

Rare event simulation

Josselin Garnier

ETICS 2020

Rare event simulation: methods and simulations

- Estimation of the probability of a rare event (such as the failure of a complex system).
- Standard methods (quadrature, Monte Carlo, reliability).
- Advanced Monte Carlo methods (different variance reduction techniques: importance sampling, control variates, with adaptive versions).
- Interacting particle systems (IPS with mutation-selection-resampling, multilevel splitting).

Uncertainty propagation

- Context: numerical code (black box) or experiment

$$Y = f(\mathbf{X})$$

with

Y = scalar output

\mathbf{X} = random input parameters, with known distribution (with pdf $p(x)$)

f = deterministic function $\mathbb{R}^d \rightarrow \mathbb{R}$ (computationally expensive).

- Goal: estimation of a quantity of the form

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

with an “error bar” and the minimal number of simulations.

Examples (for a real-valued output Y):

- $g(y) = y \rightarrow$ mean of Y , $\mathbb{E}[Y]$
- $g(y) = y^2 \rightarrow$ variance of Y , $\text{Var}(Y) = \mathbb{E}[(Y - \mathbb{E}[Y])^2] = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$
- $g(y) = \mathbf{1}_{[y_s, \infty)}(y) \rightarrow$ probability $\mathbb{P}(Y \geq y_s)$.

Analytic method

- The quantity to be estimated is a d -dimensional integral:

$$I = \mathbb{E}[g(Y)] = \mathbb{E}[h(\mathbf{X})] = \int_{\mathbb{R}^d} h(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

where $p(\mathbf{x})$ is the pdf of \mathbf{X} and $h(\mathbf{x}) = g(f(\mathbf{x}))$.

- In simple cases (when the pdf p and the function h have explicit expressions), one can sometimes evaluate the integral exactly (exceptional situation).

Quadrature method

- The quantity to be estimated is a d -dimensional integral:

$$I = \mathbb{E}[g(Y)] = \mathbb{E}[h(\mathbf{X})] = \int_{\mathbb{R}^d} h(\mathbf{x})\rho(\mathbf{x})d\mathbf{x}$$

where $\rho(\mathbf{x})$ is the pdf of \mathbf{X} and $h(\mathbf{x}) = g(f(\mathbf{x}))$.

- If $\rho(\mathbf{x}) = \prod_{i=1}^d \rho_0(x_i)$, then it is possible to apply Gaussian quadrature with a tensorized grid with n^d points:

$$\hat{I} = \sum_{j_1=1}^n \cdots \sum_{j_d=1}^n \rho_{j_1} \cdots \rho_{j_d} h(\xi_{j_1}, \dots, \xi_{j_d})$$

with the weights $(\rho_j)_{j=1, \dots, n}$ and the points $(\xi_j)_{j=1, \dots, n}$ associated to the quadrature with weighting function ρ_0 .

- There exist quadrature methods with sparse grids (cf Smolyak).
- Quadrature methods are efficient when:
 - the function $\mathbf{x} \rightarrow h(\mathbf{x})$ is smooth (and not only f),
 - the dimension d is “small” (even with sparse grids).

They require many calls.

Monte Carlo method

- Principle: replace the statistical expectation

$$I = \mathbb{E}[g(Y)] = \mathbb{E}[h(\mathbf{X})]$$

by an empirical mean.

- There are different probabilistic representations, that give different simulation methods.

Monte Carlo method

For a given y_s we want to estimate

$$\rho_s = \mathbb{P}(f(\mathbf{X}) \geq y_s)$$

The quantity of interest is an expectation:

$$\rho_s = \mathbb{E}[\mathbf{1}_{[y_s, \infty)}(f(\mathbf{X}))]$$

• Monte Carlo method:

1) Let $(\mathbf{X}^{(k)})_{k=1}^n$ be a n -sample of \mathbf{X} .

2) Compute

$$Z^{(k)} = \mathbf{1}_{[y_s, \infty)}(f(\mathbf{X}^{(k)}))$$

3) Define the empirical estimator of ρ_s :

$$\hat{\rho}_n := \frac{1}{n} \sum_{k=1}^n Z^{(k)}$$

- Empirical estimator of p_s :

$$\hat{P}_n := \frac{1}{n} \sum_{k=1}^n Z^{(k)}$$

- The estimator is **unbiased**:

$$\mathbb{E} \left[\hat{P}_n \right] = \mathbb{E} \left[\frac{1}{n} \sum_{k=1}^n Z^{(k)} \right] = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[Z^{(k)}] = \mathbb{E}[Z^{(1)}] = p_s$$

- The **law of large numbers** shows that the estimator is **convergent**:

$$\hat{P}_n = \frac{1}{n} \sum_{k=1}^n Z^{(k)} \xrightarrow{n \rightarrow \infty} \mathbb{E}[Z^{(1)}] = p_s$$

because the $Z^{(k)}$'s are independent and identically distributed (i.i.d.).

- Empirical estimator of p_s :

$$\hat{P}_n := \frac{1}{n} \sum_{k=1}^n Z^{(k)}$$

- Mean square error:

$$\begin{aligned} \mathbb{E} \left[(\hat{P}_n - p_s)^2 \right] &= \text{Var}(\hat{P}_n) = \frac{1}{n} \text{Var}(Z^{(1)}) \\ &= \frac{1}{n} (p_s - p_s^2) \end{aligned}$$

- The relative error is therefore:

$$\text{Error} = \frac{\sqrt{\text{Var}(\hat{P}_n)}}{\mathbb{E}[\hat{P}_n]} = \frac{\sqrt{\text{Var}(\hat{P}_n)}}{p_s} = \frac{1}{\sqrt{n}} \sqrt{\frac{1}{p_s} - 1} \stackrel{p_s \ll 1}{\approx} \frac{1}{\sqrt{np_s}}$$

\Leftrightarrow If $p_s \ll 1$, then we need $np_s > 1$ so that the relative error is smaller than 1 (not surprising) !

- *Question:* The estimator \hat{P}_n gives an approximate value of p_s , all the better as n is larger. How to quantify the error ?
- *Answer:* We build a confidence interval at the level 0.95, i.e. an empirical interval $[\hat{a}_n, \hat{b}_n]$ such that

$$\mathbb{P} \left(p_s \in [\hat{a}_n, \hat{b}_n] \right) \geq 0.95$$

Construction based on Central Limit Theorem:

$$\sqrt{n}(\hat{P}_n - p_s) = \sqrt{n} \left(\frac{1}{n} \sum_{k=1}^n Z^{(k)} - p_s \right) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, p_s - p_s^2) \text{ in distribution}$$

Therefore

$$\mathbb{P} \left(\left| \hat{P}_n - p_s \right| < c \frac{\sqrt{p_s - p_s^2}}{\sqrt{n}} \right) \xrightarrow{n \rightarrow \infty} \frac{2}{\sqrt{2\pi}} \int_0^c e^{-x^2/2} dx$$

$$\mathbb{P} \left(p_s \in \left[\hat{P}_n - 1.96 \frac{\sqrt{p_s - p_s^2}}{\sqrt{n}}, \hat{P}_n + 1.96 \frac{\sqrt{p_s - p_s^2}}{\sqrt{n}} \right] \right) \simeq 0.95$$

The unknown parameter p_s is still in the bounds of the interval !

