

Lecture Notes for STAT546: Computational Statistics

—Lecture 5: Monte Carlo

Faming Liang

Purdue University

August 28, 2024

Slice sampler:

Suppose that one is interested in sampling from a density $f(x)$, $x \in \mathcal{X}$. Recall that sampling $x \sim f(x)$ is equivalent to sampling uniformly from the area under the graph $f(x)$:

$$\mathcal{A} = \{(x, u) : 0 \leq u \leq f(x)\},$$

which is the basis of the acceptance-rejection algorithm described in §???. To achieve this goal, one can augment the target distribution by an auxiliary variable U , which, conditional on x , is uniformly distributed on the interval $[0, f(x)]$. Therefore, the joint density function of (X, U) is

$$f(x, u) = f(x)f(u|x) \propto \mathbf{1}_{(x,u) \in \mathcal{A}},$$

which can be sampled using the Gibbs sampler as follows:

Slice Sampler

- ▶ Draw $u_{t+1} \sim \text{Uniform}[0, f(x_t)]$.
- ▶ Draw x_{t+1} uniformly from the region $\{x : f(x) \geq u_{t+1}\}$.

This sampler is called the slice sampler by Higdon (1998), which potentially can be more efficient than the simple MH algorithm for the multimodal distributions, due to the free between-mode-transitions within a slice (as illustrated by Figure 1).

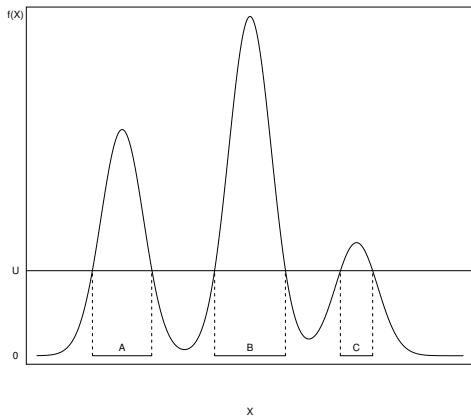


Figure 1: The slice sampler for a multimodal distribution: Draw $X|U$ uniformly from the sets labeled by A , B and C .

Edwards and Sokal (1988) noted that when $f(x)$ can be decomposed into a product of k distribution functions, i.e., $f(x) \propto f_1(x) \times f_2(x) \times \cdots \times f_k(x)$, the slice sampler can be easily implemented. To sample from such a distribution, Edwards and Sokal introduced k auxiliary variables, U_1, \dots, U_k , and proposed the following algorithm:

- ▶ Draw $u_{t+1}^{(i)} \sim \text{Uniform}[0, f_i(x_t)]$, $i = 1, \dots, k$.
- ▶ Draw x_{t+1} uniformly from the region $\mathcal{S} = \cap_{i=1}^k \{x : f_i(x) \geq u_{t+1}^{(i)}\}$.

This algorithm is a generalization of the Swendsen-Wang algorithm (Swendsen and Wang, 1987).

Consider a 2-D Ising model with the Boltzmann density

$$f(\mathbf{x}) \propto \exp \left\{ K \sum_{i \sim j} x_i x_j \right\}, \quad (1)$$

where the spins $x_i = \pm 1$, K is the inverse temperature, and $i \sim j$ represents the nearest neighbors on the lattice.

When the temperature is high, this model can be easily simulated using the Gibbs sampler (Geman and Geman, 1984): Iteratively reset the value of each spin according to the conditional distribution

$$P(x_i = 1 | x_j, j \in n(i)) = \frac{1}{1 + \exp\{-2K \sum_{j \in n(i)}\}}, \quad (2)$$
$$P(x_i = -1 | x_j, j \in n(i)) = 1 - P(x_i = 1 | x_j, j \in n(i)),$$

where $n(i)$ denotes the set of neighbors of spin i . However, the Gibbs sampler slows down rapidly when the temperature is approaching or below the critical temperature. This is the so-called “critical slowing down”.

