

Lecture notes

Stochastic Simulation

Fabio Nobile

A.Y. 2025-2026

Last update: December 18, 2025

EPFL

Contents

1	Uniform Pseudo Random Number Generation	7
1.1	Some common uniform Pseudo-RNG	8
1.2	Empirical tests for RNG	10
1.2.1	Non-parametric Goodness-of-Fit Tests	11
1.2.2	Empirical tests for independence	13
2	Random Variable Generation	15
2.1	Inverse-transform method	15
2.2	Composition method	17
2.3	Alias method	17
2.4	Acceptance-Rejection method	18
2.4.1	Squeezing	21
2.4.2	Adaptive AR for log-concave densities	21
2.5	Ad Hoc methods	22
2.5.1	Box-Muller method	22
2.6	Multivariate Random Variable Generation	23
2.6.1	Independent components	23
2.6.2	Generation from conditional distributions	24
2.6.3	Generation by transformation using copulas	24
3	Generation of Gaussain processes	27
3.1	Generation of multivariate Gaussian random variables	27
3.2	Generation from conditional Gaussian distribution	28
3.3	Gaussian process generation	30
3.3.1	Wiener process (Brownian motion)	31
3.3.2	Brownian bridge	32
3.4	Stationary Gaussian processes / random fields	33
4	Generation of Markov processes	37
4.1	Discrete time / discrete state Markov chains	37
4.2	Discrete time / continuous state Markov chains	38
4.3	Continuous time / discrete state Markov chains	39
4.4	Poisson process	40
4.5	Non-homogeneous Poisson process	42

4.6	Compound Poisson process	43
4.7	General continuous time / discrete space Markov process	43
5	Monte Carlo method	47
5.1	Confidence intervals	47
5.2	Implementation aspects	50
5.3	Non asymptotic error bounds	51
5.4	Vector valued output	53
5.5	Smooth functions of expectations and delta method	53
5.6	Monte Carlo to compute integrals	54
6	Variance Reduction Techniques	57
6.1	Antithetic Variables	58
6.2	Importance Sampling	61
6.2.1	On the choice of the importance sampling distribution g	63
6.2.2	Weighted importance sampling	66
6.2.3	Importance sampling for stochastic processes	67
6.3	Control variates	71
6.3.1	Multiple control variates	73
6.4	Stratification	74
6.4.1	Proportional allocation	75
6.4.2	Optimal allocation	77
6.5	Latin Hypercube Sampling	78
7	Quasi Monte Carlo methods	81
7.1	Low discrepancy sequences and point sets	84
7.2	Randomized QMC formulas	87
8	Markov Chain Monte Carlo	91
8.1	Markov Chains on discrete state spaces (review)	92
8.1.1	Metropolis-Hastings algorithm in discrete state spaces	96
8.1.2	Convergence results	98
8.2	Markov chains on a general state space	107
8.3	Metropolis-Hastings algorithm in general state space	111
8.3.1	Independence sampler	113
8.3.2	Random walk Metropolis	115
8.3.3	One Variable at a time Metropolis-Hastings	115
8.3.4	Gibbs sampler	118
8.3.5	Metropolis Adjusted Langevin Algorithm (MALA)	119
8.4	Convergence diagnostics	120
8.4.1	Estimating the asymptotic variance by covariance methods	123
8.4.2	Estimating the asymptotic variance by the batch means method	123

9	Sensitivities and Stochastic Optimization	125
9.1	Computation of sensitivities	126
9.2	Stochastic optimization	129
9.2.1	Stochastic Approximation and SGD	131
9.2.2	SGD with fixed step size and increasing sample size	133
9.2.3	SGD with fixed sample size and decreasing step size	137

Chapter 1

Uniform Pseudo Random Number Generation

At the heart of any Monte Carlo method, is a *Random Number Generator* (RNG), i.e. a procedure that produces an infinite stream of random variables $U_1, U_2, \dots \stackrel{iid}{\sim} \mu$ that are independent and identically distributed (i.i.d.) according to some probability distribution μ . In particular, if μ is the *uniform* distribution on $[0, 1]$, i.e. $\mu = \mathcal{U}([0, 1])$, the generator is called a *Uniform Random Number Generator*.

Although generators based on physical devices that exploit universal background radiation or quantum mechanics effects exist, the vast majority of current random number generators are based on algorithms that can be implemented on a computer. As such, these algorithms produce a *purely deterministic* stream of numbers U_1, U_2, \dots , which, however, resembles a stream of iid random variables in the sense that the stream is indistinguishable from a random one according to a number of statistical tests. Algorithmic generators are called *Pseudo-Random Number Generators* (Pseudo-RNG).

Pseudo-RNG have the general structure, illustrated in Algorithm 1.1, where \mathcal{S} is a *finite* state space, \mathcal{U} the output space, $f : \mathcal{S} \rightarrow \mathcal{S}$ and $g : \mathcal{S} \rightarrow \mathcal{U}$ two given functions.

Algorithm 1.1: General structure of a Pseudo-RNG

```
1 take  $X_0 \in \mathcal{S}$  ; // seed
2 for  $k = 1, 2, \dots$  do
3   |  $X_k = f(X_{k-1})$  ; // recursion on state variable  $X_k \in \mathcal{S}$ 
4   |  $U_k = g(X_k)$  ; // output  $U_k \in \mathcal{U}$ 
5 end
```

Few remarks are in order:

- The initial state X_0 is called the **seed**. A *Pseudo-RNG* starting from a given seed will always produce the same sequence U_1, U_2, \dots . This is actually a convenient feature when testing or debugging a code.
- Since the state space \mathcal{S} is finite, the generator eventually will repeat itself (i.e. it will revisit an already visited state). All *Pseudo-RNGs* are *periodic*.

We call **period** the largest number of steps ℓ taken before visiting an already visited state. The *maximal period* that a generator can have is $\ell = |\mathcal{S}|$ (where $|\mathcal{S}|$ denotes the cardinality of the state space).

A good uniform Pseudo-RNG should possibly:

1. *Have a large period*: if we need to run a Monte Carlo analysis using M (pseudo) random variables, the period ℓ of the generator should be $\ell \gg M$ (otherwise the property of independent samples is clearly broken).
2. *Pass a battery of statistical tests* for uniformity and independence.
3. *Be fast and efficient*: many Monte Carlo techniques require the generation of billions of random variables. In certain fields (e.g. finance) the generation time is a big issue.
4. *Be reproducible*: in certain cases it is important to be able to reproduce a stream U_1, U_2, \dots without the need of storing it (debugging purposes, advanced MC variance reduction techniques etc.)
5. *Have the possibility to generate multiple streams*. This is important when running a Monte Carlo analysis in a parallel environment: each processor should use a stream not overlapping with the ones used by the other processors.
6. *Avoid producing the numbers 0 and 1*. The value zero might produce undesirable results as “division by zero”. Since the event “ $U = 0$ ” has zero probability, the Pseudo-RNG should never produce the value zero.

1.1 Some common uniform Pseudo-RNG

The most commonly used generators are based on *linear* recurrences. We present hereafter some examples.

Linear Congruential Generator (LCG)

It is characterized by a state space $\mathcal{S} = \{0, 1, \dots, m - 1\}$ (m is called the modulus), two natural numbers $a, b \in \mathbb{N}$ and the following recurrence and output

$$X_k = (aX_{k-1} + b) \bmod m, \quad U_k = \frac{X_k}{m}, \quad k \geq 1.$$

LCG have been popular for many years but are now somewhat outdated (e.g. Matlab versions up to 5 were using one of those). LCG can generate any number in $\{0, \frac{1}{m}, \dots, \frac{m-1}{m}\}$ and m^{-1} should be chosen of the order of the floating point machine precision (ε -machine).

A popular choice is the Lewis-Goodman-Miller LCG with $a = 7^5 = 16807$, $b = 0$, $m = 2^{31} - 1 \approx 2 \cdot 10^9$, which has a maximal period of $m - 1 \approx 4 \cdot 10^9$, too small for today’s applications.

Multiple recursive generator (MRG) of order q

For natural numbers $a_1, \dots, a_q \in \mathbb{N}$ and seeds $X_0, X_{-1}, \dots, X_{-q+1} \in \{0, \dots, m-1\}$, it is defined by the recurrence and output

$$X_k = (a_1 X_{k-1} + a_2 X_{k-2} + \dots + a_q X_{k-q}) \pmod{m}, \quad U_k = \frac{X_k}{m}, \quad k \geq 1. \quad (1.1)$$

A MRG can be written in the general form of Algorithm 1.1 by introducing the vector $\mathbf{X}^{(k)} = (X_{k-q+1}, \dots, X_k)^\top$ and the integer matrix $A \in \mathbb{N}^{q \times q}$,

$$A = \begin{pmatrix} 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ a_q & a_{q-1} & \dots & a_1 \end{pmatrix}.$$

As such, the recurrence (1.1) can be written equivalently as

$$\mathbf{X}^{(k)} = A\mathbf{X}^{(k-1)} \pmod{m}, \quad U_k = \frac{(\mathbf{X}^{(k)})_q}{m}, \quad k \geq 1, \quad (1.2)$$

for which the state space is $\mathcal{S} = \{0, 1, \dots, m-1\}^q$ and the maximal period can be up to $m^q - 1$. For a more general integer valued invertible matrix A , a generator of the form (1.2) is called *Matrix Congruential Generator of order q* . (The \pmod{m} operation in (1.1) has to be interpreted componentwise.)

Combined Generators

Here, the idea is to combine the output of several generators which, individually, may be of poor quality, to make a superior quality generator.

Example 1.1 (Wichman-Hill). *This combines 3 LCGs*

$$\begin{aligned} X_k &= (171X_{k-1}) \pmod{m_1} & (m_1 &= 30269) \\ Y_k &= (172Y_{k-1}) \pmod{m_2} & (m_2 &= 30307) \\ Z_k &= (170Z_{k-1}) \pmod{m_3} & (m_3 &= 30323) \end{aligned} \quad (1.3)$$

with

$$U_k = \frac{X_k}{m_1} + \frac{Y_k}{m_2} + \frac{Z_k}{m_3} \pmod{1}.$$

It has a period of $\ell \approx 6.95 \cdot 10^{12}$ (which is not very large for today's applications) and performs quite well in simple statistical tests.

Example 1.2 (MRG32k3a). *This is a combination of 2 MRGs:*

$$\begin{aligned} X_k &= (a_2 X_{k-2} + a_3 X_{k-3}) \pmod{m_1} \\ Y_k &= (b_1 Y_{k-1} + b_3 X_{k-3}) \pmod{m_2}, \quad \text{with } m_2 < m_1 \end{aligned}$$

with

$$U_k = \begin{cases} \frac{X_k - Y_k + m_1}{m_1 + 1}, & \text{if } X_k \leq Y_k \\ \frac{X_k - Y_k}{m_1 + 1}, & \text{if } X_k > Y_k \end{cases}$$

and suitable values of $a_2, a_3, b_1, b_3, m_1, m_2$. This has a period of $\ell \approx 3 \cdot 10^{57}$ and passes all statistical tests. It has been implemented in many packages including Matlab, Mathematica, Intel's MKL library etc.

Modulo 2 Linear Generators

These are Matrix Congruential Generators with modulus $m = 2$. Since binary operations are in general faster than integer operations, these generators are usually fast. To have long periods, the order q has to be large (the maximal period is $2^q - 1$). Among these generators a popular one is the **Linear Feedback Shift Register (LFSR) Generator** also called the Tausworthe generator. The recurrence formula is in the form of a MRG (1.1) with $m = 2$, whereas the output is given by

$$U_k = \sum_{\ell=1}^w X_{kw+\ell-1} 2^{-\ell},$$

where each word of w bits $(X_0, \dots, X_{w-1}), (X_w, \dots, X_{2w-1}), \dots$ is interpreted as a binary representation of a number in $[0, 1]$. For fast generation, most of the a_j are zero. In many cases there is only one non-zero multiplier a_r apart from a_q , and the operation in the recurrence corresponds to a (modulo 2) bit addition $X_k = X_{k-r} \oplus X_{k-q}$. Generalizations of the LFSR generator include the *Mersenne Twister* generator that is now the default generator in Matlab, and R. It has a period of $2^{19937} - 1$, is very fast and passes all practical statistical tests. The default generator in Python (numpy) is instead a *Permuted Congruential Generator* (PCG). It uses a "medium quality" LCG with $m = 2^{128}$ (unsigned long long integers represented by 128 bits) and improves its performance by performing a state dependent permutation on the 128 bit and outputting only the first 64 of them. It has period of 2^{128} , excellent statistical properties and is very fast with jump ahead and multiple straming possibilities.

1.2 Empirical tests for RNG

Several statistical tests have been proposed to assess the quality of a RNG. Today's most comprehensive test suite is *TestU01* developed by L'Ecuyer and Simard [4]. In the next section we review some non-parametric Goodness-of-Fit tests that can be used to assess the uniformity of the sequence U_1, U_2, \dots produced by a Pseudo-RNG. For generality purposes, we present these tests assuming that U_j has a general cumulative distribution function F not necessarily uniform. Then, in Section 1.2.2, we discuss some tests to assess the independence of the sequence.

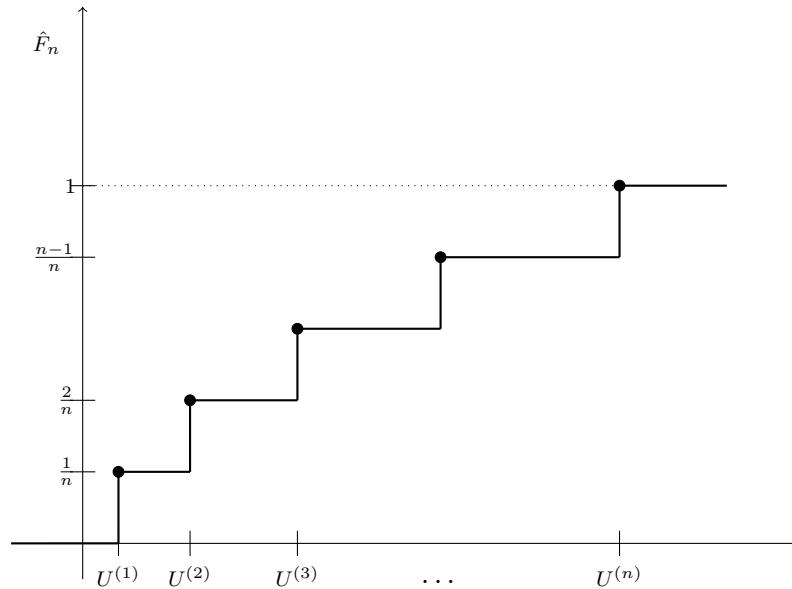


Figure 1.1: Empirical cumulative distribution function.

1.2.1 Non-parametric Goodness-of-Fit Tests

Let U be a random variable with values in a certain interval $I \subset \mathbb{R}$, and cumulative distribution function (CDF) $F(x) = \mathbb{P}(U \leq x)$. We will assume that F is absolutely continuous so that a probability density function $f : I \rightarrow \mathbb{R}_+$ exists, such that $\int_{[a,b] \subset I} f(x) dx = F(b) - F(a)$.

Let $\mathbf{U} = (U_1, \dots, U_n)$ be a random sample and denote by $\hat{F}_n(x)$ the *empirical distribution function*

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{U_i \leq x\}} = \frac{\#\{U_i \leq x, i = 1, \dots, N\}}{n}.$$

See Figure 1.1 for an illustration. In the figure, $(U^{(1)}, U^{(2)}, \dots, U^{(n)})$ denote the ordered sample \mathbf{U} . We want to test the hypothesis H_0 that \mathbf{U} has been drawn independently from the distribution F .

Q-Q plot

A first simple graphical test to see if the sample \mathbf{U} has been drawn from the distribution F is to plot the quantiles of \hat{F}_n versus the corresponding quantiles of F . We recall that the t -quantile of F is defined as

$$q_t = \operatorname{argmin}_x \{F(x) \geq t\},$$

and similarly for empirical distribution $\hat{q}_t = \operatorname{argmin}_x \{\hat{F}_n(x) \geq t\}$, which leads to $\hat{q}_{\frac{j}{n}} = U^{(j)}$, $\forall j = 1, \dots, n$, i.e. the $\frac{j}{n}$ quantile of the empirical distribution is the j -th value in the ordered sample \mathbf{U} . A better quantile estimator is actually given by $\hat{q}_{\frac{j}{n+1}} = U^{(j)}$.

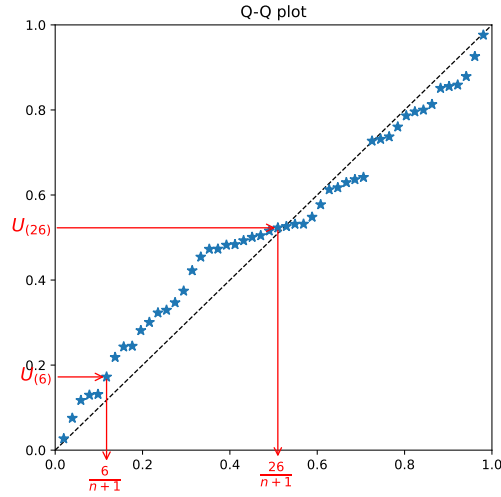


Figure 1.2: Q-Q plot of a sample from a uniform $\mathcal{U}([0, 1])$ distribution against the quantiles of the same uniform distribution

If the sample \mathbf{U} is indeed drawn from the distribution F independently, the empirical quantiles $\hat{q}_{\frac{j}{n+1}}$, when plotted against the corresponding true quantiles $q_{\frac{j}{n+1}}$, should be well aligned on the diagonal, as in Figure 1.2.

Kolmogorov-Smirnov Test

This is a more quantitative test that compares the empirical distribution \hat{F}_n with the true one F (see Figure 1.3). Let $D_n = \sup_x |\hat{F}_n(x) - F(x)|$ (which is a random variable as it depends on the random sample \mathbf{U}). For a continuous distribution F , and under the null hypothesis H_0 , it is known that

$$\sqrt{n}D_n \xrightarrow{d} K \quad \text{independently of } F$$

where \xrightarrow{d} denotes convergence in distribution and K is a random variable with Kolmogorov distribution

$$F_K(x) = \mathbb{P}(K \leq x) = \left(1 + 2 \sum_{j=1}^{\infty} (-1)^j e^{-2j^2 x^2} \right) \mathbf{1}_{\{x>0\}}.$$

The Kolmogorov distribution corresponds to the distribution of $\max_{t \in [0,1]} |B(t)|$ where $B(t)$ is a Brownian bridge in $[0, 1]$. This result shows that, under H_0 , $\hat{F}_n \rightarrow F$ uniformly at a rate $O(1/\sqrt{n})$ in a probabilistic sense. Based on this result, we can reject H_0 at level α if $\sqrt{n}D_n > K_\alpha$ with K_α the α -quantile of K : $\mathbb{P}(K \leq K_\alpha) = 1 - \alpha$. The quantiles K_α are tabulated.

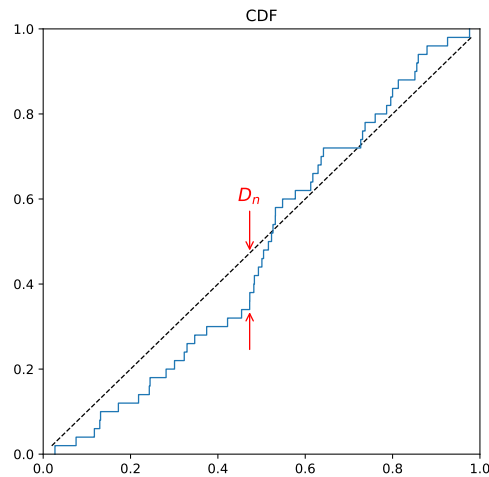


Figure 1.3: Kolmogorov-Smirnov test for a sample from the uniform $\mathcal{U}([0, 1])$ distribution. Empirical CDF \hat{F}_n in blue; exact CDF F in black; $D_n = \sup |\hat{F}_n - F|$.

χ^2 Test

We split I in $m + 1$ non-overlapping subintervals (classes) I_j , $j = 1, \dots, m + 1$ such that $\bigcup_{j=1}^{m+1} I_j = I$. For each j , let $p_j = \mathbb{P}(U \in I_j)$ be the probability that U is in I_j and define

$$N_j = \sum_{i=1}^n \mathbb{1}_{\{U_i \in I_j\}} = \#\{U_i \text{ that fall in } I_j\},$$

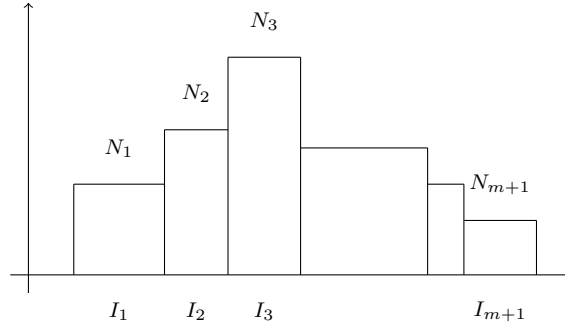
Then, under H_0 , we have $\mathbb{E}[N_j] = np_j$. We define then the statistics

$$\hat{Q}_m = \sum_{j=1}^{m+1} \frac{(N_j - np_j)^2}{np_j}$$

which has an asymptotic $\chi^2(m)$ distribution with m degrees of freedom ($m = \# \text{ classes} - 1$). We can then reject the null hypothesis H_0 at level α if $\hat{Q}_m > q_{1-\alpha}$ where $q_{1-\alpha}$ is the $1 - \alpha$ quantile of the $\chi^2(m)$ distribution. Notice that $\{(I_j, N_j), j = 1, \dots, m + 1\}$ defines a histogram of the sample and \hat{Q}_m estimates the deviation from the “true” histogram $\{(I_j, np_j), j = 1, \dots, m + 1\}$, as in Figure 1.4.

1.2.2 Empirical tests for independence

We consider here a sample $\mathbf{U} = (U_1, U_2, \dots, U_n)$ produced by a uniform Pseudo-RNG and present two statistical tests that can be used to test the null hypothesis H_0 that $\{U_i\}_i$ are mutually independent and uniformly distributed in $(0, 1)$.

Figure 1.4: χ^2 test

Serial Test

We test whether groups of variables are jointly uniformly distributed. Namely we group \mathbf{U} in groups of length d : $\mathbf{U}_1 = (U_0, \dots, U_{d-1})$, $\mathbf{U}_2 = (U_d, \dots, U_{2d-1})$, \dots and test whether $\{\mathbf{U}_j, j = 1, \dots, \frac{n}{d}\}$ are drawn independently from a multivariate uniform distribution $\mathcal{U}([0, 1]^d)$, using for instance a χ^2 test on the partition $I_{j_1 \dots j_d} = [\frac{j_1-1}{m}, \frac{j_1}{m}] \times \dots \times [\frac{j_d-1}{m}, \frac{j_d}{m}]$, $(j_1, \dots, j_d) \in \{1, \dots, m\}^d$. Of course, n/d should be sufficiently large compared to m^d so that each class has enough samples and one can apply the asymptotic result.

Gap Test

Let T_1, T_2, \dots denote the times when the process $\{U_i\}_{i=1}^n$ visits a given interval $(\alpha, \beta) \subset [0, 1]$, namely T_j is such that $U_{T_j} \in (\alpha, \beta)$ and $U_K \notin (\alpha, \beta)$, $K \notin \{T_1, T_2, \dots\}$. Let $Z_i = T_i - T_{i-1} - 1$ be the gap length between two consecutive visits (here $T_0 = 0$), i.e. if $Z_i = j$ the process has stayed j steps after T_i outside of (α, β) before entering again (α, β) at T_{i+1} . Under H_0 , Z_i are iid with a geometric distribution with parameter $p = \beta - \alpha$, i.e.

$$\mathbb{P}(Z = j) = p(1 - p)^j, \quad j = 0, 1, 2, \dots$$

One can use a $\chi^2(m)$ test to test whether the $\{Z_i\}_i$ have the correct geometric distribution, using the classes $Z = 0, Z = 1, \dots, Z = m - 1, Z \geq m$.

Chapter 2

Random Variable Generation

From a uniform (pseudo) random number generator one can construct (pseudo) random generators for many other distributions. We discuss hereafter a few approaches.

2.1 Inverse-transform method

The inverse transform method is probably the most straightforward method to generate a random variable with a given distribution and relies on the possibility to invert the cumulative distribution function. We present it separately in the case of a discrete and a continuous random variable.

Discrete random variable

Consider a discrete random variable X , which can take the values $x_1 < x_2 < \dots < x_n$ with probability mass function (pmf) $p_i = \mathbb{P}(X = x_i)$. Let $F_i = \sum_{j=1}^i p_j = \mathbb{P}(X \leq x_i)$, $i = 1, \dots, n$ and $F_0 = 0$ be the cumulative probabilities. Then X can be generated starting from a uniform random variable $U \sim \mathcal{U}([0, 1])$ by the following

Algorithm 2.1: Discrete inverse-transform.

Input: Values $\{x_i\}_{i=1}^n$, cumulative probabilities $F_i = \sum_{j=1}^i p_j$, $i = 1, \dots, n$

- 1 Generate $U \sim \mathcal{U}([0, 1])$
 - 2 Set $X = x_i$ if $F_{i-1} < U \leq F_i$
-

That this algorithm generates the correct random variable is easily seen since $\mathbb{P}(X = x_i) = \mathbb{P}(F_{i-1} < U \leq F_i) = \mathbb{P}(U \in (F_{i-1}, F_i]) = p_i$. Figure 2.1 gives a graphical illustration of the method.

Example 2.1 (Bernoulli). *Let $X \sim \text{Be}(p)$ be a Bernoulli random variable that satisfies $\mathbb{P}(X = 0) = 1 - p$, $\mathbb{P}(X = 1) = p$. Given $U \sim \mathcal{U}([0, 1])$, one sets $X = 1$ if $U > 1 - p$ and $X = 0$ otherwise.*

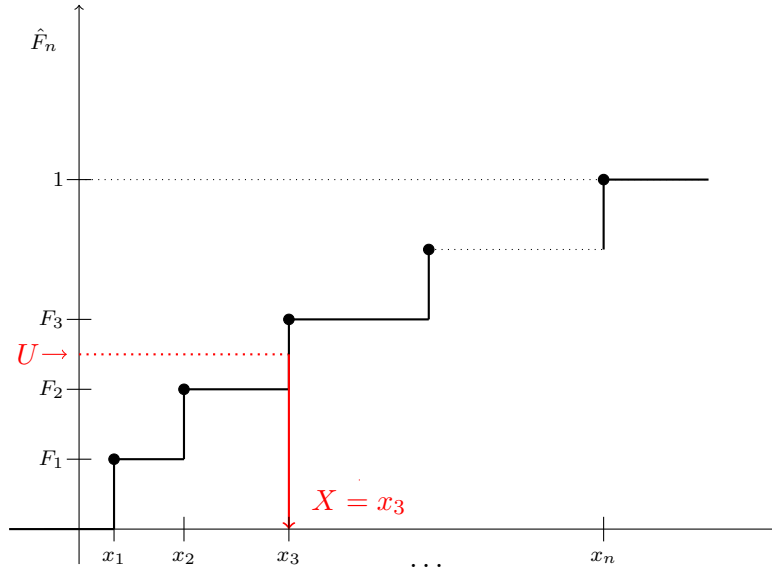


Figure 2.1: Discrete inverse transform method.

Continuous random variable

Consider a continuous random variable X taking values in an interval $[a, b]$ with continuous and strictly increasing cumulative distribution function (cdf) $F : [a, b] \rightarrow [0, 1]$, $F(x) = \mathbb{P}(X \leq x)$, with $F(a) = 0$ and $F(b) = 1$. In this case the inverse function $F^{-1} : [0, 1] \rightarrow [a, b]$ is uniquely defined and X can be generated starting from a uniform random variable $U \sim \mathcal{U}([0, 1])$ by the following

Algorithm 2.2: Continuous inverse-transform

Input: Inverse CDF F^{-1}

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

Again, one verifies easily that this algorithm generates a random variable with the correct distribution. Indeed $\mathbb{P}(X \leq x) = \mathbb{P}(F^{-1}(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x)$. Figure 2.2 gives a graphical interpretation of the method.

Example 2.2 (Exponential). Let $X \sim \text{Exp}(\lambda)$ be an exponential random variable with pdf $f(x) = \lambda e^{-\lambda x}$ and cdf $F(x) = 1 - e^{-\lambda x}$. Inversion gives $X = F^{-1}(U) = -\frac{1}{\lambda} \log(1 - U)$. Since $\tilde{U} = 1 - U$ has the same distribution as U , an equivalent inversion formula is $X = -\frac{1}{\lambda} \log U$ with $U \sim \mathcal{U}([0, 1])$.

Both the discrete and the continuous case can be combined together by defining a proper right inverse of F when it is not continuous or not strictly monotone. Let X be a random variable with cdf F . Its *Generalized inverse* is defined as $F^{-}(u) = \inf\{x : F(x) \geq u\}$. Actually, the infimum can be replaced by a minimum since F is right continuous. Then X can be generated as $X = F^{-}(U)$ with $U \sim \mathcal{U}([0, 1])$. Notice that with this definition of F^{-} we recover the discrete inverse-transform as a particular case.

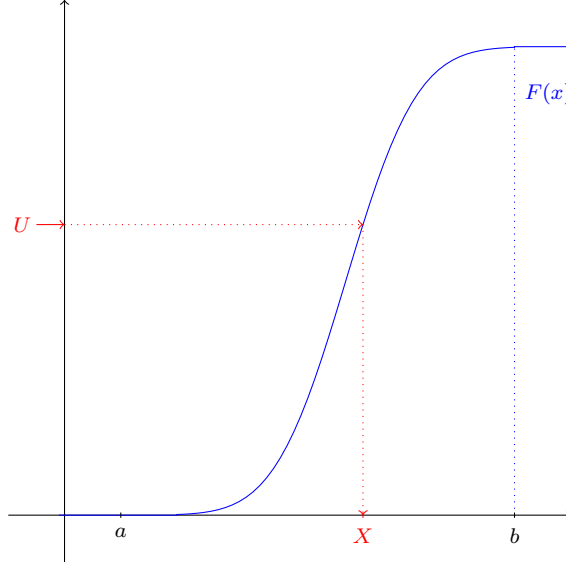


Figure 2.2: Continuous inverse transform

2.2 Composition method

Suppose that a random variable X has a mixture distribution, i.e. its cdf has the form $F(x) = \sum_{i=1}^n p_i F_i(x)$ where F_i , $i = 1, \dots, n$ are cdf functions and p_i , $i = 1, \dots, n$ are positive weights such that $\sum_{i=1}^n p_i = 1$. If the cdfs F_i are absolutely continuous with corresponding densities f_i , then X has a pdf $f(x) = \sum_{i=1}^n p_i f_i(x)$. The random variable X can be generated by the following:

Algorithm 2.3: Composition method

Input: Mixture cdf $F(x) = \sum_{i=1}^n p_i F_i(x)$

- 1 Generate discrete r.v. Y , $\mathbb{P}(Y = i) = p_i$
 - 2 Generate $X \sim F_Y$ e.g. by inversion
-

Example 2.3 (Laplace distribution). Let $X \sim \text{Lapl}(\lambda)$ with pdf

$$f(x) = \frac{\lambda}{2} e^{-\lambda|x|} = \frac{1}{2} \underbrace{\lambda e^{-\lambda x} \mathbf{1}_{\{x \geq 0\}}}_{\sim \text{Exp}(1)} + \frac{1}{2} \underbrace{\lambda e^{\lambda x} \mathbf{1}_{\{x < 0\}}}_{\sim -\text{Exp}(1)}.$$

Then, X can be generated by the composition method by first generating $B \sim \text{Be}(\frac{1}{2})$, $Y \sim \text{Exp}(\lambda)$ and then setting $X = Y$ if $B = 1$ and $X = -Y$ if $B = 0$, or, equivalently, $X = (2B - 1)Y$.

2.3 Alias method

A discrete random variable X taking the values $x_1 < x_2 < \dots < x_n$ with non-uniform probabilities $p_i = \mathbb{P}(X = x_i)$ can be generated by the discrete inverse-transform method.

However, if n is large, the search for the interval $(F_{i-1}, F_i]$ such that $U \in (F_{i-1}, F_i]$, where $F_i = \sum_{j=1}^i p_j$, might be costly. In this case, an alternative approach consists in representing the cdf $F(x)$ as a mixture distribution $F(x) = \sum_{i=1}^n \frac{1}{n} G_i(x)$ such that each G_i is a two points distribution (Bernoulli) and apply the composition method.

With little abuse of notation, we describe the algorithm using the probability “density” function which, in this case, is a linear combination of concentrated masses (delta distributions) in the points $\{x_i\}$, i.e. $f(x) = \sum_{i=1}^n p_i \delta_{x_i}(x)$. We therefore aim at rewriting it as $f(x) = \sum_{i=1}^n \frac{1}{n} g_i(x)$ where each g_i , $i = 1, \dots, n$, has the form $g_i(x) = \alpha_i \delta_{x_{\ell_i}}(x) + (1 - \alpha_i) \delta_{x_{k_i}}(x)$, with $\ell_i, k_i \in \{1, \dots, n\}$ and the distributions g_i are constructed iteratively.

- Choose ℓ_1 and k_1 such that $p_{\ell_1} < \frac{1}{n}$ and $p_{\ell_1} + p_{k_1} \geq \frac{1}{n}$ (such a choice always exists since $\{p_i\}$ are not uniform) and set $\alpha_1 = np_{\ell_1}$. Then

$$\begin{aligned} f(x) &= f^{(0)}(x) = p_{\ell_1} \delta_{x_{\ell_1}}(x) + p_{k_1} \delta_{x_{k_1}}(x) + \sum_{i \neq \ell_1, k_1} p_i \delta_{x_i}(x) \\ &= \frac{1}{n} g_1(x) + \frac{n-1}{n} f^{(1)}(x) \end{aligned}$$

with

$$\begin{aligned} g_1(x) &= \alpha_1 \delta_{x_{\ell_1}}(x) + (1 - \alpha_1) \delta_{x_{k_1}}(x), \quad \alpha_1 = np_{\ell_1} \\ f^{(1)}(x) &= \frac{n(p_{\ell_1} + p_{k_1}) - 1}{n-1} \delta_{x_{k_1}}(x) + \frac{n}{n-1} \sum_{i \neq \ell_1, k_1} p_i \delta_{x_i}(x). \end{aligned}$$

Notice that now $f^{(1)}(x)$ contains only point masses in $\{x_i, i \neq \ell_1\}$

- Iterate the procedure on $f^{(1)}, f^{(2)}, \dots$ until we reach the desired form.

We can now construct the following algorithm which does not require a search (however it requires to build in advance the table of distributions $\{g_i\}$)

Algorithm 2.4: Alias method

Input: Values $\{x_i\}_{i=1}^n$; probabilities $\{p_i\}_{i=1}^n$

- 1 Build Bernoulli distributions $g_i = \alpha_i \delta_{x_{\ell_i}} + (1 - \alpha_i) \delta_{x_{k_i}}$, $i = 1, \dots, n$
 - 2 Generate $U \sim \mathcal{U}([0, 1])$ and set $Y = \lceil nU \rceil$ // hence $Y \sim \mathcal{U}(\{1, 2, \dots, n\})$
 - 3 Generate $X \sim g_Y$ // hence $X \sim Be(\alpha_Y)$ with values $\{x_{\ell_Y}, x_{k_Y}\}$
-

2.4 Acceptance-Rejection method

Consider a continuous random variable X with pdf f and cdf F . In cases where F is difficult to invert, the inverse-transform method is not viable. Another situation which may arise is when f is known only up to a multiplicative constant, i.e. $f(x) = \kappa \tilde{f}(x)$, with $\kappa = (\int_{\mathbb{R}} \tilde{f}(x) dx)^{-1}$ and we only know \tilde{f} whereas κ is difficult or impossible to evaluate. In both cases, the acceptance-rejection method might represent a good alternative to generate X .

The idea is to find an auxiliary pdf g which is easy to sample from, and a constant $C \geq \kappa^{-1}$ such that $\tilde{f}(x) \leq Cg(x)$ for all $x \in \mathbb{R}$. Then, the acceptance-rejection algorithm reads:

Algorithm 2.5: Acceptance-Rejection (AR) algorithm

Input: \tilde{f} , g , C , such that $\tilde{f}(x) \leq Cg(x)$

- 1 Generate $Y \sim g$
 - 2 Generate $U \sim \mathcal{U}([0, 1])$ independent of Y
 - 3 If $U \leq \frac{\tilde{f}(Y)}{Cg(Y)}$ set $X = Y$, otherwise return to step 1
-

Lemma 2.1. *The acceptance-rejection Algorithm 2.5 generates a random variable X with the desired pdf $f(x) = \kappa\tilde{f}(x)$ (even without knowing κ), as long as $\tilde{f} \leq Cg$.*

Proof. Observe that the distribution of X is the distribution of Y conditional to the event $U \leq \frac{\tilde{f}(Y)}{Cg(Y)}$. Therefore $\mathbb{P}(X \leq x) = \mathbb{P}(Y \leq x \mid U \leq \frac{\tilde{f}(Y)}{Cg(Y)}) = \frac{\mathbb{P}(Y \leq x, U \leq \frac{\tilde{f}(Y)}{Cg(Y)})}{\mathbb{P}(U \leq \frac{\tilde{f}(Y)}{Cg(Y)})}$. Now,

$$\mathbb{P}\left(Y \leq x, U \leq \frac{\tilde{f}(Y)}{Cg(Y)}\right) = \int_{-\infty}^x \left(\int_0^{\frac{\tilde{f}(y)}{Cg(y)}} du \right) g(y) dy = \int_{-\infty}^x \frac{\tilde{f}(y)}{Cg(y)} g(y) dy = \frac{1}{C} \int_{-\infty}^x \tilde{f}(y) dy$$

(notice that $g = 0 \Rightarrow \tilde{f} = 0$ and we can set arbitrarily $\frac{\tilde{f}(y)}{Cg(y)} = 1$ if $g(y) = 0$) and

$$\mathbb{P}\left(U \leq \frac{\tilde{f}(Y)}{Cg(Y)}\right) = \int_{\mathbb{R}} \frac{\tilde{f}(y)}{Cg(y)} g(y) dy = \frac{1}{C} \int_{\mathbb{R}} \tilde{f}(y) dy$$

so that

$$\mathbb{P}(X \leq x) = \frac{\int_{-\infty}^x \tilde{f}(y) dy}{\int_{\mathbb{R}} \tilde{f}(y) dy} = \int_{-\infty}^x f(y) dy = F(x).$$

□

The probability of acceptance in Algorithm 2.5 is $\mathbb{P}\left(U \leq \frac{\tilde{f}(Y)}{Cg(Y)}\right) = \frac{1}{\kappa C}$ and since the trials (Y, U) are independent, the number of trials required to obtain a successful pair (X, U) has a geometric distribution $\text{Geom}(\frac{1}{\kappa C})$ with expected value κC . For the algorithm to be efficient, C should be as close as possible to κ^{-1} .

We now give a geometric interpretation of the acceptance rejection method, which is illustrated in Figure 2.3. Such interpretation is based on the following lemma.

Lemma 2.2. *Consider a non-negative integrable function $\tilde{h} : \mathbb{R} \rightarrow \mathbb{R}_+$, with $\int_{\mathbb{R}} \tilde{h} \neq 0$, the region $A_{\tilde{h}} = \{(x, u) : x \in \mathbb{R}, 0 \leq u \leq \tilde{h}(x)\}$, and the (normalized) probability density function $h(x) = \frac{\tilde{h}(x)}{\int \tilde{h}(x) dx}$ associated to \tilde{h} . A pair of random variables (X, U) is uniformly distributed in $A_{\tilde{h}}$ if and only if $X \sim h$ and $U \mid X \sim \mathcal{U}([0, \tilde{h}(X)])$.*

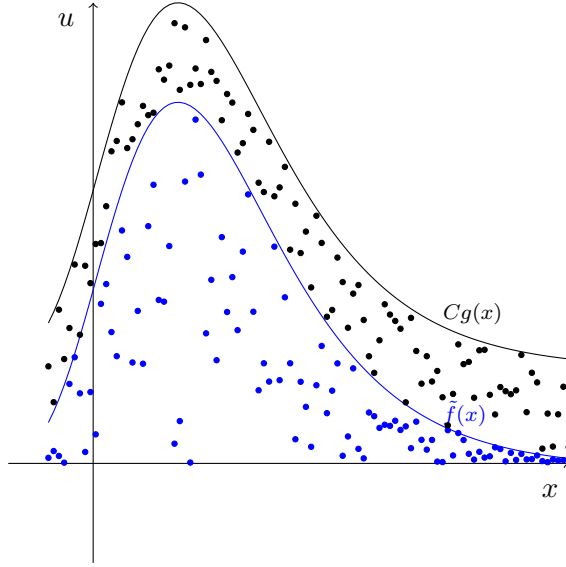


Figure 2.3: Graphical illustration of the Acceptance-Rejection method: only the blue points are retained and their abscissas are distributed according to f .

Proof. Assume first $(X, U) \sim \mathcal{U}(A_{\tilde{h}})$. Then, its probability density function is $f_{(X,U)}(x, u) = \frac{1}{|A_{\tilde{h}}|} = \frac{1}{\int_{\mathbb{R}} \tilde{h}(x) dx}$. It follows that the pdf of X is $f_X(x) = \int_0^{\tilde{h}(x)} f_{(X,U)}(x, u) du = \frac{\tilde{h}(x)}{|A_{\tilde{h}}|} = h(x)$ and the conditional probability density function of $U|X$ is $f_{U|X}(u|x) = \frac{f_{(X,U)}(x, u)}{f_X(x)} = \frac{1}{\tilde{h}(x)}$, hence $U|X \sim \mathcal{U}([0, \tilde{h}(X)])$.

Consider now the converse case, $X \sim h$ and $U|X \sim \mathcal{U}([0, \tilde{h}(X)])$. Then clearly $f_{(X,U)}(x, u) = f_{U|X}(u|x)f_X(x) = \frac{1}{\tilde{h}(x)}h(x) = \frac{1}{|A_{\tilde{h}}|}$, hence $(X, U) \sim \mathcal{U}(A_{\tilde{h}})$. \square

This observation leads to the following geometrical interpretation of the AR algorithm: in Steps 1-2 of Algorithm 2.5, one draws samples (Y, U) uniformly in the region $A_{Cg} = \{(y, u) \in \mathbb{R}^2 : 0 \leq u \leq Cg(y)\}$. In step 3, one retains only those samples that fall in the region $A_{\tilde{f}} = \{(x, u) \in \mathbb{R}^2, 0 \leq u \leq \tilde{f}(x)\}$. Hence, the abscissas of the retained points have the desired density $\frac{\tilde{f}(x)}{\int \tilde{f}(x) dx} = f(x)$.

Example 2.4. Let $Z \sim N(0, 1)$ and suppose we want to sample from $X = Z|(Z \geq 1)$, i.e. we want to sample the tail of a standard normal distribution for $Z \geq 1$. The pdf of X is $f(x) \propto e^{-x^2/2} \mathbb{1}_{\{x \geq 1\}}$. We could take as proposal distribution g an exponential $\text{Exp}(1)$ translated in 1, i.e. $g(x) = e^{-(x-1)} \mathbb{1}_{\{x \geq 1\}}$ (see Figure 2.4). We have in this case $\tilde{f}(x) = e^{-x^2/2} \mathbb{1}_{\{x \geq 1\}}$ and $\tilde{f}(x) \leq g(x) \frac{1}{\sqrt{e}}$ for all $x \geq 1$, hence we can take $C = \frac{1}{\sqrt{e}}$. The AR Algorithm reads

1. Generate $Y = 1 + \text{Exp}(1)$
2. Generate $U \sim \mathcal{U}(0, 1)$
3. If $U \leq e^{-Y^2/2 + Y - 1/2}$ set $X = Y$, otherwise return to step 1.

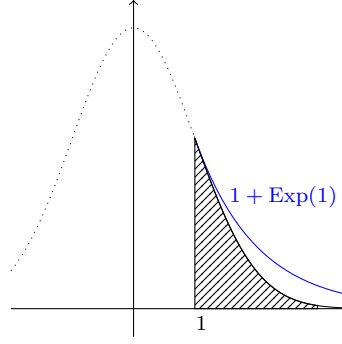


Figure 2.4: Sampling the tail of a Normal distribution by AR with an exponential proposal

The acceptance probability is $\sqrt{e} \int_1^\infty e^{-y^2/2} dy = \sqrt{2\pi e}(1 - \phi(1)) \approx 0.66$ with ϕ the cdf of a standard normal distribution. Notice that if we just sample from $N(0, 1)$ and reject all samples less than 1, we would have an acceptance rate ≈ 0.16 .

2.4.1 Squeezing

In certain cases, the expression $\tilde{f}(x)$ might be complicated and costly to evaluate, whereas $g(x)$ has generally a simple expression. To minimize the number of evaluations of \tilde{f} , one could look for another auxiliary function \hat{g} , which is also inexpensive to evaluate, such that $\hat{g}(x) \leq \tilde{f}(x) \leq Cg(x)$ for all $x \in \mathbb{R}$ and modify the AR algorithm as follows:

Algorithm 2.6: AR algorithm with squeezing.

Input: \tilde{f} , g , \hat{g} , C . such that $\hat{g} \leq \tilde{f} \leq Cg$

- 1 Generate $Y \sim g$
 - 2 Generate $U \sim \mathcal{U}([0, 1])$
 - 3 If $U \leq \frac{\hat{g}(Y)}{Cg(Y)}$ set $X = Y$, otherwise, evaluate $\tilde{f}(Y)$
 - 4 If $U \leq \frac{\tilde{f}(Y)}{Cg(Y)}$ set $X = Y$
 - 5 else reject Y and go back to 1
-

2.4.2 Adaptive AR for log-concave densities

A particularly effective adaptive AR algorithm can be set up in the case where $\log \tilde{f}(x)$ is a *concave* function. We illustrate the procedure graphically in Figure 2.5.

Let $Z_r = \{z_1, \dots, z_r\}$ be an initial set of points. Thanks to the log-concavity of \tilde{f} , we have $e^{\hat{s}(x)} \leq \tilde{f}(x) \leq e^{s(x)}$ for all $x \in \mathbb{R}$, with $s(x)$ and $\hat{s}(x)$ as in the figure. Setting now $C = \int_{\mathbb{R}} e^{s(x)} dx$, $g(x) = C^{-1}e^{s(x)}$, $\hat{g}(x) = e^{\hat{s}(x)}$, we can apply the AR algorithm with squeezing. Notice that $g(x)$ is a piecewise exponential function and can be sampled effectively by the composition method. Moreover, once a new sample X has been generated, it can be added to the set $Z_r \rightarrow Z_{r+1} = Z_r \cup \{X\}$ so that the squeezing becomes more and more effective the more variables we generate.

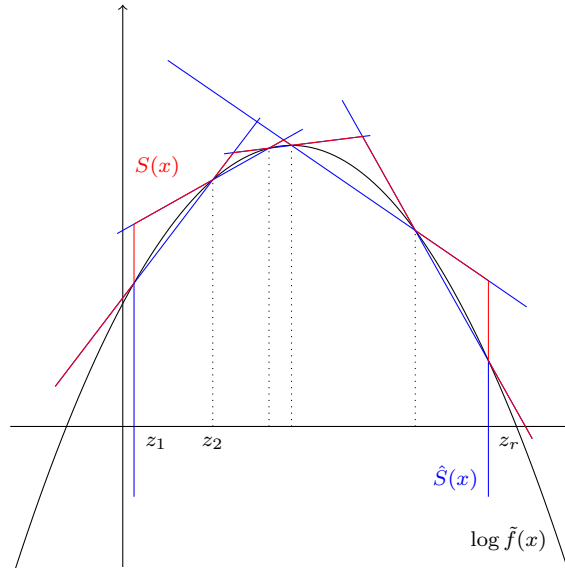


Figure 2.5: Graphical illustration of log concave density

2.5 Ad Hoc methods

The methods illustrated above are ‘general purpose’ methods, applicable to any distribution. However, for specific distributions such as Normal, Gamma, Poisson, Binomial etc, there are often much more efficient methods for random variable generation, which exploit the special structure and probabilistic interpretation of the underlying distribution. (See [3, Chapter 4]). We mention only one possible algorithm to generate variables from the Normal distribution $N(0, 1)$.

2.5.1 Box-Muller method

Let $X, Y \sim N(0, 1)$ be independent standard normal random variables, and (ρ, θ) their representation in polar coordinates. Since $X^2 + Y^2 \sim \chi_2^2 = \text{Exp}(\frac{1}{2})$ (χ_2^2 is a chi-square distribution with 2 degrees of freedom, which coincides with an exponential of parameter $\frac{1}{2}$), it follows that $\rho^2 \sim \text{Exp}(\frac{1}{2})$. Moreover, by the radial symmetry of the bivariate normal distribution $N(0, I_2)$, the distribution of (X, Y) given $\rho^2 = X^2 + Y^2$ is uniform in $[0, 2\pi)$. From these considerations, an algorithm to generate $(X, Y) \sim N(0, I_2)$ is:

Algorithm 2.7: Box-Muller method.

- | | | |
|---|---|--|
| 1 | Generate $U \sim \mathcal{U}(0, 1)$ and set $\rho = \sqrt{-2 \log U}$ | // hence $\rho^2 \sim \text{Exp}(\frac{1}{2})$ |
| 2 | Generate $V \sim \mathcal{U}(0, 1)$ and set $\Theta = 2\pi V$ | // hence $\Theta \sim \mathcal{U}([0, 2\pi])$ |
| 3 | Set $X = \rho \cos \Theta, Y = \rho \sin \Theta$. | |
-

2.6 Multivariate Random Variable Generation

We consider now the problem of generating from a multivariate distribution. Let $\mathbf{X} = (X_1, \dots, X_n)^\top \in \mathbb{R}^n$ be a vector of random variables with joint cumulative distribution function $F(\mathbf{z}) = F(z_1, \dots, z_n) = \mathbb{P}(X_1 \leq z_1, \dots, X_n \leq z_n)$ and probability density function $f(\mathbf{x}) = f(x_1, \dots, x_n)$, if it exists, such that

$$F(z_1, \dots, z_n) = \int_{-\infty}^{z_1} \cdots \int_{-\infty}^{z_n} f(x_1, \dots, x_n) dx_1 \dots dx_n.$$

The inverse transform method is not (directly) applicable in this case since the cumulative distribution function $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is not invertible. The acceptance-rejection method, on the other hand, generalizes straightforwardly to the multivariate case. However, it is in general not an easy task to find an auxiliary function $g(\mathbf{x})$ and a constant $C > 1$ such that $f(\mathbf{x}) \leq Cg(\mathbf{x})$, $\forall \mathbf{x} \in \mathbb{R}^n$, leading to reasonable acceptance rates.

In general, the problem of generating from a multivariate distribution can be very hard. We mention, hereafter, few cases where generation is relatively easy.

2.6.1 Independent components

The simplest case is when the components X_1, \dots, X_n of \mathbf{X} are independent, each with cdf $F_i : \mathbb{R} \rightarrow [0, 1]$, so that $F(\mathbf{z}) = F_1(z_1) \cdots F_n(z_n)$. (Similarly, if each X_i has a density f_i , there holds $f(\mathbf{z}) = f_1(z_1) \cdots f_n(z_n)$.) In this case, each component X_i can be generated independently of the others by using any of the techniques described for univariate functions.

Example 2.5. *Suppose we would like to draw a point $\mathbf{X} = (X_1, X_2)$ that is uniformly distributed in the unit cube $(0, 1)^2$. Since $F(\mathbf{z}) = \mathbb{P}(X_1 \leq z_1, X_2 \leq z_2) = z_1 z_2$, we conclude that X_1, X_2 are independent and $X_i \sim F_i(z) = z$, i.e. each X_i is uniformly distributed in $(0, 1)$. We can then draw $X_1, X_2 \sim \mathcal{U}(0, 1)$ independently and set $\mathbf{X} = (U_1, U_2)$.*

Example 2.6. *Suppose now that we would like to draw a point $\mathbf{X} = (X_1, X_2)$ that is uniformly distributed on the unit ball $\mathcal{B} = \{(x, y) : x^2 + y^2 \leq 1\}$. One possibility is to use an acceptance-rejection method. For instance, we could draw \mathbf{Y} uniformly on the cube $(-1, 1)^2$ and accept it by setting $\mathbf{X} = \mathbf{Y}$ only if $\mathbf{Y} \in \mathcal{B}$. The acceptance rate is $\frac{\pi}{4} \approx 78.5\%$. (Try to do it now in dimension $n \gg 2$ and see what happens ...)*

Alternatively, we could try to generate directly a point \mathbf{X} with the correct distribution, without the acceptance-rejection step. For this, let us consider a transformation in polar coordinates, $X_1 = R \cos \Theta$, $X_2 = R \sin \Theta$. Then, if $f_{(X_1, X_2)}$ denotes the joint density of (X_1, X_2) and $f_{(R, \Theta)}$ the joint density of (R, Θ) , we have $f_{(X_1, X_2)}(x, y) = \frac{1}{\pi}$, $(x, y) \in \mathcal{B}$ and

$$f_{(R, \Theta)}(\rho, \theta) = f_{(X_1, X_2)}(\rho \cos \theta, \rho \sin \theta) \left| \frac{\partial(x_1, x_2)}{\partial(\rho, \theta)} \right| = \frac{\rho}{\pi} = 2\rho \frac{1}{2\pi}.$$

We see then that (R, Θ) are independent with $\Theta \sim \mathcal{U}(0, 2\pi)$ and R having pdf $f_R(\rho) = 2\rho$ and cdf $F_R(\rho) = \rho^2$, which can be easily inverted. Therefore, starting from $U_1, U_2 \sim \mathcal{U}(0, 1)$ independent, we set $R = \sqrt{U_1}$ and $\Theta = 2\pi U_2$ so that $\mathbf{X} = (R \cos \Theta, R \sin \Theta)$ is a uniformly distributed point in the unit circle.