Two solutions:

- $p_s \in [0, 1]$, therefore $\sqrt{p_s - p_s^2} < 1/2$ and

$$\mathbb{P} \left(p_s \in \left[\hat{P}_n - 0.98 \frac{1}{\sqrt{n}}, \hat{P}_n + 0.98 \frac{1}{\sqrt{n}} \right] \right) \geq 0.95$$

- asymptotically, we can replace p_s in the bounds by \hat{P}_n (OK if $np_s > 10$ and $n(1 - p_s) > 10$):

$$\mathbb{P} \left(p_s \in \left[\hat{P}_n - 1.96 \frac{\sqrt{\hat{P}_n - \hat{P}_n^2}}{\sqrt{n}}, \hat{P}_n + 1.96 \frac{\sqrt{\hat{P}_n - \hat{P}_n^2}}{\sqrt{n}} \right] \right) \simeq 0.95$$

[Proof: consistency of \hat{P}_n and Slutsky's theorem].

Monte Carlo estimation: general model

- Black box model (numerical code)

$$Y = f(\mathbf{X})$$

We want to estimate

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

for some function $g : \mathbb{R} \rightarrow \mathbb{R}$.

For instance $g(y) = \mathbf{1}_{[y_s, \infty)}(y)$.

- Empirical estimator:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n g(f(\mathbf{X}^{(k)}))$$

where $(\mathbf{X}^{(k)})_{k=1}^n$ is a n -sample of \mathbf{X} .

This is the empirical mean of a sequence of i.i.d. random variables.

- The estimator \hat{I}_n is unbiased: $\mathbb{E}[\hat{I}_n] = I$.
- The law of large numbers gives the **convergence** of the estimator:

$$\hat{I}_n \xrightarrow{n \rightarrow \infty} I \quad \text{with probability 1}$$

- Error:

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}(g(Y))$$

Proof: the variance of a sum of i.i.d. random variables is the sum of the variances.

- Asymptotic confidence interval:

$$\mathbb{P} \left(I \in \left[\hat{I}_n - 1.96 \frac{\hat{\sigma}_n}{\sqrt{n}}, \hat{I}_n + 1.96 \frac{\hat{\sigma}_n}{\sqrt{n}} \right] \right) \simeq 0.95$$

where

$$\hat{\sigma}_n = \left(\frac{1}{n} \sum_{k=1}^n g(f(\mathbf{X}^{(k)}))^2 - \hat{I}_n^2 \right)^{1/2}$$

- Advantages of the MC method:

- 1) no regularity condition for f , g (condition: $\mathbb{E}[g(f(\mathbf{X}))^2] < \infty$).
- 2) convergence rate $1/\sqrt{n}$ in any dimension.
- 3) can be applied for any quantity that can be expressed as an expectation.
- 4) embarrassingly parallel.

- One needs to simulate many samples of \mathbf{X} and to call many times f .

Uncertainty propagation by metamodels

The complex code/experiment f is replaced by a metamodel (reduced model) f_r and one of the previous techniques is applied with f_r (analytic, quadrature, Monte Carlo).

- It is possible to call many times the metamodel.
- The choice of the metamodel is critical.
- The error control is not simple.

Taylor expansions

- The output $Y = f(\mathbf{X})$ is approximated by a Taylor series expansion $Y_r = f_r(\mathbf{X})$.
 - Example (sandwich method):
 - We want to estimate $\mathbb{E}[Y]$ and $\text{Var}(Y)$ for $Y = f(\mathbf{X})$ given $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$ and $\mathbf{C} = \text{Cov}(\mathbf{X})$.
 - We approximate $Y = f(\mathbf{X})$ by $Y_r = f_r(\mathbf{X}) = f(\boldsymbol{\mu}) + \nabla f(\boldsymbol{\mu}) \cdot (\mathbf{X} - \boldsymbol{\mu})$.
- We find:

$$\mathbb{E}[Y] \simeq \mathbb{E}[Y_r] = f(\boldsymbol{\mu}), \quad \text{Var}(Y) \simeq \text{Var}(Y_r) = \nabla f(\boldsymbol{\mu})^T \mathbf{C} \nabla f(\boldsymbol{\mu})$$

We just need to compute $f(\boldsymbol{\mu})$ and $\nabla f(\boldsymbol{\mu})$ (evaluation of the gradient by finite differences, about $d + 1$ calls to f , or adjoint method).

- Rapid, analytic, allows to evaluate approximately central trends of the output (mean, variance).
- Suitable for small variations of the input parameters and a smooth model (that can be linearized).

Reliability methods

- We consider the output $Y = f(\mathbf{X})$, with $\mathbf{X} = (X_i)_{i=1}^d$ a random vector with known pdf p .
- We want to evaluate

$$p_s = \mathbb{P}(Y \geq y_s) = \mathbb{P}(\mathbf{X} \in F) = \int_{\mathbb{R}^d} \mathbf{1}_F(\mathbf{x})p(\mathbf{x})d\mathbf{x} = \int_F p(\mathbf{x})d\mathbf{x},$$

where

$$F = \{\mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \geq y_s\}$$

is called the failure domain.

- FORM-SORM method (First or Second-Order Reliability Method):
Let us assume that the X_i 's are *i.i.d.* with distribution $\mathcal{N}(0, 1)$ (see the slides on *isoprobabilist transform*):

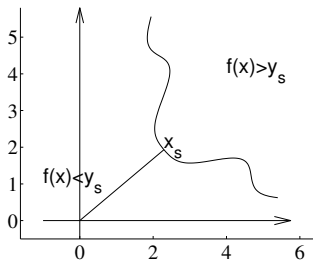
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right)$$

[Cf. O. Ditlevsen et H. O. Madsen, Structural Reliability Methods, Wiley, 1996.]

$$\rho_s = \int_F p(\mathbf{x}) d\mathbf{x}, \quad F = \{\mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \geq y_s\}$$

- we find by constrained optimization the design point \mathbf{x}_s , i.e. the point with the smallest norm s.t. $f(\mathbf{x}_s) = y_s$ (assuming it exists and is unique).

$$\mathbf{x}_s = \operatorname{argmin}_{\mathbf{x} \in F} \|\mathbf{x}\|^2$$



- the failure domain F is approximated by a half-space with smooth and simple boundary \hat{F} going through \mathbf{x}_s , which makes it possible to compute

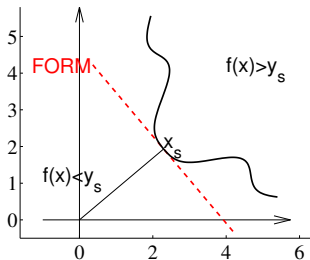
$$\hat{\rho}_s = \int_{\hat{F}} p(\mathbf{x}) d\mathbf{x}$$

$$\hat{\rho}_s = \int_{\hat{F}} p(x) dx$$

- The half-space \hat{F} is determined by a **hyperplane** by the **FORM** method; it goes through the point x_s and is orthogonal to the vector x_s . We get:

$$\hat{\rho}_s = \Phi(-\|x_s\|),$$

where Φ is the cdf of the distribution $\mathcal{N}(0, 1)$.



