

Computational statistics 2021.2

School of Applied Mathematics, Fundação Getulio Vargas
Professor Luiz Max de Carvalho

Problem sheet 2

Lucas Machado Moschen

Exercício 1. (Monte Carlo for Gaussian)

Let us consider the normal multivariate density on \mathbb{R}^d with identity covariance, that is

$$\pi(x) = \frac{1}{(2\pi)^{d/2}} \exp \left\{ -\frac{1}{2} x^T x \right\}.$$

1. (Cameron-Martin formula). Show that for any $\theta \in \mathbb{R}^d$ and function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\mathbb{E}[\phi(X)] = \mathbb{E} \left[\phi(X + \theta) \exp \left(-\frac{1}{2} \theta^T \theta - \theta^T X \right) \right].$$

Let ϕ be any measurable function and $\theta \in \mathbb{R}^d$. Denote I_2 the quantity in the right of the equation. Then,

$$\begin{aligned} I_2 &= \int_{\mathbb{R}^d} \phi(x + \theta) \exp \left(-\frac{1}{2} \theta^T \theta - \theta^T x \right) \pi(x) dx \\ &= \int_{\mathbb{R}^d} \phi(x + \theta) \exp \left(-\frac{1}{2} \theta^T \theta - \theta^T x \right) \pi(x) dx \\ &= \int_{\mathbb{R}^d} \frac{1}{(2\pi)^{d/2}} \phi(x + \theta) \exp \left(-\frac{1}{2} (x + \theta)^T (x + \theta) \right) dx \\ &= \int_{\mathbb{R}^d} \frac{1}{(2\pi)^{d/2}} \phi(y) \exp \left(-\frac{1}{2} y^T y \right) dy \\ &= \mathbb{E}[\phi(X)]. \end{aligned}$$

2. It follows directly from the Cameron-Martin formula and the strong law of large numbers that, for independent $X_1, \dots, X_n \sim \pi$, the estimator

$$\hat{I}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \phi(X_i + \theta) \exp \left(-\frac{1}{2} \theta^T \theta - \theta^T X_i \right)$$

of $\mathbb{E}[\phi(X)]$ is strongly consistent for any $\theta \in \mathbb{R}^d$ such that

$$\mathbb{E} \left[\left| \phi(X + \theta) \exp \left(-\frac{1}{2} \theta^T \theta - \theta^T X \right) \right| \right] < +\infty.$$

The case $\theta = 0$ corresponds to the usual Monte Carlo estimate. The variance of $\hat{I}_n(\theta)$ is given by $\sigma^2(\theta)/n$ where

$$\sigma^2(\theta) = \text{Var} \left(\phi(X + \theta) \exp \left(-\frac{1}{2} \theta^T \theta - \theta^T X \right) \right).$$

We assume in the sequel that $\sigma^2(\theta) < \infty$ for any θ . Show that

$$\sigma^2(\theta) = \mathbb{E} \left[\phi^2(X) \exp \left(-\frac{1}{2} X^T X + \frac{1}{2} (X - \theta)^T (X - \theta) \right) \right] - (\mathbb{E}[\phi(X)])^2$$

Let $\sigma^2(\theta) = \text{Var}(Y) = \mathbb{E}[Y^2] - (\mathbb{E}[Y])^2$ to simplify the writing. We already know that $\mathbb{E}[Y] = \mathbb{E}[\phi(X)]$ by the last exercise. Therefore, it remains to prove that

$$\mathbb{E}[Y^2] = \mathbb{E} \left[\phi^2(X) \exp \left(-\frac{1}{2} X^T X + \frac{1}{2} (X - \theta)^T (X - \theta) \right) \right].$$

For that,

$$\begin{aligned} \mathbb{E}[Y^2] &= \int_{\mathbb{R}^d} \phi^2(x + \theta) \exp(-\theta^T \theta - 2\theta^T x) \pi(x) dx \\ &= \int_{\mathbb{R}^d} \frac{1}{(2\pi)^{d/2}} \phi^2(x + \theta) \exp \left(-\theta^T \theta - 2\theta^T x - \frac{1}{2} x^T x \right) dx \\ &= \int_{\mathbb{R}^d} \frac{1}{(2\pi)^{d/2}} \phi^2(x + \theta) \exp \left(-(x + \theta)^T (x + \theta) + \frac{1}{2} x^T x \right) dx \\ &= \int_{\mathbb{R}^d} \frac{1}{(2\pi)^{d/2}} \phi^2(y) \exp \left(-y^T y + \frac{1}{2} (y - \theta)^T (y - \theta) \right) dy \\ &= \mathbb{E} \left[\phi^2(X) \exp \left(-X^T X + \frac{1}{2} (X - \theta)^T (X - \theta) \right) \right], \end{aligned}$$

as we wanted to prove.