2.6.2 Generation from conditional distributions

Another situation which may lead to a relatively easy generation algorithm is when the marginal and univariate conditional distributions of \mathbf{X} are easily accessible. For instance, let us assume that the conditional density of $X_j|X_{1:j-1}$,

$$f_{X_j|X_{1:j-1}}(z_j|z_1, \dots, z_{j-1}) = \frac{\int_{\mathbb{R}^{n-j}} f(z_1, \dots, z_j, z_{j+1}, \dots, z_n) dz_{j+1} \dots dz_n}{\int_{\mathbb{R}^{n-j+1}} f(z_1, \dots, z_j, z_{j+1}, \dots, z_n) dz_j dz_{j+1} \dots dz_n}$$

with $X_{1:j}$ a shorthand notation for (X_1, \dots, X_j) , is known explicitly for any $j = 1, \dots, n$. Assume, moreover, that we know how to generate a variable from the density $f_{X_j|X_{1:j-1}}(\cdot | z_{1:j-1})$, for any $z_{1:j-1} \in \mathbb{R}^{j-1}$. We can then generate \mathbf{X} with the following iterative Algorithm:

Algorithm 2.8: Generation from conditional distributions.

Input: conditional densities $f_{X_j|X_{1:j-1}}$ (or cumulative distributions $F_{X_j|X_{1:j-1}}$)

- 1 Generate $X_1 \sim f_{X_1}$
 - 2 For $i = 2, \dots, n$,
 - 3 Generate $X_i \sim f_{X_i|X_{1:i-1}}(\cdot | X_{1:i-1})$
-

Again, the generation of X_i from $f_{X_i|X_{1:i-1}}(\cdot | X_{1:i-1})$ can be done using any of the techniques available for univariate variables.

Example 2.7 (Generating order statistics). *Let $\mathbf{X} = (X_1, \dots, X_n) \sim \mathcal{U}((0, 1)^n)$ and denote by $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ the ordered sample (order statistics). To generate $X_{(1)}, \dots, X_{(n)}$, one can simply generate $\mathbf{X} = (X_1, \dots, X_n)$ with $X_i \stackrel{iid}{\sim} \mathcal{U}(0, 1)$ and then order the components of the vector. However, if n is large, the ‘sort’ operation might become costly so one may prefer to generate directly $X_{(1)}, \dots, X_{(n)}$ from the distribution of the order statistics.*

Observe that $X_{(n)} = \max_{i=1, \dots, n} X_i \sim F_{X_{(n)}}(z) = z^n$ so that $X_{(n)}$ can be generated easily by inversion as $X_{(n)} = (U_n)^{1/n}$ with $U_n \sim \mathcal{U}(0, 1)$. Moreover, it can be shown (exercise) that for all $j < n$,

$$\begin{aligned} F_{X_{(j)}|X_{(j+1)}, \dots, X_{(n)}}(z | x_{j+1:n}) &= \mathbb{P}(X_{(j)} \leq z | X_{(j+1)} = x_{j+1}, \dots, X_{(n)} = x_n) \\ &= \mathbb{P}(X_{(j)} \leq z | X_{(j+1)} = x_{j+1}) = \left(\frac{z}{x_{j+1}}\right)^j, \quad z \leq x_{j+1}, \end{aligned}$$

where $F_{X_{(j)}|X_{(j+1)}, \dots, X_{(n)}}(z | x_{j+1:n}) = \int_0^z f_{X_{(j)}|X_{(j+1)}, \dots, X_{(n)}}(t | x_{j+1:n}) dt$ is the cumulative conditional distribution, which can be easily inverted. Hence, we can generate $X_{(j)}$ as $X_{(j)} = (U_j)^{\frac{1}{j}} X_{(j+1)}$ with $U_j \sim \mathcal{U}(0, 1)$ independent of the previously generated ones.

2.6.3 Generation by transformation using copulas

Consider a vector $\mathbf{X} = (X_1, \dots, X_n)$ with dependent components X_i , $i = 1, \dots, n$, with marginal distributions $F_i : \mathbb{R} \rightarrow [0, 1]$. Often the dependency structure is described in terms of *copulas*.

Definition 2.1. A copula is a cdf $C : [0, 1]^n \rightarrow [0, 1]$ of n dependent uniform random variables $U_1, \dots, U_n \sim \mathcal{U}(0, 1)$

$$C(u_1, \dots, u_n) = \mathbb{P}(U_1 \leq u_1, \dots, U_n \leq u_n).$$

We say that the dependency structure of \mathbf{X} is described by the copula C and marginal distributions F_i , $i = 1, \dots, n$, if

$$F_{\mathbf{X}}(x_1, \dots, x_n) = \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n) = C(F_1(x_1), \dots, F_n(x_n)),$$

i.e. the transformed variables $U_i = F_i(X_i)$ have a uniform marginal distribution and a joint cdf given by the copula C . In such a case, an algorithm to generate \mathbf{X} is

Algorithm 2.9: Generation of dependent components via copulas

Input: marginals $\{F_i\}_{i=1}^n$; copula C

1 Generate $\mathbf{U} = (U_1, \dots, U_n) \sim C$

2 Output $\mathbf{X} = (X_1, \dots, X_n) = (F_1^{-1}(U_1), \dots, F_n^{-1}(U_n))$.

Clearly the implementability/efficiency of this algorithm depends on the possibility to generate $\mathbf{U} \sim C$. A typical example is the case of a *Gaussian copula*. In this case, let $\mathbf{Y} = (Y_1, \dots, Y_n) \sim N(\mathbf{0}, \Sigma)$ with $\Sigma \in \mathbb{R}^{n \times n}$ symmetric and positive definite and denote $\sigma_i = \sqrt{\text{Var}(Y_i)}$ the standard deviation of Y_i . Set, moreover, $\mathbf{U} = (U_1, \dots, U_n) = \left(\Phi\left(\frac{Y_1}{\sigma_1}\right), \dots, \Phi\left(\frac{Y_n}{\sigma_n}\right) \right)$ with Φ the cdf of a standard normal random variable and $C_G^{\Sigma}(u_1, \dots, u_n) = \mathbb{P}(U_1 \leq u_1, \dots, U_n \leq u_n)$. Such a copula is called a Gaussian copula with covariance matrix Σ . To generate a vector \mathbf{X} with copula C_G^{Σ} and marginals F_i , we can therefore use the following algorithm.

Algorithm 2.10: Generation of samples from a Gaussian copula

Input: marginals $\{F_i\}_{i=1}^n$; covariance matrix Σ

1 Generate $\mathbf{Y} \sim N(\mathbf{0}, \Sigma)$

2 Compute $\mathbf{U} = \left(\Phi\left(\frac{Y_1}{\sigma_1}\right), \dots, \Phi\left(\frac{Y_n}{\sigma_n}\right) \right)$

3 Compute $\mathbf{X} = (F_1^{-1}(U_1), \dots, F_n^{-1}(U_n))$

A way to generate $\mathbf{Y} \sim N(\mathbf{0}, \Sigma)$ is discussed in the next Chapter.

Chapter 3

Generation of Gaussian random variables and processes

3.1 Generation of multivariate Gaussian random variables

A multivariate Gaussian random variable $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$ with mean $\boldsymbol{\mu} \in \mathbb{R}^n$ and covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$ symmetric and positive definite has joint pdf

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad \mathbf{x} \in \mathbb{R}^n$$

and characteristic function $\phi(\mathbf{t}) = \mathbb{E}\left[e^{i\mathbf{t}^\top \mathbf{X}}\right] = \exp\left(i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \Sigma \mathbf{t}\right)$. Notice that the characteristic function is well defined also in the case of a singular covariance matrix Σ , whereas the pdf is not.

The standard algorithm to generate \mathbf{X} relies on explicit factorization of the covariance matrix as $\Sigma = AA^\top$, with $A \in \mathbb{R}^{n \times n}$, which can always be done since Σ is symmetric and positive definite. There are two common ways to compute the factor A :

- *Cholesky factorization.* It is applicable if Σ is strictly positive definite (hence invertible) and leads to a lower triangular factor A .

If Σ is nearly singular, a more stable procedure is given by the *pivoted* Cholesky factorisation, which can also be used in the singular case.

- *Spectral decomposition.* It is based on diagonalization of the covariance matrix as $\Sigma = VD V^\top$ with $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ the matrix of eigenvalues and V the orthonormal matrix of eigenvectors. We set then $A = VD^{1/2}$, with $D^{1/2} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$. It can be used also in the singular case.

Using either factorization, \mathbf{X} can be generated by the following algorithm.

Algorithm 3.1: Multivariate Gaussian generator

Given: $\boldsymbol{\mu} \in \mathbb{R}^n$ and $\Sigma = AA^\top \in \mathbb{R}^{n \times n}$ (spd)

- 1 Generate $\mathbf{Y} \sim N(\mathbf{0}, I_{n \times n})$ (i.e. $\mathbf{Y} = (Y_1, \dots, Y_n)$, $Y_i \stackrel{\text{iid}}{\sim} N(0, 1)$)
 - 2 Compute $\mathbf{X} = \boldsymbol{\mu} + A\mathbf{Y}$
-

It is easy to check that \mathbf{X} has the correct distribution. Indeed, \mathbf{X} is Gaussian being an affine transformation of a standard normal vector. Moreover, $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$ and

$$\text{Cov}[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top] = \mathbb{E}[\mathbf{A}\mathbf{Y}\mathbf{Y}^\top\mathbf{A}^\top] = \mathbf{A}\mathbb{E}[\mathbf{Y}\mathbf{Y}^\top]\mathbf{A}^\top = \Sigma.$$

The algorithm can be easily modified in the case the *precision matrix* $\Lambda = \Sigma^{-1}$ is given, instead of Σ . (Of course we assume here that Σ is invertible.)

Algorithm 3.2: Multivariate Gaussian generator from precision matrix

Given: $\boldsymbol{\mu} \in \mathbb{R}^n$ and $\Lambda = \Sigma^{-1}$

- 1 Compute the Cholesky factorisation $\Lambda = \mathbf{L}\mathbf{L}^\top$
 - 2 Generate $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$ // n independent standard normals
 - 3 Solve the linear system $\mathbf{L}^\top\mathbf{Y} = \mathbf{Z}$ // upper triangular linear system
 - 4 Output $\mathbf{X} = \boldsymbol{\mu} + \mathbf{Y}$
-

Again it is easy to verify that \mathbf{X} has the right distribution. Indeed $\mathbb{E}[\mathbf{Y}] = \mathbf{L}^{-\top}\mathbb{E}[\mathbf{Z}] = \mathbf{0}$ which implies $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$ and

$$\mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top] = \mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \mathbb{E}[\mathbf{L}^{-\top}\mathbf{Z}\mathbf{Z}^\top\mathbf{L}^{-1}] = \mathbf{L}^{-\top}\mathbf{L}^{-1} = \Lambda^{-1} = \Sigma.$$

3.2 Generation from conditional Gaussian distribution

Consider a multivariate Gaussian random variable $\mathbf{X} \in \mathbb{R}^n$, $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$, which we split into two components, $\mathbf{X} = (Y_1, \dots, Y_{n-k}, Z_1, \dots, Z_k) = (\mathbf{Y}, \mathbf{Z})$, which we suppose unobservable and observable, respectively. The mean $\boldsymbol{\mu}$ and covariance Σ also split accordingly as

$$\mathbb{E}[\mathbf{X}] = \begin{pmatrix} \boldsymbol{\mu}_\mathbf{Y} \\ \boldsymbol{\mu}_\mathbf{Z} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{\mathbf{Y}\mathbf{Y}} & \Sigma_{\mathbf{Y}\mathbf{Z}} \\ \Sigma_{\mathbf{Y}\mathbf{Z}}^\top & \Sigma_{\mathbf{Z}\mathbf{Z}} \end{pmatrix}$$

with $\boldsymbol{\mu}_\mathbf{Y} = \mathbb{E}[\mathbf{Y}]$, $\boldsymbol{\mu}_\mathbf{Z} = \mathbb{E}[\mathbf{Z}]$, $\Sigma_{\mathbf{Y}\mathbf{Z}} = \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}_\mathbf{Y})(\mathbf{Z} - \boldsymbol{\mu}_\mathbf{Z})^\top] = \Sigma_{\mathbf{Z}\mathbf{Y}}^\top$, and similarly for $\Sigma_{\mathbf{Y}\mathbf{Y}}$, $\Sigma_{\mathbf{Z}\mathbf{Z}}$.

We are interested to generate the conditional random variable $\mathbf{Y}|\mathbf{Z}$. It is well known that the conditional distribution of \mathbf{Y} given \mathbf{Z} is again a multivariate Gaussian $N(\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{Z}}, \Sigma_{\mathbf{Y}|\mathbf{Z}})$ with

$$\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{Z}} = \boldsymbol{\mu}_\mathbf{Y} + \Sigma_{\mathbf{Y}\mathbf{Z}}\Sigma_{\mathbf{Z}\mathbf{Z}}^{-1}(\mathbf{Z} - \boldsymbol{\mu}_\mathbf{Z}) \quad (3.1)$$

$$\Sigma_{\mathbf{Y}|\mathbf{Z}} = \Sigma_{\mathbf{Y}\mathbf{Y}} - \Sigma_{\mathbf{Y}\mathbf{Z}}\Sigma_{\mathbf{Z}\mathbf{Z}}^{-1}\Sigma_{\mathbf{Z}\mathbf{Y}}. \quad (3.2)$$

Assuming we have observed the realization \mathbf{z} of \mathbf{Z} , to generate $\mathbf{Y} | \mathbf{Z} = \mathbf{z}$ one can, of course, factorize the covariance matrix $\Sigma_{\mathbf{Y}|\mathbf{Z}} = \mathbf{A}\mathbf{A}^\top$ (by Cholesky or spectral decomposition). This, however, can be expensive or cumbersome if the size of \mathbf{Y} is big. In particular, there might be cases (typically in the generation of stationary Gaussian random fields) in which generating $\mathbf{X} = (\mathbf{Y}, \mathbf{Z})$ is easy, even for very large dimensions,

but generating $\mathbf{Y} \mid \mathbf{Z} = \mathbf{z}$ e.g. by Cholesky factorization, would be very costly. Here is an alternative algorithm to do so.

Algorithm 3.3: Generation from conditional Gaussian distribution - I

Given: $\boldsymbol{\mu}, \Sigma$ and $\mathbf{z} \in \mathbb{R}^k$

- 1 Generate $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$
 - 2 Set $\mathcal{Y} = (X_1, \dots, X_{n-k})$ and $\mathcal{Z} = (X_{n-k+1}, \dots, X_n)$
 - 3 Output $\mathbf{Y} = \mathcal{Y} + \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}(\mathbf{z} - \mathcal{Z})$
-

Again, we can easily verify that \mathbf{Y} has the correct distribution. Indeed,

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[\mathcal{Y}] + \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}(\mathbf{z} - \mathbb{E}[\mathcal{Z}]) = \boldsymbol{\mu}_{\mathcal{Y}} + \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}(\mathbf{z} - \boldsymbol{\mu}_{\mathcal{Z}}) = \boldsymbol{\mu}_{\mathcal{Y} \mid \mathbf{Z}=\mathbf{z}}.$$

Moreover, setting $\mathbf{Y}' = \mathbf{Y} - \mathbb{E}[\mathbf{Y}] = (\mathcal{Y} - \boldsymbol{\mu}_{\mathcal{Y}}) - \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}(\mathcal{Z} - \boldsymbol{\mu}_{\mathcal{Z}})$, we have

$$\begin{aligned} \text{Cov}(\mathbf{Y}) &= \mathbb{E}[\mathbf{Y}'\mathbf{Y}'^\top] = \mathbb{E}[\mathcal{Y}'\mathcal{Y}'^\top] - \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}\mathbb{E}[\mathcal{Z}'\mathcal{Y}'^\top] \\ &\quad - \mathbb{E}[\mathcal{Y}'\mathcal{Z}'^\top]\Sigma_{\mathcal{Z}\mathcal{Z}}^{-\top}\Sigma_{\mathcal{Y}\mathcal{Z}}^\top + \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}\mathbb{E}[\mathcal{Z}'\mathcal{Z}'^\top]\Sigma_{\mathcal{Z}\mathcal{Z}}^{-\top}\Sigma_{\mathcal{Y}\mathcal{Z}}^\top \\ &= \Sigma_{\mathcal{Y}\mathcal{Y}} - \Sigma_{\mathcal{Y}\mathcal{Z}}\Sigma_{\mathcal{Z}\mathcal{Z}}^{-1}\Sigma_{\mathcal{Z}\mathcal{Y}}. \end{aligned}$$

This construction can be generalized to the case where, instead of observing a subset \mathbf{Z} of components of \mathbf{X} , one observes k linear combinations of the components of \mathbf{X} . Let $H \in \mathbb{R}^{k \times n}$ be an ‘‘observation’’ operator and assume one has noisy observations

$$\mathbf{Z} = H\mathbf{X} + \boldsymbol{\eta}$$

with $\boldsymbol{\eta}$ a zero mean Gaussian vector, independent of \mathbf{X} , with covariance $\Gamma \in \mathbb{R}^{k \times k}$, hence $\mathbf{Z} \sim N(H\boldsymbol{\mu}, S)$, with $S = H\Sigma H^\top + \Gamma$. Notice that the previous setting can be recovered by choosing $\boldsymbol{\eta} = 0$ and $H_{ij} = 1$ whenever $j = i$ and i corresponds to an observed component of \mathbf{X} , and $H_{ij} = 0$ otherwise.

We can now define the extended vector $\tilde{\mathbf{X}} = (\mathbf{X}, \mathbf{Z}) \in \mathbb{R}^{n+k}$ which has

$$\mathbb{E}[\tilde{\mathbf{X}}] = \begin{pmatrix} \boldsymbol{\mu} \\ H\boldsymbol{\mu} \end{pmatrix}, \quad \text{Cov}(\tilde{\mathbf{X}}) = \begin{pmatrix} \Sigma & \Sigma H^\top \\ H\Sigma & S \end{pmatrix}$$

and use the previous formulas to obtain $\mathbf{X} \mid \mathbf{Z} \sim N(\boldsymbol{\mu}_{\mathbf{X} \mid \mathbf{Z}}, \Sigma_{\mathbf{X} \mid \mathbf{Z}})$ with

$$\boldsymbol{\mu}_{\mathbf{X} \mid \mathbf{Z}} = \boldsymbol{\mu} + \Sigma H^\top S^{-1}(\mathbf{Z} - H\boldsymbol{\mu}), \quad (3.3)$$

$$\Sigma_{\mathbf{X} \mid \mathbf{Z}} = \Sigma - \Sigma H^\top S^{-1}H\Sigma. \quad (3.4)$$

The matrix $K = \Sigma H^\top S^{-1}$ is called the *Kalman gain* whereas the vector $\mathbf{d} = \mathbf{Z} - H\boldsymbol{\mu}$ the *innovation*. Generating $\mathbf{X} \mid \mathbf{Z} = \mathbf{z}$ can be done with the following algorithm.

Algorithm 3.4: Generation from conditional Gaussian distribution - II

Given: $\boldsymbol{\mu}, \Sigma$ and $\mathbf{z} \in \mathbb{R}^k$

- 1 Generate $\mathcal{X} \sim N(\boldsymbol{\mu}, \Sigma)$
 - 2 Generate $\boldsymbol{\eta} \sim N(0, \Gamma)$
 - 3 Compute perturbed innovation $\tilde{\mathbf{d}} = \mathbf{z} - H\mathcal{X} + \boldsymbol{\eta}$
 - 4 Compute Kalman gain $K = \Sigma H^\top (H\Sigma H^\top + \Gamma)^{-1}$
 - 5 Output $\mathbf{X} = \mathcal{X} + K\tilde{\mathbf{d}}$
-

It is easy to check that the output \mathbf{X} of the algorithm has the right distribution. Indeed, clearly \mathbf{X} has a Gaussian distribution. Moreover

$$\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu} + K(\mathbf{z} - H\boldsymbol{\mu}) = \boldsymbol{\mu}_{\mathbf{X}|\mathbf{z}=\mathbf{z}}$$

and

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= \mathbb{E} \left[((I - KH)(\boldsymbol{\mathcal{X}} - \boldsymbol{\mu}) + K\boldsymbol{\eta}) ((I - KH)(\boldsymbol{\mathcal{X}} - \boldsymbol{\mu}) + K\boldsymbol{\eta})^\top \right] \\ &= (I - KH)\Sigma(I - KH)^\top + K\Gamma K^\top \\ &= \Sigma - KH\Sigma - \Sigma H^\top K^\top + K(H\Sigma H^\top + \Gamma)K^\top = \Sigma_{\mathbf{X}|\mathbf{z}}. \end{aligned}$$

3.3 Gaussian process generation

Let $I \subset \mathbb{R}^d$. A collection of random variables $\{X_t, t \in I\}$ indexed by $t \in I$ is called a *stochastic process* when $d = 1$ (usually t denotes time) or a *random field* if $d \geq 1$ and t denotes the space variable.

Definition 3.1 (Gaussian process). *A Gaussian process (or Gaussian random field) is a stochastic process (random field) for which all finite dimensional distributions are Gaussian, i.e. for all $n \in \mathbb{N}$ and $t_1, \dots, t_n \in I$, the random vector $\mathbf{X} = (X_{t_1}, \dots, X_{t_n})$ has a multivariate Gaussian distribution. Equivalently, any linear combination $Y_{\mathbf{b}} = \sum_{i=1}^n b_i X_{t_i}$ has a Gaussian distribution.*

Given a Gaussian process $\{X_t, t \in I\}$, we can define

Mean function : $\mu_X : I \rightarrow \mathbb{R}, \quad \mu_X(t) = \mathbb{E}[X_t], t \in I,$

Covariance funct. : $C_X : I \times I \rightarrow \mathbb{R}, \quad C_X(t, s) = \mathbb{E}[(X_t - \mu_X(t))(X_s - \mu_X(s))], t, s \in I.$

If we now take a set of points $t_1, \dots, t_n \in I$ and consider the Gaussian random vector $\mathbf{X} = (X_{t_1}, \dots, X_{t_n})$, it clearly holds $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$, with $\boldsymbol{\mu} = (\mu_X(t_1), \dots, \mu_X(t_n))$ and $\Sigma_{ij} = C_X(t_i, t_j)$. As such, the matrix Σ has to be symmetric and non negative definite. This poses restrictions to the class of functions that can be covariance functions of a stochastic process.

Definition 3.2. *A function $C : I \times I \rightarrow \mathbb{R}$ is positive (semi-)definite if, for all n and $t_1, \dots, t_n \in I$, the matrix $\Sigma \in \mathbb{R}^{n \times n}$, $\Sigma_{ij} = C(t_i, t_j)$ is positive (semi-)definite.*

Proposition 3.1. *A Gaussian process $\{X_t, t \in I\}$ is uniquely determined by the mean function $\mu_X : I \rightarrow \mathbb{R}$ and a symmetric and positive (semi-)definite covariance function $C_X : I \times I \rightarrow \mathbb{R}$.*

We use the notation $X \sim N(\mu_X, C_X)$ to denote a Gaussian process $\{X_t, t \in I\}$ with mean function μ_X and covariance function C_X .

A Gaussian process $X \sim N(\mu_X, C_X)$ can be generated exactly in a set of points $t_1, \dots, t_n \in I$ by generating the corresponding random vector $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$, with $\boldsymbol{\mu} = (\mu_X(t_1), \dots, \mu_X(t_n))$ and $\Sigma_{ij} = C_X(t_i, t_j)$. This can be done by either Cholesky or spectral factorization of the matrix Σ .

Similarly, assume that we have generated already $\mathbf{Z} = (X_{t_1}, \dots, X_{t_n})$ and we want to generate new values $\mathbf{Y} = (X_{t_{n+1}}, \dots, X_{t_m})$ conditional to the previously generated ones. This can be done by the Algorithm 3.3 illustrated in the previous section. Figure 3.1 gives a graphical interpretation of the procedure.

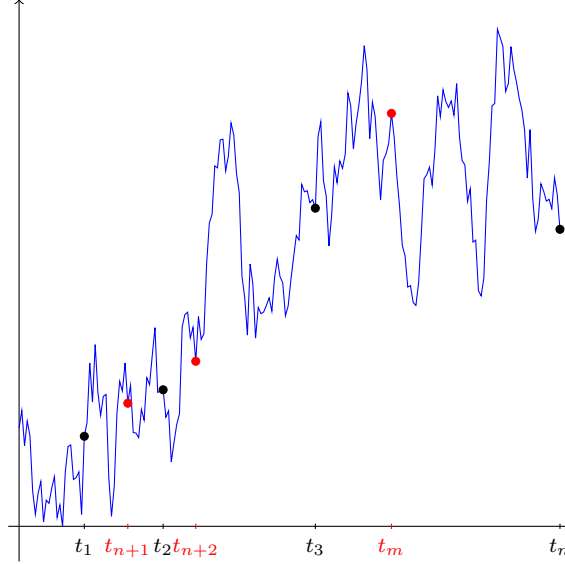


Figure 3.1: Conditioned Gaussian Process.

3.3.1 Wiener process (Brownian motion)

Definition 3.3. *The Wiener process is a Gaussian stochastic process $\{W_t, t \geq 0\}$ with the following properties:*

- $W_0 = 0$,
- *Independent increments: for all $0 < t_1 < t_2 \leq t_3 < t_4$, $(W_{t_2} - W_{t_1})$ and $(W_{t_4} - W_{t_3})$ are independent random variables*
- *Gaussian stationary increments: for all $0 \leq t_1 \leq t_2$, $W_{t_2} - W_{t_1} \sim N(0, t_2 - t_1)$*

The Wiener process is a Gaussian process with mean function $\mu_W(t) = 0$ and covariance function $\text{Cov}_W(s, t) = \min\{s, t\}$. Indeed, if $t \geq s$, $\text{Cov}_W(s, t) = \mathbb{E}[W_s W_t] = \mathbb{E}[W_s(W_t - W_s)] + \mathbb{E}[W_s^2] = s$. Similarly, if $t \leq s$ then $\text{Cov}_W(s, t) = t$. It can be shown that almost every realization $t \mapsto W_t$ is a continuous function (or, more precisely, is almost everywhere equal to a continuous function) in t . This property is referred to as pathwise continuity.

To generate $\{W_t, t \geq 0\}$ on a set of points (t_1, \dots, t_n) , one could either compute a Cholesky/spectral decomposition of the covariance matrix $\Sigma_{ij} = \min\{t_i, t_j\}$, or, more efficiently, rely on the property of independent Gaussian increments as the following Algorithm shows.

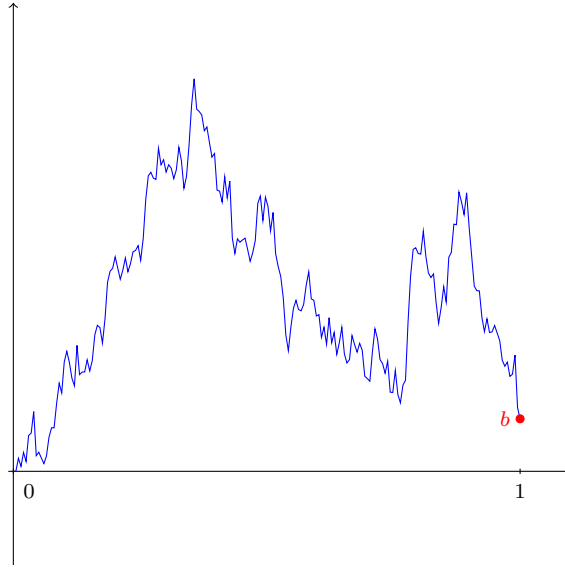


Figure 3.2: Brownian bridge.

Algorithm 3.5: Wiener process generation

- 1 Set $t_0 = 0$ and $W_{t_0} = 0$
 - 2 **for** $k = 1, \dots, n$ **do**
 - 3 Generate $\Delta W_k \sim N(0, t_k - t_{k-1})$
 - 4 Set $W_{t_k} = W_{t_{k-1}} + \Delta W_k$
 - 5 **end**
-

A Brownian motion with drift ν and diffusion coefficient σ^2 , is the solution of the stochastic differential equation

$$dB_t = \nu dt + \sigma dW_t, \quad B_0 = 0$$

that is, $B_t = \nu t + \sigma W_t$, $t \geq 0$. Hence, it can be easily generated on a set of points (t_1, \dots, t_n) as an affine transformation of the Wiener process.

3.3.2 Brownian bridge

A Brownian bridge process $\{X_t, t \in [0, 1]\}$ is a Wiener process $\{W_t, t \in [0, 1]\}$ conditioned upon $W_1 = b$. See figure 3.2 for a realization of a Brownian bridge.

The conditional mean and covariance function of a Brownian bridge can be calculated using the standard formulas for conditioned multivariate Gaussian variables. Indeed, let us first calculate the conditional mean. For that, we set $Y = W_t$, $t \in (0, 1)$ and $Z = W_1$, for which we have $\Sigma_{YY} = \Sigma_{YZ} = t$ and $\Sigma_{ZZ} = 1$. Therefore, using formula (3.1) for the conditional mean, we conclude that

$$\mu_X(t) = \mathbb{E}[X_t] = \mathbb{E}[W_t \mid W_1 = b] = \mu_W(t) + \Sigma_{YZ} \Sigma_{ZZ}^{-1} (b - \mu_W(1)) = tb.$$

An analogous procedure can be followed to compute the conditional covariance. Take this time $Y = (W_s, W_t)$, $s, t \in (0, 1)$ and $Z = W_1$, so that

$$\Sigma_{YY} = \begin{pmatrix} s & \min\{s, t\} \\ \min\{s, t\} & t \end{pmatrix}, \quad \Sigma_{YZ} = \begin{pmatrix} s \\ t \end{pmatrix}, \quad \Sigma_{ZZ} = 1.$$

Therefore

$$\begin{aligned} \Sigma_{Y|Z} &= \Sigma_{YY} - \Sigma_{YZ}\Sigma_{ZZ}^{-1}\Sigma_{ZY} \\ &= \begin{pmatrix} s & \min\{s, t\} \\ \min\{s, t\} & t \end{pmatrix} - \begin{pmatrix} s \\ t \end{pmatrix} \begin{pmatrix} s & t \end{pmatrix} = \begin{pmatrix} s - s^2 & \min\{s, t\} - st \\ \min\{s, t\} - st & t - t^2 \end{pmatrix}. \end{aligned}$$

and

$$\text{Cov}_X(s, t) = \text{Cov}(W_s, W_t | W_1 = b) = (\Sigma_{Y|Z})_{12} = \min\{s, t\} - st.$$

To generate a Brownian bridge in a set of points $0 < t_1 < \dots < t_n < t_{n+1} = 1$ one can first generate $(W_{t_1}, \dots, W_{t_n}, W_{t_{n+1}})$ from a standard Wiener process and then use Algorithm 3.3 with $\mathcal{Y} = (W_{t_1}, \dots, W_{t_n})$ and $\mathcal{Z} = W_{t_{n+1}}$. This leads to the following procedure.

Algorithm 3.6: Brownian bridge generation.

Given: $0 < t_1 < \dots < t_n < t_{n+1} = 1$ and b

- 1 Generate W_{t_i} , $i = 1, \dots, n + 1$ from standard Wiener process
 - 2 Output $X_{t_i} = W_{t_i} + t_i(b - W_{t_{n+1}})$, $i = 1, \dots, n$.
-

3.4 Stationary Gaussian processes / random fields

Definition 3.4. A Gaussian process $\{X_t, t \in \mathbb{R}\}$ is weakly stationary if $C_X(t, s)$ depends only on $(s - t)$ and is (strongly) stationary if it is weakly stationary and $\mu_X(t)$ does not depend on t .

A weakly stationary Gaussian process can be generated very efficiently on a uniform grid $\{t_j = t_0 + jh, j = 0, \dots, n\}$ with the use of FFT. This avoids the costly step of computing the Cholesky or spectral decomposition of the covariance matrix. We denote by $\mathbf{X} = (X_{t_0}, \dots, X_{t_n})$ the discrete Gaussian process on the uniform grid. Since $\{t_j, j = 0, \dots, n\}$ is a uniform grid, it follows that the corresponding covariance matrix

$$\Sigma_{ij} = C_X(t_i, t_j) = C_X(0, (j - i)h)$$

depends only on $j - i$, hence is a symmetric Toeplitz matrix and we have to store only the vector $(\sigma_0, \dots, \sigma_n)$ with $\sigma_i = C_X(t_0, t_0 + ih)$:

$$\Sigma = \begin{pmatrix} \sigma_0 & \sigma_1 & \sigma_2 & \dots & \sigma_n \\ \sigma_1 & \sigma_0 & \sigma_1 & \dots & \sigma_{n-1} \\ \sigma_2 & \sigma_1 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ \sigma_n & \sigma_{n-1} & \dots & \sigma_1 & \sigma_0 \end{pmatrix}$$

Consider now the following circulant embedding of Σ :

$$\tilde{\Sigma} = \begin{pmatrix} \sigma_0 & \sigma_1 & \sigma_2 & \dots & \sigma_{n-1} & \sigma_n & \sigma_{n-1} & \sigma_{n-2} & \dots & \sigma_1 \\ \sigma_1 & \sigma_0 & \sigma_1 & \dots & & \sigma_{n-1} & \sigma_n & \sigma_{n-1} & \dots & \sigma_2 \\ \sigma_2 & \sigma_1 & \ddots & \ddots & & \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & & \vdots & \vdots & & \ddots & \sigma_{n-1} \\ \vdots & \vdots & & & \ddots & \sigma_1 & \sigma_2 & & & \sigma_n \\ \sigma_n & \sigma_{n-1} & \dots & & \sigma_1 & \sigma_0 & \sigma_1 & \dots & \dots & \sigma_{n-1} \\ \hline \sigma_{n-1} & \sigma_n & \sigma_{n-1} & \dots & \dots & \sigma_1 & \sigma_0 & \sigma_1 & \dots & \sigma_{n-2} \\ \sigma_{n-2} & \sigma_{n-1} & \ddots & \ddots & & \sigma_2 & \sigma_1 & \sigma_0 & \dots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots \\ \sigma_1 & \dots & \dots & \sigma_{n-1} & \sigma_n & \sigma_{n-1} & \sigma_{n-2} & & \sigma_1 & \sigma_0 \end{pmatrix} \in \mathbb{R}^{2n \times 2n}$$

and the generating vector $\boldsymbol{\alpha} = (\sigma_0, \sigma_1, \dots, \sigma_n, \sigma_{n-1}, \dots, \sigma_1)$ given by the first column. We write compactly $\tilde{\Sigma} = \text{circ}(\boldsymbol{\alpha})$ and assume that $\tilde{\Sigma}$ is also non-negative definite.

Lemma 3.2. *Let $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{2n}) \in \mathbb{R}^{2n}$ and $\tilde{\Sigma} = \text{circ}(\boldsymbol{\alpha})$. Then the vectors $\mathbf{v}^{(\ell)}$, $\ell = 1, \dots, 2n$, $v_k^{(\ell)} = e^{2\pi i(\ell-1)(k-1)/2n}$ are eigenvectors of $\tilde{\Sigma}$ with corresponding eigenvalues $\lambda_\ell = \sum_{k=1}^{2n} \alpha_k e^{-2\pi i(\ell-1)(k-1)/2n}$, which are real and non-negative if $\tilde{\Sigma}$ is semi positive definite.*

Proof. It is enough to verify that $\sum_{k=1}^{2n} \tilde{\Sigma}_{jk} v_k^{(\ell)} = \lambda_\ell v_j^{(\ell)}$, for all $j = 1, \dots, 2n$. Notice that $\tilde{\Sigma}$ can be written as $\tilde{\Sigma}_{jk} = \alpha_{\{(2n+j-k+1) \bmod 2n\}}$, where we set $\alpha_0 = \alpha_{2n}$. Then

$$\begin{aligned} \sum_{k=1}^{2n} \tilde{\Sigma}_{jk} v_k^{(\ell)} &= \sum_{k=1}^{2n} \alpha_{\{(2n+j-k+1) \bmod 2n\}} e^{2\pi i(\ell-1)(k-1)/2n} \\ &= \sum_{k=1}^{2n} \alpha_{\{(2n+j-k+1) \bmod 2n\}} e^{2\pi i(\ell-1)(k-j-2n)/2n} e^{2\pi i(\ell-1)(2n+j-1)/2n} \\ &= \left(\sum_{k=1}^{2n} \alpha_k e^{-2\pi i(\ell-1)(k-1)/2n} \right) v_j^{(\ell)}. \end{aligned}$$

□

It follows from this Lemma that the matrix $\tilde{\Sigma}$ can be diagonalized as $\tilde{\Sigma} F^* = F^* \Lambda$ where $F_{k\ell} = e^{-2\pi i(\ell-1)(k-1)/2n}$ corresponds to the FFT matrix as defined in Matlab (and numpy up to a shift of indices $\ell - 1 \rightarrow \ell$, $k - 1 \rightarrow k$) and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_{2n})$. Moreover, the vector of eigenvalues corresponds to $\boldsymbol{\lambda} = F \boldsymbol{\alpha} = \text{FFT}(\boldsymbol{\alpha})$. Observe that $F^* F = F F^* = 2n I_{2n}$ so that $\tilde{\Sigma} = \frac{1}{2n} F^* \Lambda F$ and can be factorized as $\tilde{\Sigma} = A A^*$ with $A = \frac{1}{\sqrt{2n}} F^* \Lambda^{1/2}$. We consider now a vector $\mathbf{Y} = (Y_1, \dots, Y_{2n})$ of complex standard normal r.v.s, i.e. $\mathbf{Y} = \mathbf{Y}_R + i\mathbf{Y}_I$ with $\mathbf{Y}_R, \mathbf{Y}_I \stackrel{\text{iid}}{\sim} N(0, I_{2n})$, and set

$$\tilde{\mathbf{X}} = A \mathbf{Y} = \frac{1}{\sqrt{2n}} F^* \Lambda^{1/2} \mathbf{Y} = \text{iFFT}(\sqrt{2\pi} \Lambda^{1/2} \mathbf{Y}),$$

where the iFFT matrix is given by $\frac{1}{2n}F^*$. The following holds:

- $\mathbb{E}[\mathbf{Y}\mathbf{Y}^*] = 2I_{2n} = \mathbb{E}[\bar{\mathbf{Y}}\bar{\mathbf{Y}}^*]$, $\mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \mathbb{E}[\bar{\mathbf{Y}}\bar{\mathbf{Y}}^\top] = 0$,
- $\mathbb{E}[\tilde{\mathbf{X}}\tilde{\mathbf{X}}^*] = \mathbb{E}[A\mathbf{Y}\mathbf{Y}^*A^*] = 2\tilde{\Sigma} = \mathbb{E}[\tilde{\tilde{\mathbf{X}}}\tilde{\tilde{\mathbf{X}}}^*]$, $\mathbb{E}[\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top] = \mathbb{E}[\tilde{\tilde{\mathbf{X}}}\tilde{\tilde{\mathbf{X}}}^\top] = 0$,
- $\mathbb{E}[\operatorname{Re}(\tilde{\mathbf{X}})\operatorname{Re}(\tilde{\mathbf{X}})^\top] = \mathbb{E}\left[\frac{\tilde{\mathbf{X}}+\tilde{\tilde{\mathbf{X}}}}{2}\left(\frac{\tilde{\mathbf{X}}+\tilde{\tilde{\mathbf{X}}}}{2}\right)^\top\right] = \tilde{\Sigma} = \mathbb{E}[\operatorname{Im}(\tilde{\mathbf{X}})\operatorname{Im}(\tilde{\mathbf{X}})^\top]$,
- $\mathbb{E}[\operatorname{Re}(\tilde{\mathbf{X}})\operatorname{Im}(\tilde{\mathbf{X}})^\top] = \mathbb{E}\left[\frac{\tilde{\mathbf{X}}+\tilde{\tilde{\mathbf{X}}}}{2}\left(\frac{\tilde{\mathbf{X}}-\tilde{\tilde{\mathbf{X}}}}{2i}\right)^\top\right] = 0$.

Hence $\operatorname{Re}(\tilde{\mathbf{X}}), \operatorname{Im}(\tilde{\mathbf{X}}) \sim N(0, \tilde{\Sigma})$ and are independent and $\operatorname{Re}(\tilde{\mathbf{X}}_{1:n+1}), \operatorname{Im}(\tilde{\mathbf{X}}_{1:n+1}) \stackrel{\text{iid}}{\sim} N(0, \Sigma)$. This suggests the following algorithm to generate $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$.

Algorithm 3.7: Circulant embedding.

Given: $\boldsymbol{\mu} \in \mathbb{R}^n$ and $\Sigma = \begin{pmatrix} \sigma_0 & \sigma_1 & \dots & \sigma_n \\ \sigma_1 & \sigma_0 & \dots & \sigma_{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_n & \dots & \dots & \sigma_0 \end{pmatrix} \in \mathbb{R}^{n+1 \times n+1}$

- 1 Generate the vector $\boldsymbol{\alpha} = (\sigma_0, \sigma_1, \dots, \sigma_n, \sigma_{n-1}, \dots, \sigma_1) \in \mathbb{R}^{2n}$
 - 2 Compute $\boldsymbol{\lambda} = \text{FFT}(\boldsymbol{\alpha})$
 - 3 Generate $\mathbf{Y} = \mathbf{Y}_R + i\mathbf{Y}_I$ with $\mathbf{Y}_R, \mathbf{Y}_I \stackrel{\text{iid}}{\sim} N(0, I_{2n})$
 - 4 Compute $\tilde{\mathbf{X}} = \text{iFFT}(\sqrt{2n} \operatorname{diag}(\sqrt{\boldsymbol{\lambda}})\mathbf{Y})$
 - 5 Output $\mathbf{X}^{(1)} = \boldsymbol{\mu} + \operatorname{Re}(\tilde{\mathbf{X}}_{1:n+1})$ and $\mathbf{X}^{(2)} = \boldsymbol{\mu} + \operatorname{Im}(\tilde{\mathbf{X}}_{1:n+1})$
-

One may encounter the problem that the matrix $\tilde{\Sigma}$ might not be semi positive definite, even if Σ is. In such a case, one could try to enlarge the circulant embedding

$$\boldsymbol{\alpha} = (\sigma_0, \sigma_1, \dots, \sigma_n, \sigma_{n+1}^*, \dots, \sigma_m^*, \sigma_{m-1}^*, \dots, \sigma_{n+1}^*, \sigma_n, \dots, \sigma_1)$$

where $m > n$ and $\sigma_j^*, j = n+1, \dots, m$ are chosen such that $\tilde{\Sigma} = \text{circ}(\boldsymbol{\alpha})$ is semi positive definite. A typical choice is to take $\sigma_j^* = \sigma_j = C_X(0, jh)$ and m large enough.

Chapter 4

Generation of Markov processes

4.1 Discrete time / discrete state Markov chains

Let us consider a stochastic process $\{X_n, n \in \mathbb{N}_0\}$ defined on the countable set $\mathbb{N}_0 = \{0, 1, \dots\}$ and taking values in a countable set $\mathcal{X} = \{y_1, y_2, \dots\}$ i.e. $X_n \in \mathcal{X}$ for all $n \in \mathbb{N}_0$.

Definition 4.1. A stochastic process $\{X_n \in \mathcal{X}, n \in \mathbb{N}_0\}$ is a Markov chain if it satisfies the Markov property

$$\mathbb{P}(X_{n+1} = y_{n+1} \mid X_n = y_n, X_{n-1} = y_{n-1}, \dots, X_0 = y_0) = \mathbb{P}(X_{n+1} = y_{n+1} \mid X_n = y_n)$$

with $y_0, \dots, y_{n+1} \in \mathcal{X}$.

The process is therefore entirely defined by the distribution λ of the initial state X_0 and the transition matrices

$$P(n) = (P_{ij}(n))_{ij}, \quad \text{with } P_{ij}(n) = \mathbb{P}(X_n = y_j \mid X_{n-1} = y_i).$$

which are, in particular, *stochastic matrices* i.e. they satisfy

$$\sum_j P_{ij}(n) = 1, \quad \forall i = 1, 2, \dots, \quad \forall n \in \mathbb{N}.$$

A Markov chain is time-homogeneous if $P(n)$ does not depend on n . Generating a discrete time / discrete state Markov chain is rather straightforward.

Algorithm 4.1: Generation of discrete time / discrete space Markov process.

Given: λ and $P(n), n \in \mathbb{N}$

1 Generate $X_0 \sim \lambda$

2 For $n = 1, 2, \dots$,

3 Generate $X_n \sim P_{X_{n-1}, \cdot}(n)$ // pmf given by X_{n-1} -th row of $P(n)$

Exercise 4.1 (Random walk on a lattice). *A random walk on the integers $\{X_n \in \mathbb{Z}, n \in \mathbb{N}_0\}$, starting at $X_0 = 0$ is a Markov chain defined by the following transition probabilities*

$$\begin{aligned}\mathbb{P}(X_{n+1} = j \mid X_n = j - 1) &= \mathbb{P}(X_{n+1} = j \mid X_n = j + 1) = a \in (0, 1), \\ \mathbb{P}(X_{n+1} = j \mid X_n = j) &= 1 - 2a, \\ \mathbb{P}(X_{n+1} = j \mid X_n = i) &= 0, \quad i \neq j, j - 1, j + 1.\end{aligned}$$

Figure 4.1 shows a graph representation of the Markov chain. An arrow between two states denotes a connection, i.e. a non zero probability of moving from the base to the head of the arrow.

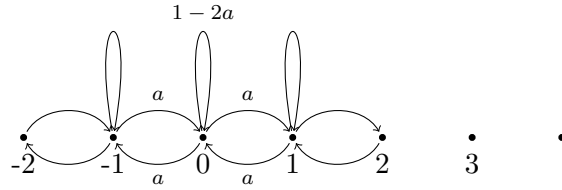


Figure 4.1: Random walk on lattice.

4.2 Discrete time / continuous state Markov chains

Consider now a stochastic process $\{X_n, n \in \mathbb{N}_0\}$ defined on $\mathbb{N}_0 = \{0, 1, \dots\}$ and taking values on a continuous set $\mathcal{X} \subset \mathbb{R}^d$. We denote by $\mathcal{B}(\mathcal{X})$ the Borel σ -algebra on \mathcal{X} so that $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ is a measurable space.

Definition 4.2. *A Markov transition kernel on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ is a function $P : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ such that*

- for all $y \in \mathcal{X}$, $P(y, \cdot)$ is a probability measure on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$;
- for all $A \in \mathcal{B}(\mathcal{X})$, $P(\cdot, A)$ is a measurable function on \mathcal{X} .

Often, the transition kernel is defined starting from a density function $p : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ such that for all $x \in \mathcal{X}$, $A \in \mathcal{B}(\mathcal{X})$, $P(x, A) = \int_A p(x, y) dy$.

Definition 4.3. *Given a Markov transition kernel $P : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$, a stochastic process $\{X_n, n \in \mathbb{N}_0\}$ with values in \mathcal{X} is a homogeneous Markov chain with kernel P and initial distribution $X_0 \sim \lambda$, denoted $\{X_n\} \sim \text{Markov}(\lambda, P)$, if for any $n \in \mathbb{N}_0$, $A \in \mathcal{B}(\mathcal{X})$,*

$$\mathbb{P}(X_{n+1} \in A \mid X_n = y_n, \dots, X_0 = y_0) = \mathbb{P}(X_{n+1} \in A \mid X_n = y_n) = P(y_n, A).$$

Again, generating a discrete time / continuous state Markov chain is rather straightforward, provided we know how to generate random variables from the probability measure $P(y, \cdot)$ (resp. probability density function $p(y, \cdot)$) for all $y \in \mathcal{X}$.

Algorithm 4.2: Generation of discrete time / continuous space Markov process.

Given: λ and P

- 1 Generate $X_0 \sim \lambda$
 - 2 For $n = 1, 2, \dots$,
 - 3 Generate $X_n \sim P(X_{n-1}, \cdot)$
-

Exercise 4.2 (Random walk in 2D). Let $\mathcal{X} = \mathbb{R}^2$ and consider the stochastic process $\{\mathbf{X}_n \in \mathcal{X}, n \in \mathbb{N}_0\}$ starting at $\mathbf{X}_0 = (0, 0)$, defined by

$$\mathbf{X}_{n+1} = \mathbf{X}_n + \boldsymbol{\xi}_n, \quad \boldsymbol{\xi}_n \stackrel{iid}{\sim} N(0, \sigma^2 I_2).$$

This is clearly a homogeneous discrete time / continuous state Markov chain with transition kernel

$$P(\mathbf{y}, A) = \mathbb{P}(\mathbf{X}_{n+1} \in A \mid \mathbf{X}_n = \mathbf{y}) = \mathbb{P}(\boldsymbol{\xi}_n + \mathbf{y} \in A) = \frac{1}{2\pi\sigma^2} \int_A e^{-\frac{\|\boldsymbol{\xi} - \mathbf{y}\|^2}{2\sigma^2}} d\boldsymbol{\xi}$$

and transition density function $p(\mathbf{x}, \mathbf{y}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{(y_1 - x_1)^2 + (y_2 - x_2)^2}{2\sigma^2}\right)$.

4.3 Continuous time / discrete state Markov chains

Let $\mathcal{X} = \{y_1, y_2, \dots\}$ be a discrete (finite or countable) set and $\{X_t, t \geq 0\}$ a stochastic process taking values in \mathcal{X} . The process is said to be right continuous if each path is so, i.e. for any realization ω ,

$$\lim_{h \rightarrow 0^+} X_{t+h}(\omega) = X_t(\omega).$$

Since the process takes only discrete values, the right continuity property implies that if $X_t = y_i$ at some t , it will stay in state y_i for a certain amount of time, i.e. there exists a (random) $\varepsilon > 0$ s.t. $X_s = y_i$, for all $t \leq s < t + \varepsilon$. We denote J_n the n -th jump time

$$J_0 = 0, \quad J_n = \inf\{t \geq J_{n-1} : X_t \neq X_{J_{n-1}}\}, \quad n > 0$$

and S_n the n -th holding time

$$S_n = \begin{cases} J_n - J_{n-1}, & \text{if } J_{n-1} < \infty, \\ \infty, & \text{otherwise.} \end{cases} \quad n = 1, 2, \dots$$

The discrete time process $\{Y_n = X_{J_n}, n \in \mathbb{N}_0\}$ is called the *jump process* (or jump chain) of $\{X_t, t \geq 0\}$. The process is therefore completely characterized by the sequence $\{J_n\}_n$ of jump times (equivalently the sequence $\{S_n\}_n$ of holding times) as well as the sequence $\{Y_n\}_n$ of visited states, i.e. the jump chain. Figure 4.2 gives an illustration of a continuous time / discrete state Markov process.

The (first) explosion time T^* is defined as $T^* = \sup_n J_n = \sum_{n=1}^{\infty} S_n$. If $T^* < +\infty$, we consider only the process $\{X_t, t \in [0, T^*)\}$ or, equivalently, we set $X_t = \infty$ for $t \geq T^*$. This is called the *minimal process*.

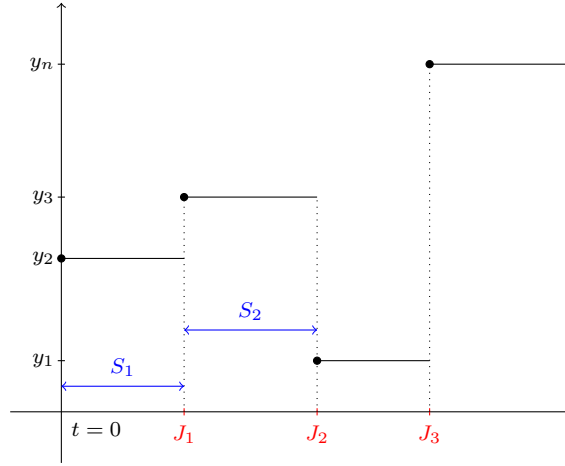


Figure 4.2: Continuous time / discrete state Markov process and associated jump and holding times.

4.4 Poisson process

The Poisson process is the simplest example of a continuous time / discrete state *Markov* process.

Definition 4.4. A Poisson process $\{N_t \in \mathbb{N}_0, t \geq 0\}$ with initial state $N_0 = 0$ and parameter $0 < \lambda < \infty$, is a non decreasing, right-continuous, integer valued process which satisfies the following properties.

1. *Independent increments:* for all $0 < t_1 < t_2 \leq t_3 < t_4$, $N_{t_2} - N_{t_1}$ and $N_{t_4} - N_{t_3}$ are independent;
2. *Poisson stationary increments:* for all $0 < s < t$, $N_t - N_s \sim \text{Pois}(\lambda(t - s))$ i.e.

$$\mathbb{P}(N_t - N_s = j) = \frac{(\lambda(t - s))^j}{j!} e^{-\lambda(t - s)}.$$

It follows, in particular, that $N_t \sim \text{Pois}(\lambda t)$. Moreover, N_t satisfies the Markov property: for any $s \geq 0$, $\tilde{N}_t = N_{s+t} - N_s$, $t \geq 0$ is also a Poisson process of rate λ , independent of $\{N_t, t \leq s\}$, as well as the strong Markov property where s is replaced by a stopping time T . (T is a stopping time if the event $\{T \leq t\}$ is measurable with respect to the σ -algebra \mathcal{F}_t generated by $\{N_s, s \leq t\}$). The following are two equivalent characterizations of a Poisson process:

- a. For any $t > 0$ and $h \rightarrow 0^+$, uniformly in t it holds

$$\mathbb{P}(N_{t+h} - N_t = 0) = 1 - \lambda h + o(h),$$

$$\mathbb{P}(N_{t+h} - N_t = 1) = \lambda h + o(h),$$

$$\mathbb{P}(N_{t+h} - N_t > 1) = o(h).$$

The last condition is actually a consequence of the first two.

- b. The holding times S_1, S_2, \dots are independent exponential random variables $\text{Exp}(\lambda)$ and the jump chain is $Y_n = N_{J_n} = n$.