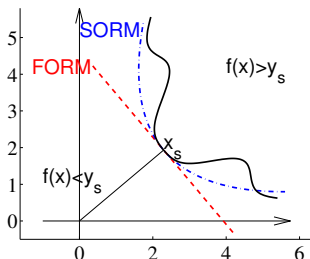
Proof. Introduce an orthonormal basis of \mathbb{R}^d whose first vector is $x_s/\|x_s\|$. Carry out the change of variable $x' = \mathbf{O}x$, where \mathbf{O} is the orthogonal matrix of basis change. □

$$\hat{\rho}_s = \int_{\hat{F}} p(x) dx$$

- The half-space \hat{F} is determined by a **quadratic surface** by the **SORM** method. We get (Breitung formula):

$$\hat{\rho}_s \simeq \Phi(-\|\mathbf{x}_s\|) \prod_{i=1}^{d-1} \frac{1}{\sqrt{1 + \|\mathbf{x}_s\| \kappa_i}} \left(1 + o_{\|\mathbf{x}_s\| \rightarrow \infty} (1) \right),$$

where the κ_i are the curvatures of the failure surface at \mathbf{x}_s (computed from the gradient and Hessian of f at \mathbf{x}_s).



Isoprobabilist transform

- Let the distribution of $\mathbf{X} = (X_i)_{i=1}^d$ be given.
How do we transform the problem into the “standard” form where X_i are i.i.d. with distribution $\mathcal{N}(0, 1)$?
- Assume that $(X_i)_{i=1}^d$ are *independent* with continuous cdf $(F_i)_{i=1}^d$.
- Let

$$\phi(\mathbf{x}) = (\Phi^{-1}(F_1(x_1)), \dots, \Phi^{-1}(F_d(x_d))),$$

where Φ is the cdf of the distribution $\mathcal{N}(0, 1)$.

- The vector

$$\tilde{\mathbf{X}} = \phi(\mathbf{X})$$

has independent coordinates which satisfy

$$\mathbb{P}(\tilde{X}_i \leq x) = \mathbb{P}(\Phi^{-1}(F_i(X_i)) \leq x) = \mathbb{P}(F_i(X_i) \leq \Phi(x)) = \Phi(x),$$

because $F_i(X_i)$ has distribution $\mathcal{U}(0, 1)$.

→ The coordinates of $\tilde{\mathbf{X}}$ are i.i.d. with distribution $\mathcal{N}(0, 1)$.

→ $p = \mathbb{P}(f(\mathbf{X}) \geq y_s) = \mathbb{P}(\tilde{f}(\tilde{\mathbf{X}}) \geq y_s)$ with $\tilde{f} = f \circ \phi^{-1}$.

- The Rosenblatt transform: If \mathbf{X} is a random vector with arbitrary pdf, then there exists a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that

$$\tilde{\mathbf{X}} = \phi(\mathbf{X})$$

has distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

- In two steps:

$$\phi = \phi^{(2)} \circ \phi^{(1)}$$

where $\phi^{(1)}$ and $\phi^{(2)}$ are the functions from \mathbb{R}^d to \mathbb{R}^d defined by:

1) $\phi_i^{(1)}(\mathbf{x}) = F_{i|1,\dots,i-1}(x_i|x_1, \dots, x_{i-1})$, $i = 1, \dots, d$,

2) $\phi^{(2)}(\mathbf{z}) = (\Phi^{-1}(z_1), \dots, \Phi^{-1}(z_d))$.

Here the function $F_{i|1,\dots,i-1}(x_i|x_1, \dots, x_{i-1})$ is the cdf of the variable X_i given $\{X_1 = x_1, \dots, X_{i-1} = x_{i-1}\}$.

- The coordinates of $\phi^{(1)}(\mathbf{X})$ are i.i.d. with distribution $\mathcal{U}(0, 1)$.
 - The coordinates of $\phi^{(2)} \circ \phi^{(1)}(\mathbf{X})$ are i.i.d. with distribution $\mathcal{N}(0, 1)$.
- ↔ Not easy to implement and to interpret (e.g., dependence w.r.t. order).

Proof.

Let $\mathbf{Z} = \phi^{(1)}(\mathbf{X})$. We want to show that \mathbf{Z} has the distribution $\mathcal{U}([0, 1]^d)$. Let g be a test function ($d = 2$).

$$\mathbb{E}[g(\mathbf{Z})] = \mathbb{E}[g(\phi^{(1)}(\mathbf{X}))] = \iint_{\mathbb{R}^2} g(\phi^{(1)}(\mathbf{x}))p(\mathbf{x})d\mathbf{x}$$

We make the change of variable $\mathbf{z} = \phi^{(1)}(\mathbf{x}) = (F_1(x_1), F_{2|1}(x_2|x_1))$, whose Jacobian is:

$$\mathbf{J} = \begin{pmatrix} \partial_{x_1} F_1(x_1) & \partial_{x_2} F_1(x_1) \\ \partial_{x_1} F_{2|1}(x_2|x_1) & \partial_{x_2} F_{2|1}(x_2|x_1) \end{pmatrix} = \begin{pmatrix} p_1(x_1) & 0 \\ * & p_{2|1}(x_2|x_1) \end{pmatrix}$$

and the determinant of the Jacobian is:

$$\text{Det}(\mathbf{J}) = p_1(x_1)p_{2|1}(x_2|x_1) = p(x_1, x_2)$$

Therefore

$$\mathbb{E}[g(\mathbf{Z})] = \iint_{[0,1]^2} g(\mathbf{z})d\mathbf{z}$$



Variance reduction techniques

Goal: reduce the variance of the Monte Carlo estimator:

$$\mathbb{E}[(\hat{I}_n - I)^2] = \frac{1}{n} \text{Var}(h(\mathbf{X}))$$

where $h(\mathbf{x}) = g(f(\mathbf{x}))$, $I = \mathbb{E}[h(\mathbf{X})]$, $\hat{I}_n = \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)})$.

- The methods

- Importance sampling
- Control variates
- Stratification

reduce the constant without changing $1/n$, stay close to the Monte Carlo method (parallelizable).

- The methods

- Quasi-Monte Carlo

aim at changing $1/n$.

- The methods

- Interacting particle systems

are different from Monte Carlo (sequential).

Importance sampling

- The goal is to estimate $I = \mathbb{E}[h(\mathbf{X})]$ for \mathbf{X} a random vector and $h(\mathbf{x}) = g(f(\mathbf{x}))$ a deterministic function.
- Observation: the representation of I as an expectation is not unique:

$$I = \mathbb{E}_p[h(\mathbf{X})] = \int h(\mathbf{x})p(\mathbf{x})d\mathbf{x} = \int \frac{h(\mathbf{x})p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} = \mathbb{E}_q\left[\frac{h(\mathbf{X})p(\mathbf{X})}{q(\mathbf{X})}\right]$$

The choice of the pdf q depends on the user.