3. A twice differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is strictly convex if $\nabla^2 f(\theta)$ (called the Hessian of f) is a positive definite matrix for any $\theta \in \mathbb{R}^d$. Deduce from the expression of $\sigma^2(\theta)$ given in (2) that the function $\theta \rightarrow \sigma^2(\theta)$ is strictly convex.

For that, we will use the derived expression in the last exercise and we differentiate under the expected value using the Leibniz Rule. Then,

$$\nabla_{\theta} \sigma^2(\theta) = \mathbb{E} \left[\phi^2(X) (\theta - X) \exp \left(-X^T X + \frac{1}{2} (X - \theta)^T (X - \theta) \right) \right]$$

and

$$\nabla_{\theta}^2 \sigma^2(\theta) = \mathbb{E} \left[\phi^2(X) \exp \left(-X^T X + \frac{1}{2} (X - \theta)^T (X - \theta) \right) ((\theta - X)^T (\theta - X) + 1) \right],$$

which is clearly positive definite since $(\theta - X)^T (\theta - X)$ is semi definite positive.

4. Show that the minimum of $\theta \rightarrow \sigma^2(\theta)$ is reached at θ^* such that

$$\mathbb{E}[\phi^2(X)(\theta^* - X) \exp(-\theta^{*T} X)] = 0.$$

Since $\sigma^2(\theta)$ is differentiable, its critical points are the solution of $\nabla_{\theta} \sigma^2(\theta) = 0$,

$$\begin{aligned} \mathbb{E} \left[\phi^2(X) (\theta - X) \exp \left(-X^T X + \frac{1}{2} (X - \theta)^T (X - \theta) \right) \right] &= 0 \\ \implies \mathbb{E} \left[\phi^2(X) (\theta - X) \exp \left(-\frac{1}{2} X^T X - \theta^T X \right) \right] &= 0, \end{aligned}$$

since $e^{\theta^T \theta/2}$ is a positive constant. Since the function is strictly convex, we already know that there is only one minimal and it occurs when the above expression is zero.

Exercício 2. (Metropolis-Hastings) Let \mathbb{X} be a finite state-space. Consider the following Markov transition kernel

$$T(x, y) = \alpha(x, y)q(x, y) + \left(1 - \sum_{z \in \mathbb{X}} \alpha(x, z)q(x, z)\right) \delta_x(y)$$

where $q(x, y) \geq 0$, $\sum_{y \in \mathbb{X}} q(x, y) = 1$ and $0 \leq \alpha(x, y) \leq 1$ for any $x, y \in \mathbb{X}$. $\delta_x(y)$ is the Kronecker symbol.

1. Explain how you would simulate a Markov chain with transition kernel T .

Let $x^{(0)} \in \mathbb{X}$ be a initial point. Then, we sample y from $q(x^{(j-1)}, \cdot)$ and $u \sim \text{Unif}[0, 1]$. Then we set $x^{(j)} = x^{(j-1)}$ if $u > \alpha(x, y)$ and $x^{(j)} = y$, otherwise. Notice that

$$\mathbb{P}(X^{(t)} = y \mid X^{(t-1)} = x) = \begin{cases} \alpha(x, y)q(x, y), & \text{if } y \neq x \\ 1 - \sum_{z \in \mathbb{X}/\{x\}} \alpha(x, z)q(x, z), & \text{if } y = x, \end{cases}$$

which is exactly the formula of T .

2. Let π be a probability mass function on \mathbb{X} . Show that if

$$\alpha(x, y) = \frac{\gamma(x, y)}{\pi(x)q(x, y)}$$

where $\gamma(x, y) = \gamma(y, x)$ and $\gamma(x, y)$ is chosen such that $0 \leq \alpha(x, y) \leq 1$ for any $x, y \in \mathbb{X}$ then T is π -reversible.

By Proposition 2.3 from the notes, we have to show that π satisfies detailed balance with respect to T , that is,

$$\pi(x)T(x, y) = \pi(y)T(y, x).$$

If $x = y$, this is clearly true. If $x \neq y$, we have that

$$\begin{aligned} \pi(x)T(x, y) &= \pi(x)\alpha(x, y)q(x, y) \\ &= \gamma(x, y) \\ &= \gamma(y, x) \\ &= \frac{\pi(y)q(y, x)}{\pi(y)q(y, x)}\gamma(y, x) \\ &= \pi(y)\alpha(y, x)q(y, x) \\ &= \pi(y)T(y, x). \end{aligned}$$

3. Show that the Metropolis-Hastings algorithm corresponds to a particular choice of $\gamma(x, y)$.

Setting $\gamma(x, y) = \min\{\pi(x)q(x, y), \pi(y)q(y, x)\}$, we will have that

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

exactly as in the Metropolis-Hastings algorithm. Besides that, it is clearly that $\gamma(x, y) = \gamma(y, x)$.