The first property follows immediately from the Poisson distribution of the increments. For the second property, observe that $\mathbb{P}(S_1 > t) = \mathbb{P}(N_t = 0) = e^{-\lambda t}$ hence $S_1 \sim \text{Exp}(\lambda)$. Similarly, $\mathbb{P}(S_{n+1} > t) = \mathbb{P}(N_{J_n+t} - N_{J_n} = 0) = e^{-\lambda t}$ so $S_{n+1} \sim \text{Exp}(\lambda)$ and independent of S_1, \dots, S_n by the property of independent increments of N_t . The second property suggests an easy algorithm to generate a Poisson process with parameter λ .

Algorithm 4.3: Homogeneous Poisson process I.

- 1 Set $N_0 = 0, J_0 = 0, Y_0 = 0$
 - 2 For $n = 1, 2, \dots$,
 - 3 Generate $S_n \sim \text{Exp}(\lambda)$ and set $J_n = J_{n-1} + S_n$
 - 4 Set $N_t = N_{J_{n-1}}, t \in [J_{n-1}, J_n)$ and $N_{J_n} = N_{J_{n-1}} + 1$.
-

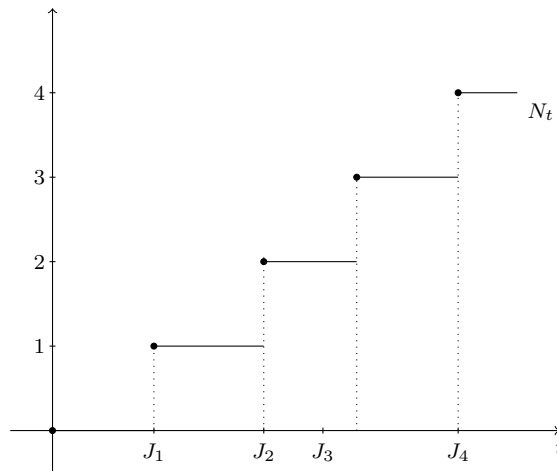


Figure 4.3: Homogeneous Poisson process.

Figure 4.3 shows a realization of a Poisson process. Another useful property of the Poisson process is the following.

- c. Conditional on $N_t = n$, the n jump times are uniformly distributed in $(0, t)$, i.e. J_1, \dots, J_n have the same distribution of the order statistics $U_{(1)}, \dots, U_{(n)}$ with $U_i \stackrel{\text{iid}}{\sim} \mathcal{U}(0, t)$.

Property c. suggests an alternative algorithm to generate a Poisson process of rate λ on $[0, T]$.

Algorithm 4.4: Homogeneous Poisson process II.

- 1 Generate $N_T \sim \text{Pois}(\lambda T)$
 - 2 Generate $U_1, \dots, U_{N_T} \stackrel{\text{iid}}{\sim} \mathcal{U}(0, T)$
 - 3 Order the sample $U_{(1)} < \dots < U_{(N_T)}$
 - 4 Set $J_0 = 0, J_n = U_{(n)}$, and $N_t = n, t \in [J_n, J_{n+1}), n = 1, \dots, N_T$
-

Finally we mention that a Poisson process $\{N_t, t \geq 0\}$ can also be thought of as a random counting measure. For a given interval $A = (t_1, t_2)$, $\mu(A) = \sum_{k=1}^{\infty} \mathbb{1}_{\{J_k \in A\}}$ counts the number of jumps that occurred in A (which is clearly a random number). Thus $d\mu(t) = \sum_{k=1}^{\infty} \delta_{J_k}(t)$ and it holds $N_t = N_0 + \int_0^t d\mu(t)$.

4.5 Non-homogeneous Poisson process

A non-homogeneous Poisson process with rate $\lambda(t)$ varying over time can be defined by extending the property b. of a Poisson process.

Definition 4.5. $\{N_t, t \geq 0, N_0 = 0\}$ is a non-homogeneous Poisson process with rate $\lambda : [0, \infty) \rightarrow \mathbb{R}_+$ if it is a right-continuous process with independent increments, such that

$$\begin{aligned}\mathbb{P}(N_{t+h} - N_t = 0) &= 1 - \lambda(t)h + o(h), \\ \mathbb{P}(N_{t+h} - N_t = 1) &= \lambda(t)h + o(h).\end{aligned}$$

Therefore, the non-homogeneous rate $\lambda(t)$ can be characterized by the following limits

$$\lambda(t) = \lim_{h \rightarrow 0^+} \frac{1 - \mathbb{P}(N_{t+h} - N_t = 0)}{h} = \lim_{h \rightarrow 0^+} \frac{\mathbb{P}(N_{t+h} - N_t = 1)}{h}.$$

To be able to generate a non-homogeneous Poisson process we need to derive the distribution of the holding times. This is shown in the next lemma.

Lemma 4.1. Let $\{N_t, t \geq 0, N_0 = 0\}$ is a non-homogeneous Poisson process with rate $\lambda : [0, \infty) \rightarrow \mathbb{R}_+$ and denote by F the cdf of the $n+1$ holding time $F_{n+1}(t) = \mathbb{P}(S_{n+1} \leq t)$. It holds

$$F_{n+1}(t) = 1 - \exp \left\{ - \int_{J_n}^{J_n+t} \lambda(s) ds \right\}.$$

Proof. We have

$$\begin{aligned}F'_{n+1}(t) &= \lim_{h \rightarrow 0} \frac{F_{n+1}(t+h) - F_{n+1}(t)}{h} \\ &= \lim_{h \rightarrow 0} \frac{\mathbb{P}(t < S_{n+1} \leq t+h)}{h} = \lim_{h \rightarrow 0} \frac{\mathbb{P}(S_{n+1} \leq t+h \mid S_{n+1} > t)}{h} (1 - F_{n+1}(t)) \\ &= \lim_{h \rightarrow 0} \frac{\mathbb{P}(N_{J_n+t+h} > n \mid N_{J_n+t} = n)}{h} (1 - F_{n+1}(t)) \\ &= \lim_{h \rightarrow 0} \frac{1 - \mathbb{P}(N_{J_n+t+h} = n \mid N_{J_n+t} = n)}{h} (1 - F_{n+1}(t)) \\ &= \lambda(J_n + t)(1 - F_{n+1}(t)).\end{aligned}$$

Solving this differential equation with initial condition $F_{n+1}(0) = 0$ leads to the desired result. \square

Hence, a non-homogeneous Poisson process with rate function $\lambda(t)$ can be generated by the following Algorithm.

Algorithm 4.5: Non-homogeneous Poisson process.

- 1 Set $N_0 = 0, J_0 = 0, Y_0 = 0$
 - 2 For $n = 1, 2, \dots$
 - 3 Generate $S_n \sim F_n(t) = 1 - \exp \left\{ - \int_{J_{n-1}}^{J_{n-1}+t} \lambda(s) ds \right\}$
 - 4 Set $J_n = J_{n-1} + S_n,$
 - 5 Set $N_t = N_{J_{n-1}}, t \in [J_{n-1}, J_n),$
 - 6 Set $N_{J_n} = N_{J_{n-1}} + 1$
-

If we define the function $\Lambda(t) = \int_0^t \lambda(s) ds,$ and let \tilde{N}_t be a homogeneous Poisson process with rate 1, it can also be shown (exercise) that the non homogeneous Poisson process N_t with rate function $\lambda(t)$ can be obtained as $N_t = \tilde{N}_t \circ \Lambda = \tilde{N}_{\Lambda(t)}.$ In particular the increments are independent, i.e. $\forall v \leq r \leq s \leq t, N_t - N_s$ is independent of $N_r - N_v,$ and $N_t - N_s$ has Poisson distribution $\text{Pois}(\Lambda(t) - \Lambda(s)).$

4.6 Compound Poisson process

A compound Poisson process $\{X_t, t \geq 0, X_0 = 0\}$ is a Poisson process with variable jump intensity. Let $\nu(dy)$ be a probability measure on \mathbb{R} and $\{N_t, t \geq 0\}$ a homogeneous Poisson process with rate $\lambda > 0.$ Then, the compound Poisson process with jump measure $\lambda\nu(dy)dt$ is given by

$$X_t = \sum_{i=1}^{N_t} Z_i, \quad Z_i \stackrel{\text{iid}}{\sim} \nu.$$

Algorithm 4.6: Compound Poisson process.

- 1 Set $N_0 = 0, J_0 = 0, Y_0 = 0$
 - 2 For $n = 1, 2, \dots$
 - 3 Generate $S_n \sim \text{Exp}(\lambda)$ and set $J_n = J_{n-1} + S_n,$
 - 4 Generate $Z_n \sim \nu,$
 - 5 Set $X_t = X_{J_{n-1}}, t \in [J_{n-1}, J_n)$ and $X_{J_n} = X_{J_{n-1}} + Z_n$
-

4.7 General continuous time / discrete space Markov process

Let $\mathcal{X} = \{y_1, y_2, \dots\}$ be a discrete (finite or countable) set and let $\mu = \{\mu_i\}_i$ be a probability mass function on $\mathcal{X},$ i.e. $\mu_i \geq 0, \forall i$ and $\sum_i \mu_i = 1.$ A continuous time Markov chain $\{X_t \in \mathcal{X}, t \geq 0\}$ with initial state $X_0 \sim \mu,$ is fully characterized by the

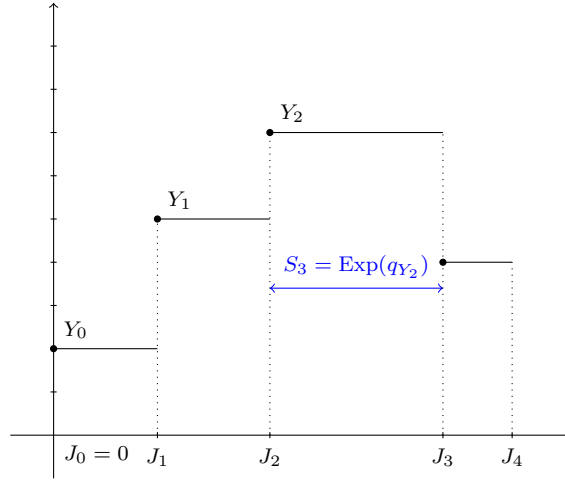


Figure 4.4: Homogeneous continuous time Markov process.

transition probabilities

$$q_{ij}(t) = \lim_{h \rightarrow 0^+} \frac{\mathbb{P}(X_{t+h} = y_j \mid X_t = y_i)}{h}$$

$$q_i(t) = \lim_{h \rightarrow 0^+} \frac{1 - \mathbb{P}(X_{t+h} = y_i \mid X_t = y_i)}{h}$$

If $q_i(t)$ and $q_{ij}(t)$ do not depend on t , the Markov chain is homogeneous. The (possibly infinite) matrix $Q = (Q_{ij})_{ij}$ given by

$$Q_{ij} = \begin{cases} q_{ij} & i \neq j \\ -q_i & i = j \end{cases}$$

is called the generator of the Markov chain. We assume here that Q is *stable*, i.e. $q_i < \infty$ for all i and *conservative*, i.e. $\sum_{j \neq i} q_{ij} = q_i$.

Definition 4.6. A homogeneous continuous time Markov chain $\{X_t \in \mathcal{X}, t \geq 0\}$ with initial state $X_0 \sim \mu$ and (stable and conservative) generator matrix Q , is a right-continuous, piecewise constant process denoted Markov (μ, Q) s.t.

- the jump process $\{Y_n = X_{J_n}, n \in \mathbb{N}_0\}$ is a discrete time Markov chain with transition probability

$$\pi_{ij} = \frac{q_{ij}}{q_i}, \quad i \neq j, \quad \pi_{ii} = 0, \quad \text{if } q_i \neq 0$$

$$\pi_{ij} = 0, \quad i \neq j, \quad \pi_{ii} = 1, \quad \text{if } q_i = 0.$$

- conditional on Y_0, Y_1, \dots, Y_{n-1} , the holding times S_1, \dots, S_n are independent random variables, $S_i \sim \text{Exp}(q_{Y_{i-1}})$, $i = 1, \dots, n$.

Notice that in this case, the holding time S_n depends on the current state of the chain $S_n \sim \text{Exp}(q_{X_{J_{n-1}}})$ and the chain can jump to any other state j with transition probability $\pi_{X_{J_{n-1}},j}$. An algorithm to generate such a process Markov (μ, Q) is given next.

Algorithm 4.7: Markov (μ, Q) .

- 1 Generate $X_0 \sim \mu$ and set $J_0 = 0, Y_0 = X_0$
 - 2 For $n = 1, 2, \dots$
 - 3 Generate $S_n \sim \text{Exp}(-Q_{Y_n Y_n})$ and set $J_n = J_{n-1} + S_n$,
 - 4 Generate $Y_{n+1} \sim \pi_{Y_n}$.
 - 5 Set $X_t = Y_n, t \in [J_{n-1}, J_n)$, and $X_{J_{n+1}} = Y_{n+1}$
-

The importance of the matrix Q is illustrated by the following calculation. Let us denote $p_i(t) = \mathbb{P}(X_t = y_i)$ the probability of finding the chain in state y_i at time t and $\mathbf{p}(t) = (p_1(t), p_2(t), \dots)$ the (row) vector of such probabilities. Then

$$\begin{aligned} \frac{dp_j}{dt}(t) &= \lim_{h \rightarrow 0^+} \frac{p_j(t+h) - p_j(t)}{h} = \lim_{h \rightarrow 0^+} \frac{1}{h} (\mathbb{P}(X_{t+h} = y_j) - p_j(t)) \\ &= \lim_{h \rightarrow 0^+} \frac{1}{h} \left(\sum_{i \neq j} \mathbb{P}(X_{t+h} = y_j \mid X_t = y_i) p_i(t) + \mathbb{P}(X_{t+h} = y_j \mid X_t = y_j) p_j(t) - p_j(t) \right) \\ &= \sum_{i \neq j} q_{ij}(t) p_i(t) - q_j(t) p_j(t) = \sum_i p_i(t) Q_{ij}(t) \end{aligned}$$

from which we deduce the following differential equation

$$\frac{d}{dt} \mathbf{p}(t) = \mathbf{p}(t) Q(t)$$

for the evolution of the probability vector \mathbf{p} .

Exercise 4.3. The Poisson process $\{N_t, t \geq 0, N_0 = 0\}$ with rate $\lambda > 0$ is a continuous time Markov chain Markov (δ_0, Q) with Q -matrix

$$Q = \begin{bmatrix} -\lambda & \lambda & 0 & \dots \\ 0 & -\lambda & \lambda & \dots \\ 0 & \ddots & \ddots & \ddots \end{bmatrix}$$

since

$$\begin{aligned} Q_{ii} &= -\lambda = \lim_{h \rightarrow 0^+} \frac{\mathbb{P}(N_{t+h} = i \mid N_t = i) - 1}{h} \\ Q_{i,i+1} &= \lambda = \lim_{h \rightarrow 0^+} \frac{\mathbb{P}(N_{t+h} = i+1 \mid N_t = i)}{h} \\ Q_{i,j} &= 0, \quad j \neq i, i+1. \end{aligned}$$

Exercise 4.4 (Birth process). Let $X_t \in \mathbb{N}$ be the size of a population at time t . Births of new individuals arrive after exponential time with rate λX_t proportional to the actual

size of the population. Hence, the birth process is characterised by the Q -matrix

$$Q = \begin{bmatrix} -\lambda & \lambda & & & \\ & -2\lambda & 2\lambda & & \\ & & -3\lambda & 3\lambda & \\ & & & \ddots & \ddots \\ & & & & \ddots & \ddots \end{bmatrix}.$$

Chapter 5

Monte Carlo method

Let us consider a random variable Z that is the output quantity of a stochastic model and the goal of computing its expectation: $\mu = \mathbb{E}[Z]$. We assume that the probability distribution of Z is not known analytically, but Z can be simulated.

Example 5.1. Consider a continuous time stochastic process $\{X_t, t \geq 0\}$ with values in a subset $\mathcal{X} \subset \mathbb{R}^d$ and the goal of computing the expectation of X_T at a given time $T \geq 0$, i.e. $Z = X_T$, or the expectation of a stopping time $Z = \inf\{t \geq 0 : X_t \in A \subset \mathcal{X}\}$ for some measurable set A .

The Monte Carlo method consists simply in generating N i.i.d. replicas $Z^{(1)}, \dots, Z^{(N)}$ of Z and estimating $\mu = \mathbb{E}[Z]$ by a *sample mean estimator*

$$\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N Z^{(i)}.$$

We assume here that $\text{Var}(Z) = \sigma^2 < +\infty$.

5.1 Properties of the Monte Carlo estimator and confidence intervals

The sample mean estimator $\hat{\mu}_N$, which we will call also the *Monte Carlo estimator*, has the following properties.

1. $\hat{\mu}_N$ is *unbiased*, i.e. $\mathbb{E}[\hat{\mu}_N] = \mu$.

The expectation here is taken with respect to the distribution of the sample $(Z^{(1)}, \dots, Z^{(N)})$.

2. $\text{Var}(\hat{\mu}_N) = \frac{\sigma^2}{N}$.

Indeed,

$$\begin{aligned}\text{Var}(\hat{\mu}_N) &= \mathbb{E}[(\hat{\mu}_N - \mathbb{E}[\hat{\mu}_N])^2] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{i=1}^N (Z^{(i)} - \mu)\right)^2\right] \\ &= \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{E}[(Z^{(i)} - \mu)(Z^{(j)} - \mu)] \\ &= \frac{1}{N^2} \sum_{i=1}^N \underbrace{\mathbb{E}[(Z^{(i)} - \mu)^2]}_{=\sigma^2 \forall i \text{ since } Z^{(i)} \text{ are iid}} + \frac{1}{N^2} \sum_{i \neq j} \underbrace{\mathbb{E}[(Z^{(i)} - \mu)(Z^{(j)} - \mu)]}_{=0 \text{ since } Z^{(i)}, Z^{(j)} \text{ are indept.}} = \frac{\sigma^2}{N}.\end{aligned}$$

3. *Almost sure convergence*: $\hat{\mu}_N \xrightarrow{N \rightarrow \infty} \mu$ a.s.

This comes from the Strong Law of Large Numbers (SLLN), since $\mathbb{E}[Z] = \mu < \infty$.

4. *Asymptotic normality*

$$\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\sigma} \xrightarrow{d} \mathcal{N}(0, 1)$$

where \xrightarrow{d} means convergence in distribution. This comes from the Central Limit Theorem (CLT), since $\text{Var}(Z) < +\infty$.

Using the CLT, we can construct an asymptotic $1 - \alpha$ confidence interval (i.e. an interval with coverage probability $1 - \alpha$)

$$I_{\alpha, N} = \left[\hat{\mu}_N - c_{1-\alpha/2} \frac{\sigma}{\sqrt{N}}, \hat{\mu}_N + c_{1-\alpha/2} \frac{\sigma}{\sqrt{N}} \right]$$

with c_α the α -quantile of the normal distribution satisfying $\Phi(c_\alpha) = \alpha$ and Φ the cdf of a standard normal random variable. This means that $\mathbb{P}(\mu \in I_{\alpha, N}) \xrightarrow{N \rightarrow \infty} 1 - \alpha$. See Figure 5.1 for an illustration. Equivalently, the error $|\mu - \hat{\mu}_N|$ satisfies

$$|\mu - \hat{\mu}_N| \leq c_{1-\alpha/2} \frac{\sigma}{\sqrt{N}} \quad \text{with probability } 1 - \alpha, \text{ asymptotically.}$$

The CLT shows that the Monte Carlo error is of order $N^{-1/2}$, which is a very slow convergence rate (to reduce the error by a factor of 10, one has to multiply N by a factor of 100) and is peculiar of Monte Carlo estimates, generally not improvable. On the other hand, it holds under quite weak assumptions ($\text{Var}(Z) < +\infty$).

The previous error estimate and confidence interval is not practical as it involves the, usually unknown, variance $\sigma^2 = \text{Var}(Z)$. We can replace it by the *sample variance estimator* computed using the same sample $(Z^{(1)}, \dots, Z^{(N)})$,

$$\hat{\sigma}_N^2 = \frac{1}{N-1} \sum_{i=1}^N (Z^{(i)} - \hat{\mu}_N)^2.$$

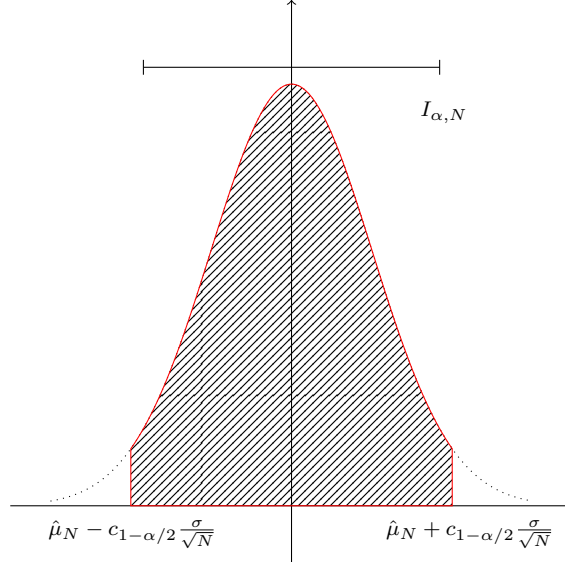


Figure 5.1: Asymptotic confidence interval for the sample mean estimator.

which is also an unbiased estimator and converges almost surely to σ^2 . It follows that $\frac{\sigma}{\hat{\sigma}_N} \rightarrow 1$ a.s. and

$$\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\hat{\sigma}_N} = \underbrace{\frac{\sigma}{\hat{\sigma}_N}}_{\rightarrow 1 \text{ a.s.}} \underbrace{\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\sigma}}_{\xrightarrow{d} \mathcal{N}(0,1)} \xrightarrow{d} \mathcal{N}(0,1).$$

Then, a *computable* asymptotic confidence interval is given by

$$\hat{I}_{\alpha, N} = \left[\hat{\mu}_N - c_{1-\alpha/2} \frac{\hat{\sigma}_N}{\sqrt{N}}, \hat{\mu}_N + c_{1-\alpha/2} \frac{\hat{\sigma}_N}{\sqrt{N}} \right] \quad (5.1)$$

which leads to $\mathbb{P}(\mu \in \hat{I}_{\alpha, N}) \xrightarrow{N \rightarrow \infty} 1 - \alpha$.

If the random variable Z is itself normally distributed, the rescaled random variable $\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\hat{\sigma}_N}$ has a Student's t distribution with $N - 1$ degrees of freedom¹, therefore an exact confidence interval can be constructed using the quantile of the Student's t distribution. In the case of Gaussian random variables, whenever the sample size N is relatively small, it is advisable to use the confidence interval based on the Student's t distribution, which is exact, rather than the one based on the CLT, which is only asymptotically exact and may not be reliable for small N .

For a generic random variable Z , one could still construct a computable *asymptotic* confidence interval based on the quantiles of the Student's t distribution

$$\hat{I}_{\alpha, N}^t = \left[\hat{\mu}_N - t_{1-\alpha/2}^{(N-1)} \frac{\hat{\sigma}_N}{\sqrt{N}}, \hat{\mu}_N + t_{1-\alpha/2}^{(N-1)} \frac{\hat{\sigma}_N}{\sqrt{N}} \right] \quad (5.2)$$

¹We recall that the probability density function of a Student's t distribution with m degrees of freedom is $f(z) = \frac{\Gamma(\frac{m+1}{2})}{\sqrt{\pi m} \Gamma(\frac{m}{2})} \left(1 + \frac{z^2}{m}\right)^{-\frac{m+1}{2}}$ where Γ is the Gamma function.

where $t_{1-\alpha/2}^{(N-1)}$ is the $1 - \frac{\alpha}{2}$ quantile of the Student's t distribution with $N - 1$ degrees of freedom. Since the Student's t distribution converges to the standard normal distribution as $N \rightarrow \infty$ the confidence interval $\hat{I}_{\alpha,N}^t$ is asymptotically exact, i.e. $\mathbb{P}\left(\mu \in \hat{I}_{\alpha,N}^t\right) \xrightarrow{N \rightarrow \infty} 1 - \alpha$. However, its use for a small sample size is less justified as it doesn't necessarily provide a more robust confidence interval than CLT, if the random variable Z has a distribution far from Gaussian. More robust confidence intervals for small sample sizes are discussed in Section 5.3.

5.2 Implementation aspects

As an output of a Monte Carlo simulation, one should *always* provide, beside the point estimate $\hat{\mu}_N$, also an estimate of the error, quantified by e.g. the $1 - \alpha$ asymptotic confidence interval $\hat{I}_{\alpha,N}$.

In practice, one would also like to choose N so as to achieve a prescribed tolerance tol . Again, this could be for instance in terms of the length of the $1 - \alpha$ confidence interval:

$$\text{choose } N: \quad |\hat{I}_{\alpha,N}| \leq 2 \text{ tol}.$$

This can be done by a two (or more) stages procedure as illustrated by the Algorithm 5.1.

Algorithm 5.1: Two stages Monte Carlo.

Given: tol, α

1 Do a pilot run with \bar{N} replicas $(Z^{(1)}, \dots, Z^{(\bar{N})})$ and compute

$$\hat{\mu}_{\bar{N}} = \frac{1}{\bar{N}} \sum_{i=1}^{\bar{N}} Z^{(i)}, \quad \hat{\sigma}_{\bar{N}}^2 = \frac{1}{\bar{N}-1} \sum_{i=1}^{\bar{N}} (Z^{(i)} - \hat{\mu}_{\bar{N}})^2$$

2 Based on the previously estimated variance, fix

$$N = \frac{c_{1-\alpha/2}^2 \hat{\sigma}_{\bar{N}}^2}{tol^2}.$$

3 Generate a new sample $(Z^{(1)}, \dots, Z^{(N)})$ and compute $\hat{\mu}_N$ and $\hat{\sigma}_N^2$

4 **if** $\hat{\sigma}_N^2 > \hat{\sigma}_{\bar{N}}^2$ **then**

5 | Set $\bar{N} = N$ and go back to 2.

6 **else**

7 | Output $\hat{\mu}_N$ and $\hat{I}_{\alpha,N}$.

8 **end**

Alternatively, one can adopt a *sequential* procedure as illustrated by Algorithm 5.2.

Algorithm 5.2: Sequential Monte Carlo.

Given: tol, α

1 Do a pilot run with \bar{N} replicas $(Z^{(1)}, \dots, Z^{(\bar{N})})$ and compute

$$\hat{\sigma}_{\bar{N}}^2 = \frac{1}{\bar{N} - 1} \sum_{i=1}^{\bar{N}} (Z^{(i)} - \hat{\mu}_{\bar{N}})^2.$$

2 Set $N = \bar{N}$, $\hat{\mu}_N = \hat{\mu}_{\bar{N}}$, $\hat{\sigma}_N = \hat{\sigma}_{\bar{N}}$.

3 **while** $\frac{\hat{\sigma}_N c_{1-\alpha/2}}{\sqrt{N}} > tol$ **do**

4 set $N = N + 1$

5 generate $Z^{(N)}$ independent of $Z^{(i)}$, $i < N$

6 recompute $\hat{\mu}_N, \hat{\sigma}_N^2$

7 **end**

An efficient implementation of Algorithm 5.2 requires stable update formulas for $\hat{\mu}_N$ and $\hat{\sigma}_N$. Two such formulas are the following:

$$\begin{aligned} \hat{\mu}_{N+1} &= \frac{N}{N+1} \hat{\mu}_N + \frac{1}{N+1} Z^{(N+1)} \\ \hat{\sigma}_{N+1}^2 &= \frac{N-1}{N} \hat{\sigma}_N^2 + \frac{1}{N+1} \left(Z^{(N+1)} - \hat{\mu}_N \right)^2. \end{aligned}$$

If $N(tol)$ denotes the sample size at the end of the while loop, which is a random variable, and $\hat{\mu}_{N(tol)}$ the corresponding sample mean estimator, it can be shown [2] under the sole assumption that $\text{Var}(Z) < +\infty$, that

$$\lim_{tol \rightarrow 0} \mathbb{P}(|\hat{\mu}_{N(tol)} - \mu| \leq tol) = 1 - \alpha \quad (\text{asymptotic consistency})$$

and

$$\frac{N(tol)tol^2}{\sigma^2 c_{1-\alpha/2}^2} \xrightarrow{\text{a.s.}} 1 \quad \text{as } tol \rightarrow 0.$$

The drawback of this algorithm is that if \bar{N} is chosen too small so that the estimator $\hat{\sigma}_N$ is unreliable, this may cause the algorithm to terminate too early, leading to a poor estimation of $\mathbb{E}[Z]$.

5.3 Non asymptotic error bounds

The confidence interval (5.1) for the sample mean estimator, derived in Section 5.1 is based on the CLT and is only valid asymptotically. Sometimes, if the distribution of the random variable Z is far from being Gaussian and the sample size is small, the distribution of the rescaled sample mean estimator $\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\hat{\sigma}_N}$ may still be far from the asymptotic Normal one and the corresponding confidence interval $\hat{I}_{\alpha, N}$ will be unreliable. Other more robust confidence intervals could be used instead, in this case, which however often lead to very conservative bounds. We mention:

- Bound based on Chebyshev inequality $\mathbb{P}(|Z - \mathbb{E}[Z]| > a) \leq \frac{\text{Var}(Z)}{a^2}$ which implies

$$\mathbb{P}\left(|\hat{\mu}_N - \mu| \leq \frac{\sigma}{\sqrt{N\alpha}}\right) \geq 1 - \alpha$$

from which we can compute the approximate confidence interval

$$\hat{I}_{\alpha,N}^{Cheb} = \left[\hat{\mu}_N - \frac{\hat{\sigma}_N}{\sqrt{\alpha N}}, \hat{\mu}_N + \frac{\hat{\sigma}_N}{\sqrt{\alpha N}}\right].$$

This should be compared with the CLT result $\mathbb{P}\left(|\hat{\mu}_N - \mu| \leq \frac{\sigma c_{1-\alpha/2}}{\sqrt{N}}\right) \geq 1 - \alpha$. For α small, one has typically $c_{1-\alpha/2} \ll \frac{1}{\sqrt{\alpha}}$. E.g. for $\alpha = 0.05$ we have $c_{.975} = 1.96$ whereas $\frac{1}{\sqrt{\alpha}} = 4.47$ and for $\alpha = 0.01$ we have $c_{.995} = 2.576$ whereas $\frac{1}{\sqrt{\alpha}} = 10$.

One could also use higher moments to mitigate the $\frac{1}{\sqrt{\alpha}}$ factor. Consider the generalized Chebyshev inequality $\mathbb{P}(|Z - \mathbb{E}[Z]| > a) \leq \frac{\mathbb{E}[|Z - \mathbb{E}[Z]|^p]}{a^p}$ and an estimator $\hat{\gamma}_{p,N} \approx \mathbb{E}[|Z - \mathbb{E}[Z]|^p]$ of the p^{th} moment based on the sample $Z^{(1)}, \dots, Z^{(N)}$. Then an approximate confidence interval is

$$\hat{I}_{\alpha,N}^{Cheb-p} = \left[\hat{\mu}_N - \frac{\sqrt[p]{\hat{\gamma}_{p,N}}}{\sqrt[p]{\alpha}\sqrt{N}}, \hat{\mu}_N + \frac{\sqrt[p]{\hat{\gamma}_{p,N}}}{\sqrt[p]{\alpha}\sqrt{N}}\right].$$

- Bound based on Berry-Essén (for random iid variables with bounded third moment)

$$\sup_x \left| \mathbb{P}\left(\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\sigma} \leq x\right) - \Phi(x) \right| \leq k \frac{\mathbb{E}[|Z - \mu|^3]}{\sqrt{N}\sigma^3}, \quad (k \approx 0.4748)$$

which quantifies how far the distribution of $\frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\sigma}$ is from the standard Normal cdf based on 3rd moment estimates.

Calling $\hat{Y} = \frac{\sqrt{N}(\hat{\mu}_N - \mu)}{\sigma}$ and $R = k \frac{\mathbb{E}[|Z - \mu|^3]}{\sqrt{N}\sigma^3}$ we have for any $x \in \mathbb{R}$

$$\begin{aligned} \mathbb{P}\left(|\hat{Y}| > x\right) &= \mathbb{P}\left(\hat{Y} > x\right) + \mathbb{P}\left(\hat{Y} < -x\right) = 1 - \underbrace{\mathbb{P}\left(\hat{Y} \leq x\right)}_{\geq \Phi(x) - R} + \underbrace{\mathbb{P}\left(\hat{Y} < -x\right)}_{\leq \Phi(-x) + R} \\ &\leq 1 - \Phi(x) + R + \Phi(-x) + R = 2 - 2\Phi(x) + 2R. \end{aligned}$$

having used that $\Phi(-x) = 1 - \Phi(x)$ and we can look for x_α such that $2 - 2\Phi(x) + 2R = \alpha$. In practice such modified quantile is not computable since σ and the third moment $\gamma_3 = \mathbb{E}[|Z - \mu|^3]$ are not known. However, they can be estimated from the sample $Z^{(1)}, \dots, Z^{(N)}$.

Given an estimate $\hat{\gamma}_{3,N} \approx \mathbb{E}[|Z - \mu|^3]$ one can then solve the problem:

$$\text{find } \hat{x}_\alpha : \quad \Phi(\hat{x}_\alpha) = 1 - \frac{\alpha}{2} + k \frac{\hat{\gamma}_{3,N}}{\sqrt{N}\hat{\sigma}_N^3}$$

and output the confidence interval $\hat{I}_{\alpha,N}^{BE} = \left[\hat{\mu}_N - \hat{x}_\alpha \frac{\hat{\sigma}_N}{\sqrt{N}}, \hat{\mu}_N + \hat{x}_\alpha \frac{\hat{\sigma}_N}{\sqrt{N}}\right]$.

5.4 Vector valued output

The Monte Carlo method extends trivially to a vector valued output $\mathbf{Z} = (Z_1, \dots, Z_m)$ and the estimation of its expected value $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Z}]$. In this case, we generate iid replicas $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(N)}$ and set $\hat{\boldsymbol{\mu}}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{Z}^{(i)}$.

We can also estimate from the same sample the covariance matrix $\hat{C}_N = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{Z}^{(i)} - \hat{\boldsymbol{\mu}}_N)(\mathbf{Z}^{(i)} - \hat{\boldsymbol{\mu}}_N)^\top$. The considerations on asymptotic confidence intervals based on the CLT extend to the case of a vector valued output as well. Indeed we have $\hat{\boldsymbol{\mu}}_N \xrightarrow{\text{a.s.}} \boldsymbol{\mu}$, $\hat{C}_N \xrightarrow{\text{a.s.}} C$ and

$$N(\hat{\boldsymbol{\mu}}_N - \boldsymbol{\mu})^\top \hat{C}_N^{-1} (\hat{\boldsymbol{\mu}}_N - \boldsymbol{\mu}) \xrightarrow{\text{d}} \chi_m^2$$

where χ_m^2 denotes the χ^2 distribution with m degrees of freedom. Based on this asymptotic result, a computable $1 - \alpha$ asymptotic confidence region is

$$\hat{I}_{\alpha, N} = \{\mathbf{y} \in \mathbb{R}^m : (\hat{\boldsymbol{\mu}}_N - \mathbf{y})^\top \hat{C}_N^{-1} (\hat{\boldsymbol{\mu}}_N - \mathbf{y}) \leq \frac{\chi_{m; 1-\alpha}^2}{N}\}$$

where $\chi_{m; 1-\alpha}^2$ is the $1 - \alpha$ quantile of the χ_m^2 distribution, so that $\mathbb{P}(\boldsymbol{\mu} \in \hat{I}_{\alpha, N}) \xrightarrow{N \rightarrow \infty} 1 - \alpha$.

5.5 Smooth functions of expectations and delta method

Consider, as in the previous section, a vector of output quantities $\mathbf{Z} = (Z_1, \dots, Z_m)$ of a stochastic model. However, now, we wish to compute a nonlinear function of the expectation of \mathbf{Z} , i.e.

$$\zeta = f(\mathbb{E}[Z_1], \dots, \mathbb{E}[Z_m])$$

with $f : \mathbb{R}^m \rightarrow \mathbb{R}$ smooth. If we denote $\mu_i = \mathbb{E}[Z_i]$ and $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m) \in \mathbb{R}^m$, the natural Monte Carlo estimator for ζ is

$$\hat{\zeta}_N = f(\hat{\mu}_{1, N}, \dots, \hat{\mu}_{m, N}) \quad \text{with} \quad \hat{\mu}_{i, N} = \frac{1}{N} \sum_{k=1}^N Z_i^{(k)}$$

with $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(N)}$ iid replicas of \mathbf{Z} . If f is continuous at $\boldsymbol{\mu}$, then $\hat{\zeta}_N \xrightarrow{\text{a.s.}} \zeta$ i.e. $\hat{\zeta}_N$ is a consistent estimator of ζ .

The question is now how to estimate the error on $\hat{\zeta}_N$ and provide a confidence interval. One way of doing this is provided by the so called *delta method*, based on a first order Taylor expansion of f around $\boldsymbol{\mu}$:

$$\hat{\zeta}_N - \zeta = f(\hat{\boldsymbol{\mu}}_N) - f(\boldsymbol{\mu}) = \nabla f(\boldsymbol{\mu})(\hat{\boldsymbol{\mu}}_N - \boldsymbol{\mu}) + o(\|\hat{\boldsymbol{\mu}}_N - \boldsymbol{\mu}\|).$$

where we have used the convention that ∇f is a row vector. Let $C = \text{Cov}(\mathbf{Z}) = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^\top]$. Then

$$\sqrt{N}(\hat{\zeta}_N - \zeta) \xrightarrow{\text{d}} \mathcal{N}(0, \nabla f(\boldsymbol{\mu})C\nabla f(\boldsymbol{\mu})^\top).$$

A computable $1 - \alpha$ confidence interval can then be constructed by replacing C with $\hat{C}_N = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{Z}^{(i)} - \hat{\boldsymbol{\mu}}_N)(\mathbf{Z}^{(i)} - \hat{\boldsymbol{\mu}}_N)^\top$ and $\nabla f(\boldsymbol{\mu})$ with $\nabla f(\hat{\boldsymbol{\mu}}_N)$ as

$$\hat{I}_{\alpha,N} = [\hat{\zeta}_N - \Delta_N, \hat{\zeta}_N + \Delta_N], \quad \Delta_N = \frac{c_{1-\alpha/2}}{\sqrt{N}} \sqrt{\nabla f(\hat{\boldsymbol{\mu}}_N) \hat{C}_N \nabla f(\hat{\boldsymbol{\mu}}_N)^\top}.$$

Obverse that the estimator $\hat{\zeta}_N$ is *biased* in general.

Example 5.2. Let Z be a scalar random variable, output of a stochastic model. Suppose we want to estimate the coefficient of variation of Z

$$\zeta = \frac{\sigma(Z)}{\mu(Z)} = \frac{\sqrt{\mathbb{E}[Z^2] - \mathbb{E}[Z]^2}}{\mathbb{E}[Z]}.$$

Setting $\mathbf{Z} = (Z_1, Z_2) = (Z, Z^2)$ and $f(x, y) = \sqrt{\frac{y}{x^2} - 1}$, then $\zeta = f(\mathbb{E}[Z_1], \mathbb{E}[Z_2])$. If $\hat{\zeta}_N$ denotes a Monte Carlo estimator for ζ , the delta method can be used to produce a $1 - \alpha$ asymptotic confidence interval. Explicit calculations are left as an exercise.

5.6 Monte Carlo to compute integrals

Consider a simple stochastic model $Z = \psi(X_1, \dots, X_d)$ with $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$ bounded and $\mathbf{X} = (X_1, \dots, X_d)$ a random vector with joint probability density function $f : \mathbb{R}^d \rightarrow \mathbb{R}_+$. Then,

$$\mathbb{E}[Z] = \int_{\mathbb{R}^d} \psi(x_1, \dots, x_d) f(x_1, \dots, x_d) dx_1 \dots dx_d = \int_{\mathbb{R}^d} \psi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}.$$

A Monte Carlo approximation of $\mu = \mathbb{E}[Z]$ consists of:

- generating N iid replicas of $\mathbf{X}^{(i)} \sim f$
- computing $\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)}) \approx \int_{\mathbb{R}^d} \psi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$.

Hence, the formula $\frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)})$ can be seen as a quadrature formula to approximate the integral $\int_{\mathbb{R}^d} \psi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$.

Conversely let us consider the problem of computing an integral $I = \int_{\mathbb{R}^d} \psi(\mathbf{x}) w(\mathbf{x}) d\mathbf{x}$ where $w : \mathbb{R}^d \rightarrow \mathbb{R}_+$ a non negative weight such that $\int_{\mathbb{R}^d} w = 1$. Then we can estimate the integral by a Monte Carlo formula

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)})$$

with $\mathbf{X}^{(i)} \stackrel{\text{iid}}{\sim} w$.

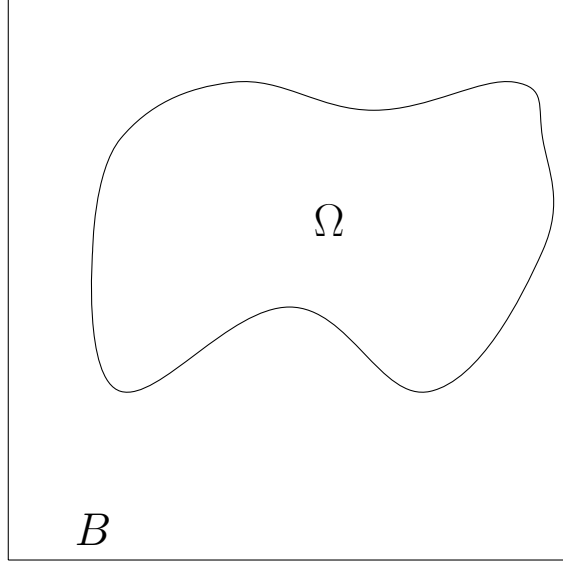


Figure 5.2: Monte Carlo to estimate the volume of Ω .

Example 5.3. Let $\Omega \subset \mathbb{R}^d$ be a bounded domain. We want to compute its volume $|\Omega|$. Let B be a rectangular domain containing Ω . Then $I = |\Omega| = |B| \int_B \mathbf{1}_\Omega(\mathbf{x}) \frac{1}{|B|} d\mathbf{x}$ and its Monte Carlo approximation is

$$\hat{I} = \frac{|B|}{N} \sum_{i=1}^N \mathbf{1}_\Omega(\mathbf{X}^{(i)}) = \frac{\#\{\mathbf{X}^{(i)} \in \Omega\}}{N} |B|,$$

with $\mathbf{X}^{(i)} \stackrel{iid}{\sim} \mathcal{U}(B)$ i.e. we draw independently uniform points in B and count how many fall in Ω . See Figure 5.2 for a graphical illustration.

The error in the Monte Carlo approximation is

$$|I - \hat{I}| \leq c_{1-\alpha/2} \frac{\sigma}{\sqrt{N}}$$

with probability $1 - \alpha$ asymptotically, where

$$\sigma^2 = \int_{\mathbb{R}^d} \psi^2(\mathbf{x}) w(\mathbf{x}) d\mathbf{x} - I^2 \leq \int_{\mathbb{R}^d} \psi^2(\mathbf{x}) w(\mathbf{x}) d\mathbf{x}.$$

Hence, the rate of convergence is $O(N^{-1/2})$ and is achieved under the sole condition $\int_{\mathbb{R}^d} \psi^2(\mathbf{x}) w(\mathbf{x}) d\mathbf{x} < +\infty$. Observe, in particular, that this rate is *independent of the dimension d !* (assuming that the variance σ^2 remains bounded when we increase the dimension of the problem). Although Monte Carlo has a very poor convergence rate $O(N^{-1/2})$, its use is still very appealing for high dimensional problems.

As a term of comparison, consider the problem of computing an integral $I = \int_{[0,1]^d} \psi(\mathbf{x}) d\mathbf{x}$ on the unit hypercube by a tensor quadrature formula, e.g. tensor mid-point rule

$$I^{mp} = \sum_{i=1}^N \psi(\mathbf{X}^{(i)}) h^d = \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)})$$

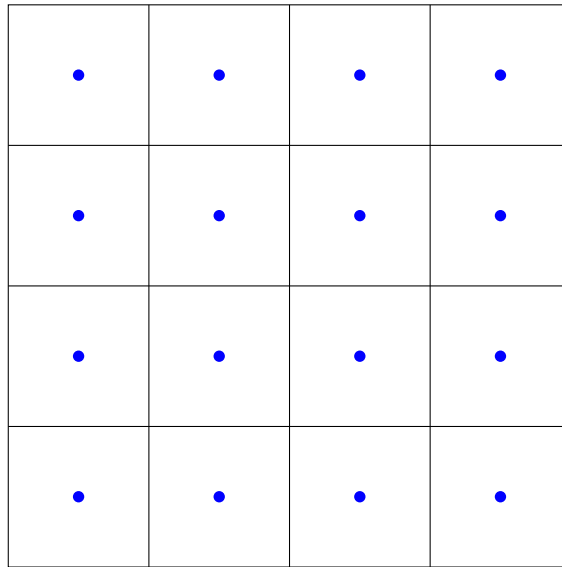


Figure 5.3: Uniform grid to estimate the integral of ψ .

where $\mathbf{X}^{(i)}$ are the centres of the cells and $h = N^{-1/d}$ is the length of each cell (see Figure 5.3).

The error of the quadrature formula can be bounded as:

$$|I - I^{mp}| \leq Ch^2 \|\psi\|_{C^2([0,1]^d)} = CN^{-2/d} \|\psi\|_{C^2([0,1]^d)}.$$

Therefore, such a formula achieves a rate $N^{-2/d}$, with respect to the number of points used, provided $\psi \in C^2([0,1]^d)$ (hence regularity is required on the integrand to achieve such rate) and already for $d > 4$ the rate will be worse than Monte Carlo (this effect is usually referred to as the “curse of dimensionality”).

Chapter 6

Variance Reduction Techniques

As in the previous chapter, let Z be a random variable, output of a stochastic model, and consider the goal of computing the expected value $\mu = \mathbb{E}[Z]$. It will be useful to assume that Z can be written as $Z = \psi(X)$ with $X = (X_1, \dots, X_d)$ a random vector with pdf $f : \mathbb{R}^d \rightarrow \mathbb{R}_+$, so that

$$\mu = \mathbb{E}[Z] = \int_{\mathbb{R}^d} \psi(x) f(x) dx.$$

The Monte Carlo approach (hereafter called “Crude Monte Carlo”) to approximate μ consists in generating N iid replicas $Z^{(1)}, \dots, Z^{(N)}$, with $Z^{(i)} = \psi(X^{(i)})$, $X^{(i)} \stackrel{\text{iid}}{\sim} f$ and computing

$$\hat{\mu}_{\text{CMC}} = \frac{1}{N} \sum_{i=1}^N Z^{(i)}.$$

As we have seen in Chapter 5, by the CLT we have that

$$|\mu - \hat{\mu}_{\text{CMC}}| \leq c_{1-\alpha/2} \frac{\sqrt{\text{Var}(Z)}}{\sqrt{N}}$$

with probability $1 - \alpha$, asymptotically as $N \rightarrow \infty$.

The techniques of variance reduction aim at improving the performance of a Crude Monte Carlo approximation by reducing the constant $\sqrt{\text{Var}(Z)}$, hence the name “variance reduction”. The idea is simple: instead of applying the sample mean estimator $\hat{\mu} = \hat{\mu}(Z)$ to the variable Z , one applies it to a cleverly modified version \tilde{Z} which satisfies

$$\mathbb{E}[\tilde{Z}] = \mathbb{E}[Z] = \mu \quad \text{and} \quad \text{Var}(\tilde{Z}) \ll \text{Var}(Z).$$

Hence, a Monte Carlo approximation with variance reduction will look like

$$\hat{\mu}_{\text{VR}} = \frac{1}{N} \sum_{i=1}^N \tilde{Z}^{(i)}$$

with $\tilde{Z}^{(i)} \stackrel{\text{iid}}{\sim} \tilde{Z}$.

6.1 Antithetic Variables

Suppose N even. Instead of generating N iid replicas of Z , the underlying idea of antithetic sampling is to generate $N/2$ iid pairs of negatively correlated random variables

$$(Z^{(1)}, Z_a^{(1)}), (Z^{(2)}, Z_a^{(2)}), \dots, (Z^{(N/2)}, Z_a^{(N/2)}),$$

where all $Z^{(i)}, Z_a^{(i)}$ have the same distribution as Z but $\text{Cov}(Z^{(i)}, Z_a^{(i)}) \leq 0, i = 0, \dots, N/2$. If we now consider the estimator

$$\hat{\mu}_{AV} = \frac{1}{N/2} \sum_{i=1}^{N/2} \frac{Z^{(i)} + Z_a^{(i)}}{2}$$

it follows immediately that

$$\mathbb{E}[\hat{\mu}_{AV}] = \mathbb{E}[Z] = \mu$$

and

$$\begin{aligned} \text{Var}(\hat{\mu}_{AV}) &= \frac{4}{N^2} \sum_{i=1}^{N/2} \text{Var}\left(\frac{Z^{(i)} + Z_a^{(i)}}{2}\right) = \frac{1}{2N} \text{Var}(Z^{(1)} + Z_a^{(1)}) \\ &= \frac{1}{2N} \left(\text{Var}(Z^{(1)}) + \text{Var}(Z_a^{(1)}) + 2 \text{Cov}(Z^{(1)}, Z_a^{(1)}) \right) \\ &= \frac{\text{Var}(Z) + \text{Cov}(Z^{(1)}, Z_a^{(1)})}{N} \leq \text{Var}(\hat{\mu}_{CMC}) \end{aligned}$$

since, by assumption, $\text{Cov}(Z^{(1)}, Z_a^{(1)}) \leq 0$. The estimator $\hat{\mu}_{AV}$ has therefore a smaller variance than the Curde Monte Carlo estimator $\hat{\mu}_{CMC}$ at the same computational cost (provided the generation of $Z_a^{(i)}$ has the same cost as the generation of $Z^{(i)}$).

The question is now how to generate pairs of negatively correlated variables $(Z^{(i)}, Z_a^{(i)})$. The following proposition presents a situation in which variance reduction can be achieved by a rather simple construction of antithetic sampling.

Proposition 6.1. *Assume that the random variable Z has the expression $Z = \psi(X)$, with $X = (X_1, \dots, X_d)$ a random vector with independent components, such that*

- X has a symmetric distribution around its mean, i.e. $2\mathbb{E}[X] - X \sim X$
- ψ is a monotonic function in each of its arguments (either non-increasing or non-decreasing).

Then $Z = \psi(X)$ and $Z_a = \psi(2\mathbb{E}[X] - X)$ satisfy

$$\mathbb{E}[Z] = \mathbb{E}[Z_a] \quad \text{and} \quad \text{Cov}(Z, Z_a) \leq 0.$$

Under the assumptions of the previous proposition, a Monte Carlo approximation of $\mu = \mathbb{E}[Z]$ with antithetic variables can be constructed by the following algorithm.

Algorithm 6.1: Antithetic variables.

- 1 Generate $N/2$ iid replicas $X^{(1)}, \dots, X^{(N/2)}$ of X ;
- 2 For each $X^{(i)}$ compute $Z^{(i)} = \psi(X^{(i)})$ and $Z_a^{(i)} = \psi(2\mathbb{E}[X] - X^{(i)})$;
- 3 Compute $\hat{\mu}_{AV} = \frac{1}{N} \sum_{i=1}^{N/2} (Z^{(i)} + Z_a^{(i)})$.
- 4 Estimate $\hat{\sigma}_{AV}^2 = \frac{1}{N/2-1} \sum_{i=1}^{N/2} \left(\frac{Z^{(i)} + Z_a^{(i)}}{2} - \hat{\mu}_{AV} \right)^2$
- 5 Output $\hat{\mu}_{AV}$ and a (asymptotic) $1 - \alpha$ confidence interval

$$\hat{I}_{\alpha, N} = \left[\hat{\mu}_{AV} - c_{1-\alpha/2} \frac{\hat{\sigma}_{AV}}{\sqrt{N/2}}, \hat{\mu}_{AV} + c_{1-\alpha/2} \frac{\hat{\sigma}_{AV}}{\sqrt{N/2}} \right]$$

The proof of Proposition 6.1 relies on the following Chebyshev Covariance inequality.

Lemma 6.2 (Chebyshev Covariance inequality). *Let X be a real-valued random variable with pdf $f : \mathbb{R} \rightarrow \mathbb{R}_+$ and let $g, h : \mathbb{R} \rightarrow \mathbb{R}$ be functions that are both non-decreasing or both non-increasing, such that $\mathbb{E}[|g(X)|], \mathbb{E}[|h(X)|], \mathbb{E}[|g(x)h(x)|] < +\infty$. Then $\text{Cov}(g(X), h(X)) \geq 0$.*

Proof. We consider the case of g, h both non-decreasing. The other case can be proven analogously. Let $\tilde{g}(x) = g(x) - \mathbb{E}[g(x)]$ and $\tilde{h}(x) = h(x) - \mathbb{E}[h(x)]$. Observe first that

$$\begin{aligned} & \frac{1}{2} \iint (g(x) - g(y))(h(x) - h(y))f(x)f(y) dx dy = \\ & \quad \frac{1}{2} \iint (\tilde{g}(x) - \tilde{g}(y))(\tilde{h}(x) - \tilde{h}(y))f(x)f(y) dx dy \\ & = \underbrace{\int \tilde{g}(x)\tilde{h}(x)f(x) dx}_{=0} - \underbrace{\left(\int \tilde{g}(x)f(x) dx \right) \left(\int \tilde{h}(y)f(y) dy \right)}_{=0} = \text{Cov}(g(X), h(X)) \end{aligned}$$

Hence

$$\begin{aligned} \text{Cov}(g(X), h(X)) &= \frac{1}{2} \int_{x \geq y} \underbrace{(g(x) - g(y))}_{\geq 0} \underbrace{(h(x) - h(y))}_{\geq 0} f(x)f(y) dx dy \\ & \quad + \frac{1}{2} \int_{x < y} \underbrace{(g(x) - g(y))}_{\leq 0} \underbrace{(h(x) - h(y))}_{\leq 0} f(x)f(y) dx dy \geq 0. \end{aligned}$$

□

The previous inequality generalizes to the multivariate case, whose proof is left as exercise.