- Idea: when we know that $h(\mathbf{X})$ is sensitive to certain values of \mathbf{X} , instead of sampling $\mathbf{X}^{(k)}$ with the original pdf $p(\mathbf{x})$ of \mathbf{X} , a biased pdf $q(\mathbf{x})$ is used that makes more likely the “important” realizations.
- Using the representation

$$I = \mathbb{E}_p[h(\mathbf{X})] = \mathbb{E}_q\left[h(\mathbf{X})\frac{p(\mathbf{X})}{q(\mathbf{X})}\right]$$

we can propose the estimator:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) \frac{p(\mathbf{X}^{(k)})}{q(\mathbf{X}^{(k)})}, \quad \mathbf{X}^{(k)} \text{ i.i.d. with pdf } q.$$

- The estimator is unbiased (provided $\text{supp}(\rho) \subset \text{supp}(q)$):

$$\begin{aligned}\mathbb{E}_q[\hat{I}_n] &= \frac{1}{n} \sum_{k=1}^n \mathbb{E}_q \left[h(\mathbf{X}^{(k)}) \frac{\rho(\mathbf{X}^{(k)})}{q(\mathbf{X}^{(k)})} \right] = \mathbb{E}_q \left[h(\mathbf{X}) \frac{\rho(\mathbf{X})}{q(\mathbf{X})} \right] \\ &= \int h(\mathbf{x}) \frac{\rho(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) d\mathbf{x} = \int h(\mathbf{x}) \rho(\mathbf{x}) d\mathbf{x} = \mathbb{E}_\rho [h(\mathbf{X})] = I\end{aligned}$$

- The estimator is convergent:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) \frac{\rho(\mathbf{X}^{(k)})}{q(\mathbf{X}^{(k)})} \xrightarrow{n \rightarrow \infty} \mathbb{E}_q \left[h(\mathbf{X}) \frac{\rho(\mathbf{X})}{q(\mathbf{X})} \right] = \mathbb{E}_\rho [h(\mathbf{X})] = I$$

- The variance of the estimator is:

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}_q \left(h(\mathbf{X}) \frac{\rho(\mathbf{X})}{q(\mathbf{X})} \right) = \frac{1}{n} \left(\mathbb{E}_\rho \left[h(\mathbf{X})^2 \frac{\rho(\mathbf{X})}{q(\mathbf{X})} \right] - \mathbb{E}_\rho [h(\mathbf{X})]^2 \right)$$

By a judicious choice of q the variance can be dramatically reduced.

- Critical points: it is necessary to know the likelihood ratio $\frac{\rho(\mathbf{x})}{q(\mathbf{x})}$ and to know how to simulate \mathbf{X} with the pdf q .

- The estimator is asymptotically normal

$$\sqrt{n}(\hat{I}_n - I) \xrightarrow{n \rightarrow +\infty} \mathcal{N}\left(0, \text{Var}_q\left(h(\mathbf{X}) \frac{p(\mathbf{X})}{q(\mathbf{X})}\right)\right)$$

- It is (in theory) possible to construct asymptotic confidence intervals. The empirical estimator of the asymptotic variance is

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{k=1}^n h^2(\mathbf{X}^{(k)}) \frac{p^2(\mathbf{X}^{(k)})}{q^2(\mathbf{X}^{(k)})} - \hat{I}_n^2, \quad \mathbf{X}^{(k)} \text{ i.i.d. with pdf } q.$$

The estimator $\hat{\sigma}_n^2$ is consistent (with standard moment conditions) but it may be strongly fluctuating.

↪ the construction of confidence interval is not easy.

- **Optimal importance sampling.**

The best biased distribution is the one that minimizes $\text{Var}(\hat{I}_n)$.

↪ solution of the minimization problem: find the pdf $q(\mathbf{x})$ minimizing

$$\mathbb{E}_p \left[h(\mathbf{X})^2 \frac{p(\mathbf{X})}{q(\mathbf{X})} \right] = \int h(\mathbf{x})^2 \frac{p^2(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

Solution (when h is nonnegative-valued):

$$q_{\text{opt}}(\mathbf{x}) = \frac{h(\mathbf{x})p(\mathbf{x})}{\int h(\mathbf{x}')p(\mathbf{x}')d\mathbf{x}'}$$

We then find

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \left(\mathbb{E}_p \left[h(\mathbf{X})^2 \frac{p(\mathbf{X})}{q_{\text{opt}}(\mathbf{X})} \right] - \mathbb{E}_p [h(\mathbf{X})]^2 \right) = 0 !$$

Practically: the denominator of $q_{\text{opt}}(\mathbf{x})$ is the desired quantity $\mathbb{E}[h(\mathbf{X})]$, which is unknown; we do not know how to sample q_{opt} .

Therefore the optimal IS method cannot be implemented.

But it is the principle for an adaptive method.

- Example: We want to estimate

$$I = \mathbb{P}(X \geq 4) = \mathbb{E}[h(X)]$$

with $X \sim \mathcal{N}(0, 1)$ and $h(x) = \mathbf{1}_{[4, \infty)}(x)$.

$$I = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \mathbf{1}_{[4, \infty)}(x) e^{-\frac{x^2}{2}} dx = \Phi(-4) = \frac{1}{2} \operatorname{erfc}\left(\frac{4}{\sqrt{2}}\right) \simeq 3.17 \cdot 10^{-5}$$

Monte Carlo: Let $(X^{(k)})_{k=1}^n$ be i.i.d. with the original distribution $\mathcal{N}(0, 1)$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X^{(k)} \geq 4}$$

We have $\operatorname{Var}(\hat{I}_n) = \frac{1}{n} 3.17 \cdot 10^{-5}$.

Importance Sampling: Let $(X^{(k)})_{k=1}^n$ be i.i.d. with the distribution $\mathcal{N}(4, 1)$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X^{(k)} \geq 4} \frac{e^{-\frac{(X^{(k)})^2}{2}}}{e^{-\frac{(X^{(k)}-4)^2}{2}}} = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X^{(k)} \geq 4} e^{-4X^{(k)}+8}$$

We have $\operatorname{Var}(\hat{I}_n) = \frac{1}{n} 5.53 \cdot 10^{-8}$.

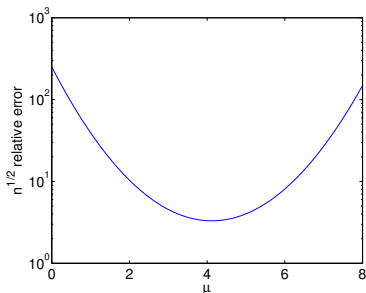
IS needs 600 times less simulations than MC to reach the same precision !

Warning: we should not bias too much.

Importance Sampling: Let $(X^{(k)})_{k=1}^n$ be i.i.d. with the distribution $\mathcal{N}(\mu, 1)$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X^{(k)} \geq 4} \frac{e^{-\frac{(X^{(k)})^2}{2}}}{e^{-\frac{(X^{(k)} - \mu)^2}{2}}} = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X^{(k)} \geq 4} e^{-\mu X^{(k)} + \frac{\mu^2}{2}}$$

$\hookrightarrow \text{Var}(\hat{I}_n) = \frac{1}{n} \frac{e^{\mu^2}}{2} \text{erfc}\left(\frac{4+\mu}{\sqrt{2}}\right) - \frac{1}{n} I^2$, which gives the normalized relative error $\sqrt{n} \mathbb{E}[(\hat{I}_n - I)^2]^{1/2} / I$:



If the bias is large, the fluctuations of the likelihood ratios become large.

