
TP 4: Improve the Metropolis-Hastings algorithm

Exercise 1: Adaptive Metropolis-Hastings within Gibbs sampler

MCMC samplers, such as the Metropolis-Hastings algorithm or Gibbs sampler, require that the user specifies a transition kernel with a given invariant distribution (the target distribution). These transition kernels usually depend on parameters which are to be given and tuned by the user. In practice, it is often difficult (if not impossible) to find the best parameters for such algorithms given a target distribution. Moreover, if the parameters are not carefully chosen, it may result in a MCMC algorithm performing poorly as in part A. *Adaptive MCMC algorithms* is a class of MCMC algorithms which addresses the problem of parameter tuning by updating automatically some of (if not all) the parameters.

1.A – Metropolis-Hastings within Gibbs sampler

We aim to sample on \mathbb{R}^2 the target distribution π given by

$$d\pi(x, y) \propto \exp\left(-\frac{x^2}{a^2} - y^2 - \frac{1}{4}\left(\frac{x^2}{a^2} - y^2\right)^2\right) dx dy$$

where $a > 0$. We consider a Markov transition kernel P defined by

$$P = \frac{1}{2}(P_1 + P_2)$$

where $P_i((x, y), dx' dy')$ for $i = 1, 2$ is the Markov transition kernel which only updates the i -th component: this update follows a symmetric random walk proposal mechanism and uses a Gaussian distribution with variance σ_i^2 .

1. Implement the MCMC algorithm with kernel P and Hastings-Metropolis filter.
2. Run the algorithm with $a = 10$, see Figure , and standard deviations of the proposal distributions chosen as follows: $(\sigma_1, \sigma_2) = (3, 3)$. Plot (i) the trajectory of the Markov chain along the coordinates as well as (ii) the auto-correlation diagram on each coordinate (see `pandas.DataFrame.corr`). Discuss the performance of the algorithm in this situation.
3. How could the performance of the above algorithm be improved ? Propose two methods.

1.B – Adaptive Metropolis-Hastings within Gibbs sampler

Let π be a density defined on an open set \mathcal{U} of \mathbb{R}^d , $d \geq 2$. We consider here a Metropolis-Hastings within Gibbs algorithm to sample from the target density π . More precisely, the HM-step is a symmetric random walk one and the proposal distribution is a Gaussian distribution centered at the current state.

As usual, for $i \in \llbracket 1, d \rrbracket$, let π_i denote the i -th full conditional of π , which is given by:

$$x_{-i} = \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d\} \quad ; \quad \pi_i(x_i | x_{-i}) \propto \pi(x)$$

and σ_i^2 the variance of the corresponding proposal distribution.

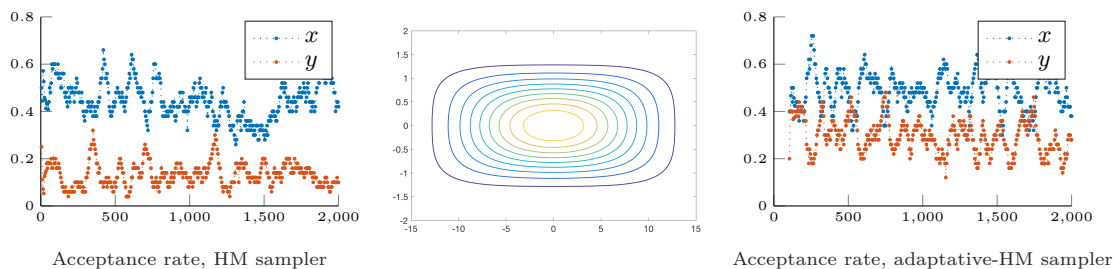


Figure 1: Mean acceptance rate and contour plot of the density – $a = 10$

Algorithm 1: Metropolis-Hastings (symmetric random walk) within Gibbs Sampler

```

1 Given  $x^{(k)} = (x_1^{(k)}, \dots, x_d^{(k)})$ 
2 for  $i = 1$  to  $d$  do
3   HM to sample from the target  $x_i^{(k+1)} \sim \pi_i(x_i | x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_{i+1}^{(k)}, \dots, x_d^{(k)})$ :
4   Proposal:  $x_i^* \sim \mathcal{N}(x_i^{(k)}, \sigma_i^2)$ 
5   Acceptance ratio  $\alpha(x_i^*, x_i^{(k)}) = \frac{\pi_i(x_i^* | x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_{i+1}^{(k)}, \dots, x_d^{(k)})}{\pi_i(x_i^{(k)} | x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_{i+1}^{(k)}, \dots, x_d^{(k)})} \wedge 1$ 
6 end
    
```