Lemma 6.3. *Let $X = (X_1, \dots, X_d) \in \mathbb{R}^d$ be a random vector with independent components and let $g, h : \mathbb{R}^d \rightarrow \mathbb{R}$ be functions whose dependence on each argument is either non-decreasing or non-increasing for both of them. Then $\text{Cov}(g(X), h(X)) \geq 0$.*

Proof of Proposition 6.1. Since $2\mathbb{E}[X] - X \sim X$ we have $Z_a \sim Z$, hence $\mathbb{E}[Z_a] = \mathbb{E}[Z]$. Moreover, observe that if $\psi(X_1, \dots, X_d)$ is e.g. non decreasing in the i -th argument, so is the function $-\psi(2\mathbb{E}[X_1] - X_1, \dots, 2\mathbb{E}[X_d] - X_d)$ and, from the previous Lemma, we have $\text{Cov}(\psi(X), -\psi(2\mathbb{E}[X] - X)) \geq 0$, hence $\text{Cov}(Z, Z_a) \leq 0$. \square

Example 6.1. *Let $Z \sim \text{Exp}(\lambda)$. Then $Z = -\frac{1}{\lambda} \log X = \psi(X)$ with $X \sim \mathcal{U}(0, 1)$ and ψ monotonic (decreasing). It follows that $\psi(X)$ and $\psi(1 - X)$ are negatively correlated and a Monte Carlo estimator with antithetic variables for the computation of $\mu = \frac{1}{\lambda} = \mathbb{E}[Z]$ is $\hat{\mu}_{AV} = \frac{1}{N} \sum_{i=1}^{N/2} (-\frac{1}{\lambda} \log(X^{(i)}) - \frac{1}{\lambda} \log(1 - X^{(i)}))$, with $X^{(i)} \stackrel{iid}{\sim} \mathcal{U}(0, 1)$.*

Example 6.2. *Consider the problem of pricing a European option $\mu = \mathbb{E}[Z]$ with $Z = \psi(S_T) = e^{-rT}(S_T - K)_+$ where S_t is the solution of the stochastic differential equation*

$$dS_t = rS_t dt + \sigma S_t dW_t$$

with W_t a standard Wiener process and S_0 given. It can be shown that $X_t = \log(S_t/S_0)$ satisfies the stochastic differential equation with constant coefficients

$$dX_t = (r - \sigma^2/2) dt + \sigma dW_t, \quad X_0 = 0$$

whose solution is $X_t = (r - \sigma^2/2)t + \sigma W_t \sim N((r - \sigma^2/2)t, \sigma^2 t)$. Hence $S_T = S_0 e^{X_T}$ has a log-normal distribution with $X_T \sim N((r - \sigma^2/2)T, \sigma^2 T)$ and $\mathbb{E}[S_T] = S_0 e^{rT}$. Observe that ψ is a non decreasing function of S_T , which, on its turn, is an increasing function of X_T whose distribution is symmetric about its mean. Hence $\tilde{\psi}(X_T) = \psi(S_0 e^{X_T})$ is non decreasing in X_T and an antithetic variable estimator

$$\hat{\mu}_{AV} = \frac{1}{N} \sum_{i=1}^{N/2} \left(\tilde{\psi}(X_T^{(i)}) + \tilde{\psi}((2r - \sigma^2)T - X_T^{(i)}) \right), \quad X_T^{(i)} \stackrel{iid}{\sim} N\left((r - \frac{\sigma^2}{2})T, \sigma^2 T\right)$$

will lead to variance reduction.

Example 6.3. *Consider a random walk on the integers: $Z_{n+1} = Z_n + X_{n+1}$ with X_i iid such that $\mathbb{P}(X_i = 1) = \mathbb{P}(X_i = -1) = 1/2$ and $Z_0 = 0$. We want to estimate by Monte Carlo $\mu = \mathbb{P}(Z_N \geq s)$ with $s \in \mathbb{N}$. Denote*

$$\psi(Z_N) = \mathbf{1}_{\{Z_N \geq s\}} = \mathbf{1}_{\{\sum_{n=1}^N X_n \geq s\}} = \tilde{\psi}(X_1, \dots, X_N).$$

Then $\mu = \mathbb{E}[\psi(Z_N)] = \mathbb{E}[\tilde{\psi}(X_1, \dots, X_N)]$. Since $\tilde{\psi}$ is a non decreasing function in each X_n , and each X_n has a symmetric distribution around its mean $\mathbb{E}[X_n] = 0$, a MC estimator with antithetic variables will lead to variance reduction. It consists in generating $N/2$ iid paths $Z_{n+1}^{(i)} = Z_n^{(i)} + X_n^{(i)}$ as well as the antithetic paths $\tilde{Z}_{n+1}^{(i)} = \tilde{Z}_n^{(i)} - X_n^{(i)}$, and build the estimator $\hat{\mu}_{AV} = \frac{1}{N} \sum_{i=1}^{N/2} (\psi(Z_N^{(i)}) + \psi(\tilde{Z}_N^{(i)}))$.

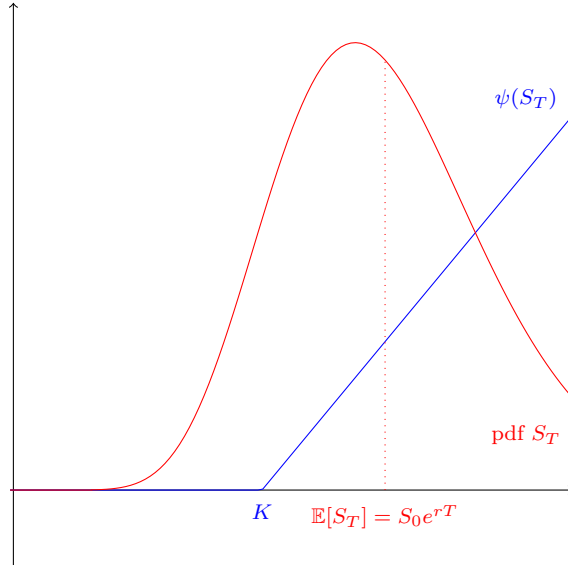


Figure 6.1: European option.

6.2 Importance Sampling

Let $X \in \mathbb{R}^d$ be a random vector with pdf $f : \mathbb{R}^d \rightarrow \mathbb{R}_+$ and $Z = \psi(X)$ with $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$. Then, computing the expected value of Z corresponds to computing the multidimensional integral

$$\mu = \mathbb{E}[Z] = \int_{\mathbb{R}^d} \psi(x) f(x) dx$$

Let now $g : \mathbb{R}^d \rightarrow \mathbb{R}_+$ be an auxiliary pdf such that $g(x) = 0$ only if $\psi(x)f(x) = 0$. Then, the integral can be rewritten as

$$\mu = \mathbb{E}[Z] = \int_{\mathbb{R}^d} \left(\frac{\psi(x)f(x)}{g(x)} \right) g(x) dx = \mathbb{E}_g \left[\frac{\psi f}{g} \right].$$

where \mathbb{E}_g denotes expectation under the measure $g(x) dx$. It follows that in a Monte Carlo approach, instead of generating iid replicas of X to estimate $\mu = \mathbb{E}_f[\psi(X)]$, we could generate iid replicas of \tilde{X} having pdf g , and estimate $\mu = \mathbb{E}_g \left[\frac{\psi(\tilde{X})f(\tilde{X})}{g(\tilde{X})} \right]$. This technique is known as *importance sampling*. The auxiliary distribution g is called the *importance sampling* or *dominating* distribution and the correcting factor $w(x) = \frac{f(x)}{g(x)}$ is often called the *likelihood ratio*.

In more general terms, if X has measure ν_X and ν^* is another probability measure that dominates ν_X , i.e. ν_X is absolutely continuous with respect to ν^* , then there exists a density $\rho = \frac{d\nu_X}{d\nu^*}$ (Radon-Nicodym derivative), and $\mathbb{E}[Z]$ can be rewritten as

$$\mu = \mathbb{E}[Z] = \int \psi(x) d\nu_X(x) = \int \psi(x) \rho(x) d\nu^*(x) = \mathbb{E}_*[\psi \rho]$$

Algorithm 6.2: Importance sampling

- 1 Generate N iid replicas $\tilde{X}^{(1)}, \dots, \tilde{X}^{(N)} \sim g$
- 2 Compute $\hat{\mu}_{\text{IS}} = \frac{1}{N} \sum_{i=1}^N \frac{\psi(\tilde{X}^{(i)})f(\tilde{X}^{(i)})}{g(\tilde{X}^{(i)})}$
- 3 Estimate $\hat{\sigma}_{\text{IS}}^2 = \frac{1}{N-1} \sum_{i=1}^N \left(\frac{\psi(\tilde{X}^{(i)})f(\tilde{X}^{(i)})}{g(\tilde{X}^{(i)})} - \hat{\mu}_{\text{IS}} \right)^2$
- 4 Output $\hat{\mu}_{\text{IS}}$ and a (asymptotic) $1 - \alpha$ confidence interval

$$\hat{I}_{\alpha, N} = \left[\hat{\mu}_{\text{IS}} - c_{1-\alpha/2} \frac{\hat{\sigma}_{\text{IS}}}{\sqrt{N}}, \hat{\mu}_{\text{IS}} + c_{1-\alpha/2} \frac{\hat{\sigma}_{\text{IS}}}{\sqrt{N}} \right]$$

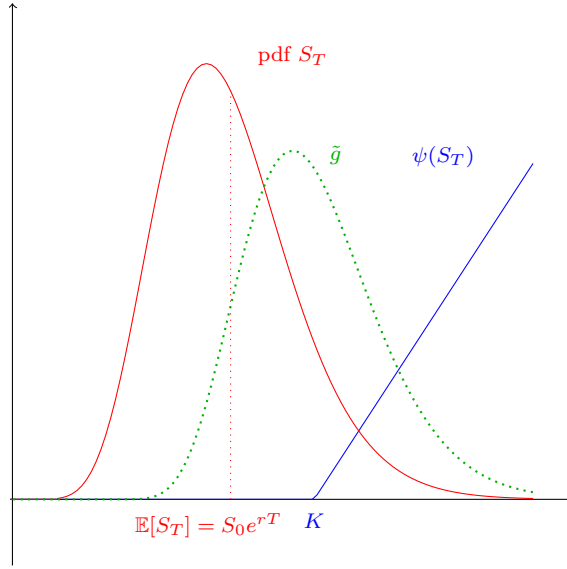


Figure 6.2: European option.

and an importance sampling strategy consists in generating iid replicas $\tilde{X}^{(i)} \stackrel{\text{iid}}{\sim} \nu^*$, $i = 1, \dots, N$ and estimating the empirical mean $\hat{\mu}_{\text{IS}} = \frac{1}{N} \sum_{i=1}^N \psi(\tilde{X}^{(i)})\rho(\tilde{X}^{(i)})$.

Example 6.4. Let us consider again the option pricing problem of computing $\mu = \mathbb{E}[Z]$, $Z = \psi(S_T) = e^{-rT}(S_T - K)_+$ with $S_T = S_0 \exp(X_T)$ and $X_T \sim N((r - \sigma^2/2)T, \sigma^2 T)$. If $K \gg \mathbb{E}[S_T] = S_0 e^{rT}$, most of the the mass of S_T falls in the region where $\psi(S_T) = 0$. Hence a crude Monte Carlo estimator will be very ineffective as only few replicas of S_T will fall in the “interesting” region $S_T > K$. The idea would then be to “artificially” push the distribution to the right. This can be achieved, for instance, by increasing the drift parameter r in the dynamics of S_t . We can therefore simulate a geometric Brownian motion

$$d\tilde{S}_t = \tilde{r}\tilde{S}_t dt + \sigma\tilde{S}_t dW_t$$

with an increased drift rate $\tilde{r} > r$. Let $X_T = \log(S_T/S_0) \sim N((r - \sigma^2/2)T, \sigma^2 T)$ and $\tilde{X}_T = \log(\tilde{S}_T/S_0) \sim N((\tilde{r} - \sigma^2/2)T, \sigma^2 T)$, and denote by f_{X_T} and $f_{\tilde{X}_T}$ the pdfs of X_T

and \tilde{X}_T , respectively. It follows that

$$\mu = \int_{\mathbb{R}} \psi(S_0 e^x) f_{X_T}(x) dx = \int_{\mathbb{R}} \psi(S_0 e^x) w(x) f_{\tilde{X}_T}(x) dx$$

with likelihood ratio

$$w(x) = \frac{f_{X_T}(x)}{f_{\tilde{X}_T}(x)} = \exp \left\{ \frac{(\tilde{r} - r)((\tilde{r} + r - \sigma^2)T - 2x)}{2\sigma^2} \right\} = (e^x)^{-\frac{\tilde{r}-r}{\sigma^2}} e^{\frac{(\tilde{r}-r)(\tilde{r}+r-\sigma^2)T}{2\sigma^2}}$$

and an importance sampling estimator is

$$\hat{\mu}_{IS} = \frac{1}{N} \sum_{i=1}^N \psi(\tilde{S}_T^{(i)}) \left(\frac{\tilde{S}_T^{(i)}}{S_0} \right)^{-\frac{\tilde{r}-r}{\sigma^2}} e^{\frac{(\tilde{r}-r)(\tilde{r}+r-\sigma^2)T}{2\sigma^2}}$$

with

$$\log \left(\tilde{S}_T^{(i)} / S_0 \right) \stackrel{iid}{\sim} N((\tilde{r} - \sigma^2/2)T, \sigma^2 T).$$

6.2.1 On the choice of the importance sampling distribution g

The importance sampling estimator

$$\hat{\mu}_{IS} = \frac{1}{N} \sum_{i=1}^N \frac{\psi(\tilde{X}^{(i)}) f(\tilde{X}^{(i)})}{g(\tilde{X}^{(i)})}, \quad \tilde{X}^{(i)} \stackrel{iid}{\sim} g$$

is unbiased and has variance

$$\mathbb{V}\text{ar}(\hat{\mu}_{IS}) = \frac{1}{N} \mathbb{V}\text{ar}_g \left(\frac{\psi f}{g} \right) = \frac{1}{N} \left(\int_{\mathbb{R}^d} \frac{\psi^2(x) f^2(x)}{g^2(x)} g(x) dx - \mu^2 \right) = \frac{1}{N} \left(\mathbb{E}_f \left[\psi^2 \frac{f}{g} \right] - \mu^2 \right).$$

Therefore, the optimal choice of g is the one that minimizes $\mathbb{V}\text{ar}(\hat{\mu}_{IS})$, i.e. it minimizes the term $\int_{\mathbb{R}^d} \psi^2 \frac{f^2}{g} dx$, under the conditions $\int_{\mathbb{R}^d} g dx = 1$ and $g \geq 0$. It is clear that the optimal distribution should vanish outside $\Gamma = \text{supp}(\psi^2 f^2)$. Moreover, introducing the Lagrangian function

$$\mathcal{L}(g, \lambda) = \int_{\Gamma} \frac{\psi^2 f^2}{g} dx + \lambda \left(\int_{\Gamma} g - 1 \right)$$

and taking variations, the (necessary) optimality condition reads

$$\frac{\partial \mathcal{L}}{\partial g}(\delta g) = - \int_{\Gamma} \left(\psi^2 \frac{f^2}{g^2} - \lambda \right) \delta g dx = 0, \quad \forall \delta g$$

which implies $g^2 = \psi^2 \frac{f^2}{\lambda}$. We see that the optimal g is given by

$$g^* = \frac{|\psi| f}{\mathbb{E}_f[|\psi|]}.$$

With such optimal g^* , the variance of the importance sampling estimator is $\mathbb{V}\text{ar}(\hat{\mu}_{IS}^*) = \mathbb{E}[|\psi|^2] - \mathbb{E}[\psi]^2$. In particular, if $\psi \geq 0$, we have $\mathbb{V}\text{ar}(\hat{\mu}_{IS}^*) = 0$! However, working with

g^* is clearly not practical as the normalizing constant $\mathbb{E}_f[|\psi|]$ is, in general, as difficult to compute as the original quantity $\mu = \mathbb{E}[\psi]$, and we need to know it explicitly to compute the likelihood ratio.

Although the optimal distribution g^* can not be used in practice, this optimization argument shows that the dominating density g should resemble as much as possible to $|\psi|f$ while still being easily simulatable and with explicit expression.

Often, this optimization is performed over a parametric family of pdfs $\{f(\cdot, \theta), \theta \in \Theta\}$. Assuming that the original pdf also belongs to the family, with parameter θ_0 , i.e. $f = f(\cdot, \theta_0)$ and that the support $\text{supp}(f(\cdot, \theta))$ of each distribution in the family is the same, we can take as dominating distribution

$$g(\cdot) = f(\cdot, \theta^*), \quad \text{with } \theta^* = \underset{\theta \in \Theta}{\text{argmin}} \mathbb{E}_\theta \left[\frac{\psi^2 f^2(\cdot, \theta_0)}{f^2(\cdot, \theta)} \right] = \underset{\theta \in \Theta}{\text{argmin}} \mathbb{E}_{\theta_0} \left[\frac{\psi^2 f(\cdot, \theta_0)}{f(\cdot, \theta)} \right].$$

A typical case is when $\{f(\cdot, \theta)\}$ is an exponential family $f(x, \theta) \propto \exp(\theta^\top x - k(\theta))$, for which the likelihood ratio $\frac{f(x, \theta_0)}{f(x, \theta)}$ takes a simple form. The optimization above can be performed numerically replacing the exact expectation with a sample average over a pilot run.

Algorithm 6.3: Importance sampling with variance minimization

- 1 Generate \bar{N} iid replicas $Y^{(1)}, \dots, Y^{(\bar{N})} \sim f(\cdot, \theta_0)$
- 2 Solve the minimization problem

$$\hat{\theta}_{\mathbf{Y}}^* = \underset{\theta \in \Theta}{\text{argmin}} \frac{1}{\bar{N}} \sum_{i=1}^{\bar{N}} \psi^2(Y^{(i)}) \frac{f(Y^{(i)}, \theta_0)}{f(Y^{(i)}, \theta)}$$

- 3 Generate N iid replicas $X^{(1)}, \dots, X^{(N)} \sim f(\cdot, \hat{\theta}_{\mathbf{Y}}^*)$
 - 4 Compute $\hat{\mu}_{\text{IS}} = \frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}) \frac{f(X^{(i)}, \theta_0)}{f(X^{(i)}, \hat{\theta}_{\mathbf{Y}}^*)}$.
-

The estimator $\hat{\mu}_{\text{IS}}$ of Algorithm 6.3 is unbiased. Indeed, if we denote by $\mathbf{X} = (X^{(1)}, \dots, X^{(N)})$ and $\mathbf{Y} = (Y^{(1)}, \dots, Y^{(\bar{N})})$, and use the tower property, we have

$$\mathbb{E}[\hat{\mu}_{\text{IS}}] = \mathbb{E}_{\mathbf{Y}}[\mathbb{E}_{\mathbf{X}}[\hat{\mu}_{\text{IS}} \mid \mathbf{Y}]] = \mathbb{E}_{\mathbf{Y}} \left[\int \psi(x) \frac{f(x, \theta_0)}{f(x, \hat{\theta}_{\mathbf{Y}}^*)} f(x, \hat{\theta}_{\mathbf{Y}}^*) dx \right] = \mu.$$

We can write as well an adaptive version of this algorithm. Notice that at step 3 we generate a sample from $\hat{\theta} = \hat{\theta}_{\mathbf{Y}}^*$. On the other hand, our functional to minimize can be written as

$$J(\theta) = \mathbb{E}_{\theta_0} \left[\frac{\psi^2 f(\cdot, \theta_0)}{f(\cdot, \theta)} \right] = \mathbb{E}_{\hat{\theta}} \left[\frac{\psi^2 f^2(\cdot, \theta_0)}{f(\cdot, \theta) f(\cdot, \hat{\theta})} \right].$$

This suggests the following adaptive Algorithm 6.4

Algorithm 6.4: Adaptive importance sampling with variance minimization

Given: $tol, \alpha, \theta_0, \bar{N} > 1, \gamma > 1$

- 1 Set $N = \bar{N}, \hat{\theta} = \theta_0, \hat{\sigma} = \infty$
- 2 **while** $\frac{\hat{\sigma} c_{1-\alpha/2}}{\sqrt{N}} > tol$ **do**
- 3 Generate N iid replicas $Y^{(1)}, \dots, Y^{(N)} \sim f(\cdot, \hat{\theta})$
- 4 Compute

$$\hat{\mu}_{IS} = \frac{1}{N} \sum_{i=1}^N \psi(Y^{(i)}) \frac{f(Y^{(i)}, \theta_0)}{f(Y^{(i)}, \hat{\theta})}, \quad \hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N \left(\psi(Y^{(i)}) \frac{f(Y^{(i)}, \theta_0)}{f(Y^{(i)}, \hat{\theta})} - \hat{\mu}_{IS} \right)^2$$
- 5 Solve the minimization problem

$$\hat{\theta}_{new} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \psi^2(Y^{(i)}) \frac{f^2(Y^{(i)}, \theta)}{f(Y^{(i)}, \theta) f(Y^{(i)}, \hat{\theta})}$$

Set $\hat{\theta} = \hat{\theta}_{new}$ and $N = \gamma N$
- 6 **end**
- 7 Output $\hat{\mu}_{IS}$

An alternative approach to determine an optimal parameter θ^* consists in minimizing over θ the *Kullback-Leibler divergence* (or cross entropy) between the candidate distribution $f(\cdot, \theta)$ and the optimal importance sampling distribution $g^*(x) = \frac{|\psi(x)|f(x, \theta_0)}{\mathbb{E}_{\theta_0}[|\psi|]}$.

Definition 6.1. The *Kullback-Leibler divergence* $D_{KL}(g|f)$ between a target pdf g and a candidate pdf f is defined as

$$D_{KL}(g|f) = \mathbb{E}_g[\log \frac{g}{f}] = \int g(x) \log g(x) dx - \int g(x) \log f(x) dx$$

In our setting, with $g = g^*$ and $f = f(\cdot, \theta)$ we have

$$\begin{aligned} D_{KL}(g^*|f(\cdot, \theta)) &= \mathbb{E}_{g^*}[\log g^*] - \mathbb{E}_{g^*}[\log f(\cdot, \theta)] \\ &= \mathbb{E}_{g^*}[\log g^*] - \frac{1}{\mathbb{E}_{\theta_0}[|\psi|]} \int |\psi(x)| f(x, \theta_0) \log f(x, \theta) dx. \end{aligned}$$

Notice that θ^* minimizes $D_{KL}(g^*|f(\cdot, \theta))$ if and only if it maximizes the quantity

$$J(\theta) = \int |\psi(x)| f(x, \theta_0) \log f(x, \theta) dx = \mathbb{E}_{\hat{\theta}} \left[|\psi(\cdot)| \frac{f(\cdot, \theta_0)}{f(\cdot, \hat{\theta})} \log f(\cdot, \theta) \right]$$

which can be approximated by an empirical mean from a sample drawn from the distribution $f(\cdot, \hat{\theta})$. Moreover, for certain families of distributions, such as the exponential one, the function $\theta \mapsto J(\theta)$ is concave, which makes the maximization problem easy to solve.

A cross-entropy adaptive Importance Sampling algorithm can be easily constructed as in Algorithm 6.4, replacing the optimization step 5 by cross entropy minimization as illustrated in Algorithm 6.5.

Algorithm 6.5: Adaptive importance sampling with cross entropy minimization

Given: tol , α , θ_0 , $\bar{N} > 1$, $\gamma > 1$

- 1 Set $N = \bar{N}$, $\hat{\theta} = \theta_0$, $\hat{\sigma} = \infty$
- 2 **while** $\frac{\hat{\sigma} c_{1-\alpha/2}}{\sqrt{N}} > tol$ **do**
- 3 Generate N iid replicas $Y^{(1)}, \dots, Y^{(N)} \sim f(\cdot, \hat{\theta})$
- 4 Compute

$$\hat{\mu}_{IS} = \frac{1}{N} \sum_{i=1}^N \psi(Y^{(i)}) \frac{f(Y^{(i)}, \theta_0)}{f(Y^{(i)}, \hat{\theta})} \quad \hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N \left(\psi(Y^{(i)}) \frac{f(Y^{(i)}, \theta_0)}{f(Y^{(i)}, \hat{\theta})} - \hat{\mu}_{IS} \right)^2$$
- 5 Solve the minimization problem

$$\hat{\theta}_{new} = \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{i=1}^N \psi(Y^{(i)}) \frac{f(Y^{(i)}, \theta_0)}{f(Y^{(i)}, \theta)} \log f(Y^{(i)}, \theta)$$
- 6 Set $\hat{\theta} = \hat{\theta}_{new}$ and $N = \gamma N$
- 7 **end**
- 8 Output $\hat{\mu}_{IS}$

6.2.2 Weighted importance sampling

In certain cases, the pdf f and/or the dominating pdf g , are known only up to a normalizing constant. (We assume, however, that we can still generate $X \sim g$ e.g. by Acceptance-Rejection). Let $f = c_f \tilde{f}$ and $g = c_g \tilde{g}$, with $c_f = (\int \tilde{f})^{-1}$ and $c_g = (\int \tilde{g})^{-1}$.

A modified (self-normalized) version of the importance sampling estimator, which does not require the explicit knowledge of the normalizing constants (c_f, c_g) is

$$\hat{\mu}_{IS}^W = \frac{\sum_{i=1}^N \psi(X^{(i)}) w(X^{(i)})}{\sum_{i=1}^N w(X^{(i)})}$$

with $w(x) = \frac{\tilde{f}(x)}{\tilde{g}(x)}$ and $X^{(i)} \stackrel{\text{iid}}{\sim} g$. Calling $\tilde{w}_i = \frac{w(X^{(i)})}{\sum_{i=1}^N w(X^{(i)})}$, the estimator $\hat{\mu}_{IS}^W$ can be written as a weighted average

$$\hat{\mu}_{IS}^W = \sum_{i=1}^N \psi(X^{(i)}) \tilde{w}_i.$$

To see that $\hat{\mu}_{IS}^W$ is a consistent estimator, observe that

$$\frac{1}{N} \sum_{i=1}^N w(X^{(i)}) \xrightarrow{\text{a.s.}} \int \frac{\tilde{f}(x)}{\tilde{g}(x)} g(x) dx = \frac{c_g}{c_f}$$

by the strong law of large numbers (SLLN) and

$$\frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}) w(X^{(i)}) \xrightarrow{\text{a.s.}} \int \psi \frac{\tilde{f}}{\tilde{g}} g dx = \frac{c_g}{c_f} \mu$$

again by SLLN. This estimator is biased, although the bias is usually small. Observe that this weighted version of the importance sampling estimator requires the stronger condition $f(x) = 0$ if $g(x) = 0$ (as opposed to the condition $\psi(x)f(x) = 0$ if $g(x) = 0$ of the standard estimator).

6.2.3 Importance sampling for stochastic processes

Discrete time Markov Chains.

Consider a homogeneous discrete time Markov chain in \mathbb{R}^d , $\{X_n, n \in \mathbb{N}_0\} \sim \text{Markov}(p_0, P)$, with Markov transition kernel defined by a density function $p : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$:

$$P(x, A) = \mathbb{P}(X_{n+1} \in A \mid X_n = x) = \int_A p(x, y) dy, \quad A \in \mathcal{B}(\mathbb{R}^d),$$

and initial probability p_0 , i.e. $X_0 \sim p_0$. We are interested in computing

$$\mu = \mathbb{E}[Z] = \mathbb{E}[\psi(X_0, \dots, X_m)]$$

for some finite horizon $m \in \mathbb{N}$. Importance sampling in this case can be done by replacing the transition kernel P by another kernel Q with density function $q : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ which dominates p , i.e. $q(x, y) = 0 \implies p(x, y) = 0$, and the initial density p_0 by a dominating one q_0 . We will use the shorthand notation $X_{0:m}$ to denote the path (X_0, \dots, X_m) and a subscript p_0, P to denote a Markov process $\text{Markov}(p_0, P)$. By successive conditioning, we have

$$\begin{aligned} \mu &= \mathbb{E}[\psi(X_0, \dots, X_m)] = \mathbb{E}_{X_{0:m-1} \sim p_{0,P}} [\mathbb{E}_{X_m \sim P(X_{m-1}, \cdot)} [\psi(X_{0:m}) \mid X_{m-1}, \dots, X_0]] \\ &= \mathbb{E}_{X_{0:m-1} \sim p_{0,P}} \left[\int \psi(X_{0:m-1}, z) p(X_{m-1}, z) dz \right] \\ &= \mathbb{E}_{X_{0:m-1} \sim p_{0,P}} \left[\int \psi(X_{0:m-1}, z) \frac{p(X_{m-1}, z)}{q(X_{m-1}, z)} q(X_{m-1}, z) dz \right] \\ &= \mathbb{E}_{X_{0:m-1} \sim p_{0,P}} \left[\mathbb{E}_{X_m \sim Q(X_{m-1}, \cdot)} [\psi(X_{0:m}) \frac{p(X_{m-1}, X_m)}{q(X_{m-1}, X_m)} \mid X_{m-1}, \dots, X_0] \right] \\ &= \mathbb{E}_{X_{0:m} \sim q_{0,Q}} \left[\psi(X_{0:m}) \frac{p_0(X_0)}{q_0(X_0)} \prod_{j=1}^m \frac{p(X_{j-1}, X_j)}{q(X_{j-1}, X_j)} \right] \\ &= \mathbb{E}_{X_{0:m} \sim q_{0,Q}} [\psi(X_{0:m}) w(X_{0:m})] \end{aligned}$$

with likelihood ratio

$$w(X_{0:m}) = \frac{p_0(X_0)}{q_0(X_0)} \prod_{j=1}^m \frac{p(X_{j-1}, X_j)}{q(X_{j-1}, X_j)}.$$

The previous argument can also be adapted to the case in which the process is stopped at some stopping time τ , e.g. $\tau = \inf\{n : X_n \in B \in \mathcal{B}(\mathbb{R}^d)\}$. Suppose we want to compute the quantity

$$\mu = \mathbb{E}[Z] = \mathbb{E}[\psi_\tau(X_0, \dots, X_\tau) \mathbb{1}_{\{\tau < +\infty\}}].$$

When doing importance sampling with the dominating density q_0 and transition kernel Q , we require that $\mathbb{P}_{q_0, Q}(\tau < +\infty) = 1$ to have a finite computational cost. Then, it can be shown that (exercise)

$$\mu = \mathbb{E}[Z] = \mathbb{E}_{p_0, P}[\psi_\tau(X_{0:\tau}) \mathbb{1}_{\{\tau < +\infty\}}] = \mathbb{E}_{q_0, Q}[\psi_\tau(X_0, \dots, X_\tau) w(X_{0:\tau})]$$

and μ can be estimated by the following algorithm:

Algorithm 6.6: Importance sampling for Markov processes.

- 1 Generate N iid paths $X_{0:\tau^{(i)}}^{(i)} = (X_0^{(i)}, \dots, X_{\tau^{(i)}}^{(i)})$, $i = 1, \dots, N$, each one up to the stopping time $\tau^{(i)}$, of the Markov chain with transition probability $q : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ and initial probability $q_0 : \mathbb{R}^d \rightarrow \mathbb{R}_+$
 - 2 Compute likelihood ratio $w(X_{0:\tau^{(i)}}^{(i)}) = \frac{p_0(X_0^{(i)})}{q_0(X_0^{(i)})} \prod_{k=1}^{\tau^{(i)}} \frac{p(X_{k-1}^{(i)}, X_k^{(i)})}{q(X_{k-1}^{(i)}, X_k^{(i)})}$
 - 3 Compute $\hat{\mu}_{\text{IS}} = \frac{1}{N} \sum_{i=1}^N \psi_{\tau^{(i)}}(X_{0:\tau^{(i)}}^{(i)}) w(X_{0:\tau^{(i)}}^{(i)})$
 - 4 Output $\hat{\mu}_{\text{IS}}$ and a confidence interval based on $\hat{\sigma}_{\text{IS}}$.
-

Discretized stochastic differential equations

Consider a stochastic differential equation in \mathbb{R}^d

$$dX_t = b(X_t, t)dt + \sigma(X_t, t)dW_t, \quad t > 0, \quad \text{with } X_0 \text{ given,} \quad (6.1)$$

where W_t is a d -dimensional Brownian motion (i.e. each component is a Brownian motion and the components are independent), and $b : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^d$ and $\sigma : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^{d \times d}$ are assumed sufficiently smooth so that a unique strong solution exists. We assume throughout that the matrix $\sigma(x, t)$ has full rank for any $(x, t) \in \mathbb{R}^{d+1}$. We aim at computing

$$\mu = \mathbb{E}[Z] = \mathbb{E}[\psi(\{X_t\}_{0 \leq t \leq T})]$$

where ψ is a function of the path $\{X_t\}_{0 \leq t \leq T}$ as, for example, $\psi(\{X_t\}_{0 \leq t \leq T}) = \int_0^T \phi(X_s) ds$ or $\psi(\{X_t\}_{0 \leq t \leq T}) = \phi(X_T)$. For this, we introduce a discretization of the stochastic differential equation (6.1) by the Euler-Maruyama scheme, on a grid $\{t_n = n\Delta t, n = 0, \dots, m = \frac{T}{\Delta t}\}$:

$$X_{n+1} = X_n + b(X_n, t_n)\Delta t + \sigma(X_n, t_n)\xi_n, \quad \xi_n \sim N(0, I_{d \times d}\Delta t).$$

It follows that $X_{n+1}|X_n \sim N(X_n + b(X_n, t_n)\Delta t, \Sigma(X_n, t_n))$ with $\Sigma(x, t) = \sigma(x, t)\sigma(x, t)^T \Delta t$. Then, $Z = \psi(\{X_t\}_{0 \leq t \leq T})$ can be approximated as $Z_{\Delta t} = \psi_{\Delta t}(X_0, \dots, X_m) = \psi(\xi_0, \dots, \xi_{m-1})$ and

$$\mu \approx \mu_{\Delta t} = \mathbb{E}_{\xi_0, \dots, \xi_{m-1}}[\hat{\psi}(\xi_0, \dots, \xi_{m-1})]$$

Importance sampling in this case, can be done, for instance, by changing the drift $b(x, t)$ to a new one $\tilde{b}(x, t)$. This can be achieved in the Euler-Maruyama scheme by changing the distribution of the Gaussian increments to $\tilde{\xi}_n \sim N(\phi(X_n, t_n)\Delta t, I_{d \times d}\Delta t)$ with

$$\phi(x, t) = \sigma^{-1}(x, t)(\tilde{b}(x, t) - b(x, t)).$$

so that the discretized path

$$X_{n+1} = X_n + b(X_n, t_n)\Delta t + \sigma(X_n, t_n)\tilde{\xi}_n,$$

has now conditional distribution $X_{n+1}|X_n \sim N(X_n + \tilde{b}(X_n, t_n)\Delta t, \Sigma(X_n, t_n))$ with the desired modified drift. If we denote by $z \mapsto p(z; \mu, \Sigma)$ the joint probability density function of a Gaussian vector with mean μ and covariance matrix Σ , then we have

$$\mu_{\Delta t} = \mathbb{E}_{\xi_{0:m-1}}[\hat{\psi}(\xi_{0:m-1})] = \mathbb{E}_{\tilde{\xi}_{0:m-1}}\left[\hat{\psi}(\tilde{\xi}_{0:m-1})w(\tilde{\xi}_{0:m-1})\right]$$

with likelihood ratio

$$\begin{aligned} w(\tilde{\xi}_{0:m-1}) &= \prod_{i=0}^{m-1} \frac{p(\tilde{\xi}_i; 0, I_{d \times d}\Delta t)}{p(\tilde{\xi}_i; \phi(X_i, t_i)\Delta t, I_{d \times d}\Delta t)} \\ &= \prod_{i=0}^{m-1} \exp\left(-\frac{1}{2\Delta t}\|\tilde{\xi}_i\|^2 + \frac{1}{2\Delta t}\|\tilde{\xi}_i - \phi(X_i, t_i)\Delta t\|^2\right) \\ &= \prod_{i=0}^{m-1} \exp\left(\frac{\Delta t}{2}\|\phi(X_i, t_i)\|^2 - \phi(X_i, t_i)^T \tilde{\xi}_i\right) \\ &= \exp\left(\frac{1}{2}\sum_{i=0}^{m-1} \Delta t\|\phi(X_i, t_i)\|^2 - \sum_{i=0}^{m-1} \phi(X_i, t_i)^T \tilde{\xi}_i\right) \end{aligned} \quad (6.2)$$

An importance sampling algorithm then reads

Algorithm 6.7: Importance sampling for SDEs.

- 1 Generate N iid paths $X_{0:m}^{(i)}$, $i = 1, \dots, N$ with modified drift

$$X_{n+1}^{(i)} = X_n^{(i)} + b(X_n^{(i)}, t_n)\Delta t + \sigma(X_n^{(i)}, t_n)\tilde{\xi}_n^{(i)}, \quad \tilde{\xi}_n^{(i)} \sim N(\phi(X_n, t_n)\Delta t, I_{d \times d}\Delta t)$$

- 2 Compute likelihood ratio

$$w(\tilde{\xi}_{0:m-1}^{(i)}) = \exp\left(\frac{1}{2}\sum_{n=0}^{m-1} \Delta t\|\phi(X_n^{(i)}, t_n)\|^2 - \sum_{n=0}^{m-1} \phi(X_n, t_n)^T \tilde{\xi}_n^{(i)}\right)$$

- 3 Compute $\hat{\mu}_{\text{IS}} = \frac{1}{N}\sum_{i=1}^N \hat{\psi}(\tilde{\xi}_{0:m-1}^{(i)})w(\tilde{\xi}_{0:m-1}^{(i)})$
 - 4 Output $\hat{\mu}_{\text{IS}}$ and a confidence interval based on $\hat{\sigma}_{\text{IS}}$.
-

As a matter of fact, what we have done is to change the distribution of the brownian increments. In the limit $\Delta t \rightarrow 0$ this corresponds to defining a drifted Brownian motion

\tilde{W}_t which satisfies $d\tilde{W}_t = \phi(X_t, t)dt + dW_t$. Then, the likelihood ratio represents the “ratio between the (joint) densities of W_t and \tilde{W}_t ” which we denote as $\frac{d\mathbb{P}_{W_t}}{d\mathbb{P}_{\tilde{W}_t}}$. Notice that in the limit $\Delta t \rightarrow 0$ the likelihood ratio (6.2) becomes

$$w(\{\tilde{W}_t\}_{0 \leq t \leq T}) = \exp\left(\frac{1}{2} \int_0^T \|\phi(X_t, t)\|^2 dt - \int_0^T \phi(X_t, t) \cdot d\tilde{W}_t\right) \quad (6.3)$$

and we have

$$\mu = \mathbb{E}_{W_t}[\psi(\{X_t\}_{0 \leq t \leq T})] = \mathbb{E}_{\tilde{W}_t} \left[\psi(\{X_t\}_{0 \leq t \leq T}) \frac{d\mathbb{P}_{W_t}}{d\mathbb{P}_{\tilde{W}_t}}(\tilde{W}_t) \right].$$

This is a well known result, known as Girsanov’s theorem, which says that, given a standard Brownian motion W_t and a “drifted” one $\tilde{W}_t = W_t + \int_0^t Z_s ds$ where $\{Z_t\}_t$ is an adapted process with enough integrability, e.g. $\mathbb{E} \left[\exp\left(\frac{1}{2} \int_0^T \|Z_s\|^2 ds\right) \right] < \infty$, then the likelihood ratio (more technically, the Radon Nikodym derivative $\frac{d\mathbb{P}_{W_t}}{d\mathbb{P}_{\tilde{W}_t}}$ of the original process with respect to the drifted one) is given by

$$\frac{d\mathbb{P}_{W_t}}{d\mathbb{P}_{\tilde{W}_t}}(B_t) = \exp\left(-\int_0^T Z_t \cdot dB_t + \frac{1}{2} \int_0^T \|Z_t\|^2 dt\right).$$

Continuous time discrete space Markov processes.

Consider a continuous time Markov process $\{X_t \in \mathcal{X}, t \geq 0\}$ taking values in the discrete space $\mathcal{X} = \{y_1, y_2, \dots\}$, defined by the stable and conservative generator matrix $(Q_{ij})_{ij}$ (see Section 4.7) and the initial distribution $X_0 \sim \lambda = (\lambda_1, \lambda_2, \dots)$, with $\lambda_i = \mathbb{P}(X_0 = y_i)$. We aim at computing

$$\mu = \mathbb{E}[Z] = \mathbb{E}[\psi(\{X_t\}_{0 \leq t \leq T})]$$

where ψ is a function of the path $\{X_t\}_{0 \leq t \leq T}$. We can do importance sampling in this case by changing the generator matrix to \tilde{Q} , and the initial distribution to $\tilde{\lambda}$, with the conditions that $\tilde{Q}_{ij} = 0$ and $\tilde{\lambda}_k = 0$ only if $Q_{ij} = 0$ and $\lambda_k = 0$, respectively. Then, if we denote by $N(t)$ the number of jumps of $\{X_t\}$ occurred in $[0, t]$, by J_n , $n = 1, \dots, N(T)$ the jump times, by $S_n = J_n - J_{n-1}$ the holding times, by $Y_n = X_{J_n}$ the jump process, by $\pi_{ij} = \frac{Q_{ij}}{Q_i}$ (with $Q_i = -Q_{ii} = \sum_{\ell} Q_{i\ell}$) the probability of jumping from state y_i to state y_j when a jump occurs, and similarly for $\tilde{\pi}_{ij} = \frac{\tilde{Q}_{ij}}{\tilde{Q}_i}$, it can be shown that

$$\mu = \mathbb{E}_{\lambda, Q}[\psi(\{X_t\}_{0 \leq t \leq T})] = \mathbb{E}_{\tilde{\lambda}, \tilde{Q}}[\psi(\{X_t\}_{0 \leq t \leq T}) w(\{X_t\}_{0 \leq t \leq T})]$$

with likelihood ratio given by

$$\begin{aligned}
w(\{X_t\}_{0 \leq t \leq T}) &= \frac{\lambda_{X_0}}{\tilde{\lambda}_{X_0}} \left(\prod_{i=1}^{N(T)} \frac{\pi_{Y_{i-1}Y_i} Q_{Y_{i-1}} \exp\{-S_i Q_{Y_{i-1}}\}}{\tilde{\pi}_{Y_{i-1}Y_i} \tilde{Q}_{Y_{i-1}} \exp\{-S_i \tilde{Q}_{Y_{i-1}}\}} \right) \frac{\exp\{-(T - J_{N(T)})Q_{Y_{N(T)}}\}}{\exp\{-(T - J_{N(T)})\tilde{Q}_{Y_{N(T)}}\}} \\
&= \frac{\lambda_{X_0}}{\tilde{\lambda}_{X_0}} \left(\prod_{i=1}^{N(T)} \frac{Q_{Y_{i-1}Y_i} \exp\{-S_i Q_{Y_{i-1}}\}}{\tilde{Q}_{Y_{i-1}Y_i} \exp\{-S_i \tilde{Q}_{Y_{i-1}}\}} \right) \frac{\exp\{-(T - J_{N(T)})Q_{Y_{N(T)}}\}}{\exp\{-(T - J_{N(T)})\tilde{Q}_{Y_{N(T)}}\}} \\
&= \frac{\lambda_{X_0}}{\tilde{\lambda}_{X_0}} \left(\prod_{i=1}^{N(T)} \frac{Q_{X_{J_{i-1}}X_{J_i}}}{\tilde{Q}_{X_{J_{i-1}}X_{J_i}}} \right) \exp \left\{ - \int_0^T (Q_{X_s} - \tilde{Q}_{X_s}) ds \right\}.
\end{aligned}$$

6.3 Control variates

We consider again the goal of computing the expected value $\mu = \mathbb{E}[Z]$ of a random variable Z , output of a stochastic model. The idea of the control variate technique is to find an auxiliary variable Y , called *control variate*, of which we *know the mean value*, and which is strongly correlated with the variable Z . We can then construct the modified variable

$$Z_\alpha = Z - \alpha(Y - \mathbb{E}[Y])$$

with $\alpha \in \mathbb{R}$, that satisfies

$$\mathbb{E}[Z_\alpha] = \mathbb{E}[Z] = \mu$$

and

$$\text{Var}(Z_\alpha) = \text{Var}(Z) + \alpha^2 \text{Var}(Y) - 2\alpha \text{Cov}(Z, Y).$$

The latter is a quadratic expression in α and is minimized for

$$\alpha_{\text{opt}} = \frac{\text{Cov}(Z, Y)}{\text{Var}(Y)}.$$

With such optimal choice, one has

$$\text{Var}(Z_{\alpha_{\text{opt}}}) = \text{Var}(Z) - \frac{\text{Cov}(Z, Y)^2}{\text{Var}(Y)} = \text{Var}(Z) (1 - \rho_{ZY}^2), \quad \rho_{ZY}^2 = \frac{\text{Cov}(Z, Y)^2}{\text{Var}(Z) \text{Var}(Y)}$$

which is always smaller than $\text{Var}(Z)$. The amount of variance reduction increases as ρ_{ZY} approaches 1 or -1 . It is clear that the ideal control variate is $Y = \gamma Z$, $\gamma \in \mathbb{R}$ for which $\text{Var}(Z_{\alpha_{\text{opt}}}) = 0$. However, $\mathbb{E}[Y] = \gamma \mathbb{E}[Z]$ is not known in this case, and such a control variate is not a viable option. The control variate Y should be a reasonable approximation of Z , of which, however, we can compute exactly its expected value, or, more generally, a random variable highly informative on Z (hence highly correlated to Z). In practice, the optimal α is not known, but it can be estimated from a pilot run.

The estimator $\hat{\mu}_{\text{CV}}$ is unbiased and, in the case σ_Y^2 known, has variance (exercise)

$$\text{Var}(\hat{\mu}_{\text{CV}}) = \mathbb{E}[(\hat{\mu}_{\text{CV}} - \mu)^2] = \frac{1}{N} (\text{Var}(Z_{\alpha_{\text{opt}}}) + \text{Var}(\hat{\alpha}_{\text{opt}}) \sigma_Y^2).$$

Algorithm 6.8: Control variate with pilot run.

- 1 Generate \bar{N} iid replicas $(Z^{(i)}, Y^{(i)})$, $i = 1, \dots, \bar{N}$ of (Z, Y)
- 2 Estimate $\hat{\alpha}_{\text{opt}} = \frac{\hat{\sigma}_{ZY}^2}{\sigma_Y^2}$ if σ_Y^2 known, or $\hat{\alpha}_{\text{opt}} = \frac{\hat{\sigma}_{ZY}^2}{\hat{\sigma}_Y^2}$ otherwise, with

$$\hat{\sigma}_{ZY}^2 = \frac{1}{\bar{N} - 1} \sum_{i=1}^{\bar{N}} (Z^{(i)} - \hat{\mu}_Z)(Y^{(i)} - \mathbb{E}[Y]), \quad \hat{\mu}_Z = \frac{1}{\bar{N}} \sum_{i=1}^{\bar{N}} Z^{(i)}$$

- 3 Generate N iid replicas $(Z^{(i)}, Y^{(i)})$, $i = 1, \dots, N$ of (Z, Y)
 - 4 Compute $\hat{\mu}_{\text{CV}} = \frac{1}{N} \sum_{i=1}^N (Z^{(i)} - \hat{\alpha}_{\text{opt}}(Y^{(i)} - \mathbb{E}[Y]))$
 - 5 Output $\hat{\mu}_{\text{CV}}$ and a confidence interval based on $\hat{\sigma}_{\text{CV}}$.
-

where $\text{Var}(\hat{\alpha}_{\text{opt}}) = O(1/\bar{N})$ since $\hat{\alpha}_{\text{opt}}$ is a Monte Carlo estimator, hence usually small compared with the first term. Moreover, $\text{Var}(Z_{\alpha_{\text{opt}}})$ can be estimated by the estimator

$$\hat{\sigma}^2(Z_{\alpha_{\text{opt}}}) = \hat{\sigma}_Z^2 - \frac{\hat{\sigma}_{ZY}^2}{\sigma_Y^2}$$

which is unbiased if σ_Y^2 is known. Based on these observations we can construct an approximate $1 - \alpha$ confidence interval as

$$\hat{I}_{\alpha, N} = \left[\hat{\mu}_{\text{CV}} - c_{1-\alpha/2} \frac{\hat{\sigma}(Z_{\alpha_{\text{opt}}})}{\sqrt{N}}, \hat{\mu}_{\text{CV}} + c_{1-\alpha/2} \frac{\hat{\sigma}(Z_{\alpha_{\text{opt}}})}{\sqrt{N}} \right]$$

which is justified, for \bar{N} not too small, by the observation that $\sqrt{\bar{N}} \frac{\hat{\mu}_{\text{CV}} - \mu}{\hat{\sigma}^2(Z_{\alpha_{\text{opt}}})} \xrightarrow{d} N(0, 1)$ as $N, \bar{N} \rightarrow \infty$.

Alternative to the previous algorithm, which uses a pilot run to estimate α_{opt} , one may consider a “one-shot” strategy.

Algorithm 6.9: Control variate – one shot

- 1 Generate N iid replicas $(Z^{(i)}, Y^{(i)})$, $i = 1, \dots, N$ of (Z, Y)
- 2 Estimate $\hat{\alpha}_{\text{opt}} = \frac{\hat{\sigma}_{ZY}^2}{\sigma_Y^2}$, with

$$\hat{\sigma}_{ZY}^2 = \frac{1}{N - 1} \sum_{i=1}^N (Z^{(i)} - \hat{\mu}_Z)(Y^{(i)} - \mathbb{E}[Y]), \quad \hat{\mu}_Z = \frac{1}{N} \sum_{i=1}^N Z^{(i)}$$

- 3 Estimate $\hat{\mu}_{\text{CV}} = \frac{1}{N} \sum_{i=1}^N (Z^{(i)} - \hat{\alpha}_{\text{opt}}(Y^{(i)} - \mathbb{E}[Y]))$
 - 4 Output $\hat{\mu}_{\text{CV}}$ and a confidence interval based on $\hat{\sigma}_{\text{CV}}$.
-

This estimator is biased, in general, contrary to the previous one. However, a CLT result still holds (exercise) and

$$\sqrt{N} \frac{\hat{\mu}_{\text{CV}} - \mu}{\hat{\sigma}^2(Z_{\alpha_{\text{opt}}})} \xrightarrow{d} N(0, 1)$$

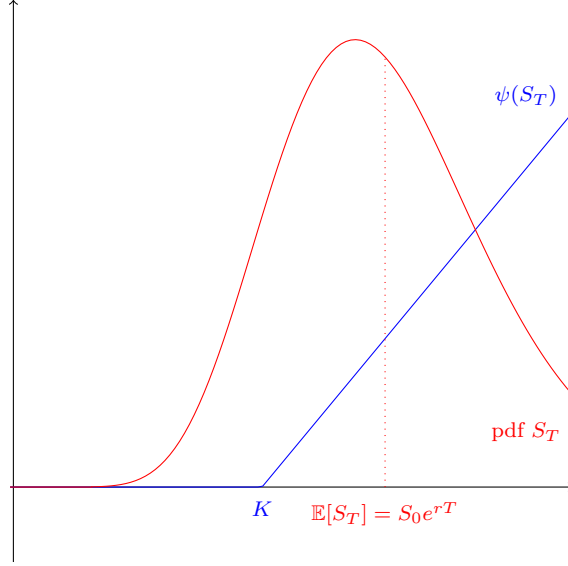


Figure 6.3: European option.

as $N \rightarrow \infty$ from which asymptotic confidence intervals can be obtained.

Example 6.5. Consider again the problem of pricing a European call option: $\mu = \mathbb{E}[Z]$, with $Z = \psi(S_T) = e^{-rT}(S_T - K)_+$, $S_T = S_0 e^{X_T}$ and $X_T \sim N((r - \sigma^2/2)T, \sigma^2 T)$. To compute $\mathbb{E}[Z]$ with Monte Carlo, we can use as a control variate the variable $Y = S_T$ whose exact mean is $\mathbb{E}[Y] = \mathbb{E}[S_T] = S_0 e^{rT}$. Observe that, since ψ is a non decreasing function of S_T , Z and Y are positively correlated, so that α should be taken positive. If the sample mean $\hat{\mu}_{S_T} = \frac{1}{N} \sum_{i=1}^N S_T^{(i)}$ is above the true mean $S_0 e^{rT}$, it is reasonable to assume that also the sample mean $\hat{\mu}_Z$ will be above the true (unknown) mean, since (Z, Y) are positively correlated, so we add a negative correctin to $\hat{\mu}_Z$ given by $-\alpha(\hat{\mu}_{S_T} - S_0 e^{rT})$, with $\alpha > 0$.