- Example: we want to estimate

$$I = \mathbb{E}[h(X)]$$

with $X \sim \mathcal{N}(0, 1)$ and $h(x) = \exp(x)$.

$$I = \frac{1}{\sqrt{2\pi}} \int e^x e^{-\frac{x^2}{2}} dx = e^{\frac{1}{2}}$$

The large values of X are important.

Importance Sampling: With a sample $(X^{(k)})_{k=1, \dots, n}$ with the distribution $\mathcal{N}(\mu, 1)$, $\mu > 0$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n h(X^{(k)}) \frac{e^{-\frac{[X^{(k)}]^2}{2}}}{e^{-\frac{[X^{(k)} - \mu]^2}{2}}} = \frac{1}{N} \sum_{k=1}^n h(X^{(k)}) e^{-\mu X^{(k)} + \frac{\mu^2}{2}}$$

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \left(e^{\mu^2 - 2\mu + 2} - e^1 \right)$$

Monte Carlo $\mu = 0$: $\text{Var}(\hat{I}_n) = \frac{1}{n} (e^2 - e^1)$

Optimal importance sampling $\mu = 1$: $\text{Var}(\hat{I}_n) = 0$.

- Gaussian random walk $X_p = X_{p-1} + \theta_p$, $X_0 = 0$, where θ_p is a sequence of i.i.d. Gaussian random variables $\mathcal{N}(0, 1)$.

$$p_s = \mathbb{P}(X_M \geq y_s) = \frac{1}{\sqrt{2\pi M}} \int_{y_s}^{\infty} \exp\left(-\frac{x^2}{2M}\right) dx$$

MC: sample n trajectories $(X_i^{(k)})_{i=0}^M$, $k = 1, \dots, n$, and estimate

$$\hat{P}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_M^{(k)} \geq y_s}, \quad \hat{e}_n = \frac{1}{\sqrt{n}} \left(\frac{1}{\hat{P}_n} - 1 \right)^{1/2}$$

IS: sample n trajectories $(X_p^{(k)})_{p=0}^M$, $k = 1, \dots, n$, with biased distribution $\theta_p \sim \mathcal{N}(a/M, 1)$ and estimate:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_M^{(k)} \geq y_s} \exp\left(\frac{y_s^2}{2M} - \frac{y_s}{M} X_M^{(k)}\right)$$

$$\hat{e}_n = \frac{1}{\sqrt{n}} \left\{ \frac{\frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_M^{(k)} \geq y_s} \exp\left(\frac{y_s^2}{M} - \frac{2y_s}{M} X_M^{(k)}\right)}{\hat{I}_n^2} - 1 \right\}^{1/2}$$

- Adaptive importance sampling: choice of a family of biased distributions

$$\{q_{\theta}, \theta \in \Theta \subset \mathbb{R}^q\}$$

- Affine transform

Change of mean $q_{\theta_{\mu}}(\mathbf{x}) = p(\mathbf{x} - \theta_{\mu})$

Change of mean and variance $q_{\theta_{\mu}, \theta_{\sigma}}(\mathbf{x}) = |\det \theta_{\sigma}|^{-1} p(\theta_{\sigma}^{-1}(\mathbf{x} - \theta_{\mu}))$

Easy to implement $\tilde{\mathbf{X}} = \theta_{\sigma} \mathbf{X} + \theta_{\mu} \sim q_{\theta_{\mu}, \theta_{\sigma}}$ iff $\mathbf{X} \sim p$.

Often used with Gaussian distributions.

- Adaptive importance sampling: choice of θ .
- Variance minimization.

Principle: minimization of the variance of the estimator:

$$\theta^* \in \operatorname{argmin}_{\theta} \mathbb{E}_p \left[h(\mathbf{X})^2 \frac{p}{q_{\theta}}(\mathbf{X}) \right]$$

Empirically

$$\hat{\theta}_n \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)})^2 \frac{p}{q_{\theta}}(\mathbf{X}^{(k)}), \quad \mathbf{X}^{(k)} \text{ i.i.d. with pdf } p$$

Studied with Gaussian distribution and change of mean and variance.
With additional hypotheses, $\hat{\theta}_n \rightarrow \hat{\theta}^*$ a.s., with central limit theorem.
But: the likelihood ratio is strongly fluctuating.

[B. Jourdain and J. Lelong. Robust adaptive importance sampling for normal random vectors. *Ann. Appl. Probab.* **19**, 1687-1718, 2009.]

- Adaptive importance sampling: choice of θ .

- Cross entropy

Let q^* be the optimal pdf $q^*(\mathbf{x}) = \frac{h(\mathbf{x})}{\mathbb{E}_p[h(\mathbf{X})]} p(\mathbf{x})$ (here $h \geq 0$).

Principle: minimization of the Kullback-Leibler distance between q^* and q_θ :

$$\theta^* \in \operatorname{argmin}_{\theta} D(q^*, q_\theta), \quad D(q^*, q_\theta) = \mathbb{E}_{q^*} \left[\log \frac{q^*}{q_\theta}(\mathbf{X}) \right]$$

In fact

$$\operatorname{argmin}_{\theta} D(q^*, q_\theta) = \operatorname{argmin}_{\theta} \mathbb{E}_p \left[h(\mathbf{X}) \log \frac{p}{q_\theta}(\mathbf{X}) \right]$$

Empirically

$$\hat{\theta}_n \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) \log \frac{p}{q_\theta}(\mathbf{X}^{(k)}), \quad \mathbf{X}^{(k)} \text{ i.i.d. with pdf } p$$

With additional hypotheses, $\hat{\theta}_n \rightarrow \theta^*$ a.s., with central limit theorem.

The log-likelihood ratio is less fluctuating than the likelihood ratio used in the variance minimization method.

- Adaptive importance sampling: choice of θ .
- Cross entropy

For any θ_0 ,

$$\operatorname{argmin}_{\theta} D(q^*, q_{\theta}) = \operatorname{argmin}_{\theta} \mathbb{E}_{q_{\theta_0}} \left[h(\mathbf{X}) \log \left(\frac{p}{q_{\theta}}(\mathbf{X}) \right) \frac{p}{q_{\theta_0}}(\mathbf{X}) \right]$$

Empirically, if $\mathbf{X}^{(k)}$ i.i.d. with pdf q_{θ_0} , then

$$\hat{\theta}_n \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) \log \left(\frac{p}{q_{\theta}}(\mathbf{X}^{(k)}) \right) \frac{p}{q_{\theta_0}}(\mathbf{X}^{(k)})$$

- Adaptive importance sampling: choice of θ .
- Recursive cross entropy

Principle: go progressively towards θ^* .

Algorithm:

Step 0: Set θ_0 (e.g. $\theta_0 = \mathbf{0}$, $q_{\theta_0} = p$).

Step $j \geq 1$: Generate $\mathbf{X}^{(k)}$ i.i.d. with pdf $q_{\theta_{j-1}}$

$$\theta_j \in \operatorname{argmin}_{\theta} \frac{1}{N_{j-1}} \sum_{k=1}^{N_{j-1}} h(\mathbf{X}^{(k)}) \log \left(\frac{p}{q_{\theta}}(\mathbf{X}^{(k)}) \right) \frac{p}{q_{\theta_{j-1}}}(\mathbf{X}^{(k)}),$$

with N_{j-1} large enough, until a stopping criterium is met, e.g.

$$|\theta_j - \theta_{j-1}| \leq \delta.$$

Problem: if $h(\mathbf{x}) = \mathbf{1}_{f(\mathbf{x}) \geq y_s}$, then one needs realizations that achieve $f(\mathbf{X}^{(k)}) \geq y_s$ during the first step (θ_0, N_0).