In [RR09], the authors propose an adaptive version of the above sampler which automatically adjusts the variances $\sigma_1^2, \dots, \sigma_d^2$ of the proposal distributions. We proceed as follows:

- For each of the variables x_i , we create an associate variable ℓ_i defined by $\ell_i := \log(\sigma_i)$;
- We initialize all ℓ_i to zero, which correspond to the unit proposal variance ;
- After the j -th ($j \in \mathbb{N}^*$) batch of 50 iterations, each variable ℓ_i is updated by adding or subtracting an amount $\delta(j)$ in order to make the acceptance rate of proposals for variable x_i as close as possible to 0.24 (proved to be optimal for one-dimensional proposals in certain settings). Specifically, if the acceptance rate for the i -th variable is greater than 0.24, ℓ_i is increased with $\delta(j)$. Otherwise, if the rate is lower than 0.24, ℓ_i is decreased by $\delta(j)$.

In practice, we take $\delta(j) := \min(0.01, j^{-1/2})$.

1. Implement the adaptive Metropolis-Hastings within Gibbs sampler and test the algorithm on the density π defined in the part A. Plot for each coordinate: the evolution of the average acceptance rate, the trajectory of the Markov chain and the auto-correlation diagrams on each coordinate. Compare the performance of the algorithm with or without adaptation.
2. Also compare the performance of this algorithm on (i) the centered d -dimensional Gaussian $\mathcal{N}(0, \Sigma)$, where $d = 20$ and Σ is given in the file <http://dept.stat.lsa.umich.edu/~yvesa/tmalaexcov.txt> and (ii) the "banana"-shaped density ($d = 20, B = 0.1$):

$$\forall x = (x_1, \dots, x_d) \in \mathbb{R}^d, \quad f_B(x) \propto \exp\left(-\frac{x_1^2}{200} - \frac{1}{2}(x_2 + Bx_1^2 - 100B)^2 - \frac{1}{2}(x_3^2 + \dots + x_d^2)\right).$$

To go further . . .

The next improvement of the Metropolis-Hastings algorithm we can make is to consider a *drift* function in the proposal distribution. Given a positive definite matrix Λ and a scale parameter $\sigma > 0$, we consider a proposal distribution of the form:

$$q_{\sigma, \Lambda}(y | x) = \frac{1}{(\sigma\sqrt{2\pi})^d} \frac{1}{\sqrt{\det(\Lambda)}} \exp\left(-\frac{1}{2\sigma^2} \left[y - x - \frac{\sigma^2}{2} \Lambda D(x)\right]^\top \Lambda^{-1} \left[y - x - \frac{\sigma^2}{2} \Lambda D(x)\right]\right).$$

$q_{\sigma, \Lambda}$ is the density (with respect to the Lebesgue measure on \mathbb{R}^d) of the d -dimensional Gaussian distribution with mean $x + \frac{\sigma^2}{2} \Lambda D(x)$ and variance-covariance matrix $\sigma^2 \Lambda$. If D vanishes everywhere ($D \equiv 0$), the corresponding algorithm is a *Metropolis-Hastings Symmetric Random Walk*. If the drift D is chosen such that:

$$\forall x \in \mathbb{R}^d, \quad D(x) = \frac{\delta}{\max(\delta, \|\nabla \log \pi(x)\|)} \nabla \log \pi(x)$$

for a constant $\delta > 0$, the corresponding algorithm is a *Metropolis Adjusted Langevin Algorithm* (MALA). In that case, the proposal distribution includes information on the gradient $\nabla \log \pi$ of the target distribution π . In [Atc06], the author proposes an adaptive version of the MALA algorithm in which the parameters σ and Λ are adjusted automatically.