6.3.1 Multiple control variates

The control variate technique can be generalized to the case in which multiple control variates Y_1, \dots, Y_p are used. We define the modified variable

$$Z_\alpha = Z - \sum_{j=1}^p \alpha_j (Y_j - \mathbb{E}[Y_j]) = Z - \boldsymbol{\alpha} \cdot (\mathbf{Y} - \mathbb{E}[\mathbf{Y}])$$

with $\mathbf{Y} = (Y_1, \dots, Y_p)$ and $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_p)$. Then

$$\begin{aligned} \text{Var}(Z_\alpha) &= \mathbb{E}[(Z - \mu - \boldsymbol{\alpha} \cdot (\mathbf{Y} - \mathbb{E}[\mathbf{Y}]))^2] \\ &= \text{Var}(Z) - 2 \text{Cov}(Z, \mathbf{Y}) \cdot \boldsymbol{\alpha} + \boldsymbol{\alpha}^\top \text{Cov}(\mathbf{Y}, \mathbf{Y}) \boldsymbol{\alpha} \end{aligned}$$

where $\text{Cov}(Z, \mathbf{Y}) = (\text{Cov}(Z, Y_i))_{i=1}^p \in \mathbb{R}^p$ and $\text{Cov}(\mathbf{Y}, \mathbf{Y}) = (\text{Cov}(Y_i, Y_j))_{i,j=1}^p \in \mathbb{R}^{p \times p}$. Again $\text{Var}(Z_\alpha)$ is a quadratic function in $\boldsymbol{\alpha}$ and is minimized by

$$\boldsymbol{\alpha}_{\text{opt}} = \text{Cov}(\mathbf{Y}, \mathbf{Y})^{-1} \text{Cov}(Z, \mathbf{Y}).$$

Algorithm 6.10: Multiple control variates – one shot

- 1 Generate N iid replicas $(Z^{(i)}, Y_1^{(i)}, \dots, Y_p^{(i)})$ of (Z, \mathbf{Y})
- 2 Estimate

$$(\hat{\sigma}_{Z\mathbf{Y}}^2)_j = \frac{1}{N-1} \sum_{i=1}^N (Z^{(i)} - \hat{\mu}_Z)(Y_j^{(i)} - \mathbb{E}[Y_j]), \quad j = 1, \dots, p$$

and

$$(\hat{\sigma}_{\mathbf{Y}\mathbf{Y}}^2)_{jk} = \frac{1}{N} \sum_{i=1}^N (Y_j^{(i)} - \mathbb{E}[Y_j])(Y_k^{(i)} - \mathbb{E}[Y_k]).$$

- 3 Estimate the optimal $\boldsymbol{\alpha}_{\text{opt}}$ by $\hat{\boldsymbol{\alpha}}_{\text{opt}} = (\hat{\sigma}_{\mathbf{Y}\mathbf{Y}}^2)^{-1} \hat{\sigma}_{Z\mathbf{Y}}^2$
 - 4 Compute $\hat{\mu}_{\text{CV}} = \frac{1}{N} \sum_{i=1}^N (Z^{(i)} - \hat{\boldsymbol{\alpha}}_{\text{opt}} \cdot (\mathbf{Y}^{(i)} - \mathbb{E}[\mathbf{Y}])$.
 - 5 Output $\hat{\mu}_{\text{CV}}$ and a confidence interval based on $\hat{\sigma}_{\text{CV}}$.
-

6.4 Stratification

As in the previous sections, we consider the problem of computing $\mu = \mathbb{E}[Z]$ where Z is the output of a stochastic model. We assume here that $Z = \psi(X_1, \dots, X_d) = \psi(X)$ where $X \in \mathbb{R}^d$ is a random vector with pdf $f : \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}_+$ so that $\mu = \int_{\Omega} \psi(x) f(x) dx$.

The idea of stratification is to divide the sample space Ω into s non overlapping regions $\Omega_1, \dots, \Omega_s$ called *strata* such that $\mathbb{P}(X \in \Omega_j) = \int_{\Omega} \mathbf{1}_{\Omega_j}(x) f(x) dx = p_j$ is known and $\sum_{j=1}^s p_j = 1$. Assume now that we can generate X conditional upon $X \in \Omega_j$. The conditional density of X given $X \in \Omega_j$ is $f_j(x) = \frac{1}{p_j} f(x) \mathbf{1}_{\{x \in \Omega_j\}}$. Let now $X_j \sim f_j$ and $Z_j = \psi(X_j)$, $j = 1, \dots, s$. Clearly, $\mu = \mathbb{E}[Z] = \sum_{j=1}^s \mathbb{E}[Z | X \in \Omega_j] \mathbb{P}(X \in \Omega_j) = \sum_{j=1}^s p_j \mathbb{E}[Z_j]$. The idea is then to sample independently each $Z_j = \psi(X_j)$ leading to the following stratified estimator

$$\hat{\mu}_{\text{Str}} = \sum_{j=1}^s p_j \hat{\mu}_j, \quad \hat{\mu}_j = \frac{1}{N_j} \sum_{i=1}^{N_j} Z_j^{(i)}, \quad \text{with } Z_j^{(i)} \stackrel{\text{iid}}{\sim} Z_j. \quad (6.4)$$

The stratified estimator (6.4) has the following properties:

1. The estimator $\hat{\mu}_{\text{Str}}$ is *unbiased*. Indeed,

$$\mathbb{E}[\hat{\mu}_{\text{Str}}] = \sum_{j=1}^s p_j \mathbb{E}[\hat{\mu}_j] = \sum_{j=1}^s p_j \mathbb{E}[Z_j] = \mathbb{E}[Z].$$

2. The variance of the estimator satisfies

$$\text{Var}(\hat{\mu}_{\text{Str}}) = \sum_{j=1}^s p_j^2 \text{Var}(\hat{\mu}_j) = \sum_{j=1}^s p_j^2 \frac{\text{Var}(Z_j)}{N_j}$$

and can be estimated by

$$\hat{\sigma}_{Str}^2 = \sum_{j=1}^s p_j^2 \frac{\hat{\sigma}_j^2}{N_j}, \quad \hat{\sigma}_j^2 = \frac{1}{N_j - 1} \sum_{i=1}^{N_j} (Z_j^{(i)} - \hat{\mu}_j)^2.$$

3. Let $N = \sum_{j=1}^s N_j$ and choose $N_j = \phi_j(N)$, with $\sum_j \phi_j(N) = N$, such that $\lim_{N \rightarrow \infty} \frac{\phi_j(N)}{N} = c_j \in (0, 1)$ for any $j = 1, \dots, s$. Then $\lim_{N \rightarrow \infty} N \text{Var}(\hat{\mu}_{Str}) = \sum_j p_j^2 \sigma_j^2 / c_j < +\infty$ and it can be shown (exercise) that

$$\frac{\hat{\mu}_{Str} - \mu}{\sqrt{\text{Var}(\hat{\mu}_{Str})}} \xrightarrow{d} N(0, 1), \quad \text{as } N \rightarrow \infty.$$

Therefore, a computable $1 - \alpha$ asymptotic confidence interval is given by

$$\hat{I}_\alpha = [\hat{\mu}_{Str} - c_{1-\alpha/2} \hat{\sigma}_{Str}, \hat{\mu}_{Str} + c_{1-\alpha/2} \hat{\sigma}_{Str}]$$

We summarize the procedure in the following Algorithm.

Algorithm 6.11: Stratification

- 1 **for** $j = 1, \dots, s$ **do**
 - 2 Generate N_j iid replicas $Z_j^{(i)}$, $i = 1, \dots, N_j$ of Z_j
 - 3 Compute $\hat{\mu}_j = \frac{1}{N_j} \sum_{i=1}^{N_j} Z_j^{(i)}$ and $\hat{\sigma}_j^2 = \frac{1}{N_j - 1} \sum_{i=1}^{N_j} (Z_j^{(i)} - \hat{\mu}_j)^2$
 - 4 **end**
 - 5 Compute $\hat{\mu}_{Str} = \sum_{j=1}^s p_j \hat{\mu}_j$ and $\hat{\sigma}_{Str}^2 = \sum_{j=1}^s p_j^2 \frac{\hat{\sigma}_j^2}{N_j}$
 - 6 Output $\hat{\mu}_{Str}$ and a confidence interval $\hat{I}_\alpha = [\hat{\mu}_{Str} - c_{1-\alpha/2} \hat{\sigma}_{Str}, \hat{\mu}_{Str} + c_{1-\alpha/2} \hat{\sigma}_{Str}]$
-

Stratification guarantees that each stratum contains a fixed number of evaluations. It remains the question of how to choose N_j in each stratum and quantify the amount of variance reduction that we can achieve.

6.4.1 Proportional allocation

If N is the total sample size, proportional allocation simply chooses $N_j = N p_j$. With this choice, we have

$$\text{Var}(\hat{\mu}_{Str}) = \sum_{j=1}^s p_j^2 \frac{\text{Var}(Z_j)}{N_j} = \frac{1}{N} \sum_{j=1}^s p_j \text{Var}(Z_j).$$

Defining the discrete random variable $J \in \{1, \dots, s\}$, $J = j \iff \{X \in \Omega_j\}$, we can rewrite $\text{Var}(\hat{\mu}_{Str})$ as

$$\text{Var}(\hat{\mu}_{Str}) = \frac{1}{N} \sum_{j=1}^s p_j \text{Var}(Z \mid J = j) = \frac{1}{N} \mathbb{E}_J[\text{Var}(Z \mid J)]$$

and, recalling the law of total variance $\text{Var}(Z) = \text{Var}(\mathbb{E}[Z | J]) + \mathbb{E}[\text{Var}(Z | J)]$ we have

$$\text{Var}(\hat{\mu}_{\text{Str}}) = \frac{1}{N} (\text{Var}(Z) - \text{Var}(\mathbb{E}[Z | J])) \leq \frac{\text{Var}(Z)}{N} = \text{Var}(\hat{\mu}_{\text{CMC}}).$$

Hence proportional allocation always leads to variance reduction. The amount of variance reduction is given by $\gamma = \mathbb{E}[\text{Var}(Z | J)] / \text{Var}(Z)$.

Example 6.6. Let $X \sim \mathcal{U}(0, 1)$ and $Z = \psi(X)$ for some function $\psi : [0, 1] \rightarrow \mathbb{R}$. To compute $\mu = \mathbb{E}[Z] = \int_0^1 \psi(x) dx$, we could use stratification by dividing the interval $\Omega = (0, 1)$ in s subintervals of equal size, $\Omega_j = \left(\frac{j-1}{s}, \frac{j}{s}\right)$, $j = 1, \dots, s$. Then

$$\mu = \sum_{j=1}^s \int_{\frac{j-1}{s}}^{\frac{j}{s}} \psi(x) dx = \sum_{j=1}^s \frac{1}{s} \int_{\frac{j-1}{s}}^{\frac{j}{s}} \psi(x) s dx = \sum_{j=1}^s \frac{1}{s} \mathbb{E}[\psi(X_j)], \quad \text{with } X_j \sim \mathcal{U}\left(\frac{j-1}{s}, \frac{j}{s}\right)$$

and a stratified estimator reads

$$\hat{\mu}_{\text{Str}} = \sum_{j=1}^s \frac{1}{s} \frac{1}{N_j} \sum_{i=1}^{N_j} \psi(X_j^{(i)}), \quad \text{with } X_j^{(i)} \stackrel{iid}{\sim} \mathcal{U}\left(\frac{j-1}{s}, \frac{j}{s}\right).$$

Figure 6.4 gives an illustration of the stratification procedure with 7 strata and 2 replicas per stratum.

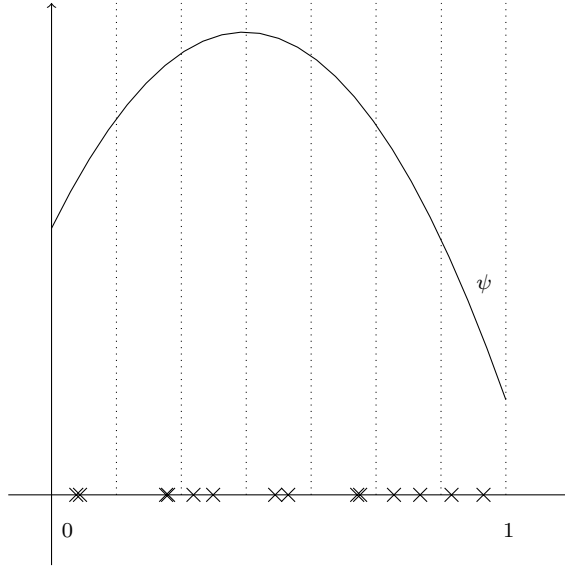


Figure 6.4: Stratification.

Figure 6.5 illustrates the variance reduction when considering proportional allocation. The variance of a crude Monte Carlo estimator is proportional to the green area in the right plot, whereas the variance of the stratified estimator is proportional to the green area in the left plot. From this graphical illustration we see that large variance reduction has to be expected when the function ψ is highly non-constant. If ψ is piecewise constant over the partition of the domain, then we even have $\text{Var}(\hat{\mu}_{\text{Str}}) = 0$.

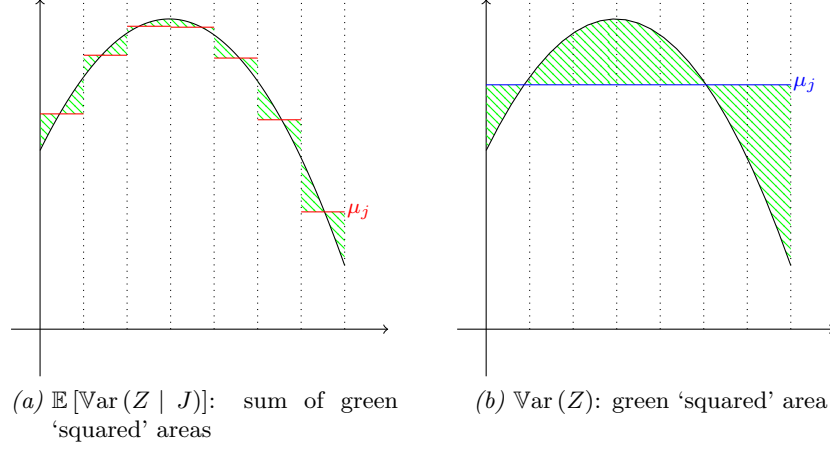


Figure 6.5: Proportional allocation

6.4.2 Optimal allocation

Instead of doing a proportional allocation, one may try to find the best choice of N_j that minimises $\text{Var}(\hat{\mu}_{\text{Str}})$:

$$\{N_j^*\} = \underset{(N_1, \dots, N_s)}{\text{argmin}} \sum_{j=1}^s p_j^2 \frac{\text{Var}(Z_j)}{N_j} \quad \text{such that} \quad \sum_{j=1}^s N_j = N.$$

Introducing a Lagrangian function $\mathcal{L}(\mathbf{N}, \lambda) = \sum_{j=1}^s p_j^2 \frac{\text{Var}(Z_j)}{N_j} + \lambda(\sum_{j=1}^s N_j - N)$, we have

$$\frac{\partial \mathcal{L}}{\partial N_j} = -p_j^2 \frac{\text{Var}(Z_j)}{N_j^2} + \lambda = 0 \quad \implies \quad N_j = p_j \sqrt{\frac{\text{Var}(Z_j)}{\lambda}}$$

and, enforcing the constraint $\sum_j N_j = N$, we obtain $\sqrt{\lambda} = \frac{\sum_{j=1}^s p_j \sqrt{\text{Var}(Z_j)}}{N}$ which leads to the optimal choice

$$N_j^* = \frac{N p_j \sigma_j}{\sum_{k=1}^s p_k \sigma_k}, \quad \sigma_j = \sqrt{\text{Var}(Z_j)}$$

and optimal variance $\text{Var}(\hat{\mu}_{\text{Str}}^*) = \frac{1}{N} \left(\sum_{j=1}^s p_j \sigma_j \right)^2$.

Since this variance is smaller than that with proportional allocation, stratification with optimal allocation will always lead to variance reduction. In practice, the σ_j are not known and can be obtained from a pilot run.

Algorithm 6.12: Stratification with optimal allocation

-
- 1 **for** $j = 1, \dots, s$ **do**
 - 2 Generate \bar{N}_j iid replicas $Z_j^{(i)}$, $i = 1, \dots, \bar{N}_j$ of Z_j
 - 3 Estimate $\hat{\sigma}_j^2 = \frac{1}{\bar{N}_j - 1} \sum_{i=1}^{\bar{N}_j} (Z_j^{(i)} - \hat{\mu}_j)^2$
 - 4 **end**
 - 5 Choose $N = (c_{1-\alpha/2} \sum_{j=1}^s p_j \hat{\sigma}_j / \text{tol})^2$ (to guarantee that $|\hat{I}_{\alpha, N}| < 2\text{tol}$)
 - 6 For $j = 1, \dots, s$, generate $N_j^* = \frac{N p_j \hat{\sigma}_j}{\sum_k p_k \hat{\sigma}_k}$ iid replicas $Z_j^{(i)}$ of Z_j
 - 7 Compute $\hat{\mu}_i = \frac{1}{N_j^*} \sum_{i=1}^{N_j^*} Z_j^{(i)}$ and $\hat{\mu}_{\text{Str}} = \sum_{j=1}^s p_j \hat{\mu}_j$
-

6.5 Latin Hypercube Sampling

Consider the problem of computing the expected value μ of $Z = \psi(X_1, \dots, X_d)$ where $X_j \in \mathbb{R}$ are independent and with pdf $f_j : \mathbb{R} \rightarrow \mathbb{R}_+$. One might want to stratify each variable X_j in s strata. However, this would lead to s^d strata which becomes unaffordable for large d . A way to overcome this problem is offered by the Latin Hypercube Sampling (LHS). For simplicity of exposition, let us assume that $X = (X_1, \dots, X_d) \sim \mathcal{U}([0, 1]^d)$. The idea of LHS is to stratify each component X_j but not the whole sampling domain $\Omega = [0, 1]^d$. In particular, N (correlated) points $X^{(i)}$, $i = 1, \dots, N$ are drawn in $[0, 1]^d$ in such a way that each component is stratified with N strata and one point per stratum. Figure 6.6 illustrates the idea.

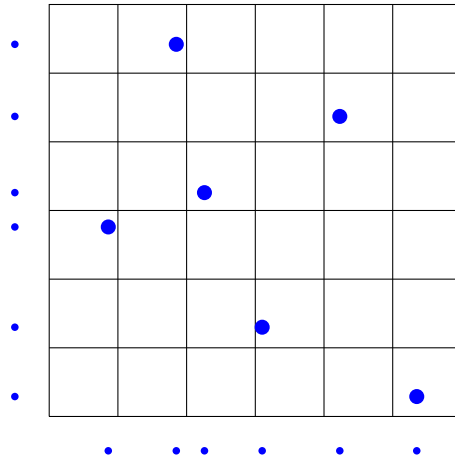


Figure 6.6: Latin hypercube.

A Latin hypercube sampling design can be generated by the following Algorithm.

Algorithm 6.13: LHS

-
- 1 Generate N iid points $\mathbf{U}^{(i)} \stackrel{\text{iid}}{\sim} \mathcal{U}((0, 1)^d)$, $i = 1, \dots, N$
 - 2 Generate d independent permutations π_j , $j = 1, \dots, d$ of $\{1, \dots, N\}$. Let
 $\boldsymbol{\pi}^{(i)} = (\pi_1(i), \pi_2(i), \dots, \pi_d(i))$
 - 3 Return $X^{(i)} = \frac{\boldsymbol{\pi}^{(i)} - \mathbf{1} + \mathbf{U}^{(i)}}{N}$, $i = 1, \dots, N$.
-

Once the LHS design generated, the LHS estimator of $\mu = \mathbb{E}[\psi(X)]$ is simply

$$\hat{\mu}_{\text{LHS}} = \frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}).$$

The following proposition illustrates the two main properties of the LHS sample and estimator.

Proposition 6.4. *Let $\{X^{(i)}, i = 1, \dots, N\}$ be a LHS design. Then*

- $X^{(i)} \sim \mathcal{U}((0, 1)^d)$ (not independent, though)
- The LHS estimator is unbiased, $\mathbb{E}[\hat{\mu}_{\text{LHS}}] = \mathbb{E}[\psi(X)]$.

Proof. By construction, each vector $X^{(i)} = \frac{\boldsymbol{\pi}^{(i)} - \mathbf{1} + \mathbf{U}^{(i)}}{N}$ has independent components. Therefore it is enough to show that each component $X_j^{(i)}$, $j = 1, \dots, d$, is uniformly distributed in $(0, 1)$. Now, $\pi_j^{(i)} = \pi_j(i)$ is the i -th component of a random permutation of $\{1, \dots, N\}$, hence $\mathbb{P}(\pi_j^{(i)} = k) = \frac{1}{N}$ for all $k = 1, \dots, N$. Moreover, the conditional cumulative distribution function of $X_j^{(i)}$ given $\pi_j^{(i)} = k$ is

$$F_{X_j^{(i)} | \pi_j^{(i)} = k}(x) = \mathbb{P}\left(X_j^{(i)} \leq x \mid \pi_j^{(i)} = k\right) = \begin{cases} 0, & x < \frac{k-1}{N} \\ Nx - k + 1, & x \in \left[\frac{k-1}{N}, \frac{k}{N}\right] \\ 1, & x > \frac{k}{N} \end{cases}$$

i.e. $X_j^{(i)} \mid \pi_j^{(i)} = k$ has distribution $\mathcal{U}\left(\frac{k-1}{N}, \frac{k}{N}\right)$ and

$$\mathbb{P}\left(X_j^{(i)} \leq x\right) = \sum_{k=1}^N \frac{1}{N} \mathbb{P}\left(X_j^{(i)} \leq x \mid \pi_j^{(i)} = k\right) = x.$$

From the uniform distribution of each $X^{(i)}$, it follows immediately that $\mathbb{E}[\hat{\mu}_{\text{LHS}}] = \mathbb{E}\left[\frac{1}{N} \sum_i \psi(X^{(i)})\right] = \mathbb{E}[\psi(X)]$. \square

Concerning the variance of the estimator $\hat{\mu}_{\text{LHS}}$, we mention the following two results.

Proposition 6.5 ([8]). *Let $Z = \psi(X)$, $X \sim \mathcal{U}((0, 1)^d)$, with $\mu = \mathbb{E}[Z] < +\infty$ and $\sigma^2 = \text{Var}(Z) < +\infty$. The LHS estimator $\hat{\mu}_{\text{LHS}}$ based on N points satisfies*

$$\text{Var}(\hat{\mu}_{\text{LHS}}) \leq \frac{\sigma^2}{N-1}.$$

This result shows that, asymptotically, $\text{Var}(\hat{\mu}_{\text{LHS}})$ is not worse than $\text{Var}(\hat{\mu}_{\text{CMC}}) = \frac{\sigma^2}{n}$ since $\lim_{N \rightarrow \infty} \text{Var}(\hat{\mu}_{\text{LHS}}) / \text{Var}(\hat{\mu}_{\text{CMC}}) \leq 1$. Moreover, LHS is very effective if the function $\psi(X)$ has an additive structure $\psi(X) = \mu + \sum_{i=1}^d \psi_j(X_j)$ as the estimator $\hat{\mu}_{\text{LHS}}$ corresponds to a stratified estimator with N strata on each function ψ_j . For a general $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$, let

$$\hat{\psi}_j(x_j) = \int_{[0,1]^{d-1}} (\psi(x_1, \dots, x_d) - \mu) dx_1 \dots dx_{j-1} dx_{j+1} \dots dx_d$$

and

$$\psi^{\text{add}}(X) = \mathbb{E}[\psi] + \sum_{i=1}^d \hat{\psi}_i(x_i).$$

The function $\hat{\psi}_j$ can be interpreted as a conditional expectation $\hat{\psi}_j(X_j) = \mathbb{E}[\psi(X) - \mu \mid X_j]$ and is often called the main effect of X_j in ψ . Then, it can be shown that:

Proposition 6.6 ([9, 5]). *For $Z = \psi(X)$, $X \sim \mathcal{U}((0, 1)^d)$ and $\mu = \mathbb{E}[Z] < +\infty$, $\sigma^2 = \text{Var}(Z) < +\infty$ and $\hat{\mu}_{\text{LHS}}$, the LHS estimator for μ based on N points satisfies*

$$\text{Var}(\hat{\mu}_{\text{LHS}}) = \frac{\text{Var}(\psi - \psi^{\text{add}})}{N} + o\left(\frac{1}{N}\right).$$

Moreover, if $\mathbb{E}[|Z|^3] < +\infty$, then $\sqrt{N}(\hat{\mu}_{\text{LHS}} - \mu) \xrightarrow{d} N(0, \text{Var}(\psi - \psi^{\text{add}}))$ as $N \rightarrow \infty$.

This result highlights the amount of variance reduction that can be achieved by the LHS estimator, compared to the CMC one. Unfortunately, the estimate of $\text{Var}(\hat{\mu}_{\text{LHS}})$ in Proposition 6.6 is not computable and can not be used to build confidence intervals for the estimator $\hat{\mu}_{\text{LHS}}$.

To control the error in the LHS estimator, we proceed in a different way by generating few independent replicas of $\hat{\mu}_{\text{LHS}}$ and estimating its variance by a sample variance estimator. Since, by proposition 6.6, $\hat{\mu}_{\text{LHS}}$ is nearly Gaussian distributed for large N , whenever Z has bounded third moments, we can build a confidence interval based on the Student's t distribution (5.2) with $K - 1$ degrees of freedom if K is the number of replicas used.

Algorithm 6.14: LHS estimator

- 1 Generate K independent LHS designs $\{X^{(i,j)}\}_{i=1}^N$ of size N , for $j = 1, \dots, K$
 - 2 For each desing compute $\hat{\mu}_{\text{LHS}}^{(j)} = \frac{1}{N} \sum_{i=1}^N \psi(X^{(i,j)})$
 - 3 Compute $\hat{\mu}_{\text{LHS}} = \frac{1}{K} \sum_{j=1}^K \hat{\mu}_{\text{LHS}}^{(j)}$ and $\hat{\sigma}_{\text{LHS}}^2 = \frac{1}{K-1} \sum_{j=1}^K \left(\hat{\mu}_{\text{LHS}}^{(j)} - \hat{\mu}_{\text{LHS}}\right)^2$
 - 4 Output $\hat{\mu}_{\text{LHS}}$ and the (Student's t based) confidence interval

$$\hat{I}_\alpha = \left[\hat{\mu}_{\text{LHS}} - t_{1-\alpha/2}^{(K-1)} \frac{\hat{\sigma}_{\text{LHS}}}{\sqrt{K}}, \hat{\mu}_{\text{LHS}} + t_{1-\alpha/2}^{(K-1)} \frac{\hat{\sigma}_{\text{LHS}}}{\sqrt{K}} \right]$$
-

Chapter 7

Quasi Monte Carlo methods

As in the previous chapter, we consider the problem of computing the expected value $\mu = \mathbb{E}[Z]$, of some random variable Z output of a stochastic model. We assume in this chapter that $Z = \psi(\mathbf{X})$, with $\mathbf{X} = (X_1, \dots, X_d) \sim \mathcal{U}([0, 1]^d)$, hence computing μ turns into computing a possibly high dimensional integral over the unit hypercube

$$\mu = \int_{[0,1]^d} \psi(x_1, \dots, x_d) dx_1 \dots dx_d.$$

A Crude Monte Carlo estimator $\hat{\mu}_{CMC}$ that uses N iid replicas of \mathbf{X} , achieves an error

$$|\mu - \hat{\mu}_{CMC}| \leq c_{1-\alpha/2} \frac{\sqrt{\text{Var}(\psi(\mathbf{X}))}}{\sqrt{N}}$$

with asymptotic confidence $1 - \alpha$. The idea of Quasi Monte Carlo (QMC) sampling, is to consider, instead, a *purely deterministic* sample $\{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$ to improve the rate $1/\sqrt{N}$, while keeping the simple structure of the sample average estimator $\hat{\mu}_{QMC} = \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)})$ with equal weights $1/N$. It relies on the observation that a random uniform sample does not seem to cover “uniformly” the hypercube and hopefully there exist better designs that achieve this goal.

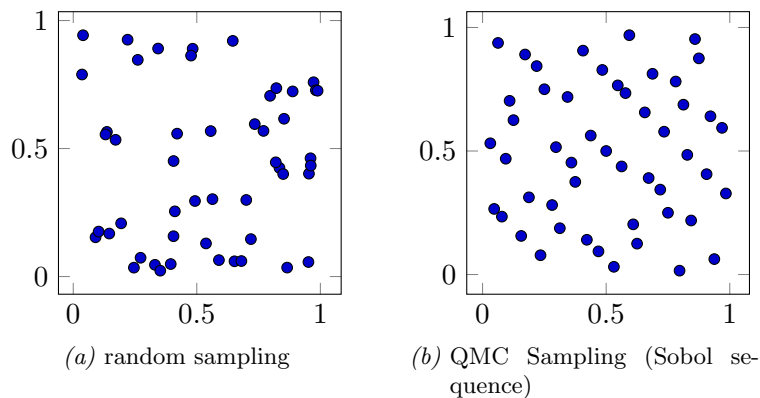


Figure 7.1: Comparing a uniform random sample (left) and a QMC sample with the same number of points on the unit hypercube.

Figure 7.1 shows a random sample and a QMC sample, with 50 points each, on the unit square.

The main notion behind QMC sampling is that of *discrepancy*. We introduce the following notation: for a point $\mathbf{y} \in [0, 1]^d$, $\mathbf{y} = (y_1, \dots, y_d)$, we denote by $[\mathbf{0}, \mathbf{y}]$ the hyperrectangle $[\mathbf{0}, \mathbf{y}] = \prod_{i=1}^d [0, y_i]$, with volume $\text{Vol}([\mathbf{0}, \mathbf{y}]) = \prod_{i=1}^d y_i$. For an arbitrary sample $\mathcal{P} = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$ of N points in $[0, 1]^d$, hereafter called a *point set*, we introduce the empirical volume estimator for $\text{Vol}([\mathbf{0}, \mathbf{y}])$, based on the point set \mathcal{P} .

$$\widehat{\text{Vol}}_{\mathcal{P}}([\mathbf{0}, \mathbf{y}]) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\mathbf{0}, \mathbf{y}]}(\mathbf{X}^{(i)}) = \frac{\#\{\mathbf{X}^{(i)} \in [\mathbf{0}, \mathbf{y}]\}}{N}.$$

Definition 7.1. We call *discrepancy function* $\Delta_{\mathcal{P}} : [0, 1]^d \rightarrow [-1, 1]$ the function

$$\Delta_{\mathcal{P}}(\mathbf{y}) = \widehat{\text{Vol}}_{\mathcal{P}}([\mathbf{0}, \mathbf{y}]) - \text{Vol}([\mathbf{0}, \mathbf{y}]) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\mathbf{0}, \mathbf{y}]}(\mathbf{X}^{(i)}) - \prod_{j=1}^d y_j.$$

From $\Delta_{\mathcal{P}}$, we define the following measures of discrepancy of a point set \mathcal{P} :

$$\begin{aligned} L_q\text{-discrepancy:} \quad D_q(\mathcal{P}) &= \|\Delta_{\mathcal{P}}\|_{L^q} = \left(\int_{[0,1]^d} |\Delta_{\mathcal{P}}(\mathbf{y})|^q d\mathbf{y} \right)^{1/q}, \quad 1 \leq q < \infty, \\ \text{Star-discrepancy:} \quad D^*(\mathcal{P}) &= \|\Delta_{\mathcal{P}}\|_{L^\infty} = \sup_{\mathbf{y} \in [0,1]^d} |\Delta_{\mathcal{P}}(\mathbf{y})|. \end{aligned}$$

Remark 7.1. There is actually nothing special in choosing only the rectangles $[\mathbf{0}, \mathbf{y}]$, so one can define also the so called extreme discrepancy

$$D(\mathcal{P}) = \sup_{\substack{\mathbf{y}, \mathbf{z} \in [0,1]^d \\ \mathbf{z} < \mathbf{y}}} |\widehat{\text{Vol}}_{\mathcal{P}}([\mathbf{z}, \mathbf{y}]) - \text{Vol}([\mathbf{z}, \mathbf{y}])|.$$

It can be easily shown that $D^*(\mathcal{P}) \leq D(\mathcal{P}) \leq 2^d D^*(\mathcal{P})$. The left inequality is obvious and the right one follows from the observation that a rectangle $[\mathbf{z}, \mathbf{y}]$ can be written as a composition (union/intersection) of 2^d rectangles of the type $[\mathbf{0}, \mathbf{z}]$. Hence, it is enough to study only the star-discrepancy.

The reason why the discrepancy plays an important role in the study of QMC quadrature formulas follows from the famous Koksma-Hlawka inequality, which we illustrate first in dimension $d = 1$. We start by deriving the following identity.

Lemma 7.1 (Zaremba's identity). Let $\psi : [0, 1] \rightarrow \mathbb{R}$ be an absolutely continuous function with integrable derivative and let $\mathcal{P} = \{X^{(1)}, \dots, X^{(N)}\}$ be any point set in $[0, 1]$. Then

$$\begin{aligned} \int_0^1 \psi(x) dx - \frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}) &= \int_0^1 \psi'(y) \Delta_{\mathcal{P}}(y) dy \\ &= \int_0^1 \psi'(y) \Delta_{\mathcal{P}}(y) dy - \Delta_{\mathcal{P}}(1) \psi(1). \end{aligned} \tag{7.1}$$

Proof. Using the identity $\psi(x) = \psi(1) - \int_x^1 \psi'(y) dy$ in the left hand side of (7.1), we have

$$\begin{aligned} \int_0^1 \psi(x) dx - \frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}) &= \psi(1) - \underbrace{\int_0^1 \int_x^1 \psi'(y) dy dx}_{=\int_0^1 \int_0^y \psi'(y) dx dy} - \frac{1}{N} \sum_{i=1}^N \psi(1) + \frac{1}{N} \sum_{i=1}^N \underbrace{\int_{X^{(i)}}^1 \psi'(y) dy}_{=\int_0^1 \psi'(y) \mathbf{1}_{[0,y]}(X^{(i)}) dy} \\ &= \int_0^1 \psi'(y) \left[\frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[0,y]}(X^{(i)}) - y \right] dy \\ &= \int_0^1 \psi'(y) \Delta_{\mathcal{P}}(y) dy. \end{aligned}$$

The second inequality follows immediately by observing that $\Delta_{\mathcal{P}}(1) = 0$ for any point set \mathcal{P} . \square

From the Zaremba's identity, we derive easily the **Koksma-Hlawka** inequality:

$$\left| \int_0^1 \psi(x) dx - \frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}) \right| \leq \|\psi'\|_{L_p} \|\Delta_{\mathcal{P}}\|_{L_q}, \quad \forall p, q \in [1, \infty], \quad \frac{1}{p} + \frac{1}{q} = 1. \quad (7.2)$$

Inequality (7.2) shows that the quadrature error is proportional to the discrepancy measure $\|\Delta_{\mathcal{P}}\|_{L_q}$, provided that $\psi' \in L_p(0, 1)$, i.e. $\psi \in W^{1,p}(0, 1)$. In particular, if ψ' is integrable (or ψ has bounded total variation) then

$$\left| \int_0^1 \psi(x) dx - \frac{1}{N} \sum_{i=1}^N \psi(X^{(i)}) \right| \leq \|\psi\|_{\text{TV}} D^*(\mathcal{P}).$$

The previous analysis extends with same care to the multi-dimensional setting. We introduce the following notation: let $\mathbf{u} = \{u_1, \dots, u_k\} \subset \{1, \dots, d\}$ be a subset of dimensions (without repetition) and set $|\mathbf{u}| = k$. For $\mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d$, we denote by $\mathbf{x}_{\mathbf{u}} = (x_{u_1}, \dots, x_{u_k}) \in [0, 1]^k$ and $\mathbf{z} = (\mathbf{x}_{\mathbf{u}}, 1)$ the vector with components $z_j = x_j$ if $j \in \mathbf{u}$ and $z_j = 1$ if $j \notin \mathbf{u}$. With this notation at hand, the Zaremba's identity generalizes to the multi-dimensional case as follows.

Lemma 7.2 (Hlawka's identity). *Let $\psi : [0, 1]^d \rightarrow \mathbb{R}$ be an integrable function with integrable mixed first order derivatives of any order, and let $\mathcal{P} = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$ be an arbitrary point set in $[0, 1]^d$. Then*

$$\frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)}) - \int_{[0,1]^d} \psi(\mathbf{x}) d\mathbf{x} = \sum_{\mathbf{u} \subset \{1, \dots, d\}} (-1)^{|\mathbf{u}|} \int_{[0,1]^{|\mathbf{u}|}} \frac{\partial^{|\mathbf{u}|} \psi}{\partial \mathbf{x}_{\mathbf{u}}}(\mathbf{x}_{\mathbf{u}}, 1) \Delta_{\mathcal{P}}(\mathbf{x}_{\mathbf{u}}, 1) d\mathbf{x}_{\mathbf{u}}$$

where $\frac{\partial^{|\mathbf{u}|} \psi}{\partial \mathbf{x}_{\mathbf{u}}} = \frac{\partial^k \psi}{\partial x_{u_1} \dots \partial x_{u_k}}$ is a mixed first order derivative.

Proof. By induction on d , one can prove the following identity

$$\psi(\mathbf{x}) = \sum_{\mathbf{u} \subset \{1, \dots, d\}} (-1)^{|\mathbf{u}|} \int_{[\mathbf{x}_{\mathbf{u}}, 1]} \frac{\partial^{|\mathbf{u}|} \psi}{\partial \mathbf{x}_{\mathbf{u}}}(\mathbf{y}_{\mathbf{u}}, 1) d\mathbf{y}_{\mathbf{u}}, \quad \forall \mathbf{x} \in [0, 1]^d, \quad (7.3)$$

which generalizes the $d = 1$ identity $\psi(x_1) = \psi(1) - \int_{x_1}^1 \frac{\partial \psi}{\partial x_1}(y) dy$ already used in the proof of Lemma 7.1. In (7.3) we have used the convention that for $\mathbf{u} = \emptyset$, $(-1)^{|\mathbf{u}|} \int_{[\mathbf{x}_{\mathbf{u}}, 1]} \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) d\mathbf{y}_{\mathbf{u}} = \psi(1, \dots, 1)$. Then

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)}) - \int_{[0,1]^d} \psi(\mathbf{x}) d\mathbf{x} \\ &= \sum_{\mathbf{u} \subset \{1, \dots, d\}} (-1)^{|\mathbf{u}|} \left(\frac{1}{N} \sum_{i=1}^N \underbrace{\int_{[\mathbf{X}_{\mathbf{u}}^{(i)}, 1]} \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) d\mathbf{y}_{\mathbf{u}}}_{=\int_{[0,1]^{|\mathbf{u}|}} \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) \mathbb{1}_{[0, \mathbf{y}_{\mathbf{u}}]}(\mathbf{X}_{\mathbf{u}}^{(i)}) d\mathbf{y}_{\mathbf{u}}} - \underbrace{\int_{[0,1]^d} \int_{[\mathbf{x}_{\mathbf{u}}, 1]} \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) d\mathbf{y}_{\mathbf{u}} d\mathbf{x}}_{=\int_{[0,1]^{|\mathbf{u}|}} \int_{[0, \mathbf{y}_{\mathbf{u}}]} \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) d\mathbf{x}_{\mathbf{u}} d\mathbf{y}_{\mathbf{u}}} \right) \\ &= \sum_{\mathbf{u} \subset \{1, \dots, d\}} (-1)^{|\mathbf{u}|} \int_{[0,1]^{|\mathbf{u}|}} \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \mathbb{1}_{[0, \mathbf{y}_{\mathbf{u}}]}(\mathbf{X}_{\mathbf{u}}^{(i)}) - \text{Vol}([0, \mathbf{y}_{\mathbf{u}}]) \right)}_{\Delta_{\mathcal{P}}(\mathbf{y}_{\mathbf{u}}, 1)} \end{aligned}$$

□

From the Hlawka's identity, the multidimensional version of the Koksma-Hlawka inequality follows. Let us define the following norm

$$\|\psi\|_{p, p'} = \left(\sum_{\mathbf{u} \subset \{1, \dots, d\}} \left(\int_{[0,1]^{|\mathbf{u}|}} \left| \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{x}_{\mathbf{u}}} \psi(\mathbf{y}_{\mathbf{u}}, 1) \right|^p d\mathbf{y}_{\mathbf{u}} \right)^{p'/p} \right)^{1/p'}$$

Then, the multidimensional **Koksma-Hlawka inequality** reads

$$\left| \int_{[0,1]^d} \psi(\mathbf{x}) d\mathbf{x} - \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)}) \right| \leq \|\psi\|_{p, p'} \|\Delta_{\mathcal{P}}\|_{q, q'}, \quad \text{with } \frac{1}{p} + \frac{1}{q} = \frac{1}{p'} + \frac{1}{q'} = 1, \quad (7.4)$$

provided $\|\psi\|_{p, p'} < +\infty$. In particular, if $\|\psi\|_{1,1} < +\infty$, then

$$\left| \int_{[0,1]^d} \psi(\mathbf{x}) d\mathbf{x} - \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)}) \right| \leq \|\psi\|_{1,1} D^*(\mathcal{P}).$$

Again, this inequality shows that the quadrature error is proportional to the star-discrepancy $D^*(\mathcal{P})$ of the point set, provided ψ has *integrable mixed first order derivatives*.

7.1 Low discrepancy sequences and point sets

There exist constructions of families of point sets $\mathcal{P} = \{\mathcal{P}_N\}_{N \in \mathbb{N}}$, with $\mathcal{P}_N = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\} \subset [0, 1]^d$, that have star-discrepancy as low as $D^*(\mathcal{P}_N) = O\left(\frac{(\log N)^{d-1}}{N}\right)$ for $\mathcal{P}_N \in \mathcal{P}$. It is widely believed that this result is sharp, i.e. there do not exist points sets that achieve a better bound. In general, these constructions do not lead to a nested sequence of points,

that is, the point set $\mathcal{P}_M = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}\} \in \mathcal{P}$ with $M > N$ does not contain $\mathcal{P}_N \in \mathcal{P}$, in general.

For nested point sets, i.e. point sets $\mathcal{S}_N = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$ that are generated as the first N points of an infinite sequence $\mathcal{S} = \{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots\}$ the lowest achievable star-discrepancy is slightly worse, namely $D^*(\mathcal{S}_N) = O\left(\frac{(\log N)^d}{N}\right)$. In view of these results, we give the following definition.

Definition 7.2. (low discrepancy sets).

- A family $\mathcal{P} = \{\mathcal{P}_N\}_{N \in \mathbb{N}}$ of non-nested point sets $\mathcal{P}_N = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\} \subset [0, 1]^d$ is called a *low discrepancy family of point sets* if $D^*(\mathcal{P}_N) = O\left(\frac{(\log N)^{d-1}}{N}\right)$;
- A point sequence $\mathcal{S} = \{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots\} \subset [0, 1]^d$ is called a *low discrepancy sequence* if $D^*(\mathcal{S}_N) = O\left(\frac{(\log N)^d}{N}\right)$.

From the above definitions and considerations, we see that a QMC quadrature formula can achieve convergence rate $1/N$ up to logarithmic terms (which however grow exponentially in the dimension!), provided the integrand function has integrable mixed first derivatives. Before presenting some common low discrepancy sequences/points sets, we give two important clarifying examples:

Example 7.1. Consider the family $\mathcal{P} = \{\mathcal{P}_N\}_{N \in \mathbb{N}}$ of point sets $\mathcal{P}_N = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$ with $\mathbf{X}^{(i)} \stackrel{iid}{\sim} \mathcal{U}([0, 1]^d)$, i.e. a random iid sample. Then, $|\Delta_{\mathcal{P}_N}(\mathbf{y})| = |\widehat{\text{Vol}}_{\mathcal{P}_N}([\mathbf{0}, \mathbf{y}]) - \text{Vol}([\mathbf{0}, \mathbf{y}])|$ is the error of the sample average estimator of $\text{Vol}([\mathbf{0}, \mathbf{y}])$, which decays as $O\left(\frac{1}{\sqrt{N}}\right)$, in a probabilistic sense. We conclude that the family of random iid point sets has not low discrepancy.

Example 7.2. Consider the family $\mathcal{P} = \{\mathcal{P}_N\}_{N \in \mathbb{N}}$ of point sets given by regular lattices (see figure 7.2)

$$\mathcal{P}_N = \left\{ \left(\frac{k_1 + 1/2}{m}, \dots, \frac{k_d + 1/2}{m} \right), 0 \leq k_j \leq m - 1, j = 1, \dots, d \right\}, \quad N = m^d$$

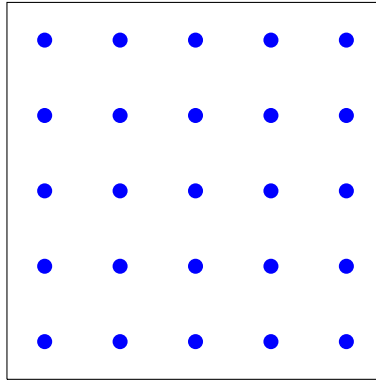


Figure 7.2: Regular lattice.

For $d = 1$, it is easy to see that $D^*(\mathcal{P}_N) = \frac{1}{2m} = \frac{1}{2N}$, hence \mathcal{P} has low discrepancy. On the other hand, in dimension $d > 1$ we have

$$D^*(\mathcal{P}_N) = \sup_{\mathbf{y} \in [0,1]^d} |\Delta_{\mathcal{P}_N}(\mathbf{y})| \geq \sup_{t \in [0,1]} |\Delta_{\mathcal{P}_N}(t, 1, \dots, 1)| = \frac{1}{2m} = \frac{1}{2N^{1/d}}.$$

We conclude then that the family of regular lattices has not low discrepancy in dimension higher than 1.

Van der Corput-Halton sequence

Let $b \geq 2$ be an integer. Any natural number $n \in \mathbb{N}_0$ can be expanded in a b -adic expansion $n = n_0 + n_1b + n_2b^2 + \dots$. The radical inverse of n is defined as

$$\varphi_b(n) = \frac{n_0}{b} + \frac{n_1}{b^2} + \dots$$

Obviously $\varphi_b : \mathbb{N}_0 \rightarrow [0, 1)$. In 1D, the b -adic Van der Corput sequence is

$$\varphi_b(0), \varphi_b(1), \varphi_b(2), \dots$$

For example, for $b = 2$, the Van der Corput sequence is $0, \frac{1}{2}, \frac{1}{4}, \frac{3}{4}, \frac{1}{8}, \frac{5}{8}, \frac{3}{8}, \frac{7}{8}, \dots$

The *Halton* sequence generalizes this construction for $d \geq 2$: Let $b_1, \dots, b_d \geq 2$ be integers pairwise relatively prime. Typically b_1, \dots, b_d are taken as the first d prime numbers. Then, the Halton sequence is

$$\mathcal{S} = \{\mathbf{X}^{(n)}, n \in \mathbb{N}_0\}, \quad \mathbf{X}^{(n)} = (\varphi_{b_1}(n), \varphi_{b_2}(n), \dots, \varphi_{b_d}(n))$$

and achieves the optimal bound on the star-discrepancy $D^*(\mathcal{S}_N) \leq c(d) \frac{(\log N)^d}{N}$.

Hammersley point set

It is derived from the Halton sequence by taking $\mathcal{P}_N = \{\mathbf{X}^{(0)}, \dots, \mathbf{X}^{(N-1)}\}$ with $\mathbf{X}^{(n)} = (\frac{n}{N}, \varphi_{b_1}(n), \dots, \varphi_{b_{d-1}}(n))$. The family $\mathcal{P} = \{\mathcal{P}_N\}_{N \in \mathbb{N}}$ of Hammersley point sets is non-nested and achieves the better bound $D^*(\mathcal{P}_N) = c(d) \frac{(\log N)^{d-1}}{N}$.

Rank-1 lattice point sets

Let $N \in \mathbb{N}$ and $\mathbf{g} \in \mathbb{N}^d$, $\mathbf{g} = (g_1, \dots, g_d)$ such that g_j has no factor in common with N . (Typically N is taken as a prime number.) Then the rank-1 N -lattice point set with generating vector \mathbf{g} is defined as

$$\mathcal{P}_N = \left\{ \frac{n\mathbf{g}}{N} \right\}_{n=0}^{N-1}$$

where $\{\cdot\}$ denotes the fractional part. Figure 7.3 shows an example of lattice point set. Good choices of \mathbf{g} lead to low discrepancy non-nested point sets.

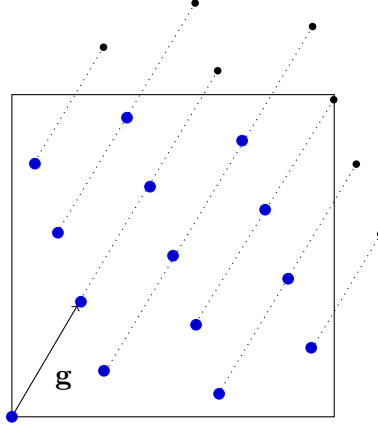


Figure 7.3: Lattice point set with $N = 14$ and $\mathbf{g} = (3, 5)$

(t-m-d)-nets and (t-d) sequences in base b

Let $0 \leq t \leq m \in \mathbb{N}$ and $b \geq 2$. A $(t-m-d)$ -net in base b is a point set \mathcal{P}_N consisting of $N = b^m$ points such that each elementary rectangle of volume b^{t-m} ,

$$R_a = \prod_{j=1}^d \left[\frac{a_j - 1}{b^{p_j}}, \frac{a_j}{b^{p_j}} \right), \quad a_j = 1, \dots, b^{p_j}$$

with $p_1 + p_2 + \dots + p_d = m - t$ contains exactly b^t points. E.g. if $t = 0$, each elementary rectangle of volume b^{-m} contains exactly 1 point.

Example 7.3. A $(0-3-2)$ -net in base $b = 2$ is a point set with $N = 2^3 = 8$ points, such that each elementary rectangle with volume $2^{-(m-t)} = 2^{-3} = 1/8$ contains exactly $2^t = 1$ point. Figure 7.4 shows graphically this property.

A $(t-d)$ sequence in base b is a sequence $\mathcal{S} = \{\mathbf{X}^{(0)}, \mathbf{X}^{(1)}, \dots\}$ such that for any $m > t$, every block of b^m points $\{\mathbf{X}^{(\ell b^m)}, \dots, \mathbf{X}^{((\ell+1)b^m-1)}\}$, $\ell \in \mathbb{N}$ is a $(t-m-d)$ -net in base b . The star-discrepancy of a $(t-m-d)$ -net satisfies $D^*(\mathcal{P}_N) = O\left(b^t \frac{(\log N)^{d-1}}{N}\right)$ and similarly for a $(t-d)$ -sequence $D^*(\mathcal{S}_N) = O\left(b^t \frac{(\log N)^d}{N}\right)$. Famous $(t-d)$ -sequences are those of Sobol, Niederreiter and Faure. For a description of their construction we refer to [3].

7.2 Randomized QMC formulas

Let us consider a point set $\mathcal{P}_N = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$ and the QMC quadrature formula

$$\hat{\mu}_{\text{QMC}} = \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{X}^{(i)}).$$

The question is how to estimate the error $|\mu - \hat{\mu}_{\text{QMC}}|$. Since the points $\mathbf{X}^{(i)}$ are not random iid, we can not use a variance estimator or a CLT as in the Monte Carlo estimator. On

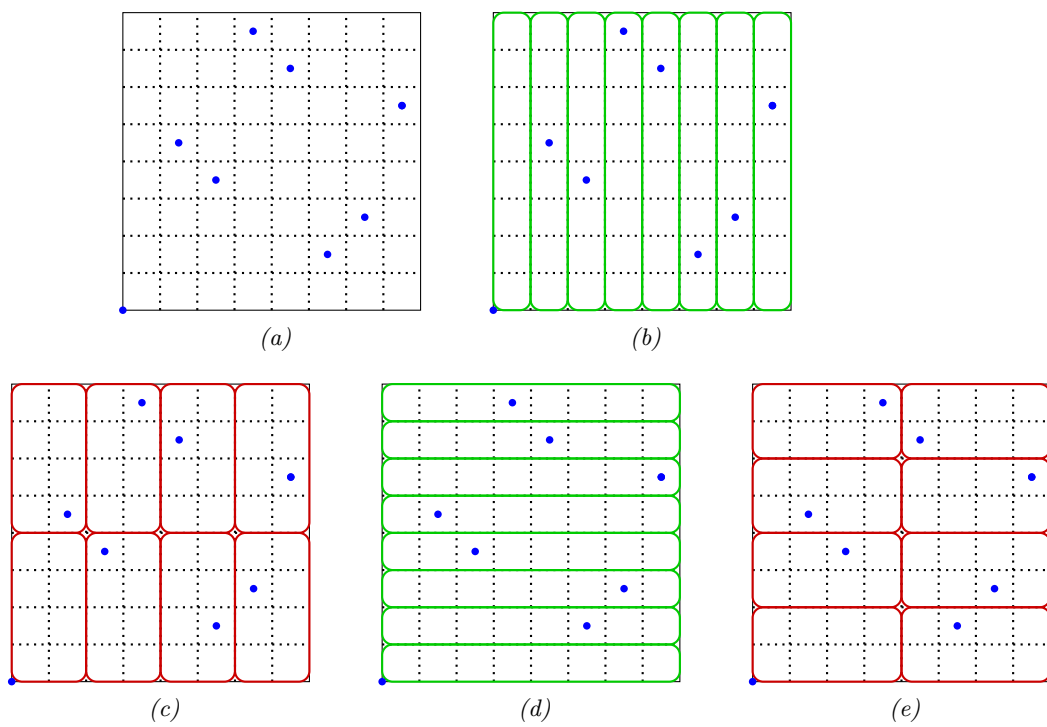


Figure 7.4: Example of a $(0,3,2)$ net in base 2. Each elementary rectangle of volume 2^{-3} contains exactly $2^0 = 1$ points.

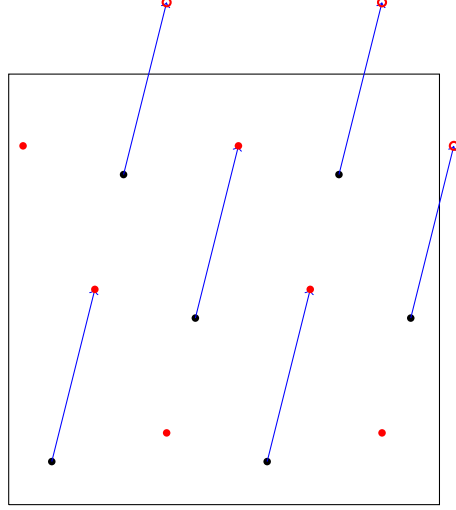


Figure 7.5: Randomized QMC.

the other hand, the error estimates in (7.4) can not be really used in practice to provide a bound on the quadrature error as they involve quantities such as the discrepancy or TV-norm of the integrand that are not known and can not be easily estimated.