4. Let π be a probability mass function on the finite space \mathbb{X} such that $\pi(x) > 0$ for any $x \in \mathbb{X}$. To sample from π , we run a Metropolis-Hastings chain $(X^{(t)})_{t \geq 1}$ with proposal $q(x, y) \geq 0$, such that $\sum_{y \in \mathbb{X}} q(x, y) = 1$ and $q(x, x) = 0$ for any $x \in \mathbb{X}$. Consider here the sequence $(Y^{(k)})_{k \geq 1}$ of accepted proposals: $Y^{(1)} = X^{(\tau_1)}$ where $\tau_1 = 1$ and, for $k \geq 2$, $Y^{(k)} = X^{(\tau_k)}$ where $\tau_k := \min\{t : t > \tau_{k-1}, X^{(t)} \neq Y^{(k-1)}\}$. Let $\phi : \mathbb{X} \rightarrow \mathbb{R}$ be a test function. Show that the estimate $\frac{1}{\tau_k - 1} \sum_{t=1}^{\tau_k-1} \phi(X^{(t)})$ can be rewritten as a function of $(Y^{(k)})_{k \geq 1}$ and $(\tau_k)_{k \geq 1}$ and prove that the sequence $(Y^{(k)})_{k \geq 1}$ is a Markov chain with transition kernel

$$K(x, y) = \frac{\alpha(x, y)q(x, y)}{\sum_{z \in \mathbb{X}} \alpha(x, z)q(x, z)}.$$

Notice that $X^{(\tau_i)} = X^{(t)}, \forall t \in [\tau_i, \tau_{i+1})$. Then,

$$\frac{1}{\tau_k - 1} \sum_{t=1}^{\tau_k-1} \phi(X^{(t)}) = \frac{1}{\tau_k - 1} \sum_{i=1}^{k-1} (\tau_{i+1} - \tau_i) \phi(Y^{(i)}).$$

Since $(X^{(t)})_{t \geq 1}$ is a Markov chain, by the Strong Markov Property, if $y \neq x$,

$$\begin{aligned} \mathbb{P}(Y^{(k)} = y \mid Y^{(k-1)} = x, \dots, Y^{(1)} = y_1) &= \mathbb{P}(X^{(\tau_k)} = y \mid X^{(\tau_{k-1})} = x, \dots, X^{(\tau_1)} = y_1) \\ &= \mathbb{P}(X^{(\tau_k)} = y \mid X^{(\tau_{k-1})} = x) \\ &= \mathbb{P}(y \text{ is accepted} \mid \text{some was accepted}, x) \\ &= K(x, y). \end{aligned}$$

5. Show that the transition kernel $K(x, y)$ of the Markov chain $(Y^{(k)})_{k \geq 1}$ is $\tilde{\pi}$ -reversible where

$$\tilde{\pi}(x) = \frac{\pi(x)m(x)}{\sum_{z \in \mathbb{X}} \pi(z)m(z)}$$

with

$$m(x) := \sum_{z \in \mathbb{X}} \alpha(x, z)q(x, z).$$

Let $Z = \sum_{z \in \mathbb{X}} \pi(z)m(z)$. We have that

$$\begin{aligned} \tilde{\pi}(x)K(x, y) &= \frac{\pi(x)m(x)}{Z} \frac{\alpha(x, y)q(x, y)}{m(x)} \\ &= \frac{\pi(x)\alpha(x, y)q(x, y)}{Z} \\ &= \frac{\pi(y)\alpha(y, x)q(y, x)}{Z} \\ &= \frac{\pi(y)m(y)}{Z} \frac{\alpha(y, x)q(y, x)}{m(y)} \\ &= \tilde{\pi}(y)K(y, x), \end{aligned}$$

what proves the detailed balance and so the reversibility.

6. Assume that for some test function $\phi : \mathbb{X} \rightarrow \mathbb{R}$ we have $\frac{1}{k} \sum_{i=1}^k \phi(Y^{(i)})$ converges to $\sum_{x \in \mathbb{X}} \phi(x) \tilde{\pi}(x)$ almost surely and additionally assume that $m(x)$ can be computed exactly for any $x \in \mathbb{X}$. Propose a strongly consistent estimate of $\sum_{x \in \mathbb{X}} \phi(x) \pi(x)$ based on the Markov chain $(Y^{(k)})_{k \geq 1}$ which does not rely on $(\tau_k)_{k \geq 1}$.