Exercise 2: Sampling from multimodal distributions

We consider a target distribution π with support $\mathcal{U} \subset \mathbb{R}^d$ ($d \in \mathbb{N}^*$). When the target distribution is multi-modal, especially with well-separated modes, classical MCMC algorithms can perform very poorly and exhibit poor mixing. Indeed, a Metropolis-Hastings algorithm with local proposal can get stuck for a long time in a local mode of the target distribution.

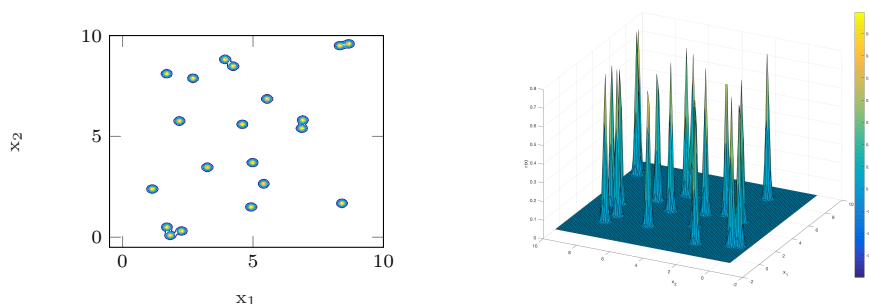


Figure 2: Mixture of 20 Gaussian distributions

2.A – A toy example

In the following, we consider a target distribution π – taken from [LW01] and plotted at Figure 2 – defined on \mathbb{R}^2 as a mixture of 20 Gaussian distributions. The target distribution writes:

$$\pi(\mathbf{x}) = \sum_{i=1}^{20} \frac{w_i}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2\sigma_i^2}\|\mathbf{x} - \boldsymbol{\mu}_i\|^2\right)$$

where, $\forall i \in \{1, \dots, 20\}$, $w_i = 0.05$ and $\sigma_i = 0.1$. The 20 means $\boldsymbol{\mu}_i$ are defined as follows:

$$\begin{aligned} (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_{20}) = & \left(\begin{pmatrix} 2.18 \\ 5.76 \end{pmatrix}, \begin{pmatrix} 8.67 \\ 9.59 \end{pmatrix}, \begin{pmatrix} 4.24 \\ 8.48 \end{pmatrix}, \begin{pmatrix} 8.41 \\ 1.68 \end{pmatrix}, \begin{pmatrix} 3.93 \\ 8.82 \end{pmatrix}, \begin{pmatrix} 3.25 \\ 3.47 \end{pmatrix}, \begin{pmatrix} 1.70 \\ 0.50 \end{pmatrix}, \right. \\ & \begin{pmatrix} 4.59 \\ 5.60 \end{pmatrix}, \begin{pmatrix} 6.91 \\ 5.81 \end{pmatrix}, \begin{pmatrix} 6.87 \\ 5.40 \end{pmatrix}, \begin{pmatrix} 5.41 \\ 2.65 \end{pmatrix}, \begin{pmatrix} 2.70 \\ 7.88 \end{pmatrix}, \begin{pmatrix} 4.98 \\ 3.70 \end{pmatrix}, \begin{pmatrix} 1.14 \\ 2.39 \end{pmatrix}, \\ & \left. \begin{pmatrix} 8.33 \\ 9.50 \end{pmatrix}, \begin{pmatrix} 4.93 \\ 1.50 \end{pmatrix}, \begin{pmatrix} 1.83 \\ 0.09 \end{pmatrix}, \begin{pmatrix} 2.26 \\ 0.31 \end{pmatrix}, \begin{pmatrix} 5.54 \\ 6.86 \end{pmatrix}, \begin{pmatrix} 1.69 \\ 8.11 \end{pmatrix} \right). \end{aligned}$$

1. Write a Metropolis-Hastings Symmetric Random Walk algorithm (you may use your code from previous tutorial classes) to sample from π .
2. Show that the Metropolis-Hastings algorithm (even the adaptive Metropolis-Hastings algorithm) fails to sample from π .