- Adaptive importance sampling: choice of θ .
- Adaptive cross entropy

Assume $h(\mathbf{x}) = \mathbf{1}_{f(\mathbf{x}) \geq y_s}$.

Principle: go progressively towards θ^* by “increasing” y_s .

Algorithm:

Step 0: Set θ_0 (e.g. $\theta_0 = \mathbf{0}$, $q_{\theta_0} = p$), y_0 , and α (e.g. $\alpha = 0.1$).

Step $j \geq 1$: Generate $\mathbf{X}^{(k)}$ i.i.d. with pdf $q_{\theta_{j-1}}$

$$\theta_j \in \underset{\theta}{\operatorname{argmin}} \frac{1}{N_{j-1}} \sum_{k=1}^{N_{j-1}} \mathbf{1}_{f(\mathbf{X}^{(k)}) \geq y_{j-1}} \log \left(\frac{p}{q_{\theta}}(\mathbf{X}^{(k)}) \right) \frac{p}{q_{\theta_{j-1}}}(\mathbf{X}^{(k)}),$$

$$y_j = (1 - \alpha) \text{--empirical quantile of } (f(\mathbf{X}^{(k)}))_{k=1}^{N_{j-1}}$$

with N_{j-1} large enough, until $y_j \geq y_s$.

Estimation of I with or without recycling.

[P.T. De Boer, D.P. Kroese, S. Mannor, and R.Y. Rubinstein. A tutorial on the cross-entropy method, *Annals of Operations Research* **134**, 19-67, 2005.]

Control variates

- The goal is to estimate $I = \mathbb{E}[h(\mathbf{X})]$ for \mathbf{X} a random vector and $h(\mathbf{x}) = g(f(\mathbf{x}))$ a deterministic function.
- Assume that we have a reduced model $f_r(\mathbf{X})$.

- Importance sampling method:

First we evaluate (we approximate) the optimal density

$$q_{\text{opt},r}(\mathbf{x}) = \frac{g(f_r(\mathbf{x}))p(\mathbf{x})}{I_r}, \text{ with } I_r = \int g(f_r(\mathbf{x}))p(\mathbf{x})d\mathbf{x}.$$

Second we use it as a biased density for estimating I .

Dangerous, use conservative version.

Control variates

- The goal is to estimate $I = \mathbb{E}[h(\mathbf{X})]$ for \mathbf{X} a random vector and $h(\mathbf{x}) = g(f(\mathbf{x}))$ a deterministic function.
- Assume that we have a reduced model $f_r(\mathbf{X})$.

- Control variates method:

We denote $h(\mathbf{x}) = g(f(\mathbf{x}))$, $h_r(\mathbf{x}) = g(f_r(\mathbf{x}))$.

We assume that we know $I_r = \mathbb{E}[h_r(\mathbf{X})]$.

By considering the representation

$$I = \mathbb{E}[h(\mathbf{X})] = I_r + \mathbb{E}[h(\mathbf{X}) - h_r(\mathbf{X})]$$

we propose the estimator:

$$\hat{I}_n = I_r + \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) - h_r(\mathbf{X}^{(k)}),$$

where $(\mathbf{X}^{(k)})_{k=1}^n$ is a n -sample (with the pdf p).

- Estimator:

$$\hat{l}_n = l_r + \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) - h_r(\mathbf{X}^{(k)})$$

- The estimator is unbiased:

$$\begin{aligned} \mathbb{E}[\hat{l}_n] &= l_r + \frac{1}{n} \sum_{k=1}^n \mathbb{E}[h(\mathbf{X}^{(k)}) - h_r(\mathbf{X}^{(k)})] = l_r + \mathbb{E}[h(\mathbf{X})] - \mathbb{E}[h_r(\mathbf{X})] \\ &= l_r + \mathbb{E}[h(\mathbf{X})] - l_r = l \end{aligned}$$

- The estimator is convergent:

$$\hat{l}_n \xrightarrow{n \rightarrow \infty} l_r + \mathbb{E}[h(\mathbf{X}) - h_r(\mathbf{X})] = l$$

- The variance of the estimator is:

$$\text{Var}(\hat{l}_n) = \frac{1}{n} \text{Var}[h(\mathbf{X}) - h_r(\mathbf{X})]$$

↔ The use of a reduced model can reduce the variance.

- Example: we want to estimate

$$I = \mathbb{E}[h(X)]$$

with $X \sim \mathcal{U}(0, 1)$, $h(x) = \exp(x)$.

Result: $I = e - 1 \simeq 1.72$.

Monte Carlo.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \exp[X^{(k)}]$$

Variance of the MC estimator = $\frac{1}{n}(2e - 1) \simeq \frac{1}{n}4.44$.

Control Variates. Reduced model: $h_r(x) = 1 + x$ (here $I_r = \frac{3}{2}$). CV estimator:

$$\hat{I}_n = I_r + \frac{1}{n} \sum_{k=1}^n \left\{ \exp[X^{(k)}] - 1 - X^{(k)} \right\}$$

Variance of the CV estimator = $\frac{1}{n}(3e - \frac{e^2}{2} - \frac{53}{12}) \simeq \frac{1}{n}0.044$.

The CV method needs 100 times less simulations to reach the same precision as MC !

- Application: estimation of

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

We have a reduced model f_r of the full code f . The ratio between the computational cost of one call of f and one call of f_r is $q > 1$.

Estimator

$$\hat{I}_n = \frac{1}{n_r} \sum_{k=1}^{n_r} h_r(\tilde{\mathbf{X}}^{(k)}) + \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) - h_r(\mathbf{X}^{(k)})$$

with $n_r > n$, $h(\mathbf{x}) = g(f(\mathbf{x}))$, $h_r(\mathbf{x}) = g(f_r(\mathbf{x}))$.

Allocation between calls to the complete code and calls to the reduced model can be optimized with the constraint $n_r/q + n(1 + 1/q) = n_{\text{tot}}$:

$$\frac{n}{n_{\text{tot}}} = \frac{q}{1+q} \frac{1}{1 + \frac{1}{\sqrt{1+q}} \frac{\sqrt{\text{Var}(h_r(\mathbf{X}))}}{\sqrt{\text{Var}((h-h_r)(\mathbf{X}))}}}$$

Classical trade-off between approximation error and estimation error. Used when $f(\mathbf{X})$ is the solution of an ODE or PDE with fine grid, while $f_r(\mathbf{X})$ is the solution obtained with a coarse grid (MultiLevel Monte Carlo).

- Optimal control variate

By considering the representation (for a fixed λ)

$$I = \mathbb{E}[h(\mathbf{X})] = \lambda I_r + \mathbb{E}[h(\mathbf{X}) - \lambda h_r(\mathbf{X})]$$

we propose the estimator:

$$\hat{I}_n = \lambda I_r + \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) - \lambda h_r(\mathbf{X}^{(k)}),$$

where $(\mathbf{X}^{(k)})_{k=1}^n$ is a n -sample (with the pdf p).