An easy idea to obtain error bounds is to randomize the QMC formula. Let $\mathbf{U} \sim \mathcal{U}([0, 1]^d)$. If $\mathcal{P} = \{\mathcal{P}_N\}_{N \in \mathbb{N}}$, with $\mathcal{P}_N = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}\}$, is a low discrepancy point set, so is

$$\mathcal{P}^{\mathbf{U}} = \{\mathcal{P}_N^{\mathbf{U}}\}_{N \in \mathbb{N}}, \quad \text{with } \mathcal{P}_N^{\mathbf{U}} = \{\{\mathbf{X}^{(1)} + \mathbf{U}\}, \{\mathbf{X}^{(2)} + \mathbf{U}\}, \dots, \{\mathbf{X}^{(N)} + \mathbf{U}\}\}$$

where the same shift is applied to all points and again $\{\cdot\}$ denotes the fractional part. $\mathcal{P}_N^{\mathbf{U}}$ is called a randomly shifted point set. We could then compute $\hat{\mu}_{\text{QMC}}^{(j)}$, $j = 1, \dots, K$, for few randomly shifted point sets and average the obtained results. The resulting randomly shifted QMC estimator is then $\hat{\mu}_{\text{QMC}} = \frac{1}{K} \sum_{j=1}^K \hat{\mu}_{\text{QMC}}^{(j)}$. Since $\mathbf{U}^{(j)} \sim \mathcal{U}([0, 1]^d)$, so is $\{\mathbf{X}^{(i)} + \mathbf{U}^{(j)}\}$ for any $i = 1, \dots, N$. It follows that $\hat{\mu}_{\text{QMC}}$ is an *unbiased* estimator of $\mu = \mathbb{E}[\psi]$. Moreover, since $\hat{\mu}_{\text{QMC}}^{(j)}$ are independent, the variance of the estimator is $\text{Var}(\hat{\mu}_{\text{QMC}}) = \frac{\sigma_{\text{QMC}}^2}{K}$ with $\sigma_{\text{QMC}}^2 = \mathbb{E}[(\hat{\mu}_{\text{QMC}}^{(j)} - \mu)^2] = O\left(\frac{(\log N)^{2(d-1)}}{N^2}\right)$ hence, very small, in general, and can be estimated by the standard sample variance estimator $\hat{\sigma}_{\text{QMC}}^2 = \frac{1}{K-1} \sum_{j=1}^K (\hat{\mu}_{\text{QMC}}^{(j)} - \hat{\mu}_{\text{QMC}})^2$.

Algorithm 7.1: Randomly shifted QMC.

- 1 Generate point set $\mathcal{P}_N = (\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)})$
- 2 Generate $\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(K)} \stackrel{\text{iid}}{\sim} \mathcal{U}([0, 1]^d)$;
- 3 For $j = 1, \dots, K$, compute $\hat{\mu}_{\text{QMC}}^{(j)} = \frac{1}{N} \sum_{i=1}^N \psi(\{\mathbf{X}^{(i)} + \mathbf{U}^{(j)}\})$;
- 4 Compute $\hat{\mu}_{\text{QMC}} = \frac{1}{K} \sum_{j=1}^K \hat{\mu}_{\text{QMC}}^{(j)}$ as well as $\hat{\sigma}_{\text{QMC}}^2 = \frac{1}{K-1} \sum_{j=1}^K (\hat{\mu}_{\text{QMC}}^{(j)} - \hat{\mu}_{\text{QMC}})^2$;
- 5 Output $\hat{\mu}_{\text{QMC}}$ as well as a $1 - \alpha$ confidence interval

$$\hat{I}_\alpha = \left[\hat{\mu}_{\text{QMC}} - c_{1-\alpha/2} \frac{\hat{\sigma}_{\text{QMC}}}{\sqrt{K}}, \hat{\mu}_{\text{QMC}} + c_{1-\alpha/2} \frac{\hat{\sigma}_{\text{QMC}}}{\sqrt{K}} \right]$$

Chapter 8

Markov Chain Monte Carlo

Let π be a given probability density function on a state space $\mathcal{X} \subset \mathbb{R}^n$ and $\psi : \mathcal{X} \rightarrow \mathbb{R}$ an integrable function with respect to π . We consider the goal of computing $\mu = \mathbb{E}_\pi[\psi] = \int_{\mathcal{X}} \psi(x) \pi(x) dx$.

If we can generate independent replicas of $Z \sim \pi$, then μ can be computed by Monte Carlo or any improved version using variance reduction techniques. Assume, however, that sampling directly from π is not viable either because the expression of π is too complicated and possibly high dimensional, or because π is known only up to a multiplicative constant and computing the normalization constant might be too expensive, if not impossible.

Example 8.1 (Bayesian statistics). *Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from a parametric density $g(\mathbf{x} \mid \theta)$. Then the joint density of \mathbf{X} given θ is $g(\mathbf{X} \mid \theta) = \prod_{i=1}^n g(X_i \mid \theta)$ and we want to estimate θ from the sample \mathbf{X} . In the Bayesian paradigm, θ is thought as a random variable itself, with prior density $\pi_0(\theta)$, which summarizes any prior information on θ in the absence of data. Then, the posterior density of θ given the data is*

$$\pi(\theta) = \frac{1}{Z(\mathbf{X})} g(\mathbf{X} \mid \theta) \pi_0(\theta)$$

with $Z(\mathbf{X}) = \int g(\mathbf{X} \mid \theta) \pi_0(\theta) d\theta$ which is often unknown and difficult to compute.

Example 8.2 (Statistical physics). *Let $x \in \mathcal{X}$ be a configuration of a physical system and \mathcal{X} the configuration space. Let $H : \mathcal{X} \rightarrow \mathbb{R}$ be an energy function and T the temperature. Then the probability density function of finding the system in a given state x is*

$$\pi(x) = \frac{1}{Z} \exp \left\{ -\frac{H(x)}{kT} \right\}$$

where k is the Boltzmann constant and $Z = \int e^{-H(x)/kT} dx$ is the partition function, often difficult to compute.

The idea of Markov chain Monte Carlo (MCMC) is to construct an ergodic Markov Chain $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ on \mathcal{X} that has π as its invariant distribution. Then we can approximate $\mu = \mathbb{E}_\pi[\psi]$ by the ergodic estimator

$$\hat{\mu}_N^{\text{MCMC}} = \frac{1}{N} \sum_{i=1}^N \psi(X_i) \tag{8.1}$$

or

$$\hat{\mu}_{N, N_0}^{\text{MCMC}} = \frac{1}{N} \sum_{i=1}^N \psi(X_{i+N_0})$$

if we want to “cut” out the first part of the chain, which might be too sensitive to the initial state $X_0 \sim \lambda$ of the chain (this operation is usually called “burn-in”). We will see that constructing a Markov Chain with a given invariant distribution is not so difficult and can be achieved by the well known and celebrated Metropolis-Hastings algorithm. Before discussing such algorithm, however, it is worth recalling some basic concepts in the theory of Markov Chains. We will do so in the finite state space case in the next section and briefly mention generalizations to general state spaces in Section 8.2.

8.1 Markov Chains on discrete state spaces (review)

Let $\mathcal{X} = \{x_1, x_2, \dots, x_d\}$ be a finite ($d < \infty$) or countably infinite ($d = \infty$) state space, $\lambda = \{\lambda_1, \lambda_2, \dots, \lambda_d\}$ a probability mass function on \mathcal{X} , with $\lambda_i \geq 0$ for all $i = 1, \dots, d$, $\sum_i \lambda_i = 1$, and $P = \{P_{ij}\}_{i,j=1}^d$ a stochastic matrix, such that $P_{ij} \geq 0$ for all $i, j = 1, \dots, d$ and $\sum_j P_{ij} = 1$ for all $i = 1, \dots, d$.

We consider hereafter a *homogeneous* Markov chain $\{X_n, n \in \mathbb{N}_0\} \sim \text{Markov}(\lambda, P)$ having initial state $X_0 \sim \lambda$ and transition matrix P independent of n (See Chapter 4 for the definition of a Markov Chain). To highlight the dependence of the chain on the initial distribution λ , we denote by $\mathbb{P}_\lambda(A)$ the probability of an event A under $X_0 \sim \lambda$. If $\lambda = \delta_{x_i}$, i.e. $\mathbb{P}(X_0 = x_i) = 1$, we use the notation \mathbb{P}_{x_i} or simply \mathbb{P}_i . We introduce also the following notion

Definition 8.1 (Stopping time). *A random variable τ is called a stopping time if the event $\{\tau \leq n\}$ depends only on X_0, \dots, X_n , i.e. the event $\{\tau \leq n\}$ is measurable with respect to the σ -algebra $\sigma(X_0, \dots, X_n)$ generated by X_0, \dots, X_n .*

Typical examples of stopping times are the following: given a subset $A \subset \mathcal{X}$

- hitting time of A : $\tau_A = \inf\{n \geq 0 : X_n \in A\}$,
- return time to A : $\sigma_A = \inf\{n > 0 : X_n \in A\}$,
- successive return times to A : $\sigma_A^{(k)} = \inf\{n > \sigma_A^{(k-1)} : X_n \in A\}$, for $k \geq 1$,

with the conventions that $\sigma_A^{(0)} = 0$, $\tau_A = +\infty$ if $X_n \notin A$ for any n , and $\sigma_A^{(k)} = +\infty$ if $X_n \notin A$ for any $n > \sigma_A^{(k-1)}$. From the definition of a homogeneous Markov chain, the following *Markov property* follows.

Lemma 8.1. *Let $\{X_n, n \in \mathbb{N}_0\} \sim \text{Markov}(\lambda, P)$.*

- (Weak Markov Property). *Conditional on $X_m = x_i$, $\{X_{m+n}, n \in \mathbb{N}_0\}$ is Markov (δ_{x_i}, P) and independent of $\{X_0, \dots, X_m\}$.*
- (Strong Markov property). *Let τ be a stopping time of $\{X_n\}$. Conditional on $\tau < +\infty$ and $X_\tau = x_i$, $\{X_{\tau+n}, n \in \mathbb{N}_0\}$ is Markov (δ_{x_i}, P) independent of X_0, \dots, X_τ .*

Given $\{X_n, n \in \mathbb{N}_0\} \sim \text{Markov}(\lambda, P)$, let $P^{(n)}$ denote the n -step transition matrix, i.e. $P_{ij}^{(n)} = \mathbb{P}(X_{m+n} = x_j \mid X_m = x_i)$. Thanks to the homogeneity of the Markov chain, $P_{ij}^{(n)}$ does not depend on m . Clearly $P^{(1)} = P$ and for $n > 1$,

$$\begin{aligned} P_{ij}^{(n)} &= \sum_{\ell} \mathbb{P}(X_{m+n} = x_j \mid X_{m+n-1} = x_{\ell}, X_m = x_i) \mathbb{P}(X_{m+n-1} = x_{\ell} \mid X_m = x_i) \\ &= \sum_{\ell} P_{\ell j} P_{i\ell}^{(n-1)}. \end{aligned}$$

Introducing the matrix multiplication $(P^2)_{ij} = \sum_{\ell} P_{i\ell} P_{\ell j}$, we see that $P^{(n)} = P^n$. More generally, $P^{(n+m)} = P^n P^m$ which is often referred to as the *Chapman Kolmogorov equation*.

We may also ask what is the probability distribution of X_n at any given $n > 0$, i.e. the probability mass function $\pi^{n,\lambda} = (\pi_1^{n,\lambda}, \dots, \pi_d^{n,\lambda})$, taken as a row vector in \mathbb{R}^d , with $\pi_i^{n,\lambda} = \mathbb{P}_{\lambda}(X_n = x_i)$. It is easy to see that

$$\pi_i^{n,\lambda} = \sum_{\ell} \mathbb{P}(X_n = x_i \mid X_{n-1} = x_{\ell}) \mathbb{P}(X_{n-1} = x_{\ell}) = \sum_{\ell} P_{\ell i} \pi_{\ell}^{n-1,\lambda}.$$

In matrix notation,

$$\pi^{n,\lambda} = \pi^{n-1,\lambda} P = \lambda P^n.$$

If we denote by $M_1(\mathcal{X}) = \{(\mu_1, \dots, \mu_d) \in \mathbb{R}^d : \mu_i \geq 0, \sum_i \mu_i = 1\}$ the set of probability mass functions on \mathcal{X} , then the transition matrix P can be interpreted as an operator $P : M_1(\mathcal{X}) \rightarrow M_1(\mathcal{X})$ acting (to the left) on probability measures. We may ask if such an operator has a fixed point.

Definition 8.2. A probability mass function $\pi \in M_1(\mathcal{X})$ is called invariant distribution for P if $\pi P = \pi$.

Hence, for a Markov chain $\{X_n\} \sim \text{Markov}(\lambda, P)$ whose initial state $X_0 \sim \pi$ is distributed as the invariant distribution π , it follows that $X_n \sim \pi$ for any n and the chain is said to be “at equilibrium” or “at stationarity”. Observe that, if an invariant distribution π exists, then it is a left eigenvector of the transition matrix P , associated to the eigenvalue $\lambda_1 = 1$.

Consider now the set $\mathcal{F}(\mathcal{X}) = \{\varphi : \mathcal{X} \rightarrow \mathbb{R}\}$ of real valued measurable functions on \mathcal{X} , which can be identified with \mathbb{R}^d . We represent any function $\varphi \in \mathcal{F}(\mathcal{X})$ as a column vector $\varphi = (\varphi_1, \dots, \varphi_d)^{\top} \in \mathbb{R}^d$. Given $\varphi \in \mathcal{F}(\mathcal{X})$, we can define the following function $g \in \mathcal{F}(\mathcal{X})$:

$$g_i = \mathbb{E}[\varphi(X_{n+1}) \mid X_n = x_i] = \mathbb{E}_{x_i}[\varphi(X_1)], \quad i = 1, \dots, d,$$

the last equality being justified thanks to the Markov property. Clearly we have $g_i = \sum_{j=1}^d \varphi(x_j) \mathbb{P}(X_{n+1} = x_j \mid X_n = x_i) = \sum_j \varphi_j P_{ij}$ which, in matrix notation gives

$$g = P\varphi.$$

Hence, the transition matrix P can also be interpreted as an operator $P : \mathcal{F}(\mathcal{X}) \rightarrow \mathcal{F}(\mathcal{X})$ acting (to the right) on functions. Observe, in particular, that the constant unit function $\varphi = (1, \dots, 1) \in \mathcal{F}(\mathcal{X})$ satisfies

$$(P\varphi)_i = \sum_j 1 \cdot P_{ij} = 1 = \varphi_i,$$

since P is a stochastic matrix, and is therefore a *right eigenvector* of P corresponding to the eigenvalue $\lambda_1 = 1$. This argument shows that $\lambda_1 = 1$ is *always* an eigenvalue of P . The eigenvalue $\lambda_1 = 1$ turns out to be the largest in absolute value.

Lemma 8.2. *Given a stochastic matrix $P \in \mathbb{R}^{d \times d}$, let (λ, v) be a left eigenpair of P , i.e. $vP = \lambda v$, with $\|v\|_{\ell^1} = \sum_j |v_j| < \infty$. Then $|\lambda| \leq 1$.*

Proof. We have

$$|\lambda v_i| = \left| \sum_j v_j P_{ji} \right| \leq \sum_j |v_j| P_{ji}.$$

Hence

$$|\lambda| \sum_i |v_i| \leq \sum_i \sum_j |v_j| P_{ji} = \sum_j |v_j| \underbrace{\sum_i P_{ji}}_{=1} = \sum_j |v_j|$$

which implies $|\lambda| \leq 1$. □

It follows that an invariant distribution π is a left eigenvector of P corresponding to the largest (in absolute value) eigenvalue. The iterates $\pi^{n,\lambda} = \lambda P^n$ correspond to power iterations so we should expect $\pi^{n,\lambda}$ to converge to π as long as $\lambda_1 = 1$ is a simple eigenvalue and there are no other eigenvalues with absolute value 1.

In practice, in MCMC algorithms, we construct a Markov Chain so that the target distribution we want to sample from corresponds to an invariant distribution of the Markov Chain. (This also guarantees existence of an invariant distribution in the infinite dimensional case). However, it remains the question whether such invariant distribution is unique ($\lambda_1 = 1$ is simple) and whether the second largest eigenvalue $\beta = \max_{i=2,\dots,d} |\lambda_i(P)|$ in absolute value is strictly smaller than one as the *spectral gap* $1 - \beta$ will dictate the speed of convergence of $\pi^{n,\lambda}$ to π . We postpone this discussion to Section 8.1.2.

We now address the important case in which the transition matrix P features some symmetry properties. This will be indeed the case for the most popular MCMC algorithms, namely the Metropolis-Hastings ones. Let P be a transition matrix with invariant distribution π , and $\{X_n\}_{n=0}^N \sim \text{Markov}(\pi, P)$ a Markov chain at equilibrium. Let us look at the chain $\{Y_n = X_{N-n}, n = 0, \dots, N\}$, called the *time-reversal* of $\{X_n, n = 0, \dots, N\}$. It is not difficult to see that $\{Y_n\}_{n=0}^N$ is also a Markov chain. Indeed, assuming that

$\mathbb{P}(X_{N-n+1} = x_{i_{n-1}}, \dots, X_N = x_0) > 0$, we have for any $n = 1, \dots, N$

$$\begin{aligned} & \mathbb{P}(Y_n = x_{i_n} \mid Y_0 = x_{i_0}, \dots, Y_{n-1} = x_{i_{n-1}}) \\ &= \mathbb{P}(X_{N-n} = x_{i_n} \mid X_N = x_{i_0}, \dots, X_{N-n+1} = x_{i_{n-1}}) \\ &= \frac{\mathbb{P}(X_{N-n} = x_{i_n}, \dots, X_N = x_{i_0})}{\mathbb{P}(X_{N-n+1} = x_{i_{n-1}}, \dots, X_N = x_{i_0})} \\ &= \frac{P_{i_1 i_0} P_{i_2 i_1} \dots P_{i_n i_{n-1}} \mathbb{P}(X_{N-n} = x_{i_n})}{P_{i_1 i_0} P_{i_2 i_1} \dots P_{i_{n-1} i_{n-2}} \mathbb{P}(X_{N-n+1} = x_{i_{n-1}})} \\ &= P_{i_n i_{n-1}} \frac{\pi_{i_n}}{\pi_{i_{n-1}}} =: \hat{P}_{i_{n-1}, i_n}. \end{aligned}$$

Hence, the probability $\mathbb{P}(Y_n = x_{i_n} \mid Y_0 = x_{i_0}, \dots, Y_{n-1} = x_{i_{n-1}})$ of Y_n given the past depends only on i_{n-1} and $\{Y_n\}_{n=0}^N$ is a Markov chain $\{Y_n\}_{n=0}^N \sim \text{Markov}(\pi, \hat{P})$ with transition matrix

$$\hat{P}_{ij} = P_{ji} \frac{\pi_j}{\pi_i}.$$

Definition 8.3. Let P be a stochastic matrix, π a distribution on \mathcal{X} and $\{X_n\} \sim \text{Markov}(\pi, P)$ a Markov chain. We say that $\{X_n\}_{n \geq 0}$ is reversible if for all $N \geq 1$, $\{X_{N-n}\}_{n=0}^N \sim \text{Markov}(\pi, P)$. (Equivalently, for any $N \geq 0$, the joint distributions of (X_0, \dots, X_N) and (X_N, \dots, X_0) are the same.)

The definition of reversibility implies, in particular, that $X_N \sim \pi$ for all $N \geq 0$, hence π has to be an invariant distribution and a necessary condition for reversibility is that the chain is at equilibrium.

Definition 8.4. A stochastic matrix P and a probability distribution λ on \mathcal{X} are said to be in detailed balance if $\lambda_i P_{ij} = \lambda_j P_{ji}$ for all i, j .

The following Lemma establishes the relation between the detailed balance condition and the reversibility of the chain.

Lemma 8.3. Let P be a stochastic matrix and π a distribution on \mathcal{X} . (P, π) are in detailed balance if and only if π is invariant for P and $\{X_n\} \sim \text{Markov}(\pi, P)$ is reversible.

Proof. Suppose first that (P, π) are in detailed balance. Then, from direct calculation

$$(\pi P)_i = \sum_j \pi_j P_{ji} = \sum_j \pi_i P_{ij} = \pi_i \sum_j P_{ij} = \pi_i.$$

Hence π is an invariant distribution. Moreover, the detailed balance condition directly implies $\hat{P} = P$, hence the chain $\{X_n\} \sim \text{Markov}(\pi, P)$ is reversible.

The opposite implication is immediate: if π is invariant for P and $\{X_n\} \sim \text{Markov}(\pi, P)$ is reversible, by definition $\hat{P} = P$ which is equivalent to the detailed balance condition. \square

The detailed balance is a useful condition to verify that a certain distribution π is invariant (often easier than verifying $\pi P = \pi$). Intuitively, it says that under π , the probability of going from i to j is the same as the probability of going from j to i . Another

way to interpret the detailed balance equation is the following. Let us define the Hilbert space $\ell_\pi^2 = \{\varphi : \mathcal{X} \rightarrow \mathbb{R} : \sum_i \varphi_i^2 \pi_i < +\infty\}$ with inner product $(\varphi, \psi)_\pi = \sum_i \varphi_i \psi_i \pi_i$, where π is in detailed balance with P . Then, the matrix P is symmetric with respect to such an inner product (the corresponding operator $P : \ell_\pi^2 \rightarrow \ell_\pi^2$ is self adjoint). Indeed,

$$(P\varphi, \psi)_\pi = \sum_i \pi_i (P\varphi)_i \psi_i = \sum_{i,j} \pi_i P_{ij} \varphi_j \psi_i = \sum_{i,j} \pi_j P_{ji} \psi_i \varphi_j = (\varphi, P\psi)_\pi.$$

Hence, if (P, π) are in detailed balance, all eigenvalues of P are real and, at least in the finite dimensional case, the matrix is diagonalizable by an ℓ_π^2 -orthonormal set of eigenvectors.

8.1.1 Metropolis-Hastings algorithm in discrete state spaces

We come back to the original goal of constructing a Markov chain $\{X_n\} \sim \text{Markov}(\lambda, P)$ on \mathcal{X} which has a given invariant distribution π . We assume $\pi_i > 0$ for all i . (If $\pi_j = 0$ for some j , we can just remove the corresponding state x_j from the state space.) The Metropolis-Hastings is probably the most popular algorithm used for this purpose. It constructs a transition matrix P which is in detailed balance with the target distribution π . The idea is the following:

- Take a stochastic matrix Q with the condition that $Q_{ij} = 0 \iff Q_{ji} = 0$. Q is called the *proposal*. In general, Q will not have π as invariant distribution so we have to “correct” it.
- For any $i, j \in \{1, \dots, d\}$, define the acceptance probability

$$\alpha(i, j) = \min \left\{ 1, \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}} \right\} \quad \text{if } Q_{ij} \neq 0, \quad \alpha(i, j) = 0, \quad \text{if } Q_{ij} = 0.$$

The Metropolis-Hastings algorithm then reads:

Algorithm 8.1: Metropolis-Hastings

Given: λ (initial distribution), Q (proposal), π (target distribution)

```

1 Generate  $X_0 \sim \lambda$  for  $n = 0, 1, \dots$ , do
2   | Generate candidate new state  $\tilde{X}_{n+1} \sim Q_{X_n, \cdot}$ 
3   | Generate  $U \sim \mathcal{U}([0, 1])$ 
4   | if  $U \leq \alpha(X_n, \tilde{X}_{n+1})$  then
5     | set  $X_{n+1} = \tilde{X}_{n+1}$  //  $\tilde{X}_n$  accepted with prob.  $\alpha(X_n, \tilde{X}_{n+1})$ 
6     | else
7     | set  $X_{n+1} = X_n$  //  $\tilde{X}_n$  rejected with prob.  $1 - \alpha(X_n, \tilde{X}_{n+1})$ 
8     | end
9 end

```

If Q is symmetric, then the acceptance probability simplifies to $\alpha(i, j) = \min \left\{ 1, \frac{\pi_j}{\pi_i} \right\}$. In this case, step 4 of the algorithm will always accept \tilde{X}_{n+1} if the probability mass of

the new state $\pi_{\tilde{X}_{n+1}}$ is higher than the probability mass of the old state π_{X_n} . In case where $\pi_{\tilde{X}_{n+1}} < \pi_{X_n}$, the new state is accepted only with probability $\pi_{\tilde{X}_{n+1}}/\pi_{X_n}$. Hence, if $\pi_{\tilde{X}_{n+1}} \ll \pi_{X_n}$, the new state has a high probability of being rejected. Notice that in Algorithm 8.1, only the ratio $\pi_{\tilde{X}_{n+1}}/\pi_{X_n}$ appears. Therefore, the algorithm is applicable also in the case of a un-normalized target distribution.

We may ask what is the transition matrix associated to the Markov chain $\{X_n\}_n$ generated by 8.1. The following Lemma answers the question.

Lemma 8.4. *Let $\alpha_j^* = \sum_j \alpha(i, j)Q_{ij}$. Then, the transition matrix of the chain produced by the Metropolis-Hastings algorithm is given by*

$$P_{ij} = \alpha(i, j)Q_{ij} + (1 - \alpha_j^*)\delta_{ij}. \quad (8.2)$$

Proof. For $j \neq i$, we have

$$\begin{aligned} P_{ij} &= \mathbb{P}(X_{n+1} = j \mid X_n = i) = \mathbb{P}(\tilde{X}_{n+1} = j, X_{n+1} = \tilde{X}_{n+1} \mid X_n = i) \\ &= \mathbb{P}(X_{n+1} = \tilde{X}_{n+1} \mid \tilde{X}_{n+1} = j, X_n = i) \mathbb{P}(\tilde{X}_{n+1} = j \mid X_n = i) \\ &= \alpha(i, j)Q_{ij}. \end{aligned}$$

On the other hand, if $j = i$,

$$\begin{aligned} P_{ii} &= \mathbb{P}(X_{n+1} = i \mid X_n = i) \\ &= \mathbb{P}(\tilde{X}_{n+1} = i, X_{n+1} = \tilde{X}_{n+1} \mid X_n = i) + \mathbb{P}(X_{n+1} \neq \tilde{X}_{n+1} \mid X_n = i) \\ &= \alpha(i, i)Q_{ii} + \sum_j \mathbb{P}(\tilde{X}_{n+1} = j, X_{n+1} \neq \tilde{X}_{n+1} \mid X_n = i) \\ &= \alpha(i, i)Q_{ii} + \sum_j (1 - \alpha(i, j))Q_{ij} \\ &= \alpha(i, i)Q_{ii} + (1 - \alpha_i^*). \end{aligned}$$

□

The quantity $\alpha_i^* = \sum_j \alpha(i, j)Q_{ij}$ represents the overall probability of accepting a new state when being in state i . If such acceptance probability is very close to 0, with high probability the chain will not move, hence the random variables $\{X_n\}_n$ will be highly correlated, which is not desirable for constructing the ergodic estimator (8.1). A very high acceptance probability might not be desirable either. Consider the two possible strategies: a) jump only to neighboring states with high acceptance rate; b) jump to far away states but with lower acceptance rate. It is not obvious which strategy is more effective in decorrelating (mixing) the chain. Rule of thumb says that the average acceptance rate should be around 0.2.

That Algorithm 8.1 produces the right chain, i.e. a chain that has invariant distribution π , is shown in the following Lemma and is a consequence of the fact that the transition matrix P in (8.2) is in detailed balance with π .

Lemma 8.5. *The transition matrix P in (8.2) is in detailed balance with π . Hence, the chain produced by Algorithm 8.1 is reversible and has π as invariant distribution.*

Proof. We have to show that $\pi_i P_{ij} = \pi_j P_{ji}$ for all i, j . This is obviously true for $i = j$. Consider then $i \neq j$. If $\pi_i P_{ij} = 0$, then $P_{ij} = 0$ which implies $Q_{ij} = Q_{ji} = 0$ so $P_{ji} = 0$ and $\pi_i P_{ij} = \pi_j P_{ji}$. If $\pi_i P_{ij} \neq 0$, then $P_{ij} \neq 0$ so $Q_{ij}, Q_{ji} \neq 0$ and

$$\begin{aligned} \pi_i P_{ij} &= \pi_i \alpha(i, j) Q_{ij} = \pi_i Q_{ij} \min \left\{ 1, \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}} \right\} \\ &= \min \{ \pi_i Q_{ij}, \pi_j Q_{ji} \} \\ &= \min \left\{ \frac{\pi_i Q_{ij}}{\pi_j Q_{ji}}, 1 \right\} \pi_j Q_{ji} = \pi_j \alpha(j, i) Q_{ji} = \pi_j P_{ji}. \end{aligned}$$

□

8.1.2 Convergence results

Let $\{X_n\} \sim \text{Markov}(\lambda, P)$ be a Markov chain with invariant distribution π . We want to understand under which conditions on (λ, P) , π is the *unique* invariant distribution and the sequence $\pi^{n, \lambda}$ converges to π as $n \rightarrow \infty$. Three concepts are key to answer this question: *irreducibility*, *reversibility* and *aperiodicity*.

Definition 8.5 (Irreducible chain). *Let P be a transition matrix on \mathcal{X} .*

- *We say that a state $x_i \in \mathcal{X}$ communicates with another state $x_j \in \mathcal{X}$ if $\mathbb{P}_i(X_n = x_j \text{ for some } n) > 0$. Equivalently, there exists $n > 0$: $P_{ij}^{(n)} > 0$.*
- *The transition matrix P is said to be irreducible if every state communicates with every other state, i.e. for all i, j , there exists $n > 0$ such that $P_{ij}^{(n)} > 0$. A Markov chain $\text{Markov}(\lambda, P)$ is irreducible if P is so.*

An equivalent definition of irreducibility is the following:

Definition 8.6 (Irreducible chain). *Let P be a transition matrix on \mathcal{X} .*

- *A state x_j is said to be accessible if $\mathbb{P}_i(\sigma_j < \infty) > 0$ for any $x_i \in \mathcal{X}$.*
- *P is irreducible if every state is accessible.*

In the definition of accessible state, σ_j is the return time to the state x_j . It is easy to see that the two definitions are equivalent. Figure 8.1 shows an example of an irreducible chain (left) and a reducible one (right).

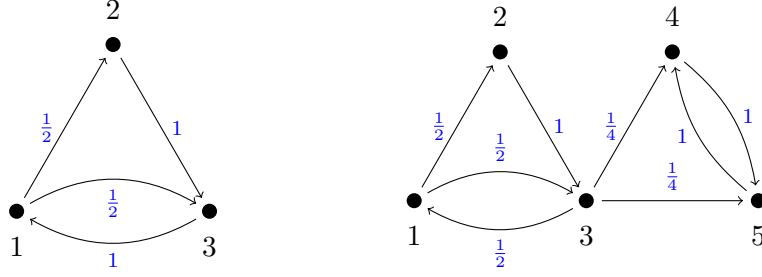


Figure 8.1: Left: irreducible chain – every state communicates with every other state. Right: reducible chain – $\{4, 5\}$ is an absorbing class and does not communicate with $\{1, 2, 3\}$.

We now turn to the notion of recurrence. Given a state $x_i \in \mathcal{X}$ we denote by $V_i = \sum_{n=0}^{\infty} \mathbb{1}_{\{X_n=x_i\}}$ the number of visits to x_i . Notice that

$$\mathbb{E}_i[V_i] = \mathbb{E}_i \sum_{n=0}^{\infty} \mathbb{1}_{\{X_n=x_i\}} = \sum_{n=0}^{\infty} \mathbb{P}_i(X_n = x_i) = \sum_{n=0}^{\infty} P_{ii}^{(n)}.$$

Definition 8.7 (Recurrent state). *A state $x_i \in \mathcal{X}$ is said to be recurrent if $\mathbb{P}_i(X_n = x_i \text{ infinitely often}) = 1$, or equivalently $\mathbb{P}_i(V_i = \infty) = 1$. It is transient if $\mathbb{P}_i(V_i = \infty) = 0$.*

An interesting fact is that a state x_i is either recurrent or transient, i.e. it can never happen that $\mathbb{P}_i(V_i = \infty) \in (0, 1)$.

Lemma 8.6. *A given state x_i is either recurrent or transient. Moreover,*

- x_i is recurrent $\iff \mathbb{P}_i(\sigma_i < \infty) = 1 \iff \mathbb{E}_i[V_i] = \infty$;
- x_i is transient $\iff \mathbb{P}_i(\sigma_i < \infty) < 1 \iff \mathbb{E}_i[V_i] < \infty$.

Proof. Let $\sigma_i^{(r)}$ be the r -return time to the state x_i , i.e. $\sigma_i^{(r)} = \inf\{n > \sigma_i^{(r-1)} : X_n = x_i\}$, with $\sigma_i^{(1)} = \sigma_i$. Then

$$\begin{aligned} \mathbb{P}_i(V_i > r + 1) &= \mathbb{P}_i(\sigma_i^{(r+1)} < \infty) \\ &= \mathbb{P}_i(\sigma_i^{(r)} < \infty, \sigma_i^{(r+1)} - \sigma_i^{(r)} < +\infty) \\ &= \mathbb{P}_i(\sigma_i^{(r+1)} - \sigma_i^{(r)} < \infty \mid \sigma_i^{(r)} < \infty) \mathbb{P}_i(\sigma_i^{(r)} < \infty) \\ &= \mathbb{P}_i(\sigma_i^{(1)} < \infty) \mathbb{P}_i(\sigma_i^{(r)} < \infty) \text{ (by the strong Markov property)} \\ &= \mathbb{P}_i(\sigma_i < \infty)^{r+1}. \end{aligned}$$

Hence

$$\mathbb{P}_i(V_i = \infty) = \lim_{r \rightarrow \infty} \mathbb{P}_i(V_i > r) = \begin{cases} 0, & \iff \mathbb{P}_i(\sigma_i < \infty) < 1, \\ 1, & \iff \mathbb{P}_i(\sigma_i < \infty) = 1. \end{cases}$$

Moreover

$$\mathbb{E}_i[V_i] = \sum_r \mathbb{P}_i(V_i > r) = \begin{cases} C < \infty, & \iff \mathbb{P}_i(\sigma_i < \infty) < 1, \\ \infty, & \iff \mathbb{P}_i(\sigma_i < \infty) = 1. \end{cases}$$

□

If all the states communicate with each other, i.e. the chain is irreducible, it is easy to see that if a chain has a recurrent state, all the states are recurrent.

Lemma 8.7. *Let $\{X_n\} \sim \text{Markov}(\lambda, P)$ with an irreducible transition matrix P . Then either all states are transient or recurrent.*

Proof. Suppose x_i is transient and take $x_j \neq x_i$. Since P is irreducible, there exist $n, m > 0$: $P_{ij}^{(n)} > 0$ and $P_{ji}^{(m)} > 0$. Then for all $r \geq 0$,

$$P_{ii}^{(n+m+r)} = \sum_{\ell, k} P_{li}^{(m)} P_{k\ell}^{(r)} P_{ik}^{(n)} \geq P_{ji}^{(m)} P_{jj}^{(r)} P_{ij}^{(n)}.$$

On the other hand, being x_i transient, we have $\mathbb{E}_i[V_i] < \infty$ and

$$\mathbb{E}_j[V_j] = \sum_{r=0}^{\infty} P_{jj}^{(r)} \leq \frac{1}{P_{ji}^{(m)} P_{ij}^{(n)}} \sum_{r=0}^{\infty} P_{ii}^{(m+n+r)} \leq \frac{1}{P_{ji}^{(m)} P_{ij}^{(n)}} \mathbb{E}_i[V_i] < \infty,$$

hence, from the previous lemma, x_j is also transient. \square

The previous result justifies the following

Definition 8.8. *An irreducible Markov chain $\{X_n\} \sim \text{Markov}(\lambda, P)$ is said to be recurrent if it has at least one recurrent state (equivalently if every state is recurrent).*

Example 8.3. *Consider a random walk on \mathbb{Z} :*

$$\mathbb{P}(X_{n+1} = i + 1 \mid X_n = i) = p, \quad \mathbb{P}(X_{n+1} = i - 1 \mid X_n = i) = q = 1 - p.$$

If we start at $X_0 = 0$, we can return to zero only after an even number of steps, say $2n$, with n moves to the right and n to the left. Hence

$$P_{00}^{(2n)} = \binom{2n}{n} p^n q^n = \frac{(2n)!}{(n!)^2} (pq)^n \sim \frac{(4pq)^n}{\sqrt{2\pi n}}$$

where we have used Stirling's formula $n! \sim \sqrt{2\pi n}(n/e)^n$. Hence

$$\mathbb{E}_0[V_0] = \sum_{n=0}^{\infty} P_{00}^{(n)} \begin{cases} = \infty, & \text{for } p = q = \frac{1}{2}, \\ < \infty, & \text{for } p \neq q. \end{cases}$$

We conclude that $\{X_n\}$ is recurrent if $p = \frac{1}{2}$ and transient otherwise. By similar calculations, one can show that a symmetric random walk on \mathbb{Z}^2 with

$$\mathbb{P}(X_{n+1} = (i \pm 1, j) \mid X_n = (i, j)) = \mathbb{P}(X_{n+1} = (i, j \pm 1) \mid X_n = (i, j)) = \frac{1}{4}$$

is also recurrent, whereas a symmetric random walk on \mathbb{Z}^3 is transient.

For an irreducible recurrent Markov chain $\{X_n, n \geq 0\}$ on an infinite countable state space \mathcal{X} , we can further distinguish two cases:

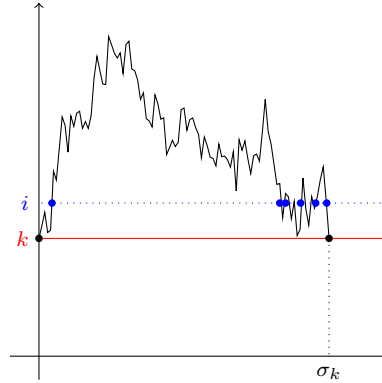


Figure 8.2: The invariant measure corresponds to the expected time spent by the chain in each state, between two consecutive visits of a fixed (recurrent) state x_k .

Definition 8.9. An irreducible Markov chain $\{X_n\} \sim \text{Markov}(\lambda, P)$ is said to be positive recurrent (or simply positive) if $\mathbb{E}_i[\sigma_i] < \infty$ for at least one state $x_i \in \mathcal{X}$ and null recurrent otherwise.

Again one can show that if there exists a state $x_i \in \mathcal{X}$ for which $\mathbb{E}_i[\sigma_i] < \infty$ and the chain is irreducible, then $\mathbb{E}_j[\sigma_j] < \infty$ for every state $x_j \in \mathcal{X}$. The property of recurrence / positive recurrence is key to obtain existence of an invariant measure as the next theorem shows.

Theorem 8.8. Let P be irreducible and recurrent. Then P has a unique invariant measure (not necessarily finite) up to a multiplicative constant.

The invariant measure can be constructed in the following way. Let us fix a (recurrent) state $x_k \in \mathcal{X}$ and consider $\tilde{\pi}_i^k = \mathbb{E}_k \left[\sum_{n=0}^{\sigma_k-1} \mathbb{1}_{\{X_n=x_i\}} \right]$ which corresponds to the expected number of visits to x_i between two consecutive visits to x_k . Then $\tilde{\pi}^k$ is an invariant measure and is unique up to a multiplicative factor. Notice that $\tilde{\pi}_k^k = 1$.

Proof of existence. Observe first that for all $i \neq k$,

$$\tilde{\pi}_i^k = \mathbb{E}_k \left[\sum_{n=0}^{\sigma_k-1} \mathbb{1}_{\{X_n=x_i\}} \right] = \mathbb{E}_k \left[\sum_{n=1}^{\sigma_k} \mathbb{1}_{\{X_n=x_i\}} \right] = \mathbb{E}_k \left[\sum_{n=0}^{\sigma_k-1} \mathbb{1}_{\{X_{n+1}=x_i\}} \right]$$

i.e. $\tilde{\pi}_i^k$ is invariant by a +1 right shift of the chain, which basically shows that $\tilde{\pi}_i^k$ is an

invariant measure. More precisely, $\tilde{\pi}_i^k$ can be equivalently written as

$$\begin{aligned}
\tilde{\pi}_i^k &= \mathbb{E}_k \left[\sum_{n=0}^{\infty} \mathbb{1}_{\{X_{n+1}=x_i, \sigma_k > n\}} \right] = \sum_{n=0}^{\infty} \mathbb{P}_k(X_{n+1} = x_i, \sigma_k > n) \\
&= \sum_j \sum_{n=0}^{\infty} \mathbb{P}_k(X_{n+1} = x_i, \sigma_k > n, X_n = x_j) \\
&= \sum_j \sum_{n=0}^{\infty} \underbrace{\mathbb{P}_k(X_{n+1} = x_i \mid X_n = x_j, \sigma_k > n)}_{=P_{ji} \text{ since } \{\sigma_k > n\} \text{ depends only on } X_0, \dots, X_n} \mathbb{P}_k(X_n = x_j, \sigma_k > n) \\
&= \sum_j P_{ji} \underbrace{\sum_{n=0}^{\infty} \mathbb{P}_k(X_n = x_j, \sigma_k > n)}_{\tilde{\pi}_j^k} = \sum_j \tilde{\pi}_j^k P_{ji}.
\end{aligned}$$

□

Proof of uniqueness. Let λ be another invariant measure. Then for $j \neq k$,

$$\begin{aligned}
\lambda_j &= \sum_{i_1} \lambda_{i_1} P_{i_1 j} \\
&= \lambda_k P_{kj} + \sum_{i_1 \neq k} \lambda_{i_1} P_{i_1 j} \\
&= \lambda_k \underbrace{P_{kj}}_{\mathbb{P}_k(X_1=x_j, \sigma_k > 0)} + \lambda_k \underbrace{\sum_{i_1 \neq k} P_{ki_1} P_{i_1 j}}_{\mathbb{P}_k(X_2=x_j, \sigma_k > 2)} + \sum_{i_1, i_2 \neq k} \lambda_{i_2} P_{i_2 i_1} P_{i_1 j} = \dots \\
&\geq \lambda_k \sum_{n=0}^{\infty} \mathbb{P}_k(X_n = x_j, \sigma_k > n) = \lambda_k \tilde{\pi}_j^k.
\end{aligned}$$

Moreover, since P is irreducible, there exists $n > 0$ such that $P_{jk}^{(n)} > 0$. Hence

$$0 = \frac{\lambda_k}{\lambda_k} - \tilde{\pi}_k^k = \sum_i \left(\frac{\lambda_i}{\lambda_k} - \tilde{\pi}_i^k \right) P_{ik}^{(n)} \geq \left(\frac{\lambda_j}{\lambda_k} - \tilde{\pi}_j^k \right) P_{jk}^{(n)} \implies \lambda_j \leq \lambda_k \tilde{\pi}_j^k.$$

It follows that $\lambda_j = \lambda_k \tilde{\pi}_j^k$ for all j , therefore $\lambda \propto \tilde{\pi}^k$.

□

The measure $\tilde{\pi}^k$ is not necessarily finite. Indeed

$$\sum_i \tilde{\pi}_i^k = \sum_i \sum_{n=0}^{\infty} \mathbb{P}_k(X_n = x_i, \sigma_k > n) = \sum_{n=0}^{\infty} \mathbb{P}_k(\sigma_k > n) = \mathbb{E}_k[\sigma_k]$$

hence we see that $\tilde{\pi}^k$ is finite if and only if P is positive recurrent.

Theorem 8.9. *Let P be irreducible, then P has an invariant distribution (invariant probability measure) π if and only if P is positive recurrent. Moreover, in this case, π is unique and is given by*

$$\pi_i = \frac{\tilde{\pi}_i^k}{\mathbb{E}_k[\sigma_k]} = \frac{1}{\mathbb{E}_i[\sigma_i]}.$$

(The last equality follows by simply taking $k = i$.)

Consider now an irreducible and positive recurrent Markov chain $\{X_n\}_n$ with invariant distribution π and an integrable function $\psi : \mathcal{X} \rightarrow \mathbb{R}$ with respect to π , i.e. $\mathbb{E}_\pi[\psi] < \infty$. We ask whether $\mu = \mathbb{E}_\pi[\psi]$ (ensemble average of ψ) can be computed as a *temporal* average over only one realization of the chain. The following result holds.

Theorem 8.10 (Ergodic theorem). *Let $\{X_n\} \sim \text{Markov}(\lambda, P)$ with P irreducible and positive recurrent, with invariant distribution π . Then, for any function $\psi : \mathcal{X} \rightarrow \mathbb{R}$, such that $\mathbb{E}_\pi[|\psi|] < \infty$, it holds*

$$\mathbb{P}_\lambda \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \psi(X_j) = \mathbb{E}_\pi[\psi] \right) = 1$$

for any $\lambda \in \mathcal{M}_1(\mathcal{X})$.

Idea of the proof. Let $\sigma_k^{(r)}$ be the r -return time to the state x_k , with $\sigma_k^{(0)} = 0$ and $V_k(n) = \sum_{j=1}^{n-1} \mathbb{1}_{\{X_j=x_k\}}$ be the number of visits to x_k before time n . Set $Y_r = \sum_{j=\sigma_k^{(r-1)}+1}^{\sigma_k^{(r)}} \psi(X_j)$, $r = 1, \dots, V_k(n)$. By the strong Markov property, $Y_r \stackrel{\text{iid}}{\sim} Y_2$ for all $r \geq 2$. Hence by the strong law of large numbers (SLLN)

$$\frac{1}{V_k(n) - 1} \sum_{r=2}^{V_k(n)} Y_r \xrightarrow{\text{a.s.}} \mathbb{E}[Y_r]$$

and

$$\frac{1}{n} \sum_{j=1}^n \psi(X_j) = \underbrace{\frac{Y_1}{n}}_{\rightarrow 0} + \underbrace{\frac{V_k(n) - 1}{n}}_{\rightarrow \pi_k} \underbrace{\frac{1}{V_k(n) - 1} \sum_{r=2}^{V_k(n)} Y_r}_{\rightarrow \mathbb{E}[Y_r]} + \underbrace{\frac{1}{n} \sum_{j=\sigma_k^{V_k(n)}+1}^n \psi(X_j)}_{\rightarrow 0} \xrightarrow{\text{a.s.}} \pi_k \mathbb{E}[Y_r].$$

Moreover, $\mathbb{E}[Y_2] = \sum_i \psi(x_i) \tilde{\pi}_i^k = \frac{1}{\pi_k} \mathbb{E}_\pi[\psi]$ hence the result. \square

The interval $[\sigma_k^{(r-1)} + 1, \sigma_k^{(r)}]$ is called a renewal cycle. The fact that the chain regenerates itself every time it visits a given state k is what allowed us to use the SLLN. Thanks to the renewal structure highlighted in the proof of the previous theorem, one can obtain also a Central Limit Theorem (CLT).

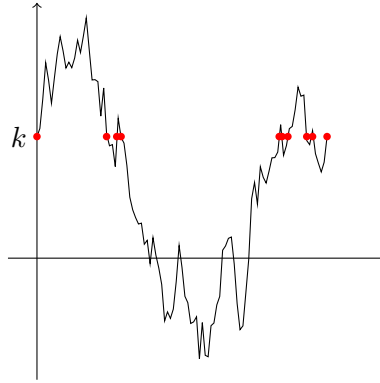


Figure 8.3: Renewal cycles for an irreducible and positive recurrent Markov chain.

Theorem 8.11 (CLT for Markov Chains). *Let $\{X_n\} \sim \text{Markov}(\lambda, P)$ with P irreducible, positive recurrent and with invariant distribution π . Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ be such that $\mathbb{E}_\pi[|\psi|] < \infty$ and $C(\psi) = \frac{1}{\pi_k} \mathbb{E}_k[(\sum_{j=1}^{\sigma_k} \psi(X_j) - \sigma_k \mathbb{E}_\pi[\psi])^2] < \infty$. Then*

$$\sqrt{n} \left(\frac{1}{n} \sum_{j=1}^n \psi(X_j) - \mathbb{E}_\pi[\psi] \right) \xrightarrow{d} N(0, C(\psi)).$$

We turn now to the stronger question whether the sequence of distributions $\pi^{n,\lambda}$ of the steps $\{x_n, n \in \mathbb{N}\}$ of a Markov chain $\text{Markov}(\lambda, P)$ converges, in a suitable sense, to the invariant distribution π as $n \rightarrow \infty$ for any choice of initial distribution $\lambda \in \mathcal{M}_1(\mathcal{X})$. The fact that this property is not always true is shown in the next example.

Example 8.4. *Consider the transition matrix $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ which is clearly irreducible and has an invariant distribution $\pi = (\frac{1}{2}, \frac{1}{2})$ (hence it is positive recurrent). However, if we take the initial distribution $\lambda = (1, 0)$, we have $\pi^{1,\lambda} = \lambda P = (0, 1)$, $\lambda^{2,\lambda} = \lambda^{1,\lambda} P = (1, 0)$ and clearly $\pi^{n,\lambda}$ does not converge to π . The problem in this example is that the chain visits periodically (with period 2) the two states.*

We need therefore to exclude such cases.

Definition 8.10. *Given a transition matrix P , we say that a state x_i is aperiodic if $P_{ii}^{(n)} > 0$ for all sufficiently large n , or equivalently if the set $\{n > 0 : P_{ii}^{(n)} > 0\}$ has no common divisor other than 1.*

Using the Chapman-Kolmogorov equation, it is easy to see that if P is irreducible and has an aperiodic state x_i , then all states $x_j \in \mathcal{X}$ are aperiodic. We will then say that P is aperiodic. The next theorem states that for an irreducible, positive recurrent, aperiodic Markov chain $\{X_n\}_n \sim \text{Markov}(\lambda, P)$, $\pi^{n,\lambda} \rightarrow \pi$ in total variation as $n \rightarrow \infty$ for any initial distribution λ .

Before stating the theorem, we recall the definition of *total variation* of a measure: given a measurable space $(\mathcal{X}, \mathcal{B})$, with \mathcal{B} a σ -algebra on \mathcal{X} , and a (signed) measure

$\mu : \mathcal{B} \rightarrow \mathbb{R}$ the total variation of μ is defined as

$$\|\mu\|_{\text{TV}} = \sup_{A \in \mathcal{B}} \mu(A) - \inf_{A \in \mathcal{B}} \mu(A) = \sup_{\substack{f: \mathcal{X} \rightarrow \mathbb{R} \text{ meas.} \\ \|f\|_{\infty} \leq 1}} \int_{\mathcal{X}} f(x) \mu(dx).$$

In the case of a discrete set \mathcal{X} the above definition reduces to the ℓ^1 -norm of the row vector $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots)$

$$\|\mu\|_{\text{TV}} = \sum_{x_i \in \mathcal{X}} |\mu(\{x_i\})| = \|\boldsymbol{\mu}\|_{\ell^1}.$$

Theorem 8.12. *Let P be irreducible, aperiodic and positive recurrent with invariant distribution π . Let λ be any distribution on \mathcal{X} and $\{X_n\} \sim \text{Markov}(\lambda, P)$. Then, for $\pi_i^{n,\lambda} = \mathbb{P}_{\lambda}(X_n = x_i)$ it holds*

$$\lim_{n \rightarrow \infty} \|\pi^{n,\lambda} - \pi\|_{\text{TV}} = \lim_{n \rightarrow \infty} \sum_i |\pi_i^{n,\lambda} - \pi_i| = 0.$$

Idea of the proof. Consider a chain $\{Y_n\} \sim \text{Markov}(\pi, P)$ at equilibrium and independent of $\{X_n\}$. The joint process $\{W_n = (X_n, Y_n)\}_n$ is also a Markov chain with transition matrix

$$\tilde{P}_{ik,j\ell} = \mathbb{P}(X_{n+1} = x_j, Y_{n+1} = x_{\ell} \mid X_n = x_i, Y_n = x_k) = P_{ij}P_{k\ell}$$

and invariant distribution $\tilde{\pi}_{ik} = \pi_i\pi_k$. Since P is aperiodic, for all i, j, k, ℓ , $\tilde{P}_{ik,j\ell}^{(n)} = P_{ij}^{(n)}P_{k\ell}^{(n)} > 0$ for sufficiently large n . Hence \tilde{P} is irreducible and positive recurrent (since it has an invariant distribution).

Consider now the return time $\sigma_k = \inf\{n > 0 : X_n = Y_n = k\}$ for which $\mathbb{P}(\sigma_k < +\infty) = 1$ since \tilde{P} is irreducible and recurrent. At time σ_k , the two chains meet. Hence for $n \geq \sigma_k$, we can follow the path $\{Y_n\}$ which is at equilibrium. More precisely, the process

$$Z_n = \begin{cases} X_n, & n < \sigma_k \\ Y_n, & n \geq \sigma_k \end{cases}$$

is also Markov (λ, P) , i.e. has the same distribution as $\{X_n\}$. Moreover

$$\begin{aligned} |\mathbb{P}(Z_n = x_j) - \pi_j| &= |\mathbb{P}(Z_n = x_j) - \mathbb{P}(Y_n = x_j)| \\ &= |\mathbb{P}(Z_n = x_j, n < \sigma_k) + \mathbb{P}(Z_n = x_j, n \geq \sigma_k) - \mathbb{P}(Y_n = x_j)| \\ &= |\mathbb{P}(X_n = x_j, n < \sigma_k) - \mathbb{P}(Y_n = x_j, n < \sigma_k)| \\ &\leq \mathbb{P}(X_n = x_j, n < \sigma_k) + \mathbb{P}(Y_n = x_j, n < \sigma_k). \end{aligned}$$

Hence

$$\begin{aligned} \sum_j |\mathbb{P}(Z_n = x_j) - \pi_j| &\leq \sum_j \mathbb{P}(X_n = x_j, n < \sigma_k) + \sum_j \mathbb{P}(Y_n = x_j, n < \sigma_k) \\ &\leq 2\mathbb{P}(n < \sigma_k) \rightarrow 0 \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Therefore,

$$\sum_j |\pi_j^{\lambda,n} - \pi_j| = \sum_j |\mathbb{P}(Z_n = x_j) - \pi_j| \rightarrow 0.$$

□

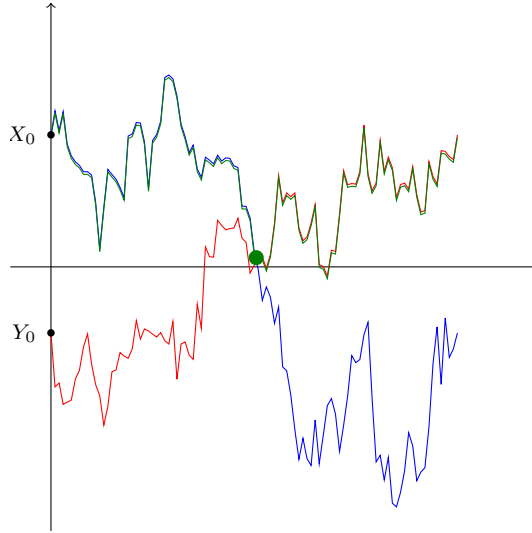


Figure 8.4: Paths to follow to equilibrium.