Notice that

$$\sum_{x \in \mathbb{X}} \phi(x) \pi(x) = \frac{\sum_{x \in \mathbb{X}} \phi(x) \frac{\pi(x)}{\tilde{\pi}(x)} \tilde{\pi}(x)}{\sum_{x \in \mathbb{X}} \frac{\pi(x)}{\tilde{\pi}(x)} \tilde{\pi}(x)},$$

since π is a probability distribution. Notice that this can be turned to a Importance Sampling. We have that $\tilde{\pi}$ is the stationary distribution of $(Y^{(k)})_{k \geq 1}$, which we already know how to sample using a Metropolis-Hastings algorithm. It remains to calculate the weights $\pi(x)/\tilde{\pi}(x)$, which can be calculated from

$$\frac{\pi(x)}{\tilde{\pi}(x)} = \frac{Z}{m(x)},$$

and

$$\sum_{x \in \mathbb{X}} \phi(x) \pi(x) = \frac{\sum_{x \in \mathbb{X}} \phi(x) \frac{Z}{m(x)} \tilde{\pi}(x)}{\sum_{x \in \mathbb{X}} \frac{Z}{m(x)} \tilde{\pi}(x)} = \frac{\sum_{x \in \mathbb{X}} \phi(x) \frac{1}{m(x)} \tilde{\pi}(x)}{\sum_{x \in \mathbb{X}} \frac{1}{m(x)} \tilde{\pi}(x)}.$$

Hence, a strongly consistent estimate is

$$\frac{\sum_{i=1}^k \phi(Y^{(i)}) \frac{1}{m(Y^{(i)})}}{\sum_{i=1}^k \frac{1}{m(Y^{(i)})}}.$$

Exercício 3.

Exercício 4. (Gibbs Sampler) Suppose that we wish to use the Gibbs sampler on

$$\pi(x, y) \propto \exp \left(-\frac{1}{2}(x-1)^2(y-2)^2 \right).$$

1. Write down the two “full” conditional distributions associated to $\pi(x, y)$.

$$\begin{aligned} \pi(x) &= \int_{y \in \mathbb{R}} \pi(x, y) dy \\ &= \int_{y \in \mathbb{R}} c \exp \left(-\frac{1}{2}(x-1)^2(y-2)^2 \right) dy \\ &= c \int_{y \in \mathbb{R}} \exp \left(-\frac{1}{2}(x-1)^2(y-2)^2 \right) dy \\ &= c \int_{y \in \mathbb{R}} \exp \left(-\frac{1}{2/(x-1)^2}(y-2)^2 \right) dy \\ &= c \sqrt{2\pi} \frac{1}{|x-1|} \end{aligned}$$

since setting $\sigma^2 = 1/(x-1)^2$, the integrand is the kernel of a normal distribution with mean 2 and variance σ^2 . Therefore, its integral is the normalization constant. With the same reasoning, we have that

$$\pi(y) = c \sqrt{2\pi} \frac{1}{|y-2|}.$$

Then, we have that

$$\pi(x|y) = \frac{\pi(x, y)}{\pi(y)} = \frac{1}{\sqrt{2\pi}} |y - 2| \exp\left(-\frac{1}{2}(x - 1)^2(y - 2)^2\right)$$

and

$$\pi(y|x) = \frac{\pi(x, y)}{\pi(x)} = \frac{1}{\sqrt{2\pi}} |x - 1| \exp\left(-\frac{1}{2}(x - 1)^2(y - 2)^2\right),$$

which implies that

$$X | Y = y \sim \text{Normal}(1, |y - 2|^2)$$

and

$$Y | X = x \sim \text{Normal}(2, |x - 1|^2).$$

2. Does the resulting Gibbs sampler make any sense?

The problem with that Gibbs sampler is that the samples (x^n, y^n) converges to $(1, 2)$ when n is sufficiently high. This happens because the variances are very low in the region of greater mass. Therefore, even if the initial points are far from the mode, the sampling will explore it and when it happens, it will get stuck.

Exercício 5. (Gibbs Sampler) For $i = 1, \dots, T$ consider $Z_i = X_i + Y_i$ with independent X_i, Y_i such that

$$X_i \sim \text{Bin}(m_i, \theta_1), Y_i \sim \text{Bin}(n_i, \theta_2).$$

1. We assume $0 \leq z_i \leq m_i + n_i$ for $i = 1, \dots, T$. We observe z_i for $i = 1, \dots, T$ and the n_i, m_i for $i = 1, \dots, T$ are given. Give the expression of the likelihood function $p(z_1, \dots, z_T | \theta_1, \theta_2)$.