The estimator is unbiased and consistent for any λ . The variance is

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}(h(\mathbf{X}) - \lambda h_r(\mathbf{X}))$$

The λ that minimizes the variance is

$$\lambda = \frac{\text{Cov}(h(\mathbf{X}), h_r(\mathbf{X}))}{\text{Var}(h_r(\mathbf{X}))}$$

and then

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}(h(\mathbf{X})) (1 - \rho^2), \quad \rho = \text{Corr}(h(\mathbf{X}), h_r(\mathbf{X}))$$

- Empirical optimal control variate

$$\hat{l}_n = \hat{\lambda}_n l_r + \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)}) - \hat{\lambda}_n h_r(\mathbf{X}^{(k)}),$$

with

$$\hat{\lambda}_n = \frac{\sum_{k=1}^n (h(\mathbf{X}^{(k)}) - \frac{1}{n} \sum_{k'=1}^n h(\mathbf{X}^{(k')})) (h_r(\mathbf{X}^{(k)}) - \frac{1}{n} \sum_{k'=1}^n h_r(\mathbf{X}^{(k')}))}{\sum_{k=1}^n (h_r(\mathbf{X}^{(k)}) - \frac{1}{n} \sum_{k'=1}^n h_r(\mathbf{X}^{(k')}))^2}$$

The estimator is slightly biased ($O(1/n)$).

We have the (optimal) asymptotic normality result

$$\sqrt{n}(\hat{l}_n - l) \xrightarrow{n \rightarrow +\infty} \mathcal{N}(0, \text{Var}(h(\mathbf{X})) (1 - \rho^2))$$

Stratification

Principle: The sample $(\mathbf{X}^{(k)})_{k=1}^n$ is enforced to obey prescribed proportions in some “strata”.

Method used in polls (representative sample).

Here: we want to estimate $\mathbb{E}[h(\mathbf{X})]$, \mathbf{X} with values in D .

• Two ingredients:

i) A partition of the state space $D = \bigcup_{i=1}^m D_i$. We know $p_i = \mathbb{P}(\mathbf{X} \in D_i)$.

ii) Total probability formula:

$$I = \mathbb{E}[h(\mathbf{X})] = \sum_{i=1}^m \underbrace{\mathbb{E}[h(\mathbf{X}) | \mathbf{X} \in D_i]}_{J_i} \underbrace{\mathbb{P}(\mathbf{X} \in D_i)}_{p_i}$$

• Estimation:

1) For all $i = 1, \dots, m$, J_i is estimated by MC with n_i simulations:

$$\hat{J}_{i,n_i} = \frac{1}{n_i} \sum_{k=1}^{n_i} h(\mathbf{X}^{(i,k)}), \quad \mathbf{X}^{(i,k)} \sim \mathcal{L}(\mathbf{X} | \mathbf{X} \in D_i) \text{ ind.}$$

2) The estimator is $\hat{I}_n = \sum_{i=1}^m \hat{J}_{i,n_i} p_i$

$$\hat{I}_n = \sum_{i=1}^m p_i \hat{J}_{i,n_i}, \quad \hat{J}_{i,n_i} = \frac{1}{n_i} \sum_{k=1}^{n_i} h(\mathbf{X}^{(i,k)}), \quad \mathbf{X}^{(i,k)} \sim \mathcal{L}(\mathbf{X} | \mathbf{X} \in D_i)$$

The total number of simulations is $n = \sum_{i=1}^m n_i$.

- The estimator is unbiased, convergent and its variance is

$$\text{Var}(\hat{I}_n)_S = \sum_{i=1}^m p_i^2 \text{Var}(\hat{J}_{i,n_i}) = \sum_{i=1}^m p_i^2 \frac{\sigma_i^2}{n_i}, \quad \text{with } \sigma_i^2 = \text{Var}(h(\mathbf{X}) | \mathbf{X} \in D_i)$$

The user is free to choose the n_i (with the constraint $\sum_{i=1}^m n_i = n$).

- **Proportional stratification:** $n_i = p_i n$.

$$\hat{I}_n = \sum_{i=1}^m \frac{p_i}{n_i} \sum_{k=1}^{n_i} h(\mathbf{X}^{(i,k)}) = \frac{1}{n} \sum_{i=1}^m \sum_{k=1}^{n_i} h(\mathbf{X}^{(i,k)}), \quad \mathbf{X}^{(i,k)} \sim \mathcal{L}(\mathbf{X} | \mathbf{X} \in D_i)$$

Then

$$\text{Var}(\hat{I}_n)_{SP} = \frac{1}{n} \sum_{i=1}^m p_i \sigma_i^2$$

$$\text{Var}(\hat{I}_n)_{MC} = \frac{1}{n} \text{Var}(h(\mathbf{X})) \geq \frac{1}{n} \sum_{i=1}^m p_i \sigma_i^2 = \text{Var}(\hat{I}_n)_{SP}$$

Proof: We have

$$\begin{aligned}\mathbb{E}[h(\mathbf{X})]^2 &= \left(\sum_{i=1}^m p_i \mathbb{E}[h(\mathbf{X}) | \mathbf{X} \in D_i] \right)^2 \\ &\leq \left(\sum_{i=1}^m p_i \right) \left(\sum_{i=1}^m p_i \mathbb{E}[h(\mathbf{X}) | \mathbf{X} \in D_i]^2 \right) \\ &= \sum_{i=1}^m p_i \mathbb{E}[h(\mathbf{X}) | \mathbf{X} \in D_i]^2\end{aligned}$$

Therefore

$$\begin{aligned}\sum_{i=1}^m p_i \sigma_i^2 &= \sum_{i=1}^m p_i \left(\mathbb{E}[h(\mathbf{X})^2 | \mathbf{X} \in D_i] - \mathbb{E}[h(\mathbf{X}) | \mathbf{X} \in D_i]^2 \right) \\ &= \mathbb{E}[h(\mathbf{X})^2] - \sum_{i=1}^m p_i \mathbb{E}[h(\mathbf{X}) | \mathbf{X} \in D_i]^2 \\ &\leq \mathbb{E}[h(\mathbf{X})^2] - \mathbb{E}[h(\mathbf{X})]^2 = \text{Var}(h(\mathbf{X}))\end{aligned}$$

However, the proportional allocation is not optimal !

- The optimal allocation is the one that minimizes the variance

$$\text{Var}(\widehat{I}_n)_S = \sum_{i=1}^m p_i^2 \frac{\sigma_i^2}{n_i}.$$

It is the solution of the minimization problem: find $(n_i)_{i=1,\dots,m}$ minimizing

$$\sum_{i=1}^m p_i^2 \frac{\sigma_i^2}{n_i} \text{ with the constraint } \sum_{i=1}^m n_i = n$$

Solution (optimal allocation, obtained with Lagrange multiplier method):

$$n_i = n \frac{p_i \sigma_i}{\sum_{l=1}^m p_l \sigma_l}$$

The minimal variance is

$$\text{Var}(\widehat{I}_n)_{SO} = \frac{1}{n} \left(\sum_{i=1}^m p_i \sigma_i \right)^2,$$

We have:

$$\text{Var}(\widehat{I}_n)_{SO} \leq \text{Var}(\widehat{I}_n)_{SP} \leq \text{Var}(\widehat{I}_n)_{MC}$$

Practically: the σ_i 's are unknown, therefore the optimal allocation is unknown (principle of an adaptive method).

