probability measure, 32
- Inversion of the distribution function, 18
- irreducibility, 36

- Law of Large Numbers, 11

- Markov chain, 28
- Markov kernel, 27
- Markov property, 31
- Metropolis Hastings, 33

- Optimal allocation, 45

- positive part of a measure, 36

- Proportional allocation, 45

- Quantile function, 17

- Radon Nikodym, 36
- random walk MH sampler, 35
- Rejection algorithm, 20
- reversibility, 32

- Sampling
 - By transformation, 22
 - Approximate, 23
 - By marginalisation, 22
 - Conditional, 20
 - Exact, 17
- Stratified sampling, 44
- Suites equiréparties, 46

- Weak discrepancy sequences, 46