Supposing conditionally independent samples, we have that

$$p(z_1, \dots, z_T | \theta_1, \theta_2) = \prod_{i=1}^T p(z_i | \theta_1, \theta_2)$$

Note that

$$\begin{aligned} \mathbb{P}(Z_i = X_i + Y_i = k) &= \sum_{j=0}^k \mathbb{P}(X_i = j) \mathbb{P}(Y_i = k - j) \\ &= \sum_{j=0}^k \binom{m_i}{j} \binom{n_i}{k-j} \theta_1^j (1 - \theta_1)^{m_i-j} \theta_2^{k-j} (1 - \theta_2)^{n_i-k+j}. \end{aligned}$$

Therefore,

$$p(z_1, \dots, z_T | \theta_1, \theta_2) = \prod_{i=1}^T \sum_{j=0}^{z_i} \binom{m_i}{j} \binom{n_i}{z_i-j} \theta_1^j (1 - \theta_1)^{m_i-j} \theta_2^{z_i-j} (1 - \theta_2)^{n_i-z_i+j}.$$

2. Assume we set independent uniform priors $\theta_1 \sim \text{Unif}[0, 1], \theta_2 \sim \text{Unif}[0, 1]$. Propose a Gibbs sampler to sample from $p(\theta_1, \theta_2 | z_1, \dots, z_T)$.

We know that

$$p(\theta_1, \theta_2 | z_1, \dots, z_T) \propto p(z_1, \dots, z_T | \theta_1, \theta_2).$$

and we want to specify

$$p(\theta_1^t \mid \theta_2^{t-1}, z_1, \dots, z_T)$$

and

$$p(\theta_2^t \mid \theta_1^t, z_1, \dots, z_T).$$

Since we do not know how to sample from these distributions, we have to introduce auxiliary variables to help out. Let X_1, \dots, X_T and Y_1, \dots, Y_T be the variables introduced in the exercise. We denote $X_{1:T} = X_1, \dots, X_T$. Then

$$\mathbb{P}(\theta_1 \mid \theta_2, X_{1:T}, Y_{1:T}, Z_{1:T}) = \mathbb{P}(\theta_1 \mid X_{1:T}),$$

since given $X_{1:T}$, we have that θ_1 is independent of the other random variables. Analogously,

$$\mathbb{P}(\theta_2 \mid \theta_1, X_{1:T}, Y_{1:T}, Z_{1:T}) = \mathbb{P}(\theta_2 \mid Y_{1:T}).$$

Given that $\{X_i \mid \theta_1\}_{i=1}^T$ are independent and have binomial distribution, and θ_1 has Beta distribution with parameters $\alpha = 1$ and $\beta = 1$, we know that,

$$\theta_1 \mid X_{1:T} = x_{1:T} \sim \text{Beta} \left(1 + \sum_{i=1}^T x_i, 1 + \sum_{i=1}^T m_i - x_i \right)$$

and

$$\theta_2 \mid Y_{1:T} = y_{1:T} \sim \text{Beta} \left(1 + \sum_{i=1}^T y_i, 1 + \sum_{i=1}^T n_i - y_i \right).$$

It remains to sample from

$$X_{1:T}, Y_{1:T} \mid \theta_1, \theta_2, Z_{1:T}.$$

Note that

$$\begin{aligned} \mathbb{P}(X_{1:T} = x_{1:T}, Y_{1:T} = y_{1:T} \mid \theta_1, \theta_2, Z_{1:T}) &= \mathbb{P}(X_{1:T} = x_{1:T} \mid \theta_1, \theta_2, Z_{1:T}, Y_{1:T}) \\ &\quad \times \mathbb{P}(Y_{1:T} = y_{1:T} \mid \theta_1, \theta_2, Z_{1:T}). \end{aligned}$$

The first term is $1_{\{x_{1:T}=Z_{1:T}-Y_{1:T}\}}$, while the second can be writing as follows:

$$\begin{aligned} \mathbb{P}(Y_{1:T} = y_{1:T} \mid \theta_1, \theta_2, Z_{1:T}) &= \prod_{i=1}^T \mathbb{P}(Y_i = y_i \mid \theta_1, \theta_2, Z_i) \\ &= \prod_{i=1}^T \text{Bin}(y_i; n_i, \theta_2) \cdot \text{Bin}(z_i - y_i; m_i, \theta_1). \end{aligned}$$

Since this distribution has finite support, one can calculate its constant. Therefore, we have defined our Gibbs sampler.