Swendsen-Wang algorithm

As in slice sampling, the density (1) is rewritten in the form

$$f(\mathbf{x}) \propto \prod_{i \sim j} \exp\{K(1 + x_i x_j)\} = \prod_{i \sim j} \exp\{\beta I(x_i = x_j)\}, \quad (3)$$

where $\beta = 2K$ and $I(\cdot)$ is the indicator function, as $1 + x_i x_j$ is either 0 or 2. If we introduce auxiliary variables $\mathbf{u} = (u_{i \sim j})$, where each component $u_{i \sim j}$, conditional on x_i and x_j , is uniformly distributed on $[0, \exp\{\beta I(x_i = x_j)\}]$, then

$$f(\mathbf{x}, \mathbf{u}) \propto \prod_{i \sim j} I(0 \leq u_{i \sim j} \leq \exp\{\beta I(x_i = x_j)\}).$$

In Swendsen and Wang (1987), $u_{i \sim j}$ is called a “bond variable”, which can be viewed as a variable physically sitting on the edge between spin i and spin j .

If $u_{i\sim j} > 1$, then $\exp\{\beta I(x_i = x_j)\} > 1$ and there must be $x_i = x_j$. Otherwise, there is no constraint on x_i and x_j . Let $b_{i\sim j}$ be an indicator variable for the constraint. If x_i and x_j are constrained to be equal, we set $b_{i\sim j} = 1$ and 0 otherwise. Note that for any two “like-spin” (i.e., the two spins have the same values) neighbors, they are bonded with probability $1 - \exp(-\beta)$. Based on the configurations of \mathbf{u} , we can “cluster” the spins according to whether they are connected via a “mutual bond” (i.e., $b_{i\sim j} = 1$). Then all the spins within the same cluster will have identical values, and flipping all the spins in a cluster simultaneously will not change the equilibrium of $f(\mathbf{x}, \mathbf{u})$.

The Swendsen-Wang Algorithm

- ▶ *Update bond values*: Check all “like-spin” neighbors, and set $b_{i\sim j} = 1$ with probability $1 - \exp(-\beta)$.
- ▶ *Update spin values*: Cluster spins by connecting neighboring sites with a mutual bond, and then flip each cluster with probability 0.5.

For the Ising model, the introduction of the auxiliary variable \mathbf{u} has the dependence between neighboring spins partially decoupled, and the resulting sampler can thus converge substantially faster than the single site updating algorithm. As demonstrated by Swendsen and Wang (1987), this algorithm can eliminate much of the critical slowing down.

spatial autologistic model

Spatial models, e.g., the autologistic model, the Potts model, and the autonormal model (Besag, 1974), have been used in modeling of many scientific problems. A major problem with these models is that the normalizing constant is intractable. Suppose we have a dataset \mathbf{X} generated from a statistical model with the likelihood function

$$f(\mathbf{x}|\theta) = \frac{1}{Z(\theta)} \exp\{-U(\mathbf{x}, \theta)\}, \quad \mathbf{x} \in \mathcal{X}, \quad \theta \in \Theta, \quad (4)$$

where θ is the parameter, and $Z(\theta)$ is the normalizing constant which depends on θ and is not available in closed form. Let $f(\theta)$ denote the prior density of θ . The posterior distribution of θ given \mathbf{x} is then given by

$$f(\theta|\mathbf{x}) \propto \frac{1}{Z(\theta)} \exp\{-U(\mathbf{x}, \theta)\} f(\theta). \quad (5)$$

Let \mathbf{y} denote the auxiliary variable, which shares the same state space with \mathbf{x} . Let

$$f(\theta, \mathbf{y}|\mathbf{x}) = f(\mathbf{x}|\theta)f(\theta)f(\mathbf{y}|\theta, \mathbf{x}), \quad (6)$$

denote the joint distribution of θ and \mathbf{y} conditional on \mathbf{x} , where $f(\mathbf{y}|\theta, \mathbf{x})$ is the distribution of the auxiliary variable \mathbf{y} .

