

Practical: Monte Carlo and Markov chain theory

Instructors: Kari Auranen, Elizabeth Halloran and Vladimir Minin

July 11 – July 13, 2018

Estimating the tail of the standard normal distribution

Let $Z \sim \mathcal{N}(0, 1)$. We would like to estimate the tail probability $\Pr(Z > c)$, where c is large (e.g., $c = 4.5$).

Naive Monte Carlo: simulate $Z_1, \dots, Z_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$. Then

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n 1_{\{Z_i > c\}} \approx \mathbb{E}(1_{\{Z > c\}}) = \Pr(Z > c).$$

This estimator will most likely give you 0 even for $n = 10,000$. The problem is the large variance of the integrand:

$$\text{Var}(\hat{\mu}) = \frac{1}{n} \text{Var}(1_{\{Z_1 > c\}}) = \frac{1}{n} \Pr(Z_1 > c)[1 - \Pr(Z_1 > c)] = \mathbf{3.4 \times 10^{-10}}$$
 for $n = 10,000$ and $c = 4.5$.

This variance is huge, because the quantity of interest is $\Pr(Z_1 > c) = 3.39 \times 10^{-6}$ and the standard deviation of our estimator is 1.84×10^{-5} .

Importance sampling: Simulate $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Exp}(c, 1)$ from a shifted exponential with density

$$g(y) = e^{-(y-c)} 1_{\{y > c\}}.$$

Generating such random variables is very easy: just simulate a regular exponential $\text{Exp}(1)$ and add c to the simulated value. Then the importance sampling estimator becomes

$$\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n \frac{\phi(Y_i)}{g(Y_i)} 1_{\{Y_i > c\}},$$

where $\phi(x)$ is the standard normal density. The variance of this estimator amounts to

$$\begin{aligned} \text{Var}(\tilde{\mu}) &= \frac{1}{n} \text{Var} \left[\frac{\phi(Y)}{g(Y)} 1_{\{Y > c\}} \right] = \frac{1}{n} \left\{ \mathbb{E}_g \left[\frac{\phi^2(Y)}{g^2(Y)} 1_{\{Y > c\}} \right] - \left[\mathbb{E}_g \left(\frac{\phi(Y)}{g(Y)} 1_{\{Y > c\}} \right) \right]^2 \right\} \\ &= \frac{1}{n} \left[\int_c^\infty \frac{\phi^2(y)}{g(y)} dy - \Pr(Z > c)^2 \right] = \mathbf{1.9474 \times 10^{-15}} \end{aligned}$$
 for $n = 10,000$ and $c = 4.5$.

This means that we reduced Monte Carlo variance roughly by a factor of 10^5 using importance sampling.

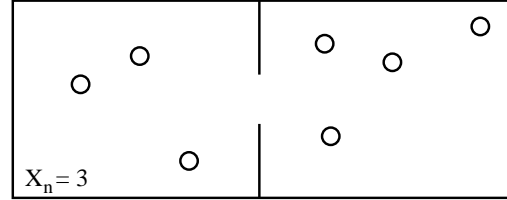
Your task

Implement naive and importance sampling Monte Carlo estimates of $\Pr(Z > 4.5)$, where $Z \sim \mathcal{N}(0, 1)$. Download ‘import_sampl_reduced.R’ from the course web page. The code has a couple of things to get you started.

Ehrenfest model of diffusion

Imagine a two dimensional rectangular box with a divider in the middle. The box contains N balls (gas molecules) distributed somehow between the two halves. The divider has a small gap, through which balls can go through one at a time. We assume that at each time step we select a ball uniformly at random and force it go through the gap to the opposite side of the divider. Letting X_n denote the total number of balls in the left half of the box, our Markov process is described by the following transition probabilities.

$$p_{ij} = \begin{cases} \frac{i}{N}, & \text{for } j = i - 1, \\ 1 - \frac{i}{N}, & \text{for } j = i + 1, \\ 0, & \text{otherwise.} \end{cases}$$



If we want to derive a stationary distribution of the system, we can solve the global balance equations $\pi^T \mathbf{P} = \pi^T$. Alternatively, we may “guess” that at equilibrium $X_n \sim \text{bin}(\frac{1}{2}, N)$ and verify this candidate stationary distribution via detailed balance. Notice we do not know whether the Ehrenfest chain is reversible, but we’ll go ahead with the detailed balance check anyway. First, notice that entries of our candidate vector are

$$\pi_i = \binom{N}{i} \left(\frac{1}{2}\right)^i \left(1 - \frac{1}{2}\right)^{N-i} = \binom{N}{i} \frac{1}{2^N}$$

Since X_n can only increase or decrease by one at each time step, we need to check detailed balance only for i and $j = i + 1$.

$$\begin{aligned} \pi_i p_{i,i+1} &= \frac{1}{2^N} \binom{N}{i} \frac{N-i}{N} = \frac{1}{2^N} \frac{N!}{i!(N-i)!} \frac{N-i}{N} = \frac{1}{2^N} \frac{N!}{(i+1)!(N-i-1)!} \frac{i+1}{N} \\ &= \binom{N}{i+1} \frac{1}{2^N} \frac{i+1}{N} = \pi_{i+1} p_{i+1,i}, \end{aligned}$$

confirming our guess.

Now, consider the Ehrenfest model with $N = 100$ gas molecules. From our derivations we know that the stationary distribution of the chain is $\text{Bin}(\frac{1}{2}, N)$. The chain is irreducible and positive recurrent (why?). The stationary variance can be computed analytically as $N \times \frac{1}{2} \times \frac{1}{2}$.

Your task

Use ergodic theorem to approximate the stationary variance and compare your estimate with the analytical result. Don’t panic! You will not have to write everything from scratch. Download ‘ehrenfest_diff_reduced.R’ file from the course web page. Follow comments in this R script to fill gaps in the code.