2.B – Parallel Tempering

The general idea of the Parallel Tempering (PT) [Gey91, ED05] algorithm is to use *tempered* versions of the distribution π and run parallel Metropolis-Hastings algorithms to sample from these tempered distributions. The tempered distributions are obtained by "warming up" the target distribution π at different *temperatures*. At each iteration of the algorithm, a *swap* between two chains (chains running at different temperature levels) is proposed. The Parallel Tempering uses the fast mixing of the chains at high temperature to improve the mixing of the chains at low temperatures.

Let K denote a positive integer. We consider a sequence of temperatures $(T_i)_{1 \leq i \leq K}$ such that:

$$T_1 > T_2 > \dots > T_K = 1.$$

In the Parallel Tempering algorithm, K chains run in parallel: for $i \in \llbracket 1, K \rrbracket$, the i -th chain targets the tempered distribution $\pi_i := \pi^{1/T_i}$; the distribution of interest corresponds to the lowest temperature, $T_K = 1$. Let $(X_n^{(i)})_{n \in \mathbb{N}}$ denote the i -th chain, sampling from the tempered distribution π_i .

At the n -th iteration of the Parallel Tempering algorithm, a candidate $Y_{n+1}^{(i)}$ for the i -th chain is proposed using the transition kernel $P^{(i)}(X_n^{(i)}, \cdot)$ of a Metropolis-Hastings algorithm. The next step consists in proposing a swap between two different chains (running at different temperatures): given $(i, j) \in \llbracket 1, K \rrbracket^2$, with $i \neq j$, a swap is proposed with probability $\alpha(i, j)$.

1. Implement the Parallel Tempering algorithm.

Algorithm 2: Parallel Tempering

```

1 For all  $i \in \llbracket 1, K \rrbracket$ , initialize  $X_0^{(i)}$  ;
2 for  $n = 1$  to  $N_{\text{iter}}$  do
3   For all  $i \in \llbracket 1, K \rrbracket$ , draw  $Y_{n+1}^{(i)}$  using the transition kernel  $P^{(i)}(X_n^{(i)}, \cdot)$  and standard MH filter;
4   Choose uniformly  $(i, j) \in \llbracket 1, K \rrbracket^2$ , with  $|i - j| = 1$  ;
5   Compute the swap acceptance probability  $\alpha(i, j) = \min \left( 1, \frac{\pi_i(Y_{n+1}^{(j)}) \pi_j(Y_{n+1}^{(i)})}{\pi_i(Y_{n+1}^{(i)}) \pi_j(Y_{n+1}^{(j)})} \right)$  ;
6   Draw  $U \sim \mathcal{U}([0, 1])$  ;
7   if  $U \leq \alpha(i, j)$  then
8     |
9     |            $X_{n+1}^{(i)} = Y_{n+1}^{(j)}$    and    $X_{n+1}^{(j)} = Y_{n+1}^{(i)}$  ;
10    |
11    | else
12    |            $X_{n+1}^{(i)} = Y_{n+1}^{(i)}$    and    $X_{n+1}^{(j)} = Y_{n+1}^{(j)}$  ;
13    | end
14   For all  $k \in \llbracket 1, K \rrbracket$ ,  $k \neq i, j$ , set  $X_{n+1}^{(k)} = Y_{n+1}^{(k)}$  .
15 end

```

2. In order to illustrate the performance of the algorithm, use your code to sample from the distribution π of Part A. Use the algorithm with $K = 5$ and with the following temperatures ladder:

$$(T_1, \dots, T_5) = (60, 21.6, 7.7, 2.8, 1).$$

For the Metropolis-Hastings step (Line 3), take as proposal distribution the bivariate Gaussian distribution centered at $X_n^{(i)}$, with variance-covariance matrix $\tau_i^2 \mathbf{I}_2$:

$$\forall i \in \llbracket 1, K \rrbracket, \quad Y_{n+1}^{(i)} \sim \mathcal{N}_{\mathbb{R}^2}(X_n^{(i)}, \tau_i^2 \mathbf{I}_2) \quad \text{where} \quad \tau_i = 0.25 \sqrt{T_i}.$$

The scale parameters τ_i are tuned to ensure a reasonable acceptance rate in the algorithm.

In practice, the performance of the Parallel Tempering algorithm strongly depends on the choice of the temperatures ladder, the number of chains and the choice of proposal kernels. For most distributions, tuning these parameters may be infeasible. In [MMV13], the authors have proposed an adaptive Parallel Tempering algorithm to address these difficulties.