Concerning the rate of convergence, we introduce the following definitions

Definition 8.11. An irreducible, positive recurrent, aperiodic Markov chain $\{X_n\}_n$ with transition matrix P and invariant distribution π is

- geometrically ergodic if there exists a function $h : \mathcal{X} \rightarrow \mathbb{R}$, with $\mathbb{E}_\pi[h] < +\infty$ and $r \in (0, 1)$ such that

$$\|\pi^{n, \delta_{x_i}} - \pi\|_{TV} \leq h(x_i)r^n \quad \text{for all } x_i \in \mathcal{X},$$

- uniformly ergodic if there exists $C > 0$ and $r \in (0, 1)$ such that

$$\|\pi^{n, \delta_{x_i}} - \pi\|_{TV} \leq Cr^n \quad \text{for all } x_i \in \mathcal{X}.$$

Establishing geometric/uniform ergodicity is in general not easy, but can be done in special cases, exploiting the structure of the transition matrix P . One such special case is that of a *finite* state space \mathcal{X} . We recall here some properties. Let \mathcal{X} be a finite set of cardinality dimension d and P an irreducible, aperiodic transition matrix. Then

- P is recurrent and positive recurrent (exercise)
- P has an eigenvalue $\lambda_1 = 1$ simple (Perrou-Frobenius theorem) and all other eigenvalues satisfy $|\lambda_i| < 1$, $i = 2, \dots, d$.
- A Markov chain $\{X_n\}_n$ with transition matrix P is always uniformly ergodic and $\|\pi^{n, \delta_{x_i}} - \pi\|_{TV} \leq C|\lambda_2|^n$ with $|\lambda_2| = \max_{|\lambda_i| < 1} |\lambda_i|$ if P is diagonalizable. If P is not diagonalizable, the estimate has to be modified as $\|\pi^{n, \delta_{x_i}} - \pi\|_{TV} \leq C(\epsilon)(|\lambda_2| + \epsilon)^n$ for $\epsilon > 0$ arbitrary.

8.2 Markov chains on a general state space

We give here a brief overview of how the theory of Markov chains generalizes to a continuous state space \mathcal{X} , typically a subset of \mathbb{R}^d with non zero Lebesgue measure.

Definition 8.12. A Markov transition kernel on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, where $\mathcal{B}(\mathcal{X})$ is the Borel σ -algebra on \mathcal{X} , is a function $P : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ s.t.

1. for all $x \in \mathcal{X}$, $P(x, \cdot)$ is a probability measure on \mathcal{X} ,
2. for all $A \in \mathcal{B}(\mathcal{X})$, $P(\cdot, A)$ is measurable.

Whenever $P(x, \cdot)$ admits a density with respect to the Lebesgue measure, we denote it by $p : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ i.e. for all $x \in \mathcal{X}$, $A \in \mathcal{B}(\mathcal{X})$,

$$P(x, A) = \int_A p(x, y) dy.$$

Definition 8.13. Given a Markov transition kernel P and a measure λ on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, a sequence of random variables $\{X_n \in \mathcal{X}, n \geq 0\}$ is a homogeneous Markov chain with transition Kernel P and initial distribution λ , in short $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ if

- $X_0 \sim \lambda$
- $\mathbb{P}(X_{n+1} \in A \mid \mathcal{F}_n) = \mathbb{P}(X_{n+1} \in A \mid X_n) = P(X_n, A)$

where $\mathcal{F}_n = \sigma(X_0, \dots, X_n)$ is the σ -algebra generated by X_0, \dots, X_n . A Markov chain $\{X_n\} \sim \text{Markov}(\lambda, P)$ satisfies the strong Markov property. Let τ be a stopping time and $\mathcal{F}_\tau = \sigma(X_0, \dots, X_\tau)$ be the σ -algebra generated by X_0, \dots, X_τ . Then it holds

$$\mathbb{E}_\lambda[h(X_{\tau+1}, X_{\tau+2}, \dots) \mathbf{1}_{\{\tau < \infty\}} \mid \mathcal{F}_\tau] = \mathbb{E}_{X_\tau}[h(X_1, X_2, \dots)]$$

for any bounded and measurable function $h : \mathcal{X}^{\mathbb{N}} \rightarrow \mathbb{R}$.

The n -step transition kernel $P^{(n)}(x, A) = \mathbb{P}(X_n \in A \mid X_0 = x)$ is given by the recursion

$$P^{(n)}(x, A) = \int_{\mathcal{X}} P^{(n-1)}(y, A) P(x, dy), \quad P^{(1)}(x, A) = P(x, A).$$

Similarly, if $p^{(n)} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ denotes the density of $P^{(n)}$ (provided it exists), then

$$p^{(n)}(x, y) = \int_{\mathcal{X}} p^{(n-1)}(z, y) p(x, z) dz, \quad p^{(1)}(x, y) = p(x, y).$$

To each Markov transition kernel P we can associate the Markov operator \mathcal{P} acting to the left on measures, $\mathcal{P} : \mathcal{M}_1(\mathcal{X}) \rightarrow \mathcal{M}_1(\mathcal{X})$, with $\mathcal{M}_1(\mathcal{X})$ the set of probability measures on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ as

$$\mu = \lambda \mathcal{P} \quad \implies \quad \mu(A) = \int_{\mathcal{X}} P(x, A) \lambda(dx), \quad \forall A \in \mathcal{B}(\mathcal{X}).$$

Notice that

$$\begin{aligned}\lambda\mathcal{P}^2(A) &= (\lambda\mathcal{P})\mathcal{P}(A) = \int_{\mathcal{X}} P(y, A)(\lambda\mathcal{P})(dy) \\ &= \int_{\mathcal{X}} \int_{\mathcal{X}} P(y, A)P(x, dy)\lambda(dx) = \int_{\mathcal{X}} P^{(2)}(x, A)\lambda(dx)\end{aligned}$$

so \mathcal{P}^2 is the operator associated to $P^{(2)}$ and more generally \mathcal{P}^n is the operator associated to $P^{(n)}$. If $\pi^{n,\lambda}$ denotes the measure of X_n , i.e. $\pi^{n,\lambda}(A) = \mathbb{P}_\lambda(X_n \in A)$, it follows that $\pi^{n,\lambda} = \lambda\mathcal{P}^n = \int_{\mathcal{X}} P^{(n)}(y, \cdot)\lambda(dy)$.

Definition 8.14. A measure π on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ is called invariant (or stationary) if $\pi = \pi\mathcal{P} = \int_{\mathcal{X}} P(y, \cdot)\pi(dy)$. If the measure π has a density $f : \mathcal{X} \rightarrow \mathbb{R}_+$ (i.e. $\pi(A) = \int_A f(y) dy$, $\forall A \in \mathcal{B}(\mathcal{X})$), and the kernel P has a density p , then $f(x) = \int_{\mathcal{X}} p(y, x)f(y) dy$.

Similarly, a Markov transition kernel P defines an operator acting on functions to the right, $\mathcal{P} : \mathcal{F}(\mathcal{X}) \rightarrow \mathcal{F}(\mathcal{X})$, where $\mathcal{F}(\mathcal{X})$ is the set of measurable functions on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, as

$$g = \mathcal{P}\varphi \quad \implies \quad g(x) = \int_{\mathcal{X}} P(x, dy)\varphi(y) = \mathbb{E}_x[\varphi(X_1)].$$

Definition 8.15. A chain $\{X_n\}_{n=0}^N \sim \text{Markov}(\pi, P)$ is reversible if, for any $N > 0$, the chain $\{Y_n = X_{N-n}\}_{n=0}^N \sim \text{Markov}(\pi, P)$.

Reversibility implies, in particular, that π is an invariant distribution for \mathcal{P} . As for discrete state spaces, $\{X_n\}_n \sim \text{Markov}(\pi, P)$ is reversible if and only if (P, π) satisfy the detailed balance condition, which in this case reads

$$\int_A P(x, B)\pi(dx) = \int_B P(y, A)\pi(dy), \quad \forall A, B \in \mathcal{B}(\mathcal{X}), \quad \pi(A), \pi(B) > 0,$$

or, equivalently (whenever \mathcal{X} is a metric space)

$$\int_{\mathcal{X} \times \mathcal{X}} g(x, y)P(x, dy)\pi(dx) = \int_{\mathcal{X} \times \mathcal{X}} g(x, y)P(y, dx)\pi(dy), \quad \forall g \in \mathcal{F}_b(\mathcal{X} \times \mathcal{X})$$

where \mathcal{F}_b denotes the space of bounded measurable functions. This condition is often written shortly as

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

If (P, π) are in detailed balance, then π is an invariant distribution for \mathcal{P} . Indeed,

$$(\pi\mathcal{P})(B) = \int_{\mathcal{X}} P(x, B)\pi(dx) = \int_B \underbrace{P(y, \mathcal{X})}_{=1} \pi(dy) = \pi(B), \quad \forall B \in \mathcal{B}(\mathcal{X}).$$

We now extend the concepts of irreducibility, recurrence and aperiodicity. In the discrete setting, we have said that x_i communicates with x_j if there exists $n > 0 : P_{ij}^{(n)} > 0$ and a chain is irreducible if any state communicates with every other state. In the general state space case, the definition is slightly more cumbersome. Indeed, if we work with continuous random variables and assume that the transition kernel P has a density, then $P^{(n)}(x, \{y\}) = \mathbb{P}_x(X_n = y) = 0$ for all n since the set $\{y\}$ is of zero measure.

Definition 8.16 (irreducibility). *We say that $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ is irreducible if there exists a (σ -finite) measure φ on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ such that for any $x \in \mathcal{X}$ and $A \in \mathcal{B}(\mathcal{X})$, with $\varphi(A) > 0$, there exists $n > 0$ for which $P^{(n)}(x, A) > 0$. In this case, φ is called an irreducibility measure.*

Recall that a set $A \in \mathcal{B}(\mathcal{X})$ is accessible if $\mathbb{P}_x(\sigma_A < \infty) > 0$ for all $x \in \mathcal{X}$, where $\sigma_A = \inf\{n > 0 : X_n \in A\}$ is the return time to the set A . The above definition of irreducibility, with irreducibility measure φ , implies that all sets $A \in \mathcal{B}(\mathcal{X})$ with non-zero φ -measure are accessible. The notion of irreducibility does not really depend on the irreducibility measure φ as shown by the next result.

Theorem 8.13 ([6, Proposition 4.2.2]). *If $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ is irreducible for some irreducibility measure φ on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, then there exists a probability measure ψ on $\mathcal{B}(\mathcal{X})$, called maximal irreducibility measure such that*

- $\{X_n\}_n$ is ψ -irreducible
- For any other measure φ' on $\mathcal{B}(\mathcal{X})$ for which $\{X_n\}$ is φ' -irreducible, one has φ' is absolutely continuous with respect to ψ (i.e. for all $A \in \mathcal{B}(\mathcal{X})$, $\psi(A) = 0 \implies \varphi'(A) = 0$)

The maximal irreducibility measure is in general not unique, but all maximal irreducibility measures have the same null sets (i.e. they are equivalent). If $\{X_n\}_n$ is irreducible for some measure φ and has an invariant distribution π , then π is a maximal irreducibility measure. This, in particular, implies that an irreducible chain has at most one invariant probability measure. In the context of Markov Chain Monte Carlo, where the target distribution π is given, we have to check that the chain π is an irreducibility measure.

The notions of aperiodicity generalizes quite straightforwardly to accessible sets.

Definition 8.17 (aperiodicity). *A Markov chain $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ is aperiodic if for any $x \in \mathcal{X}$ and any accessible set $A \in \mathcal{B}(\mathcal{X})$*

$$\exists n_0 \geq 0 : \quad P^{(n)}(x, A) > 0 \quad \forall n \geq n_0.$$

On the other hand, in general state spaces, one can give two different definitions of recurrence, which are equivalent in discrete state spaces. For a set $A \in \mathcal{B}(\mathcal{X})$ let $V_A = \sum_{n \geq 0} \mathbb{1}_{\{X_n \in A\}}$ be the number of visits to A . In the discrete setting, a state x_i is recurrent if $\mathbb{P}_i(V_i = \infty) = 1$ which happens if and only if $\mathbb{E}_i[V_i] = \infty$. This “if and only if” result is not true anymore for general state spaces, so we can give two notions of recurrence depending on whether we take the definition $\mathbb{P}_i(V_i = \infty) = 1$ or $\mathbb{E}_i[V_i] = \infty$.

Definition 8.18 (recurrence). *A set $A \in \mathcal{B}(\mathcal{X})$ is said to be*

- recurrent if $\mathbb{E}_x[V_A] = \infty$, for all $x \in A$;
- Harris recurrent if $\mathbb{P}_x(V_A = \infty) = 1$, for all $x \in A$.

A Markov chain $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ is recurrent / Harris recurrent if it is irreducible and every accessible set is recurrent / Harris recurrent.

The first notion of recurrence is weaker than the second one as it requires only that the expected number of visits to A is infinite as opposed to the Harris recurrence condition that requires that almost surely the number of visits is infinite.

As in the discrete case, we have that if $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ is irreducible and recurrent, then it has a unique non-zero invariant measure $\hat{\pi}$ (non necessarily finite) up to a multiplicative constant.

We finally generalize the notion of positive recurrence. In the discrete case, an irreducible chain $\{X_n\}_n$ is positive recurrent, or simply *positive*, if $\mathbb{E}_i[\sigma_i] < \infty$ for any i , and we have seen that this happens if and only if the chain has an invariant probability measure. The latter condition is taken as definition of positive recurrence in general state spaces.

Definition 8.19. *A φ -irreducible Markov chain $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ is positive if it admits an invariant probability measure.*

A Markov chain $\{X_n\}_n$ constructed for MCMC will always be positive as long as it is irreducible, since it is designed to have a prescribed invariant probability measure.

With the above definitions, the convergence Theorem 8.12 generalizes as

Theorem 8.14. *Let $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ be irreducible, Harris recurrent, positive and aperiodic, with (unique) invariant probability distribution π . Then*

$$\forall \lambda \in \mathcal{M}_1(\mathcal{X}), \quad \lim_{n \rightarrow \infty} \|\lambda P^n - \pi\|_{TV} = 0.$$

We also mention how the ergodic theorem generalizes:

Theorem 8.15. *Let $\{X_n\}_n \sim \text{Markov}(\delta_x, P)$ be an irreducible, positive chain with invariant probability distribution π and $\psi \in \mathcal{F}(\mathcal{X})$ a π -integrable function with $\mathbb{E}_\pi[|\psi|] < \infty$. Then, for π -a.e. $x \in \mathcal{X}$*

$$\mathbb{P}_x \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \psi(X_j) = \mathbb{E}_\pi[\psi] \right) = 1.$$

8.3 Metropolis-Hastings algorithm in general state space

We generalize here the Metropolis-Hastings algorithm, already introduced in Section 8.1.1, to the case of a general state space, as a tool to construct a Markov Chain $\{X_n\} \sim \text{Markov}(\lambda, P)$ on $\mathcal{X} \subset \mathbb{R}^d$ which has a given invariant measure π with density $f : \mathcal{X} \rightarrow \mathbb{R}_+$ with respect to the Lebesgue measure. In the following discussion, we accept that the density f may be known only up to a multiplicative constant, i.e. it does not necessarily integrate to one (in which case, the invariant density is $\tilde{f}(x) = f(x) / \int_{\mathcal{X}} f(y) dy$).

Let $Q : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ be a Markov transition kernel $Q(x, A) = \int_A q(x, y) dy$ for all $x \in \mathcal{X}$, $A \in \mathcal{B}(\mathcal{X})$ with density $q : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ satisfying $q(x, y) = 0 \Leftrightarrow q(y, x) = 0$, also called the *proposal* or instrumental *density*, and define the following acceptance rate $\alpha : \mathcal{X} \times \mathcal{X} \rightarrow [0, 1]$,

$$\alpha(x, y) = \min \left\{ \frac{f(y) q(y, x)}{f(x) q(x, y)}, 1 \right\}, \quad \text{if } q(x, y) \neq 0, \quad \alpha(x, y) = 0, \quad \text{if } q(x, y) = 0.$$

The Metropolis-Hastings algorithm then reads

Algorithm 8.2: Metropolis-Hastings.

Given: λ (initial measure), q (proposal transition density), f (target density)

```

1 Generate  $X_0 \sim \lambda$ 
2 for  $n = 0, 1, \dots$ , do
3   | Generate  $Y_{n+1} \sim q(X_n, \cdot)$  // proposal state
4   | Generate  $U \sim \mathcal{U}(0, 1)$ 
5   | if  $U \leq \alpha(X_n, Y_{n+1})$  then
6   |   | set  $X_{n+1} = Y_{n+1}$  // accept proposal
7   | else
8   |   | set  $X_{n+1} = X_n$  // reject proposal
9   | end
10 end

```

For the algorithm to work, the chain has to be able to explore the whole density f . Let us denote $D_f = \text{supp}(f) = \{x \in \mathcal{X} : f(x) > 0\}$ the support of f . Minimum requirements are:

- $X_0 \in D_f$, otherwise $\alpha(X_0, \cdot)$ is not defined. This guarantees, in particular, that $X_n \in D_f, \forall n$;
- $\bigcup_{x \in D_f} \text{supp}(q(x, \cdot)) \supset D_f$, otherwise the chain fails to visit some parts of D_f .

We derive now the transition kernel P , resp. transition density p , of the Markov chain generated by the Metropolis-Hastings algorithm. There is, in general, a non-zero probability that $X_{n+1} = X_n$, so $P(X_n, \cdot)$ has a point mass in X_n :

$$\mathbb{P}(X_{n+1} = x \mid X_n = x) = \int_{\mathcal{X}} q(x, y)(1 - \alpha(x, y)) dy = 1 - \int_{\mathcal{X}} \alpha(x, y)q(x, y) dy$$

so the transition density p is

$$p(x, y) = \alpha(x, y)q(x, y) + (1 - \alpha^*(x))\delta_x(y), \quad \alpha^*(x) = \int_{\mathcal{X}} \alpha(x, y)q(x, y) dy$$

where $\delta_x(y)$ is a Dirac mass in x . Equivalently, the transition kernel P is given by

$$P(x, A) = \int_A \alpha(x, y)q(x, y)dy + (1 - \alpha^*(x))\mathbb{1}_A(x).$$

As in the discrete state space case, we can verify that P and f are in detailed balance.

Lemma 8.16. *The transition kernel P of the Metropolis-Hastings algorithm 8.2, with density $p(x, y) = \alpha(x, y)q(x, y) + (1 - \alpha^*(x))\delta_x(y)$ is in detailed balance with the probability density f . Hence f is an invariant probability density for P .*

Proof. Observe first that

$$\begin{aligned} f(x)q(x, y)\alpha(x, y) &= f(x)q(x, y) \min \left\{ \frac{f(y)q(y, x)}{f(x)q(x, y)}, 1 \right\} \\ &= \min \{ f(y)q(y, x), f(x)q(x, y) \} = f(y)q(y, x)\alpha(y, x). \end{aligned}$$

Hence

$$\begin{aligned} \int_A P(x, B)f(x) dx &= \int_A \left(\int_B (\alpha(x, y)q(x, y) + (1 - \alpha^*(x))\delta_x(y)) dy \right) f(x) dx \\ &= \int_A \int_B f(y)\alpha(y, x)q(y, x) dy dx + \int_{A \cap B} (1 - \alpha^*(x))f(x) dx \\ &= \int_B \left(\int_A (\alpha(y, x)q(y, x) + (1 - \alpha^*(y))\delta_y(x)) dx \right) f(y) dy \\ &= \int_B P(y, A)f(y) dy. \end{aligned}$$

□

To assess the convergence to equilibrium of the chain, we should further check irreducibility and aperiodicity. In particular, irreducibility should be checked with respect to the invariant density f .

- f -irreducibility is something that should be checked every time depending on the choice of the proposal density. If it holds, then for all $\varphi : \mathcal{X} \rightarrow \mathbb{R}$, $\mathbb{E}_f[|\varphi|] < +\infty$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \varphi(X_j) = \mathbb{E}_f[\varphi] = \int_{\mathcal{X}} \varphi(x)f(x) dx.$$

- Concerning aperiodicity, observe that in general $\mathbb{P}(X_{n+1} = x \mid X_n = x) > 0$ as long as $\alpha^*(x) < 1$, since the transition kernel $P(x, \cdot)$ has an atom at x . Consider the set $C = \{x : \alpha(x) < 1\}$. This is a f -zero measure set, i.e. $\int_C f(x) dx = 0$, if and only if

(exercise) $f(x)q(x, y) = f(y)q(y, x)$ for f -almost every $x, y \in D_f$ which corresponds to the case in which the proposal q is in detailed balance with f . In this case, the acceptance-rejection step is useless and one should check the aperiodicity of q . If, on the other hand, (q, f) are *not* in detailed balance, then the chain is aperiodic. If, moreover, the chain is f -irreducible, then for any initial distribution $\lambda \ll f$ we have (denoting π the measure associated to f)

$$\lim_{n \rightarrow \infty} \|\pi^{n, \lambda} - \pi\|_{TV} = 0.$$

We describe in the next subsections few methods to choose proposal densities q .

8.3.1 Independence sampler

Let $g : \mathcal{X} \rightarrow \mathbb{R}_+$ be a probability density function such that $g(x) > 0$ whenever $f(x) > 0$ (i.e. $f \ll g$). We choose simply $q(x, y) = g(y)$ independently of the current state x (hence the name of *independence sampler*).

Algorithm 8.3: Independence sampler Metropolis-Hastings

Given: $X_0 \sim \lambda$, $\text{supp}(\lambda) \subset D_f$

- 1 **for** $n = 0, 1, \dots$, **do**
- 2 Generate $Y_{n+1} \sim g$
- 3 Compute $\alpha(X_n, Y_{n+1}) = \min \left\{ \frac{f(Y_{n+1})}{f(X_n)} \frac{g(X_n)}{g(Y_{n+1})}, 1 \right\}$
- 4 Generate $U \sim \mathcal{U}(0, 1)$ and set

$$X_{n+1} = \begin{cases} Y_{n+1}, & \text{if } U \leq \alpha(X_n, Y_{n+1}) \\ X_n, & \text{otherwise} \end{cases}$$
- 5 **end**

Concerning the convergence to equilibrium, we recall first a useful result for general state space Markov chains.

Lemma 8.17. *Let $P : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ be a Markov transition kernel with invariant measure π . If there exists $\epsilon \in (0, 1)$ and a probability measure ν on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ such that*

$$P(x, A) \geq \epsilon \nu(A), \quad \forall x \in \mathcal{X}, \quad A \in \mathcal{B}(\mathcal{X}) \quad (8.3)$$

then

$$\|\pi^{n, \lambda} - \pi\|_{TV} \leq 2(1 - \epsilon)^n.$$

More generally, if there exists $k_0 \in \mathbb{N}$ such that $P^{(k_0)}(x, A) \geq \epsilon \nu(A)$ for all $x \in \mathcal{X}$, $A \in \mathcal{B}(\mathcal{X})$, then $\|\pi^{n, \lambda} - \pi\|_{TV} \leq 2(1 - \epsilon)^{\lfloor n/k_0 \rfloor}$. The condition (8.3) is called uniform minorizing condition.

Idea of the proof. (For $k_0 = 1$): We build two coupled chains $\{X_n\} \sim \text{Markov}(\lambda, P)$ and $\{Y_n\} \sim \text{Markov}(\pi, P)$ using the following algorithm. Notice, in particular, that the chain $\{Y_n\}$ is at stationarity.

Algorithm 8.4: Coupled chains.

```

1 Let  $X_0 \sim \lambda, Y_0 \sim \pi$ 
2 for  $n = 0, 1, \dots$ , do
3   Draw  $Z_n \sim \text{Be}(\epsilon), \mathbb{P}(Z_n = 1) = \epsilon, \mathbb{P}(Z_n = 0) = 1 - \epsilon$ 
4   if  $Z_n = 1$  then
5     draw  $W \sim \nu$  and set  $X_{n+1} = Y_{n+1} = W$ 
6   else
7     draw  $X_{n+1} \sim \frac{P(X_n, \cdot) - \epsilon \nu(\cdot)}{1 - \epsilon}$  and  $Y_{n+1} \sim \frac{P(Y_n, \cdot) - \epsilon \nu(\cdot)}{1 - \epsilon}$  independently
8   end
9 end

```

It is easy to verify that indeed $\{X_n\} \sim \text{Markov}(\lambda, P)$ and $\{Y_n\} \sim \text{Markov}(\pi, P)$. Let $T = \inf\{n \geq 0 : Z_n = 1\}$. It is clear that after T , the two chains have the same distribution $X_n \sim Y_n, n > T$. Moreover, $\mathbb{P}(T \geq n) = (1 - \epsilon)^n$. Now

$$\begin{aligned}
\|\pi^{n,\lambda} - \pi\|_{\text{TV}} &= 2 \sup_{A \in \mathcal{B}(\mathcal{X})} |\mathbb{P}(X_n \in A) - \mathbb{P}(Y_n \in A)| \\
&= 2 \sup_A |\mathbb{P}(X_n \in A, T < n) + \mathbb{P}(X_n \in A, T \geq n) \\
&\quad - \mathbb{P}(Y_n \in A, T < n) - \mathbb{P}(Y_n \in A, T \geq n)| \\
&= 2 \sup_A |\mathbb{P}(X_n \in A, T \geq n) - \mathbb{P}(Y_n \in A, T \geq n)| \\
&= 2 \sup_A |\mathbb{P}(X_n \in A, Y_n \notin A, T \geq n) - \mathbb{P}(X_n \notin A, Y_n \in A, T \geq n)| \\
&\leq 2\mathbb{P}(T \geq n) \leq 2(1 - \epsilon)^n.
\end{aligned}$$

□

In the case of the independence sampler, the following result holds.

Theorem 8.18. *If there exists $M < +\infty$ such that $f(x) \leq Mg(x)$ for all $x \in \mathcal{X}$, then the chain generated by the independence sampler algorithm 8.3 is uniformly ergodic and*

$$\|\pi^{n,\lambda} - \pi\|_{\text{TV}} \leq \left(1 - \frac{\int f(x) dx}{M}\right)^n, \quad \text{for any } \lambda.$$

Proof. If f is not normalized, let $\tilde{f} = f/C, C = \int_{\mathcal{X}} f$. Notice that

$$\alpha(x, y)q(x, y) = g(y) \min \left\{ \frac{f(y)g(x)}{f(x)g(y)}, 1 \right\} = \min \left\{ f(y) \underbrace{\frac{g(x)}{f(x)}}_{\geq 1/M}, \underbrace{g(y)}_{\geq f(y)/M} \right\} \geq \frac{1}{M} f(y).$$

It follows that for any $A \in \mathcal{B}(\mathcal{X})$,

$$P(x, A) = \int_A (\alpha(x, y)q(x, y) + (1 - \alpha^*(x))\delta_x(y)) dy \geq \frac{1}{M} \int_A f(y) dy \geq \frac{C}{M} \pi(A)$$

and the result follows from Lemma 8.17. □

Under the same condition as in Theorem 8.18, it can be shown that the expected acceptance probability satisfies $\mathbb{E}[\alpha(X_n, Y_{n+1})] \geq \frac{C}{M}$ (exercise). This result has to be compared with a pure acceptance-rejection sampling strategy, for which the expected acceptance probability is $\frac{C}{M}$. Hence, independence MH sampler accepts more often than a pure acceptance-rejection sampler.

8.3.2 Random walk Metropolis

Let $g_\sigma : \mathcal{X} \rightarrow \mathbb{R}_+$ be a probability density function with zero mean, σ being a scaling parameter; a typical choice is $g_\sigma = N(0, \sigma^2)$. In the random walk Metropolis we choose $q(x, y) = g_\sigma(y - x)$, i.e. the proposal density is g_σ centred in the current state x . If we further assume $g_\sigma(\cdot)$ symmetric around the origin, the acceptance probability takes the simplified form

$$\alpha(x, y) = \min \left\{ \frac{f(y)}{f(x)}, 1 \right\}.$$

The choice of σ is rather delicate. Small σ imply small steps from the current state, hence high correlation in the chain. Large steps might lead to high rejection rate, hence the chain will stay for a long time in the given state, which also leads to high correlation in the chain. One should then expect that some “optimal” choice of σ exists.

Concerning convergence of this algorithm, one could try to verify a uniform minorizing condition

$$g_\sigma(y - x) \geq \epsilon f(y)$$

for all $x, y \in D_f$. By the same arguments as for independence sampler, this would imply $P(x, A) \geq \epsilon \pi(A)$ hence uniform ergodicity $\|\pi^{n, \lambda} - \pi\|_{\text{TV}} \leq 2(1 - \epsilon)^n$ for all initial distributions λ . However, such minorizing condition does not hold, in general for unbounded or non-compact $D_f \subset \mathcal{X}$. We mention a result by Mengersen and Tweedie ('96) showing geometric ergodicity for tail-log-concave f and $\mathcal{X} = \mathbb{R}$.

Definition 8.20. *A probability density function f on \mathbb{R} is log-concave in the tails if there exists $\alpha, M > 0$ such that $\log f(x) - \log f(y) \geq \alpha(|y| - |x|)$ for all $|y| \geq |x| \geq M$.*

Theorem 8.19. *If the invariant density f on \mathbb{R} is log concave in tails for some $\alpha, M > 0$ and $\inf_{|x| \leq R} f(x) > 0$ for all $R > 0$, then the Markov chain generated by the random walk Metropolis-Hastings algorithm with symmetric proposal $g_\sigma(\cdot)$ is geometrically ergodic.*

8.3.3 One Variable at a time Metropolis-Hastings

Suppose that a state $x \in \mathcal{X}$ has several components, $x = (x^{(1)}, \dots, x^{(d)})$, with $x^{(i)} \in \mathcal{X}^{(i)}$. One can thus construct a Metropolis-Hastings algorithm by updating one component at a time, either chosen randomly or by performing a systematic sweep over the components. Say that the i -th component has been chosen. We use the notation $x = (x^{(i)}, x^{(-i)})$ with $x^{(-i)} = (x^{(1)}, \dots, x^{(i-1)}, x^{(i+1)}, \dots, x^{(d)})$. Let $q_i : \mathcal{X} \times \mathcal{X}^{(i)} \rightarrow \mathbb{R}$ be a family of proposal density functions on $\mathcal{X}^{(i)}$, i.e. $q_i(x, \cdot)$ is a density function on $\mathcal{X}^{(i)}$ for any $x \in \mathcal{X}$. Then the one variable at a time MH algorithm with random coordinate selection reads:

Algorithm 8.5: One variable at a time MH with random selection.

```

1 Generate  $X_0 \sim \lambda$ 
2 for  $n = 0, 1, \dots$  do
3   Draw index  $i_n \sim \beta$  (p.m.f on  $\{1, \dots, d\}$ )
4   Draw  $y \sim q_{i_n}(X_n, \cdot)$  and set  $Y_{n+1} = (y, X_n^{(-i_n)})$ 
5   Compute  $\alpha_{i_n}(X_n, Y_{n+1}) = \min \left\{ \frac{f(Y_{n+1}) q_{i_n}(Y_{n+1}, X_n^{(i_n)})}{f(X_n) q_{i_n}(X_n, Y_{n+1}^{(i_n)})}, 1 \right\}$ 
6   Set  $X_{n+1} = \begin{cases} Y_{n+1} & \text{with prob. } \alpha_{i_n}(X_n, Y_{n+1}) \\ X_n & \text{otherwise} \end{cases}$ 
7 end

```

whereas the one variable at a time MH algorithm with systematic sweep over the coordinates reads:

Algorithm 8.6: One variable at a time MH with systematic sweep.

```

1 Generate  $X_0 \sim \lambda$ 
2 for  $n = 0, 1, \dots$  do
3   Set  $Y_{n+1,0} = X_n$ 
4   for  $i = 1, \dots, d$  do
5     Draw  $y \sim q_i(X_n, \cdot)$  and set  $\tilde{Y} = (y, Y_{n+1,i-1}^{(-i)})$ 
6     Set  $Y_{n+1,i} = \begin{cases} \tilde{Y}, & \text{with prob. } \alpha_i(Y_{n+1,i-1}, \tilde{Y}) \\ Y_{n+1,i-1}, & \text{otherwise} \end{cases}$ 
7   end
8    $X_{n+1} = Y_{n+1,d}$ 
9 end

```

To see that this algorithm produces a Markov chain with the correct invariant distribution we proceed as follows. Suppose index i has been selected in either Algorithm 8.5 or 8.6 and define the Markov transition density induced by the Metropolis-Hastings step on the i -th component:

$$\begin{aligned} p_i(x, y) &= p_i((x^{(i)}, x^{(-i)}), (y^{(i)}, y^{(-i)})) \\ &= \left(\alpha_i(x, y^{(i)}) q_i(x, y^{(i)}) + (1 - \alpha_i^*(x)) \delta_{x^{(i)}}(y^{(i)}) \right) \delta_{x^{(-i)}}(y^{(-i)}) \end{aligned}$$

with $\alpha_i(x, y^{(i)}) = \min \left\{ 1, \frac{f(y^{(i)}, x^{(-i)}) q_i((y^{(i)}, x^{(-i)}), x^{(i)})}{f(x^{(i)}, x^{(-i)}) q_i((x^{(i)}, x^{(-i)}), y^{(i)})} \right\}$ and $\alpha_i^*(x) = \int_{\mathcal{X}} \alpha_i(x, y^{(i)}) q_i(x, y^{(i)}) dy^{(i)}$. Let $P_i(x, A) = \int_A p_i(x, y) dy$ denote the corresponding Markov transition kernel and \mathcal{P}_i the associated operator.

We show first that (P_i, f) are in detailed balance, which implies that f is an invariant distribution for \mathcal{P}_i .

Lemma 8.20. *Let $P_i : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ be the Markov transition kernel defined by the Metropolis-Hastings step in Algorithms 8.5 or 8.6 when the i -th component is selected. Then P_i is in detailed balance with f .*

Proof. We use the characterization of detailed balance by integration against measurable and bounded functions. We aim at showing that for any $g \in \mathcal{F}_b(\mathcal{X} \times \mathcal{X})$ we have

$$\int_{\mathcal{X}^2} g(x, y) P_i(x, dy) f(x) dx = \int_{\mathcal{X}^2} g(x, y) P_i(y, dx) f(y) dy. \quad (8.4)$$

Let us split the transition density as $p_i(x, y) = p_i^1(x, y) + p_i^2(x, y)$ with

$$p_i^1(x, y) = \alpha_i(x, y^{(i)}) q_i(x, y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) \quad p_i^2(x, y) = (1 - \alpha_i^*(x)) \delta_{x^{(i)}}(y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}).$$

Then

$$\int_{\mathcal{X}^2} g(x, y) P_i(x, dy) f(x) dx = \underbrace{\int_{\mathcal{X}^2} g(x, y) p_i^1(x, y) f(x) dx dy}_A + \underbrace{\int_{\mathcal{X}^2} g(x, y) p_i^2(x, y) f(x) dx dy}_B$$

and

$$\begin{aligned} A &= \int_{\mathcal{X}^2} g(x, y) \min \left\{ 1, \frac{f(y^{(i)}, x^{(-i)}) q_i((y^{(i)}, x^{(-i)}), x^{(i)})}{f(x) q_i(x, y^{(i)})} \right\} q_i(x, y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) f(x) dx dy \\ &= \int_{\mathcal{X}^2} g(x, y) \min \left\{ 1, \frac{f(y) q_i(y, x^{(i)})}{f(x) q_i(x, y^{(i)})} \right\} q_i(x, y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) f(x) dx dy \\ &= \int_{\mathcal{X}^2} g(x, y) \min \left\{ \frac{f(x) q_i(x, y^{(i)})}{f(y) q_i(y, x^{(i)})}, 1 \right\} q_i(y, x^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) f(y) dx dy \\ &= \int_{\mathcal{X}^2} g(x, y) \min \left\{ \frac{f(x^{(i)}, y^{(-i)}) q_i((x^{(i)}, y^{(-i)}), y^{(i)})}{f(y) q_i(y, x^{(i)})}, 1 \right\} q_i(y, x^{(i)}) \delta_{y^{(-i)}}(x^{(-i)}) f(y) dx dy \\ &= \int_{\mathcal{X}^2} g(x, y) \alpha_i(y, x^{(i)}) q_i(y, x^{(i)}) \delta_{y^{(-i)}}(x^{(-i)}) f(y) dx dy. \end{aligned}$$

Likewise

$$B = \int_{\mathcal{X}^2} g(x, y) (1 - \alpha_i^*(x)) \delta_x(y) f(x) dx dy = \int_{\mathcal{X}^2} g(x, y) (1 - \alpha_i^*(y)) \delta_y(x) f(y) dx dy.$$

Summing A and B we obtain (8.4). \square

Having established this result, we can now analyze the full Algorithms 8.5 and 8.6. The Markov transition kernel of algorithm 8.5 is given by

$$\begin{aligned} P^{\text{rand}}(x, A) &= \mathbb{P}(X_{n+1} \in A \mid X_n = x) \\ &= \sum_{i=1}^d \mathbb{P}(X_{n+1} \in A \mid X_n = x, i_n = i) \mathbb{P}(i_n = 1) = \frac{1}{d} \sum_{i=1}^d P_i(x, A) \end{aligned}$$

and the associated Markov operator is $\mathcal{P}^{\text{rand}} = \frac{1}{d} \sum_{i=1}^d \mathcal{P}_i$. Since each \mathcal{P}_i is in detailed balance with f it follows that also $\mathcal{P}^{\text{rand}}$ is in detailed balance with f , which is then as invariant distribution, and, moreover, it produces a *reversible* chain.

On the other hand, the Markov transition kernel of Algorithm 8.6 is given by

$$P^{\text{sweep}}(x, A) = \mathbb{P}(X_{n+1} \in A \mid X_n = x) = \int_A \int_{\mathcal{X}^{d-1}} P_d(y_{d-1}, dy_d) \cdots P_2(y_1, dy_2) P_1(x, dy_1)$$

and the associated Markov operator is $\mathcal{P}^{\text{sweep}} = \mathcal{P}_1 \cdots \mathcal{P}_d$. Since each \mathcal{P}_i leaves f invariant, it also follows that $\mathcal{P}^{\text{sweep}}$ has f as invariant distribution. However, contrary to Algorithm 8.5, the Markov chain produced by this algorithm *is not reversible*, although each \mathcal{P}_i is. To see this, observe that

$$\begin{aligned} \int_{\mathcal{X}^2} g(x, y_d) P^{\text{sweep}}(x, dy_d) f(x) dx &= \int_{\mathcal{X}^{d+1}} g(x, y_d) P_d(y_{d-1}, dy_d) \cdots P_2(y_1, dy_2) P_1(x, dy_1) f(x) dx \\ &= \int_{\mathcal{X}^{d+1}} g(x, y_d) P_d(y_{d-1}, dy_d) \cdots P_2(y_1, dy_2) f(y_1) P_1(y_1, dx) dy_1 \\ &= \int_{\mathcal{X}^{d+1}} g(x, y_d) f(y_d) P_d(y_d, dy_{d-1}) \cdots P_2(y_2, dy_1) P_1(y_1, dx) dy_d \\ &= \int_{\mathcal{X}^2} g(x, y_d) \hat{P}^{\text{sweep}}(y_d, dx) f(y_d) dx \end{aligned}$$

with

$$\hat{P}^{\text{sweep}}(y_d, A) = \mathbb{P}(X_n \in A \mid X_{n+1} = y_d) = \int_A \int_{\mathcal{X}^{d-1}} P_1(y_1, dx) \cdots P_{d-1}(y_{d-1}, dy_{d-2}) P_d(y_d, dy_{d-1})$$

that is, the Markov operator associated to the reversed chain is $\hat{\mathcal{P}}^{\text{sweep}} = \mathcal{P}_d \cdots \mathcal{P}_1 \neq \mathcal{P}^{\text{sweep}} = \mathcal{P}_1 \cdots \mathcal{P}_d$.

To recover a reversible chain one should use a palindromic structure $\mathcal{P} = \mathcal{P}_1 \cdots \mathcal{P}_d \cdots \mathcal{P}_1$, i.e. a forward loop $i = 1, \dots, d$ followed by a backward loop $i = d - 1, \dots, 1$ over the components.

8.3.4 Gibbs sampler

The Gibbs sampler is a *one variable at a time* MH algorithm in which the component-wise proposal density is the conditional density $q_i(x, \cdot) = f_{X^{(i)} \mid X^{(-i)}}(\cdot \mid x^{(-i)})$. Observe that, in this case, the Hastings ratio for $x = (x^{(i)}, x^{(-i)})$ and $y = (y^{(i)}, y^{(-i)})$ is

$$\begin{aligned} \alpha_i(x, y) &= \min \left\{ \frac{f(y) f_{X^{(i)} \mid X^{(-i)}}(x^{(i)} \mid x^{(-i)})}{f(x) f_{X^{(i)} \mid X^{(-i)}}(y^{(i)} \mid x^{(-i)})}, 1 \right\} \\ &= \min \left\{ \frac{f(y) f(x) / f_{X^{(-i)}}(x^{(-i)})}{f(x) f(y) / f_{X^{(-i)}}(x^{(-i)})}, 1 \right\} = 1 \end{aligned}$$

i.e. the move is always accepted, or, in other words, the transition kernel Q_i which samples independently the i -th component from the conditional density $f_{X^{(i)} \mid X^{(-i)}}$ preserves the density f . The next algorithm presents the Gibbs sampler with random sweep.

Algorithm 8.7: Gibbs with random sweep.

```

1 Generate  $X_0 \sim \lambda$ 
2 for  $n = 0, 1, \dots$  do
3   | Draw  $i_n$  from a pmf  $\beta$  on  $\{1, \dots, d\}$ 
4   | Generate  $y^{(i_n)} \sim f(\cdot \mid X_n^{(-i_n)})$ 
5   | Set  $X_{n+1} = (y^{(i_n)}, X_n^{(-i_n)})$ 
6 end

```

8.3.5 Metropolis Adjusted Langevin Algorithm (MALA)

The MALA algorithm relies on the following observation: consider the stochastic differential equation (Langevin dynamics)

$$dX_t = \nabla \log f(X_t) + \sqrt{2} dW_t, \quad t > 0, \quad X_0 \sim \lambda \quad (8.5)$$

with $X_t \in \mathbb{R}^d$, where f is the target density, W_t is a standard Wiener process and λ is a probability density function on \mathbb{R}^d . At any t , let us denote by $\rho(x, t) : \mathcal{X} \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$ the probability density function of X_t , i.e.

$$\int_A \rho(x, t) dx = \mathbb{P}_\lambda(X_t \in A).$$

It is well known that ρ satisfies the so-called Fokker-Planck equation

$$\partial_t \rho + \operatorname{div}(\rho \nabla \log f) - \Delta \rho = 0, \quad \text{in } \mathbb{R}^d, \quad t > 0,$$

with $\rho(x, 0) = \lambda(x)$, from which we see that $\bar{\rho}(x, t) = f(x)$ is a stationary solution. Indeed,

$$\partial_t \bar{\rho} + \operatorname{div}(\bar{\rho} \nabla \log f) - \Delta \bar{\rho} = \sum_{i=1}^d \partial_{x_i} \left(\bar{\rho} \frac{\partial_{x_i} f}{f} \right) - \Delta \bar{\rho} = 0.$$

Under mild assumptions on f , such stationary solution is unique and $\lim_{t \rightarrow \infty} \rho(\cdot, t) = f$ (in a suitable sense) for any initial density λ . Hence, the time continuous process (8.5) has f as unique invariant distribution. The problem is that we are not able to find exact solutions of (8.5), in general, and we have to use some numerical scheme for example the Euler-Maruyama:

$$X_{n+1} = X_n + \Delta t \nabla \log f(X_n) + \sqrt{2\Delta t} \xi_n, \quad \xi_n \sim N(0, I) \quad (8.6)$$

i.e.

$$X_{n+1} \sim N(X_n + \Delta t \nabla \log f(X_n), 2\Delta t I). \quad (8.7)$$

However, after discretization, (8.6) does not have anymore f as invariant distribution. Yet (8.7) can be used as a *proposal* in a Metropolis-Hastings algorithm. This leads to the following

Algorithm 8.8: Metropolis Adjusted Langevin Algorithm (MALA).

```

1 Generate  $X_0 \sim \lambda$ 
2 for  $n = 0, 1, \dots$  do
3   Generate  $Y \sim N(X_n + \Delta t \nabla \log f(X_n), 2\Delta t I)$ 
4   Compute  $\alpha(X_n, Y) = \min \left\{ 1, \frac{f(Y)}{f(X_n)} \frac{\exp(-\|X_n - Y - \Delta t \nabla \log f(Y)\|^2 / 2\Delta t)}{\exp(-\|Y - X_n - \Delta t \nabla \log f(X_n)\|^2 / 2\Delta t)} \right\}$ 
5   Set  $X_{n+1} = \begin{cases} Y & \text{with prob. } \alpha(X_n, Y) \\ X_n & \text{otherwise} \end{cases}$ 
6 end

```

8.4 Convergence diagnostics

Let us consider an f -irreducible aperiodic Metropolis-Hastings Markov chain. Given any function $\varphi : \mathbb{E}_f[\varphi] < +\infty$ and f -a.e. initial state, by the ergodic theorem

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \varphi(X_j) = \mathbb{E}_f[\varphi].$$

Hence, to compute $\mu = \mathbb{E}_f[\phi]$, we can consider the estimator

$$\hat{\mu}_N^{\text{mcmc}} = \frac{1}{N} \sum_{j=1}^N \varphi(X_j).$$

The question is how to monitor properly the convergence of $\hat{\mu}^{\text{mcmc}}$ to μ and how to choose N .

We start by analyzing the *Bias*. The estimator $\hat{\mu}_N^{\text{mcmc}}$ is biased, in general, since $X_n \sim f$ only asymptotically as $n \rightarrow \infty$. The bias is generally of order $\frac{1}{N}$ as shown in the next lemma.

Lemma 8.21. *Let $\{X_n\} \sim \text{Markov}(\delta_x, P)$ with P a Metropolis-Hastings transition kernel with invariant distribution π and density f . If $\{X_n\}$ is geometrically ergodic, that is, there exists $\gamma > 0$ and $h : \mathcal{X} \rightarrow \mathbb{R}_+$ such that $\|\pi^{n, \delta_x} - \pi\|_{TV} \leq h(x)e^{-\gamma n}$ then for any bounded $\varphi : \mathcal{X} \rightarrow \mathbb{R}$, there exists $C_\varphi > 0$ such that*

$$|\mathbb{E}[\hat{\mu}_N^{\text{mcmc}} - \mu]| \leq \frac{C}{N}.$$

Proof.

$$\begin{aligned}
|\mathbb{E} [\hat{\mu}_N^{\text{mcmc}} - \mu]| &= \left| \frac{1}{N} \sum_{j=1}^N \mathbb{E} [\varphi(X_j) - \mu] \right| \\
&\leq \frac{1}{N} \sum_{j=1}^N \left| \int_{\mathcal{X}} \varphi(y) (\pi^{j, \delta_x}(dy) - \pi(dy)) \right| \\
&\leq \frac{1}{N} \sum_{j=1}^N \sup_{x \in \mathcal{X}} |\varphi(x)| \|\pi^{j, \delta_x} - \pi\|_{\text{TV}} \\
&\leq \frac{1}{N} \sup_{x \in \mathcal{X}} |\varphi(x)| h(x) \frac{1}{1 - e^{-\gamma}}.
\end{aligned}$$

□

Such bias can be further reduced by considering the estimator $\hat{\mu}_{N,B}^{\text{mcmc}} = \frac{1}{N} \sum_{j=B+1}^{N+B} \varphi(X_j)$ i.e. by disregarding the first B terms of the chain. The lag B is often called the *burn-in* or *warm-up* period. Under the assumptions of the previous lemma, the bias of the estimator $\hat{\mu}_{N,B}^{\text{mcmc}}$ is bounded by

$$|\mathbb{E} [\hat{\mu}_{N,B}^{\text{mcmc}}] - \mu| \leq \frac{e^{-\gamma B}}{N} \sup_{x \in \mathcal{X}} |\varphi(x)| \frac{h(x)}{1 - e^{-\gamma}}$$

and is thus reduced by a factor $e^{-\gamma B}$ with respect to the base estimator $\hat{\mu}_N^{\text{mcmc}} = \hat{\mu}_{N,0}^{\text{mcmc}}$. The quantity $\frac{1}{\gamma}$ is often called the *relaxation time* and choosing $B = \frac{m}{\gamma}$ with moderate m makes the bias negligible. Estimating the relaxation time is not easy. However, a graphical inspection of the trace plot of the chain $\{\varphi(X_n)\}$ is often sufficient to have a reasonable estimation of the time at which the chain reaches stationarity.

We focus now on the variance of the estimator $\hat{\mu}_N^{\text{mcmc}}$ (or $\hat{\mu}_{N,B}^{\text{mcmc}}$). Assuming that a sufficient burn-in period has been considered, we can reasonably assume in the analysis that follows that the chain is at stationarity when computing the estimator $\hat{\mu}_N^{\text{mcmc}}$, i.e. $\{X_n\} \sim \text{Markov}(\pi, P)$. Let us denote $c(k) = \text{Cov}_{\pi}(\varphi(X_0), \varphi(X_k)) = \text{Cov}_{\pi}(\varphi(X_j), \varphi(X_{j+k}))$ for all j thanks to the stationarity of the the chain.

Lemma 8.22. *Let $\{X_n\} \sim \text{Markov}(\pi, P)$ Then*

$$\text{Var}_{\pi}[\hat{\mu}_N^{\text{mcmc}}] = \frac{\sigma_{\text{mcmc},N}^2}{N}, \quad \text{with } \sigma_{\text{mcmc},N}^2 = c(0) + 2 \sum_{\ell=1}^{N-1} \left(1 - \frac{\ell}{N}\right) c(\ell).$$

Moreover, if $\sum_{k=0}^{\infty} |c(k)| < +\infty$, then

$$\lim_{N \rightarrow \infty} N \text{Var}(\hat{\mu}_N^{\text{mcmc}}) = \sigma_{\text{mcmc}}^2$$

with $\sigma_{\text{mcmc}}^2 = c(0) + 2 \sum_{k=1}^{\infty} c(k)$.

Proof.

$$\begin{aligned}
\text{Var}_\pi[\hat{\mu}_N^{\text{mcmc}}] &= \mathbb{E}_\pi \left[\left(\frac{1}{N} \sum_{j=1}^N \varphi(X_j) - \mu \right)^2 \right] \\
&= \frac{1}{N^2} \sum_{j=1}^N \sum_{k=1}^N \mathbb{E}_\pi [(\varphi(X_j) - \mu)(\varphi(X_k) - \mu)] \\
&= \frac{1}{N^2} \left[\sum_{j=1}^N \underbrace{\text{Var}_\pi[\varphi(X_j)]}_{c(0)} + 2 \sum_{j=1}^{N-1} \sum_{k=j+1}^N \underbrace{\text{Cov}_\pi(\varphi(X_j), \varphi(X_k))}_{c(k-j)} \right] \\
&= \frac{c(0)}{N} + \frac{2}{N^2} \sum_{j=1}^{N-1} \sum_{\ell=1}^{N-j} c(\ell) \\
&= \frac{c(0)}{N} + \frac{2}{N} \sum_{\ell=1}^{N-1} \frac{N-\ell}{N} c(\ell) \\
&= \frac{1}{N} \left(c(0) + 2 \sum_{\ell=1}^{N-1} \left(1 - \frac{\ell}{N}\right) c(\ell) \right).
\end{aligned}$$

Under the assumption $\sum_{\ell=0}^{\infty} |c(\ell)| < +\infty$, it follows that $\lim_{N \rightarrow \infty} N \text{Var}_\pi[\hat{\mu}_N^{\text{mcmc}}] = \sigma_{\text{mcmc}}^2$. \square

The quantity σ_{mcmc}^2 is called *time-average variance constant* (TAVC) or *asymptotic variance*. If $\{X_n\}_{n=1}^N$ were independent and all distributed as π , then the variance of the Crude Monte Carlo estimator $\hat{\mu}_N^{\text{MC}} = \frac{1}{N} \sum_{j=1}^N \varphi(X_j)$ would be $\text{Var}(\hat{\mu}_N^{\text{MC}}) = \frac{c(0)}{N}$. From this we see that

$$\lim_{N \rightarrow \infty} \frac{\text{Var}(\hat{\mu}_N^{\text{mcmc}})}{\text{Var}(\hat{\mu}_N^{\text{MC}})} = \frac{\sigma_{\text{mcmc}}^2}{c(0)} = 1 + 2 \sum_{k=1}^{\infty} \frac{c(k)}{c(0)}.$$

Hence $\hat{\mu}_N^{\text{mcmc}}$ is generally less effective than a pure iid sampling from π , due to the correlation in the chain. The quantity

$$ESS = N \frac{c(0)}{\sigma_{\text{mcmc}}^2}$$

is called the *effective sample size* and represents the size of an equivalent independent sample that would lead to the same variance of the estimator.