- Example: we want to estimate

$$\mathbb{E}[h(X)]$$

with $X \sim \mathcal{U}(-1, 1)$, $h(x) = \exp(x)$.

Result: $\mathbb{E}[\exp(X)] = \sinh(1) \simeq 1.18$.

MC. With a sample $X^{(1)}, \dots, X^{(n)}$ with the distribution $\mathcal{U}(-1, 1)$

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \exp[X^{(k)}]$$

Variance of the estimator = $\frac{1}{n} \left(\frac{1}{2} - \frac{e^{-2}}{2} \right) \simeq \frac{1}{n} 0.43$.

Proportional stratification. With a sample

- $X^{(1)}, \dots, X^{(n/2)}$ with the distribution $\mathcal{U}(-1, 0)$,
- $X^{(n/2+1)}, \dots, X^{(n)}$ with the distribution $\mathcal{U}(0, 1)$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^{n/2} \exp[X^{(k)}] + \frac{1}{n} \sum_{k=n/2+1}^n \exp[X^{(k)}] = \frac{1}{n} \sum_{k=1}^n \exp[X^{(k)}]$$

Variance of the PS estimator $\simeq \frac{1}{n} 0.14$.

Here PS needs 3 times less simulations to reach the same precision as MC.

Nonproportional stratification. With a sample

- $X^{(1)}, \dots, X^{(n/4)}$ with the distribution $\mathcal{U}(-1, 0)$,
- $X^{(n/4+1)}, \dots, X^{(n)}$ with the distribution $\mathcal{U}(0, 1)$.

$$\hat{I}_n = \frac{2}{n} \sum_{k=1}^{n/4} \exp[X^{(k)}] + \frac{1}{2n} \sum_{k=n/4+1}^n \exp[X^{(k)}]$$

Variance of the estimator $\simeq \frac{1}{n}0.048$.

The stratification method needs 9 times less simulations to reach the same precision as MC.

Similar to importance sampling with a stepwise constant biased pdf.

Antithetic variables

- We want to compute

$$I = \int_{[0,1]^d} h(\mathbf{x}) d\mathbf{x}$$

MC with a n -sample $(\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)})$ with the distribution $\mathcal{U}([0, 1]^d)$:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n h(\mathbf{X}^{(k)})$$

$$\mathbb{E}[(\hat{I}_n - I)^2] = \frac{1}{n} \text{Var}(h(\mathbf{X})) = \frac{1}{n} \left(\int_{[0,1]^d} h^2(\mathbf{x}) d\mathbf{x} - I^2 \right)$$

- We consider the representations

$$I = \int_{[0,1]^d} h(1 - \mathbf{x}) d\mathbf{x} \text{ and } I = \int_{[0,1]^d} \frac{h(\mathbf{x}) + h(1 - \mathbf{x})}{2} d\mathbf{x}$$

MC with a $n/2$ -sample $(\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n/2)})$ distributed as $\mathcal{U}([0, 1]^d)$:

$$\tilde{I}_n = \frac{1}{n} \sum_{k=1}^{n/2} h(\mathbf{X}^{(k)}) + h(1 - \mathbf{X}^{(k)})$$

- MC estimator with the sample $(\tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(n)}) := (\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n/2)}, 1 - \mathbf{X}^{(1)}, \dots, 1 - \mathbf{X}^{(n/2)})$ that is not i.i.d.:

$$\tilde{I}_n = \frac{1}{n} \sum_{k=1}^n h(\tilde{\mathbf{X}}^{(k)})$$

The function f is called n times.

- Error:

$$\begin{aligned} \mathbb{E}[(\tilde{I}_n - I)^2] &= \frac{1}{n} \left(\text{Var}(h(\mathbf{X})) + \text{Cov}(h(\mathbf{X}), h(1 - \mathbf{X})) \right) \\ &= \frac{1}{n} \left(\int_{[0,1]^d} h^2(\mathbf{x}) + h(\mathbf{x})h(1 - \mathbf{x})d\mathbf{x} - 2I^2 \right) \end{aligned}$$

The variance is reduced if $\text{Cov}(h(\mathbf{X}), h(1 - \mathbf{X})) < 0$.

Sufficient condition: h is monotoneous.

Proof: If \mathbf{X}, \mathbf{X}' i.i.d.

$$\begin{aligned} [h(\mathbf{X}) - h(\mathbf{X}')][-h(1 - \mathbf{X}) + h(1 - \mathbf{X}')] &\geq 0 \text{ a.s.} \\ -2\mathbb{E}[h(\mathbf{X})h(1 - \mathbf{X})] + 2\mathbb{E}[h(\mathbf{X})]^2 &\geq 0 \end{aligned}$$



- Example:

$$I = \int_0^1 \frac{1}{1+x} dx$$

Result: $I = \ln 2$.

Monte Carlo:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \frac{1}{1+X^{(k)}}$$

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \left(\int_0^1 (1+x)^{-2} dx - \ln^2 2 \right) = \frac{1}{n} \left(\frac{1}{2} - \ln^2 2 \right) \simeq \frac{1}{n} 1.95 \cdot 10^{-2}.$$

Antithetic variables:

$$\tilde{I}_n = \frac{1}{n} \sum_{k=1}^{n/2} \frac{1}{1+X^{(k)}} + \frac{1}{2-X^{(k)}}$$

$$\text{Var}(\tilde{I}_n) = \frac{2}{n} \left(\int_0^1 \left(\frac{1}{2(1+x)} + \frac{1}{2(2-x)} \right)^2 dx - \ln^2 2 \right) \simeq \frac{1}{n} 1.2 \cdot 10^{-3}.$$

The AV method requires 15 times less simulations than MC.

- More generally: one needs to find a pair $(\mathbf{X}, \tilde{\mathbf{X}})$ such that $h(\mathbf{X})$ and $h(\tilde{\mathbf{X}})$ have the same distribution and $\text{Cov}(h(\mathbf{X}), h(\tilde{\mathbf{X}})) < 0$.
- Monte Carlo with an i.i.d. sample $((\mathbf{X}^{(1)}, \tilde{\mathbf{X}}^{(1)}), \dots, (\mathbf{X}^{(n/2)}, \tilde{\mathbf{X}}^{(n/2)}))$:

$$\tilde{I}_n = \frac{1}{n} \sum_{k=1}^{n/2} h(\mathbf{X}^{(k)}) + h(\tilde{\mathbf{X}}^{(k)})$$

$$\mathbb{E}[(\tilde{I}_n - I)^2] = \frac{1}{n} \left(\text{Var}(h(\mathbf{X})) + \text{Cov}(h(\mathbf{X}), h(\tilde{\mathbf{X}})) \right)$$

- Recent application: computation of effective tensors in stochastic homogenization (the effective tensor is the expectation of a functional of the solution of an elliptic PDE with random coefficients; antithetic pairs of the realizations of the composite medium are sampled; gain by a factor 3; in fact, better results with control variates; cf C. Le Bris, F. Legoll).
- Not very useful for the estimation of probabilities of rare events.