Exercise 3: Bayesian analysis of a one-way random effects model

We recall that the density of an *Inverse-Gamma* distribution with positive parameters (a, b) is proportional to

$$x \mapsto \frac{1}{x^{a+1}} \exp\left(-\frac{b}{x}\right) \mathbf{1}_{\mathbb{R}^+}(x).$$

We recall that we can sample y from this distribution by generating x from a gamma distribution of parameters $(a, \frac{1}{b})$ and then taking $y = \frac{1}{x}$. In practice, consider `scipy.stats.invgamma` or the Statistics and Machine Learning Toolbox in Matlab.

Assume that we are provided with observations $Y = \{y_{i,j}, i \in \llbracket 1, N \rrbracket, j \in \llbracket 1, k_i \rrbracket\}$. Define $k := \sum_{i=1}^N k_i$ as the total number of observations. We consider the following random effects model:

- (i) $y_{i,j}$ is a realization of the variable $Y_{i,j}$ where $Y_{i,j} = X_i + \varepsilon_{i,j}$;
- (ii) The random effects $X = \{X_i, i \in \llbracket 1, N \rrbracket\}$ are i.i.d from a Gaussian $\mathcal{N}(\mu, \sigma^2)$ and independent of the errors $\varepsilon = \{\varepsilon_{i,j}, i \in \llbracket 1, N \rrbracket, j \in \llbracket 1, k_i \rrbracket\}$;
- (iii) The errors ε are i.i.d from the centered Gaussian $\mathcal{N}(0, \tau^2)$;

where (μ, σ, τ) **are the unknown parameters**. Bayesian analysis using this model requires specifying a prior distribution, for which we consider:

$$\pi_{prior}(\mu, \sigma^2, \tau^2) \propto \frac{1}{\sigma^{2(1+\alpha)}} \exp\left(-\frac{\beta}{\sigma^2}\right) \frac{1}{\tau^{2(1+\gamma)}} \exp\left(-\frac{\beta}{\tau^2}\right)$$

where α, β **and** γ **are known hyper-parameters** (Inverse-gamma on σ^2 , Inverse-gamma on τ^2 and no prior on μ).

1. Write the density of the *a posteriori* distribution $(X, \mu, \sigma^2, \tau^2)$ up to a normalizing constant.
2. Implement a Gibbs sampler which updates in turn $(\sigma^2, \tau^2, \mu, X)$ one at a time.
3. Implement a Block-Gibbs sampler which updates σ^2 , then τ^2 and then the block (X, μ) .
4. Discuss the theoretical performance of these two algorithms.
5. Test your code on a synthetic dataset $Y = \{y_{i,j}, i \in \llbracket 1, N \rrbracket, j \in \llbracket 1, k_i \rrbracket\}$ generated from the previous model.

References

- [Atc06] Yves F. Atchadé. An adaptive version for the metropolis adjusted langevin algorithm with a truncated drift. *Methodology and Computing in Applied Probability*, 8(2):235–254, 2006.
- [ED05] David J Earl and Michael W Deem. Parallel tempering: Theory, applications, and new perspectives. *Physical Chemistry Chemical Physics*, 7(23):3910–3916, 2005.
- [Gey91] Charles J Geyer. Markov chain monte carlo maximum likelihood. 1991.
- [KTB13] D.P. Kroese, T. Taimre, and Z.I. Botev. *Handbook of Monte Carlo Methods*. Wiley Series in Probability and Statistics. Wiley, 2013.
- [LW01] Faming Liang and Wing Hung Wong. Real-parameter evolutionary monte carlo with applications to bayesian mixture models. *Journal of the American Statistical Association*, 96(454):653–666, 2001.
- [MMV13] Błażej Miasojedow, Eric Moulines, and Matti Vihola. An adaptive parallel tempering algorithm. *Journal of Computational and Graphical Statistics*, 22(3):649–664, 2013.
- [RR09] Gareth O. Roberts and Jeffrey S. Rosenthal. Examples of adaptive mcmc. *Journal of Computational and Graphical Statistics*, 18(2):349–367, 2009.