For the estimator $\hat{\mu}_N^{\text{mcmc}}$ a CLT is also available (and more generally for aperiodic, irreducible and reversible chains with invariant distribution π).

Theorem 8.23 (CLT for Metropolis-Hastings Markov Chains). *Let $\{X_n\}$ be an f -irreducible, aperiodic Metropolis-Hastings chain, with invariant distribution π (resp. density f) and $\varphi : \mathcal{X} \rightarrow \mathbb{R}$ such that*

$$\sigma_{\text{mcmc}}^2 := \text{Var}_\pi(\varphi(X_0)) + 2 \sum_{\ell=1}^{\infty} \text{Cov}_\pi(\varphi(X_0), \varphi(X_\ell)) < +\infty,$$

then,

$$\sqrt{N}(\hat{\mu}_N^{\text{mcmc}} - \mu) \xrightarrow{d} N(0, \sigma_{\text{mcmc}}^2)$$

as $N \rightarrow \infty$.

From the CLT, asymptotic confidence intervals can be derived. The practical question, however, is how to estimate σ_{mcmc}^2 .

8.4.1 Estimating the asymptotic variance by covariance methods

We recall the formula $\sigma_{\text{mcmc}}^2 = c(0) + 2 \sum_{k=1}^{\infty} c(k)$. Given a path $\{X_n\}_{n=0}^N$, if we discard a sufficient burn-in lag B , we can reasonably assume that $\{X_n\}_{n=B+1}^{N+B}$ is (nearly) stationary, so that a sample estimator for $c(k)$ is

$$\hat{c}(k) = \frac{1}{N-k-1} \sum_{j=B+1}^{N+B-k} (\varphi(X_j) - \hat{\mu}_{N,B}^{\text{mcmc}})(\varphi(X_{j+k}) - \hat{\mu}_{N,B}^{\text{mcmc}})$$

and an estimator for σ_{mcmc}^2 is

$$\hat{\sigma}_{\text{mcmc}}^2 = \hat{c}(0) + 2 \sum_{k=1}^{N-2} \hat{c}(k).$$

However, the last terms in the sum are very unstable since these are sample averages of very few terms. It is often wiser to truncate the sum much earlier

$$\hat{\sigma}_M^2 = \hat{c}(0) + 2 \sum_{k=1}^M \hat{c}(k).$$

where $M < N - 2$. It has been shown [Geyer '92] that the sequence $\Gamma_k = c(2k) + c(2k+1)$ is strictly positive, decreasing and convex for a reversible Markov Chain. Hence a good choice is

$$M = 2 \min\{k : \hat{c}(2k) + \hat{c}(2k+1) < 0\}.$$

8.4.2 Estimating the asymptotic variance by the batch means method

An alternative idea to estimate σ_{mcmc}^2 is to split the sequence $\{X_n\}_{n=B+1}^{N+B}$ into M blocks of size $T = N/M$ (assumed to be an integer). Then we can build M different sample averages

$$\hat{\mu}^{(i)} = \frac{1}{T} \sum_{j=(i-1)T+B+1}^{iT+B} \varphi(X_j), \quad \text{and} \quad \hat{\mu}_{N,B}^{\text{mcmc}} = \frac{1}{M} \sum_{i=1}^M \hat{\mu}^{(i)}.$$

If T is sufficiently large (larger than the relaxation time), the M blocks are nearly independent so $\text{Var}(\hat{\mu}_{N,B}^{\text{mcmc}}) \approx \frac{\sigma_{\text{mcmc}}^2}{N} \approx \frac{\text{Var}(\hat{\mu}^{(1)})}{M}$ and $\text{Var}(\hat{\mu}^{(1)})$ can be estimated by a sample variance estimator

$$\text{Var}(\hat{\mu}^{(1)}) \approx \hat{\sigma}_{\hat{\mu}^{(1)}}^2 = \frac{1}{M-1} \sum_{i=1}^M (\hat{\mu}^{(i)} - \hat{\mu}_{N,B}^{\text{mcmc}})^2.$$

Finally, an estimator for σ_{mcmc}^2 is

$$\hat{\sigma}_{\text{mcmc}}^2 = \frac{N}{M} \hat{\sigma}_{\hat{\mu}^{(1)}}^2 = \frac{T}{M-1} \sum_{i=1}^M \left(\hat{\mu}^{(i)} - \hat{\mu}_{N,B}^{\text{mcmc}} \right)^2.$$

Chapter 9

Sensitivities and gradient-based Stochastic Optimization

Let Z be the output of some stochastic model, which can be expressed as a parametric function $Z = \psi(\theta, X)$ of a random vector X with known probability distribution, taking values in $\mathcal{X} \subset \mathbb{R}^d$. Here, $\theta \in \Gamma \subset \mathbb{R}^n$ is a vector of parameters characterizing the model. We are interested in studying how the expected output depends on the model parameters, i.e. we focus on the (deterministic) function

$$J(\theta) = \mathbb{E}[\psi(\theta, X)], \quad \theta \in \Gamma.$$

In particular, we may be interested in computing sensitivities of J at a given point $\theta_0 \in \Gamma$, i.e. partial derivatives $\partial_{\theta_i} J(\theta_0)$, $i = 1, \dots, n$ or the full gradient $\nabla J(\theta_0)$, or solve an optimization problem

$$\text{find} \quad J^* = \min_{\theta \in \Gamma} \mathbb{E}[\psi(\theta, X)] \quad \text{or} \quad \theta^* \in \operatorname{argmin}_{\theta \in \Gamma} \mathbb{E}[\psi(\theta, X)]. \quad (9.1)$$

Optimization problems involving the minimization of an expected cost as in (9.1) are usually referred to as *stochastic optimization* or *stochastic programming* problems. We give two examples hereafter.

Example 9.1 (News vendor problem). *A company has to decide about the quantity $\theta \geq 0$ of a certain product to order to satisfy the demand X . The unit cost for the order is $c > 0$. If the demand X is larger than θ , the company makes additional orders at unit cost $b > c$. If the demand X is smaller than θ , the company incurs holding cost of h per unit. Given a demand X and order θ , the total cost incurred is therefore*

$$\begin{aligned} \psi(\theta, X) &= c\theta + b(X - \theta)_+ + h(\theta - X)_+ \\ &= \max\{(c - b)\theta + bX, (c + h)\theta - hX\}. \end{aligned}$$

Since the order has to be placed before knowing the actual demand, one may use a probabilistic forecast of the demand, i.e. treat the demand as a random variable whose probability distribution may be estimated from historical data. Then, a reasonable strategy to decide the quantity to order is to solve the stochastic optimization problem $\theta^ \in \operatorname{argmin}_{\theta \geq 0} \mathbb{E}[\psi(\theta, X)]$.*

Example 9.2 (Portfolio optimization). *Suppose we have a capital W_0 and we want to invest an amount θ_i in asset i , $i = 1, \dots, n$ and keep $\theta_0 = W_0 - \sum_{i=1}^n \theta_i$ in cash. The return rate, per unit of time, of asset i is a random variable X_i with known probability distribution. Denote $X = (X_1, \dots, X_n)$ and $\theta = (\theta_1, \dots, \theta_n)$. Then, the total wealth after one unit of time is*

$$\begin{aligned}\psi(\theta, X) &= \theta_0 + \sum_{i=1}^n (1 + X_i)\theta_i \\ &= W_0 + \sum_{i=1}^n X_i\theta_i = W_0 + X^\top \theta\end{aligned}$$

and we may try to maximise the expected return of the investment

$$\theta^* \in \operatorname{argmax}_{\theta \geq 0, |\theta|_1 \leq W_0} \mathbb{E}[\psi(\theta, X)]$$

with $|\theta|_1 = \sum_{i=1}^n \theta_i$, or, more generally, minimize a utility function $U = U(z)$ concave and non decreasing,

$$\theta^* \in \operatorname{argmax}_{\theta \geq 0, |\theta|_1 \leq W_0} \mathbb{E}[U(\psi(\theta, X))].$$

To solve these optimization problems, a common approach is to use gradient based iterative methods, provided that the cost function J is differentiable. This, however, implies estimating $\nabla J(\theta) = \nabla \mathbb{E}[\psi(\theta, X)]$ at each iteration. It is often the case that the expectation is not computable in closed form and has to be approximated, e.g. by sampling techniques.

In the next section we focus on the problem of computing sensitivities of expectations using Monte-Carlo type estimators. Then, in Section 9.2 we will discuss gradient descent type methods for stochastic optimization problems.

9.1 Computation of sensitivities

Let X be a random variable taking values in $\mathcal{X} \subset \mathbb{R}^d$ with known pdf f and $\psi : \Gamma \times \mathcal{X} \rightarrow \mathbb{R}$ be a given function. We denote $J(\theta) = \mathbb{E}[\psi(\theta, X)]$. Our goal is to approximate $\nabla J(\theta) = \nabla \mathbb{E}[\psi(\theta, X)]$ using sampling techniques. We illustrate hereafter three approaches to achieve this goal.

Sample path approximation (or infinitesimal perturbation analysis – IPA)

Suppose that $\theta \mapsto \psi(\theta, \cdot)$ is a.s. differentiable and we can interchange the gradient and expectation

$$\nabla J(\theta) = \mathbb{E}[\nabla_{\theta} \psi(\theta, X)]. \quad (9.2)$$

This is not always possible and requires some properties of the integrand function. The next Lemma gives sufficient conditions for (9.2) to be true.

Lemma 9.1. *Assume that*

- a) $\theta \mapsto \psi(\theta, x)$ is differentiable at $\theta_0 \in \overset{\circ}{\Gamma}$ for f -almost every $x \in \mathcal{X}$;
 b) there exists a open set $\mathcal{U}(\theta_0) \subset \Gamma$ containing θ_0 and a function $M_{\theta_0} : \mathcal{X} \rightarrow \mathbb{R}_+$ with $\mathbb{E}[M_{\theta_0}(X)] < \infty$ such that

$$|\psi(\theta_1, x) - \psi(\theta_2, x)| \leq M_{\theta_0}(x) \|\theta_1 - \theta_2\|, \quad \forall \theta_1, \theta_2 \in \mathcal{U}(\theta_0), \quad \forall x \in \mathcal{X};$$

- c) $x \mapsto \psi(\theta, x)$ is measurable and $\mathbb{E}[\psi(\theta, X)] < \infty$ for every $\theta \in \mathcal{U}(\theta_0)$.

Then $\theta \mapsto \mathbb{E}[\psi(\theta, X)]$ is differentiable in θ_0 and

$$\nabla \mathbb{E}[\psi(\theta, X)]|_{\theta=\theta_0} = \mathbb{E}[\nabla_{\theta} \psi(\theta, X)|_{\theta=\theta_0}].$$

Proof. Property b) implies $\|\nabla_{\theta} \psi(\theta_0, x)\| \leq M_{\theta_0}(x)$ for f -a.e. $x \in \mathcal{X}$ and $\mathbb{E}[\nabla_{\theta} \psi(\theta_0, X)]$ is well defined and finite. For any $h \in \mathbb{R}^n$ such that $\theta_0 + h \in \mathcal{U}(\theta_0)$, let $\psi_h(x) = \frac{1}{\|h\|}[\psi(\theta_0 + h, x) - \psi(\theta_0, x) - \nabla_{\theta} \psi(\theta_0, x) \cdot h]$. Thanks to assumption b), $|\psi_h(x)| \leq 2M_{\theta_0}(x)$ and, from assumption a), $\lim_{h \rightarrow 0} \psi_h(x) = 0$, f -a.e. Hence, by dominated convergence and denoting $J(\theta) = \mathbb{E}[\psi(\theta, X)]$, we have

$$\lim_{h \rightarrow 0} \frac{1}{\|h\|} [J(\theta_0 + h) - J(\theta_0) - \mathbb{E}[\nabla_{\theta} \psi(\theta_0, X)] \cdot h] = \lim_{h \rightarrow 0} \mathbb{E}[\psi_h(x)] = \mathbb{E}\left[\lim_{h \rightarrow 0} \psi_h(x)\right] = 0,$$

which shows that J is differentiable at θ_0 and $\nabla J(\theta_0) = \mathbb{E}[\nabla_{\theta} \psi(\theta_0, X)]$. \square

Thanks to (9.2) we can then approximate $\nabla J(\theta)$ by Monte Carlo as

$$\widehat{\nabla J}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \psi(\theta, X^{(i)}), \quad X^{(i)} \stackrel{\text{iid}}{\sim} X$$

Notice that this estimator is nothing but the gradient of the Monte Carlo estimator $\widehat{J}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \psi(\theta, X^{(i)})$ for J , i.e. $\widehat{\nabla J}_N = \nabla \widehat{J}_N$. In particular, it is unbiased and has variance $\text{Var}(\widehat{\nabla J}_N) = \frac{1}{N} \text{Var}(\nabla_{\theta} \psi(\theta, X))$, where we have used the notation $\text{Var}(X) = \text{Tr}(\text{Cov}(X)) = \mathbb{E}[\|X - \mathbb{E}[X]\|^2]$ for a random vector X .

The limitation of this estimator is that it is only applicable when (9.2) holds.

Example 9.3. Consider the two random variables

$$\psi_1(\theta, X) = \max\{X, \theta\}, \quad \psi_2(\theta, X) = \mathbb{1}_{\{X > \theta\}}, \quad \theta \in \mathbb{R}.$$

where $X \in \mathbb{R}$ is a random variable with continuous cumulative distribution function $F_X : \mathbb{R} \rightarrow [0, 1]$.

For ψ_1 , property (9.2) holds. Indeed at $\theta_0 \in \mathbb{R}$

$$\begin{aligned} \frac{d}{d\theta} \max\{x, \theta\} &= \mathbb{1}_{\{x < \theta\}} \quad \text{for all } x \neq \theta \\ |\max\{x, \theta_1\} - \max\{x, \theta_2\}| &\leq |\theta_1 - \theta_2| \quad \text{for all } \theta_1, \theta_2, x \in \mathbb{R} \end{aligned}$$

Hence $\frac{d}{d\theta} \mathbb{E}[\psi_1(\theta, X)] = \mathbb{E}[\mathbb{1}_{\{X < \theta\}}]$ for all θ .

For ψ_2 , instead, property (9.2) does not hold. Indeed $\frac{d}{d\theta} \mathbb{1}_{\{x > \theta\}} = 0$ for all $\theta \neq x$ and is not defined at $\theta = x$. Then, for all θ for which $0 < F_X(\theta) < 1$ we have

$$\frac{d}{d\theta} \mathbb{E}[\mathbb{1}_{\{X > \theta\}}] = \frac{d}{d\theta} (1 - F_X(\theta)) \neq \mathbb{E}\left[\frac{d}{d\theta} \mathbb{1}_{\{X > \theta\}}\right] = 0.$$

Finite difference approximation

Another approach to estimate $\nabla J(\theta)$ is to use a finite difference approximation $\partial_{\theta_j} J(\theta) \approx \frac{J(\theta + he_j) - J(\theta)}{h}$ where $J(\theta + he_j)$ and $J(\theta)$ are estimated by Monte Carlo

$$\widehat{\partial_{\theta_j} J}_N(\theta) = \frac{\widehat{J}_N(\theta + he_j) - \widehat{J}_N(\theta)}{h} = \frac{1}{N} \sum_{i=1}^N \frac{\psi(\theta + he_j, X^{(i)}) - \psi(\theta, X^{(i)})}{h}, \quad X^{(i)} \stackrel{\text{iid}}{\sim} X.$$

It is a priori not needed to use the same sample $(X^{(1)}, \dots, X^{(N)})$ in the two Monte Carlo estimators for $J(\theta)$ and $J(\theta + he_j)$, however, it is usually advantageous to do so.

This estimator is biased since

$$\text{Bias} = \left| \mathbb{E} \left[\widehat{\partial_{\theta_j} J}_N(\theta) \right] - \partial_{\theta_j} J(\theta) \right| = \left| \frac{J(\theta + he_j) - J(\theta)}{h} - \partial_{\theta_j} J(\theta) \right| \neq 0.$$

and $\text{Bias} = O(h)$ if $\theta \mapsto J(\theta)$ is twice differentiable. On the other hand, it is applicable also to cases where (9.2) does not hold. The MSE of the estimator reads:

$$\text{MSE}(\widehat{\partial_{\theta_j} J}_N) = \frac{1}{Nh^2} [\text{Var}(\psi(\theta + he_j, X)) - \text{Var}(\psi(\theta, X))] + \text{Bias}^2.$$

If (9.2) holds with square integrable partial derivative $\partial_{\theta_j} \psi(\theta, X)$, then $\text{Var}(\psi(\theta + he_j, X)) - \text{Var}(\psi(\theta, X)) = O(h^2)$ and

$$\text{MSE}(\widehat{\partial_{\theta_j} J}_N) = O\left(\frac{1}{N}\right) + O(h^2).$$

Choosing $h \propto \frac{1}{\sqrt{N}}$ we obtain $\text{MSE} = O\left(\frac{1}{N}\right)$ with a cost of $2N$ evaluations of ψ , hence an analogous performance as that of a Monte Carlo estimator for $J(\theta)$.

If, instead, $\theta \mapsto \psi(\theta, \cdot)$ is only α -Hölder continuous $|\psi(\theta + he_j, x) - \psi(\theta, x)| \leq M_\theta(x)|h|^\alpha$ with $\mathbb{E}[M_\theta(X)^2] < \infty$, but not differentiable (still assuming J twice differentiable), then $\text{Var}(\psi(\theta + he_j, X)) - \text{Var}(\psi(\theta, X)) = O(h^{2\alpha})$ and

$$\text{MSE}(\widehat{\partial_{\theta_j} J}_N) = O\left(\frac{1}{Nh^{2-2\alpha}}\right) + O(h^2).$$

Equilibrating the two error terms leads to $h \propto N^{-\frac{1}{4-2\alpha}}$ and a mean squared error $\text{MSE} = O(N^{-\frac{1}{2-\alpha}})$, hence a reduced performance with respect to a standard Monte Carlo estimator. In the extreme case $\alpha = 0$ (discontinuous integrand), we have $\text{MSE} = O\left(\frac{1}{\sqrt{N}}\right)$, which is half the rate of a standard Monte Carlo estimator.

These rates can be improved by considering a higher order finite difference approximation, e.g. a centered one

$$\widehat{\partial_{\theta_j}^c J}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{\psi(\theta + he_j, X^{(i)}) - \psi(\theta - he_j, X^{(i)})}{2h}, \quad X^{(i)} \stackrel{\text{iid}}{\sim} X,$$

in which case one has $\text{Bias} = O(h^2)$ assuming J three times differentiable and $\text{MSE} = O(N^{-\frac{2}{3-\alpha}})$ or, for $\alpha = 0$, $\text{MSE} = O(N^{-\frac{2}{3}})$.

Likelihood ratio estimator (or score function method)

Suppose one can make a change of variable $\tilde{X} = g_\theta(X)$ so that $\psi(\theta, X) = \tilde{\psi}(\tilde{X})$ does not depend on θ anymore, when expressed in the variable \tilde{X} . Let f_θ denote the pdf of \tilde{X} (assuming it exists). Then

$$J(\theta) = \mathbb{E}[\psi(\theta, X)] = \mathbb{E}_{\tilde{X} \sim f_\theta}[\tilde{\psi}(\tilde{X})] = \int \tilde{\psi}(x) f_\theta(x) dx.$$

If we can exchange the gradient with the integral, we can write

$$\nabla J(\theta) = \int \tilde{\psi}(x) \nabla_\theta f_\theta(x) dx = \int \tilde{\psi}(x) (\nabla_\theta \log f_\theta(x)) f_\theta(x) dx = \mathbb{E}_{\tilde{X} \sim f_\theta}[\tilde{\psi}(\tilde{X}) S_\theta(\tilde{X})] \quad (9.3)$$

where $S_\theta(x) = \nabla_\theta \log f_\theta(x)$ is called the *score function*. We can now use a standard Monte Carlo estimator to approximate the last expectation

$$\widehat{\nabla J}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \tilde{\psi}(\tilde{X}^{(i)}) S_\theta(\tilde{X}^{(i)}), \quad \tilde{X}^{(i)} \stackrel{\text{iid}}{\sim} \tilde{X}.$$

The estimator is unbiased and the approach may work also in the case of a non-differentiable integrand $\theta \mapsto \psi(\theta, \cdot)$, if one can rewrite the expectation as in (9.3) with a parametric probability density function f_θ that depends smoothly on θ .

Example 9.4. We consider again the example $\psi_2(\theta, x) = \mathbf{1}_{\{x > \theta\}}$ and the goal of computing $\nabla \mathbb{E}[\psi_2(\theta, X)]$ where X is a random variable with pdf f . Defining the new random variable $\tilde{X} = X - \theta$, which has probability density function $f_\theta(x) = f(x + \theta)$ we have

$$\mathbb{E}[\psi_2(\theta, X)] = \int \mathbf{1}_{\{x > \theta\}} f(x) dx = \int \mathbf{1}_{\{\tilde{x} > 0\}} f(\tilde{x} + \theta) d\tilde{x}, \quad \tilde{X} \sim f_\theta(\cdot) = f(\cdot + \theta)$$

and

$$\frac{d}{d\theta} \mathbb{E}[\psi_2(\theta, X)] = \mathbb{E}_{\tilde{X} \sim f_\theta} \left[\mathbf{1}_{\{\tilde{X} > 0\}} \frac{f'(\tilde{X} + \theta)}{f(\tilde{X} + \theta)} \right].$$

9.2 Stochastic optimization

We focus now on the stochastic optimization problem:

$$\text{find } \theta^* \in \underset{\theta \in \Gamma}{\text{argmin}} J(\theta) := \mathbb{E}[\psi(\theta, X)]. \quad (9.4)$$

We assume hereafter that $\theta \mapsto J(\theta)$ is differentiable, Γ is closed and convex, and the set $\Sigma = \{\theta \in \Gamma : J(\theta) \leq J(\theta_0)\}$ is compact so that at least one solution of (9.4) exists and is in Σ .

A common approach to find a solution of (9.4) is given by the projected gradient descent method. In its simplest form it reads

$$\theta_{j+1} = \Pi_\Gamma(\theta_j - \tau_j \nabla J(\theta_j)), \quad j = 0, 1, \dots \quad (9.5)$$

where τ_j is the step size, which may depend on j , and Π_Γ is the projection on the admissible set Γ :

$$\forall \theta \in \mathbb{R}^n, \quad \Pi_\Gamma(\theta) = \underset{y \in \Gamma}{\operatorname{argmin}} \|\theta - y\|,$$

which is well defined when Γ is convex and closed. The projector Π_Γ is also contractive

$$\|\Pi_\Gamma(\theta_1) - \Pi_\Gamma(\theta_2)\| \leq \|\theta_1 - \theta_2\| \quad \forall \theta_1, \theta_2 \in \mathbb{R}^n.$$

However, $J(\theta_i)$ might not be accessible as it implies computing an expectation. A natural idea is then to replace it by a suitable estimator $\widehat{\nabla} J_{N_j}(\theta_j)$, as one of those discussed in the previous section, based on an iid sample of size N_j drawn from X . This leads to an approximate gradient descent algorithm

$$\theta_{j+1} = \Pi_\Gamma \left(\theta_j - \tau_j \widehat{\nabla} J_{N_j}(\theta_j) \right), \quad j = 0, 1, \dots \quad (9.6)$$

For instance, if one can exchange the gradient with the expectation as in (9.2), then an unbiased gradient estimator is $\widehat{\nabla} J_{N_j}(\theta_j) = \frac{1}{N} \sum_{i=1}^{N_j} \nabla_{\theta} \psi(\theta_j, X^{(i)})$ with $X^{(i)} \stackrel{\text{iid}}{\sim} X$.

Several options are available

- draw a large sample $\{X^{(i)}\}_{i=1}^N$ of size N once and for all and use it for all iterations of the approximate gradient descent algorithm (9.6). This approach is referred to as *batch gradient* and is related to the so called *Sample Average Approximation* (SAA). When a IPA estimator is used, it corresponds to applying a gradient descent algorithm to the discretized optimization problem

$$\operatorname{fing} \quad \hat{\theta}_N^* \in \underset{\theta \in \Gamma}{\operatorname{argmin}} \widehat{J}_N(\theta), \quad \widehat{J}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \psi(\theta, X^{(i)}).$$

The obtained minimum and minimizer will be affected by a Monte Carlo error due to the finite sample size. Under rather mild conditions on ψ one can show that $\operatorname{dist}(\hat{\theta}_N^*, S^*) \rightarrow 0$ a.s. as $N \rightarrow \infty$, where $S^* = \operatorname{argmin}_{\theta \in \Gamma} J(\theta)$ denotes the set of minimizers of J in Γ .

- Draw independently a new small sample $\{X^{(i,j)}\}_{i=1}^N$ of size N (mini-batch) at each iteration. The variance of the estimator will not be small, in general, and convergence is achieved by shrinking the step size $\tau_j \rightarrow 0$ at a suitable rate. This approach is referred to as *Stochastic Approximation* (SA) and *Stochastic Gradient Descent* (SGD).
- Draw independently a new sample $\{X^{(i,j)}\}_{i=1}^{N_j}$ at each iteration with an increasing sample size over the iterations, $N_j \rightarrow \infty$. The choice of the sequence $\{N_j\}_j$ can be made a priori or adaptively based on the estimated size of the gradient. In the limit, the variance of the estimator goes to zero and this approach recovers full gradient iterations asymptotically. This approach can be equally referred to as *Stochastic Gradient Descent of Inexact Gradient Descent*.

In the next sections we analyze and compare the last two approaches related to Stochastic Gradient Descent, in the case of a strongly convex stochastic optimization problem.

9.2.1 Stochastic Approximation and SGD

We work under the following assumptions:

Assumption 9.1.

- a) $\theta \mapsto J(\theta)$ is continuously differentiable in Γ ;
- b) ∇J is Lipschitz continuous with constant L

$$\|\nabla J(\theta_1) - \nabla J(\theta_2)\| \leq L\|\theta_1 - \theta_2\| \quad \forall \theta_1, \theta_2 \in \Gamma;$$

- c) J is strongly convex

$$(\theta_1 - \theta_2)^T (\nabla J(\theta_1) - \nabla J(\theta_2)) \geq c\|\theta_1 - \theta_2\|^2 \quad \forall \theta_1, \theta_2 \in \Gamma;$$

- d) Γ is convex and closed.

Under the above assumptions, the stochastic optimization problem (9.4) has a unique minimizer θ^* , which is characterized by the variational inequality

$$(\theta - \theta^*)^T \nabla J(\theta^*) \geq 0 \quad \forall \theta \in \Gamma. \quad (9.7)$$

This condition can be equivalently written as

$$\theta^* = \Pi_\Gamma(\theta^* - \tau \nabla J(\theta^*)) \quad \forall \tau > 0. \quad (9.8)$$

Moreover, the full gradient descent method (9.5) with fixed step size converges geometrically in the number of iterations.

Lemma 9.2. *Consider the iterates $\{\theta_j\}_j$ given by (9.5) with $\tau_j = \tau$, $\forall j$. Under Assumption 9.1, it holds for any $j \geq 0$*

$$\|\theta_j - \theta^*\|^2 \leq \rho_\tau^j \|\theta_0 - \theta^*\|^2, \quad \text{with } \rho_\tau = 1 - 2c\tau + L^2\tau^2.$$

In particular, $1 - \frac{c}{L^2} \leq \rho_\tau < 1$ whenever $0 < \tau < \frac{2c}{L^2}$.

Proof. We have

$$\begin{aligned} \|\theta_{j+1} - \theta^*\|^2 &= \|\Pi_\Gamma(\theta_j - \tau \nabla J(\theta_j)) - \Pi_\Gamma(\theta^* - \tau \nabla J(\theta^*))\|^2 \\ &\leq \|\theta_j - \theta^* - \tau(\nabla J(\theta_j) - \nabla J(\theta^*))\|^2 \quad (\text{by contractivity of } \Pi_\Gamma) \\ &= \|\theta_j - \theta^*\|^2 - 2\tau \underbrace{(\theta_j - \theta^*)^T (\nabla J(\theta_j) - \nabla J(\theta^*))}_{\geq c\|\theta_j - \theta^*\|^2 \text{ by strong conv.}} + \tau^2 \underbrace{\|\nabla J(\theta_j) - \nabla J(\theta^*)\|^2}_{\leq L^2\|\theta_j - \theta^*\|^2 \text{ by Lipsch. gradient}} \\ &\leq \rho_\tau \|\theta_j - \theta^*\|^2 \leq \rho_\tau^{j+1} \|\theta_0 - \theta^*\|^2. \end{aligned}$$

Being ρ_τ quadratic in τ , a direct calculation shows that $0 < \tau < \frac{2c}{L^2}$ implies $1 - \frac{c}{L^2} \leq \rho_\tau < 1$. \square

The full gradient descent method, however, is generally not directly implementable since $\nabla J(\theta)$ is not accessible because of the expectation. We focus then on the approximate gradient descent algorithm (9.6) and we make the following assumption on the gradient estimator $\widehat{\nabla J}_N(\theta)$:

Assumption 9.2.

- a) $\widehat{\nabla J}_N(\theta)$ is unbiased for any $\theta \in \Gamma$;
- b) the (conditional) variance grows at most quadratically in θ : there exist $V, M > 0$ s.t.

$$\text{Var}(\widehat{\nabla J}_N(\theta) \mid \theta) := \mathbb{E} \left[\|\widehat{\nabla J}_N(\theta) - \nabla J(\theta)\|^2 \mid \theta \right] \leq \frac{1}{N} (V + M \|\theta - \theta^*\|^2) \quad \forall \theta \in \Gamma.$$

The assumption that $\widehat{\nabla J}_N(\theta)$ is unbiased is very convenient. It can be relaxed, but the bias will have to decrease over the iterations to be able to converge to the exact solution. The considered Stochastic Gradient Descent (SGD) method is summarized in Algorithm 9.1

Algorithm 9.1: Stochastic Gradient Descent.

Given: $\theta_0 \in \Gamma$

- 1 **for** $j = 0, 1, \dots$, **do**
- 2 draw $X^{(1,j)}, \dots, X^{(N_j,j)} \stackrel{\text{iid}}{\sim} X$, independently of the previous iterations
- 3 estimate the gradient $\widehat{\nabla J}_{N_j}(\theta_j)$ based on the sample $\{X^{(1,j)}, \dots, X^{(N_j,j)}\}$
- 4 update $\theta_{j+1} = \Pi_\Gamma \left(\theta_j - \tau_j \widehat{\nabla J}_{N_j}(\theta_j) \right)$
- 5 **end**

In general, both τ_j and N_j may depend on the iteration. Under Assumptions 9.1 and 9.2 the following recursion on the mean square error $e_j := \mathbb{E} [\|\theta_j - \theta^*\|^2]$ can be derived.

Lemma 9.3. *Under Assumptions 9.1 and 9.2, the iterations $\{\theta_j\}_{j=1}^\infty$ produced by the SGD Algorithm 9.1 satisfy*

$$\mathbb{E} [\|\theta_{j+1} - \theta^*\|^2] \leq (1 - 2c\tau_j + L^2\tau_j^2) \mathbb{E} [\|\theta_j - \theta^*\|^2] + \tau_j^2 V_j \quad (9.9)$$

with $V_j = \mathbb{E} [\text{Var}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j)]$

Proof. We have

$$\begin{aligned} \|\theta_{j+1} - \theta^*\|^2 &= \|\Pi_\Gamma \left(\theta_j - \tau_j \widehat{\nabla J}_{N_j}(\theta_j) \right) - \Pi_\Gamma \left(\theta^* - \tau_j \nabla J(\theta^*) \right)\|^2 \\ &\leq \|\theta_j - \theta^* - \tau_j \left(\widehat{\nabla J}_{N_j}(\theta_j) - \nabla J(\theta^*) \right)\|^2 \\ &= \underbrace{\|\theta_j - \theta^*\|^2}_{(A)} - 2\tau_j \underbrace{(\theta_j - \theta^*)^T \left(\widehat{\nabla J}_{N_j}(\theta_j) - \nabla J(\theta^*) \right)}_{(B)} + \tau_j^2 \underbrace{\|\widehat{\nabla J}_{N_j}(\theta_j) - \nabla J(\theta^*)\|^2}_{(C)}. \end{aligned}$$

Let us denote by $\mathcal{F}_j = \sigma\{X^{(i,\ell)}, i = 1, \dots, N_\ell, \ell = 1, \dots, j-1\}$ the σ -algebra generated by all the random variables $\{X^{(i,j)}\}$ drawn up to iteration $j-1$. Notice that θ_j is fully determined by $\{X^{(i,\ell)}, i = 1, \dots, N_\ell, \ell = 1, \dots, j-1\}$, hence, for any measurable function $g : \Gamma \rightarrow \mathbb{R}$ we have $\mathbb{E}[g(\theta_j) \mid \mathcal{F}_j] = g(\theta_j)$. We take now the conditional expectation with respect to \mathcal{F}_j in the previous inequality and notice that

- $\mathbb{E}[(A) \mid \mathcal{F}_j] = \|\theta_j - \theta^*\|^2$
- $\mathbb{E}[(B) \mid \mathcal{F}_j] = (\theta_j - \theta^*)^T \mathbb{E} \left[\widehat{\nabla J}_{N_j}(\theta_j) - \nabla J(\theta^*) \mid \mathcal{F}_j \right]$
 $= (\theta_j - \theta^*)^T (\nabla J(\theta_j) - \nabla J(\theta^*))$ (since $\widehat{\nabla J}_{N_j}(\theta_j)$ is unbiased)
 $\geq c\|\theta_j - \theta^*\|^2$ (by strong convexity)
- $\mathbb{E}[(C) \mid \mathcal{F}_j] = \mathbb{E} \left[\underbrace{\|\widehat{\nabla J}_{N_j}(\theta_j) - \nabla J(\theta_j) + \nabla J(\theta_j) - \nabla J(\theta^*)\|^2}_{\text{zero conditional mean}} \mid \mathcal{F}_j \right]$
 $= \mathbb{E} \left[\|\widehat{\nabla J}_{N_j}(\theta_j) - \nabla J(\theta_j)\|^2 \mid \mathcal{F}_j \right] + \tau_j^2 \|\nabla J(\theta_j) - \nabla J(\theta^*)\|^2$
 $= \text{Var}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j) + \tau_j^2 L^2 \|\theta_j - \theta^*\|^2$

Putting all these estimates together, we obtain

$$\mathbb{E} [\|\theta_{j+1} - \theta^*\|^2 \mid \mathcal{F}_j] \leq (1 - 2c\tau_j + L^2\tau_j^2)\|\theta_j - \theta^*\|^2 + \tau_j^2 \text{Var}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j)$$

and, taking a further expectation, we obtain the final estimate (9.9). \square

Thanks to Assumption 9.2, the recursion (9.9) can be further bounded as

$$e_{j+1} \leq \rho_j e_j + \frac{\tau_j^2}{N_j} V_j, \quad \text{with } \rho_j = 1 - 2c\tau_j + \tau_j^2(L^2 + \frac{M}{N_j}). \quad (9.10)$$

Notice that $\rho_j < 1$ if one takes τ_j small enough, so the initial error is reduced at each iteration. On the other hand, the algorithm introduces an error $\frac{\tau_j^2}{N_j} V_j$ at each iteration due to the approximate gradient estimate. To recover to the correct solution, such error has to vanish in the limit $j \rightarrow \infty$, which can be achieved by reducing the step size τ_j and/or by increasing the sample size N_j .

9.2.2 SGD with fixed step size and increasing sample size

We analyze first the case in which the step size in Algorithm 9.1 is kept fixed, i.e. $\tau_j = \tau, \forall j$ whereas the sample size increases geometrically as $N_j = N_0 \xi^{-j}$ for some $0 < \xi < 1$. The next lemma shows that if τ and ξ are properly chosen, the SGD algorithm recovers a geometric convergence as the ideal full gradient descent method.

Lemma 9.4. *Consider the iterates $\{\theta_j\}_j$ generated by Algorithm 9.1 with $\tau_j = \tau \forall j$ with $\tau \in (0, \frac{2c}{L^2+M})$ and $N_j = N_0 \xi^{-j}$ with $\xi \in (0, 1)$. Under Assumptions 9.1 and 9.2 it holds for any $j \geq 0$*

$$\mathbb{E} [\|\theta_j - \theta^*\|^2] \leq C \bar{\rho}^j \quad \text{for some } C > 0.$$

with $\bar{\rho} = \max\{\rho_\tau, \xi\}$ and $\rho_\tau = 1 - 2c\tau + \tau^2(L^2 + M) \in [\frac{c}{L^2+M}, 1)$.

Proof. We start from the error recursion (9.10), using $\rho_j \leq \rho_\tau \forall j$, which leads to the estimate

$$\begin{aligned} \mathbb{E} [\|\theta_j - \theta^*\|^2] &\leq \rho_\tau^j \|\theta_0 - \theta^*\|^2 + \tau^2 \sum_{\ell=0}^{j-1} \rho_\tau^{j-1-\ell} \frac{V}{N_\ell} \\ &\leq \rho_\tau^j \|\theta_0 - \theta^*\|^2 + \tau^2 \frac{\rho_\tau^{j-1} V}{N_0} \sum_{\ell=0}^{j-1} \left(\frac{\xi}{\rho_\tau} \right)^\ell \\ &\leq \rho_\tau^j \|\theta_0 - \theta^*\|^2 + \tau^2 \frac{V}{N_0} \rho_\tau^{j-1} \frac{(\xi/\rho_\tau)^j - 1}{\xi/\rho_\tau - 1} \end{aligned}$$

We analyze now the last term in the two cases $\xi > \rho_\tau$ and $\xi < \rho_\tau$:

$$\begin{aligned} \xi > \rho_\tau & \quad \rho_\tau^{j-1} \frac{(\xi/\rho_\tau)^j - 1}{\xi/\rho_\tau - 1} \leq \rho_\tau^{j-1} \frac{(\xi/\rho_\tau)^j}{\xi/\rho_\tau - 1} = \frac{\xi^j}{\xi - \rho_\tau} = \frac{\bar{\rho}^j}{|\xi - \rho_\tau|}, \\ \xi < \rho_\tau & \quad \rho_\tau^{j-1} \frac{(\xi/\rho_\tau)^j - 1}{\xi/\rho_\tau - 1} \leq \rho_\tau^{j-1} \frac{1}{1 - \xi/\rho_\tau} = \frac{\rho_\tau^j}{\rho_\tau - \xi} = \frac{\bar{\rho}^j}{|\xi - \rho_\tau|}. \end{aligned}$$

Hence

$$\mathbb{E} [\|\theta_j - \theta^*\|^2] \leq \rho_\tau^j \|\theta_0 - \theta^*\|^2 + \frac{\tau^2 V}{N_0 |\xi - \rho_\tau|} \bar{\rho}^j \leq C \bar{\rho}^j$$

with $C = \|\theta_0 - \theta^*\|^2 + \frac{\tau^2 V}{N_0 |\xi - \rho_\tau|}$. \square

It is worth investigating the computational cost of the algorithm to achieve a mean square error tol^2 . From the previous lemma, we should stop the iterations at

$$j = j(tol) : \quad C \bar{\rho}^{j(tol)} = tol^2 \quad \implies \quad j(tol) = \frac{\log(Ctol^{-2})}{\log(\bar{\rho}^{-1})}.$$

On the other hand, the overall cost of the algorithm can be estimated as the sum of sample sizes used at each iteration, i.e.

$$\text{Cost} = \sum_{\ell=0}^{j(tol)} N_\ell \lesssim \sum_{\ell=0}^{j(tol)} \xi^{-\ell} \lesssim \xi^{-j(tol)} = e^{j(tol) \log \xi^{-1}} \lesssim tol^{-2 \frac{\log \xi^{-1}}{\log \bar{\rho}^{-1}}}$$

where we have used the notation $A \lesssim B$ to indicate that there exists a constant $K > 0$ such that $A \leq KB$. This calculation shows that if we choose $\xi > \rho_\tau$, then $\text{Cost} \lesssim tol^{-2}$ which is the standard complexity of a Monte Carlo estimator. The case $\xi < \rho_\tau$ leads to a higher complexity and is therefore not recommended.

Alternatively to the a-priori choice $N_j = N_0 \xi^{-j}$ which we have discussed above, one may choose N_j adaptively, still in such a way that $N_j \rightarrow \infty$ as $j \rightarrow \infty$.

Let us define the reduced gradient

$$R_\tau(\theta) = \frac{1}{\tau} (\theta - \Pi_\Gamma(\theta - \tau \nabla J(\theta))).$$

Notice that $R_\tau(\theta^*) = 0 \forall \tau$ and, in the unconstrained case, $R_\tau(\theta) = \nabla J(\theta)$. It is then sound to relate the variance of the gradient estimator with the size of the reduced gradient,

which is hopefully going to zero if the algorithm converges, thus driving also the variance of the gradient estimator to zero, and the sample size to infinity. We consider the following (ideal) adaptive sampling strategy:

$$\text{choose } N_j \text{ such that } \quad \text{Var}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j) \leq \eta \|R_\tau(\theta_j)\|^2 \quad (9.11)$$

In practice, neither $\text{Var}(\widehat{\nabla J}_{N_j}(\theta_j \mid \theta_j))$ nor $R_\tau(\theta_j)$ can be computed exactly, as they involve full expectations, however, they can be replaced by suitable estimators $\widehat{\text{Var}}(\widehat{\nabla J}_{N_j}(\theta_j \mid \theta_j))$ and $\widehat{R}_{\tau, N_j}(\theta_j) = \frac{1}{\tau} \left(\theta - \Pi_\Gamma(\theta - \tau \widehat{\nabla J}_{N_j}(\theta_j)) \right)$. A practical adaptive strategy is then given by:

$$\text{choose } N_j \text{ such that } \quad \widehat{\text{Var}}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j) \leq \eta \|\widehat{R}_\tau(\theta_j)\|^2 \quad (9.12)$$

The resulting adaptive SGD is illustrated in Algorithm 9.2.

Algorithm 9.2: Adaptive Stochastic Gradient Descent.

Given: $\theta_0 \in \Gamma$, $N_0 \in \mathbb{N}$, $\eta \in (0, 1)$, $\tau > 0$

- 1 **for** $j = 0, 1, \dots$, **do**
- 2 draw $X^{(1,j)}, \dots, X^{(N_j,j)} \stackrel{\text{iid}}{\sim} X$, independently of the previous iterations
- 3 estimate the gradient $\widehat{\nabla J}_{N_j}(\theta_j)$, the reduced gradient $\widehat{R}_{\tau, N_j}(\theta_j)$ and variance of the gradient estimator $\widehat{\text{Var}}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j)$ based on the drawn sample
- 4 **if** (9.12) *not satisfied* **then**
- 5 | increase N_j and go back to 2
- 6 **end**
- 7 update $\theta_{j+1} = \Pi_\Gamma \left(\theta_j - \tau \widehat{\nabla J}_{N_j}(\theta_j) \right)$
- 8 **end**

We analyze hereafter the convergence of this algorithm, assuming that the idealized adaptive sampling strategy (9.11) is used. The presentation follows closely [1]. For this, we need first a technical result that can be found e.g. in [7, Chapter 2]

Lemma 9.5. *Under Assumption 9.1 it holds*

$$R_\tau(\theta) \leq \frac{1}{\tau} (1 + \sqrt{1 - c\tau}) \|\theta - \theta^*\|^2, \quad \forall \theta \in \Gamma, \quad \forall \tau \in (0, \frac{1}{L}).$$

Proof. For $\theta \in \Gamma$, let $\hat{\theta} = \Pi_\Gamma(\theta - \tau \nabla J(\theta))$. Notice that $R_\tau(\theta) = \frac{\theta - \hat{\theta}}{\tau}$ and

$$(z - \hat{\theta})^T (\theta - \tau \nabla J(\theta) - \hat{\theta}) \leq 0 \quad \forall z \in \Gamma.$$

Since J is strongly convex, it holds

$$\begin{aligned}
J(\theta^*) &\geq J(\theta) + \nabla J(\theta) \cdot (\theta^* - \theta) + \frac{c}{2} \|\theta^* - \theta\|^2 \\
&\geq J(\theta) + \nabla J(\theta) \cdot (\hat{\theta} - \theta) + \underbrace{\left(\frac{\hat{\theta} - \theta}{\tau} + \nabla J(\theta) \right) \cdot (\theta^* - \hat{\theta}) - \frac{(\hat{\theta} - \theta)^T}{\tau} (\theta^* - \hat{\theta})}_{\geq 0} + \frac{c}{2} \|\theta^* - \theta\|^2 \\
&\geq J(\theta) + \nabla J(\theta) \cdot (\hat{\theta} - \theta) + \underbrace{\frac{1}{\tau} \|\hat{\theta} - \theta\|^2}_{\geq (\frac{1}{2} + \frac{1}{2\tau}) \|\hat{\theta} - \theta\|^2} - \frac{(\hat{\theta} - \theta)^T}{\tau} (\theta^* - \theta) + \frac{c}{2} \|\theta^* - \theta\|^2 \\
&\geq \underbrace{J(\theta) + \nabla J(\theta) \cdot (\hat{\theta} - \theta) + \frac{L}{2} \|\hat{\theta} - \theta\|^2}_{\geq J(\hat{\theta}) \geq J(\theta^*)} + \frac{1}{2\tau} \|\hat{\theta} - \theta\|^2 - \frac{(\hat{\theta} - \theta)^T}{\tau} (\theta^* - \theta) + \frac{c}{2} \|\theta^* - \theta\|^2 \\
&\geq J(\theta^*) + \frac{\tau}{2} \|R_\tau(\theta)\|^2 + R_\tau(\theta)^T (\theta^* - \theta) + \frac{c}{2} \|\theta^* - \theta\|^2.
\end{aligned}$$

It follows that

$$\frac{\tau}{2} \|R_\tau(\theta)\|^2 + \frac{c}{2} \|\theta^* - \theta\|^2 \leq R_\tau(\theta)^T (\theta - \theta^*) \leq \|R_\tau(\theta)\| \|\theta - \theta^*\|,$$

which implies

$$\tau \|R_\tau(\theta)\|^2 - 2 \|R_\tau(\theta)\| \|\theta - \theta^*\| + c \|\theta - \theta^*\|^2 \leq 0.$$

Since $\tau < \frac{1}{L} \leq \frac{1}{c}$, this implies $\frac{\|R_\tau(\theta)\|}{\|\theta - \theta^*\|} \in \left(\frac{1 - \sqrt{1 - c\tau}}{\tau}, \frac{1 + \sqrt{1 - c\tau}}{\tau} \right)$, hence the thesis. \square

With this lemma at hand, we can prove that the ideal adaptive SGD algorithm converges geometrically.

Lemma 9.6. *Consider the iterates $\{\theta_j\}_j$ generated by Algorithm 9.2, with the ideal adaptive strategy (9.11) instead of the practical one (9.12). Under Assumptions 9.1 and 9.2, and choosing $\tau < \frac{1}{L}$ and $\eta < \frac{2c/\tau - L^2}{(1 + \sqrt{1 - c\tau})^2}$, it holds*

$$\mathbb{E} [\|\theta_j - \theta^*\|^2] \leq \bar{\rho}^j \|\theta_0 - \theta^*\|^2, \quad \text{with } \bar{\rho} = (1 - 2c\tau + L^2\tau^2 + \tau^2\eta(1 + \sqrt{1 - c\tau})^2) \in (0, 1).$$

Proof. Starting from the recursion (??) we have

$$\begin{aligned}
\mathbb{E} [\|\theta_{j+1} - \theta^*\|^2 \mid \mathcal{F}_j] &\leq (1 - 2c\tau + L^2\tau^2) \|\theta_j - \theta^*\|^2 + \tau^2 \text{Var}(\widehat{\nabla J}_{N_j}(\theta_j) \mid \theta_j) \\
&\leq (1 - 2c\tau + L^2\tau^2) \|\theta_j - \theta^*\|^2 + \tau^2 \eta \|R_\tau(\theta_j)\|^2 \leq \bar{\rho} \|\theta_j - \theta^*\|^2.
\end{aligned}$$

Taking a further expectation leads to the thesis:

$$\mathbb{E} [\|\theta_{j+1} - \theta^*\|^2] \leq \bar{\rho} \mathbb{E} [\|\theta_j - \theta^*\|^2] \leq \bar{\rho}^{j+1} \|\theta_0 - \theta^*\|^2.$$

The assumptions $\tau < \frac{1}{L}$ and $\eta < \frac{2c/\tau - L^2}{(1 + \sqrt{1 - c\tau})^2}$ imply $\bar{\rho} \in (0, 1)$. \square

Using Assumption 9.2, the adaptive sampling condition is satisfied if

$$\frac{1}{N_j}(V + M\|\theta_j - \theta^*\|^2) \leq \eta\|\mathcal{R}_\tau(\theta_j)\|^2 \leq \eta(1 + \sqrt{1 - c\tau})^2\|\theta_j - \theta^*\|^2$$

which implies $N_j \geq C_1 + C_2\|\theta_j - \theta^*\|^{-2}$ with suitable constants $C_1, C_2 > 0$. The previous lemma shows that $\mathbb{E}[\|\theta_j - \theta^*\|^2] = O(\bar{\rho}^j)$. Assuming that this rate holds pathwise and not just in expectation, we have $N_j = O(\bar{\rho}^{-j})$ and the total cost of the adaptive SGD algorithm to achieve a mean square error of order tol^2 again scales as $O(\text{tol}^{-2})$ as in standard Monte Carlo estimation. This argument is not rigorous as Lemma 9.6 only provides a result in expectation. To make it rigorous, one would have to prove a similar result in high probability.

9.2.3 SGD with fixed sample size and decreasing step size

We now turn to the case in which the sample size is kept fixed, i.e. $N_j = N \forall j$, whereas the step size is decreased over the iterations and $\tau_j \rightarrow 0$ as $j \rightarrow \infty$. The following lemma gives a sufficient condition for convergence.

Lemma 9.7. *Consider the iterates $\{\theta_j\}_j$ generated by Algorithm 9.1 with $N_j = N \forall j$ and a sequence $\{\tau_j\}_j$ of step sizes satisfying*

$$\sum_{j=0}^{\infty} \tau_j = \infty \quad \text{and} \quad \sum_{j=0}^{\infty} \tau_j^2 < \infty.$$

Then, under Assumptions 9.1 and 9.2 it holds $\lim_{j \rightarrow \infty} \mathbb{E}[\|\theta_j - \theta^\|^2] = 0$.*

We consider now the particular choice $\tau_j = \frac{\tau_0}{j+k}$ with $\tau_0, k > 0$, which leads to the classic *Robbins-Monro* algorithm. The following result can be proved.

Lemma 9.8. *Consider the iterates $\{\theta_j\}_j$ generated by Algorithm 9.1 with $N_j = N \forall j$ and $\tau_j = \frac{\tau_0}{j+k}$ with $\tau_0 > \frac{1}{2c}$. Under Assumptions 9.1 and 9.2 it holds*

$$\mathbb{E}[\|\theta_j - \theta^*\|^2] \leq \frac{C}{j} \left(\|\theta_0 - \theta^*\|^2 + \frac{V}{N} \right)$$

for a suitable constant $C > 0$.

To achieve a mean square error of order tol^2 we should stop the iterations at $j = j(\text{tol}) = O(\text{tol}^{-2})$. Since in this case the cost per iteration is constant (generation of N random variables), we conclude that the total cost of the adaptive SGD algorithm to achieve a mean square error of order tol^2 scales as $O(\text{tol}^{-2})$ as in the version of SGD discussed in the previous section.

The restriction $\tau_0 \geq \frac{1}{2c}$ may be difficult to enforce. A more robust algorithm is provided by the so called *Polyak-Ruppert* averaging technique which consists in returning the value $\hat{\theta}_j = \frac{1}{j} \sum_{\ell=0}^{j-1} \theta_\ell$ instead of θ_j itself at the j -th iteration. It can be shown that, in the strongly convex case, $\mathbb{E}[\|\hat{\theta}_j - \theta^*\|^2] = O(j^{-1})$ for any choice $\tau_j = \tau_0 j^{-\gamma}$, $\gamma \in (0, 1)$, even when τ_0 does not satisfy the condition $\tau_0 > \frac{1}{2c}$.

Bibliography

- [1] F. Beiser, B. Keith, S. Urbainczyk, and B. Wohlmuth. Adaptive sampling strategies for risk-averse stochastic optimization with constraints. *IMA Journal of Numerical Analysis*, 43(6):3729–3765, 2023.
- [2] Y. S. Chow and Herbert Robbins. On the asymptotic theory of fixed-width sequential confidence intervals for the mean. *Ann. Math. Statist.*, 36:457–462, 1965.
- [3] D. Kroese, T. Taimre, and Z. Botev. *Handbook of Monte Carlo Methods*. Wiley, 2011.
- [4] P. L’Ecuyer and R. Simard. TestU01: A C library for empirical testing of random number generators. *ACM Transactions on Mathematical Software*, 33(Article 22), 2007.
- [5] Wei-Liem Loh. On Latin hypercube sampling. *Ann. Statist.*, 24(5):2058–2080, 1996.
- [6] S. Meyn and R.L. Tweedie. *Markov chains and stochastic stability*. Cambridge University Press, Cambridge, second edition, 2009. With a prologue by Peter W. Glynn.
- [7] Yurii Nesterov. *Lectures on convex optimization*, volume 137 of *Springer Optimization and Its Applications*. Springer, Cham, second edition, 2018.
- [8] Art B. Owen. Monte Carlo variance of scrambled net quadrature. *SIAM J. Numer. Anal.*, 34(5):1884–1910, 1997.
- [9] Michael Stein. Large sample properties of simulations using Latin hypercube sampling. *Technometrics*, 29(2):143–151, 1987.