

TP 4 : Improve the Metropolis-Hastings algorithm

MCMC samplers, such as the Metropolis-Hastings algorithm or Gibbs sampler, require that the user specify 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 even 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 exercise 1. *Adaptive MCMC algorithms* is a class of MCMC algorithms which address the problem of parameters tuning by updating automatically some of (if not all) the parameters.

Exercise 1: Metropolis-Hastings within Gibbs sampler

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

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

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

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

where $P_i((x, y); dx' \times dy')$ 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 an algorithm which samples the distribution $P_1(z; \cdot)$ where $z \in \mathbb{R}^2$; likewise for the distribution $P_2(z; \cdot)$. Then, implement an algorithm which samples a chain with kernel P .
2. Run the algorithm with $a = 10$ and standard deviations of the proposal distributions chosen as follows : $(\sigma_1, \sigma_2) = (3, 3)$. Discuss the performance of the algorithm in this situation.
3. How could the performance of the above algorithm be improved ? Propose two methods.

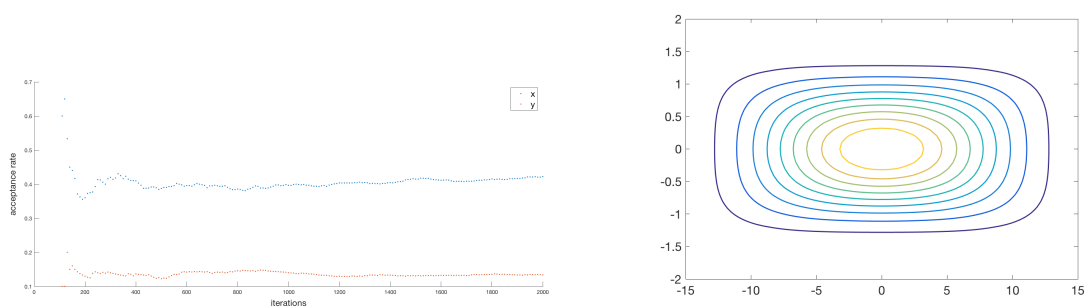


Figure 1: Mean acceptance rate and contour plot of the density – $a = 10$

Exercise 2: Adaptive Metropolis-Hastings within Gibbs sampler

Let π be a density defined on an open set \mathcal{U} of \mathbb{R} , $d \geq 2$. We consider here a Metropolis-Hastings within Gibbs algorithm to sample from the target density π . More precisely, the Metropolis-Hastings step is a symmetric random walk one — as in the population model of TP 3 — 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.

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 propose as follows :

- For each of the variables x_i , we create an associate variable ℓ_i giving for the logarithm of the standard deviation σ_i to be used when proposing a normal increment to variable : $\ell_i := \log(\sigma_i)$;
- 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 $\delta(j)$. The adapting attempts to make the acceptance rate of proposals for variable x_i as close as possible to 0.44, which is optimal for one-dimensional proposals in certain settings according to [RR09]. Specifically, if the acceptance rate for the i -th variable is greater than 0.44, ℓ_i is increased with $\delta(j)$. Otherwise, if ℓ_i is lower than 0.44, ℓ_i is decreased by $\delta(j)$.

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

1. Implement the adaptative Metropolis-Hastings within Gibbs sampler and test the algorithm on the density π defined in the exercise 1 : Using auto-correlation plots (use a built-in function), compare the performance of the algorithm with or without adaptation.
2. We can also compare the performance of our algorithm on more complicated target densities. For example centered d -dimensional Gaussian $\mathcal{N}(0, \Sigma)$ or "banana"-shaped density as in TP 2 :

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

In practice, you can choose $d = 20$ and $B = 0.1$ for the density f_B . You may also choose $d = 20$ for the second one and use the 20×20 variance-covariance matrix Σ given in the file <http://dept.stat.lsa.umich.edu/~yvesa/tmalaexcov.txt>.

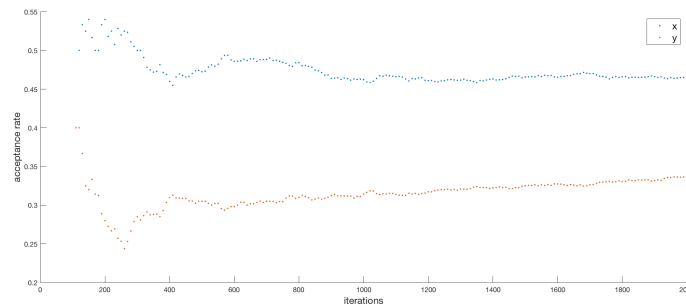


Figure 2: Mean acceptance rate – $a = 10$ – Exercise 1

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 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

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

and especially that we can sample y from the inverse gamma distribution of parameters (a, b) by generating x from a gamma distribution of parameters $(a, \frac{1}{b})$ and then taking $y = \frac{1}{x}$. We also recall that we can find a Matlab toolbox-free Gamma Generator in the Handbook of Monte Carlo Methods [KTB13] : <https://people.smp.uq.edu.au/DirkKroese/montecarlohandbook/probdist/>. Otherwise, we can use directly `scipy.stats.invgamma` or the Statistics and Machine Learning Toolbox in Matlab.

Let suppose we collect the observations $Y = \{y_{i,j}, i \in \llbracket 1, N \rrbracket, j \in \llbracket 1, k_i \rrbracket\}$ and set $k := \sum_{i=1}^N k_i$ the total number of observations. Let 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 centred 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.

1. Write the density of the *a posteriori* distribution $(X, \mu, \sigma^2, \tau^2)$ — it can be given up to a normalizing constant — *i.e* the density of the distribution $(Y, X, \mu, \sigma^2, \tau^2)$.
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 at each iterations σ^2 , then τ^2 and then the block (X, μ) .
4. Discuss the performance of these two algorithms. As previously, you can use the auto-correlation (built-in) function.
5. Test your algorithm 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.
- [KTB13] D.P. Kroese, T. Taimre, and Z.I. Botev. *Handbook of Monte Carlo Methods*. Wiley Series in Probability and Statistics. Wiley, 2013.
- [RR09] Gareth O. Roberts and Jeffrey S. Rosenthal. Examples of adaptive mcmc. *Journal of Computational and Graphical Statistics*, 18(2):349–367, 2009.