To simulate from (6) using the MH algorithm, one can use the proposal distribution

$$q(\theta', \mathbf{y}' | \theta, \mathbf{y}) = q(\theta' | \theta, \mathbf{y})q(\mathbf{y}' | \theta'), \quad (7)$$

which corresponds to the usual change on the parameter vector $\theta \rightarrow \theta'$, followed by an exact sampling (Propp and Wilson, 1996) step of drawing \mathbf{y}' from $q(\cdot | \theta')$.

If $q(\mathbf{y}'|\theta')$ is set as $f(\mathbf{y}'|\theta)$, then the MH ratio can be written as

$$r(\theta, \mathbf{y}, \theta', \mathbf{y}'|\mathbf{x}) = \frac{f(\mathbf{x}|\theta')f(\theta')f(\mathbf{y}'|\theta', \mathbf{x})q(\theta|\theta', \mathbf{y}')f(\mathbf{y}|\theta)}{f(\mathbf{x}|\theta)f(\theta)f(\mathbf{y}|\theta, \mathbf{x})q(\theta'|\theta, \mathbf{y})f(\mathbf{y}'|\theta')}, \quad (8)$$

where the unknown normalizing constant $Z(\theta)$ can be canceled. To ease computation, Møller *et al.* further suggested to set the proposal distributions $q(\theta'|\theta, \mathbf{y}) = q(\theta'|\theta)$ and $q(\theta|\theta', \mathbf{y}') = q(\theta|\theta')$, and to set the auxiliary distributions

$$f(\mathbf{y}|\theta, \mathbf{x}) = f(\mathbf{y}|\hat{\theta}), \quad f(\mathbf{y}'|\theta', \mathbf{x}) = f(\mathbf{y}'|\hat{\theta}), \quad (9)$$

where $\hat{\theta}$ denotes an estimate of θ , for example, which can be obtained by maximizing a pseudo-likelihood function.

The Møller Algorithm

- ▶ Generate θ' from the proposal distribution $q(\theta'|\theta_t)$.
- ▶ Generate an exact sample \mathbf{y}' from the distribution $f(\mathbf{y}|\theta')$.
- ▶ Accept (θ', \mathbf{y}') with probability $\min(1, r)$, where

$$r = \frac{f(\mathbf{x}|\theta')f(\theta')f(\mathbf{y}'|\hat{\theta})q(\theta_t|\theta')f(\mathbf{y}|\theta_t)}{f(\mathbf{x}|\theta_t)f(\theta_t)f(\mathbf{y}|\hat{\theta})q(\theta'|\theta_t)f(\mathbf{y}'|\theta')}$$

If it is accepted, set $(\theta_{t+1}, \mathbf{y}_{t+1}) = (\theta', \mathbf{y}')$. Otherwise, set $(\theta_{t+1}, \mathbf{y}_{t+1}) = (\theta_t, \mathbf{y}_t)$.

The exchange algorithm

- ▶ Propose $\theta' \sim q(\theta'|\theta, \mathbf{x})$.
- ▶ Generate an auxiliary variable $\mathbf{y} \sim f(\mathbf{y}|\theta')$ using an exact sampler.
- ▶ Accept θ' with probability $\min\{1, r(\theta, \theta', \mathbf{y}|\mathbf{x})\}$, where

$$r(\theta, \theta', \mathbf{y}|\mathbf{x}) = \frac{\pi(\theta')f(\mathbf{x}|\theta')f(\mathbf{y}|\theta)q(\theta|\theta')}{\pi(\theta)f(\mathbf{x}|\theta)f(\mathbf{y}|\theta')q(\theta'|\theta)}. \quad (10)$$

The exchange algorithm can be viewed as an auxiliary variable MCMC algorithm with the proposal distribution being augmented, for which the proposal distribution can be written as

$$T(\theta \rightarrow (\theta', \mathbf{y})) = q(\theta'|\theta)f(\mathbf{y}|\theta'), \quad T(\theta' \rightarrow (\theta, \mathbf{y})) = q(\theta|\theta')f(\mathbf{y}|\theta).$$

This simply validates the algorithm, following the arguments for auxiliary variable Markov chains made in previous lectures.

Monte Carlo Metropolis-Hastings Algorithm (Liang and Jin, 2010)

Let θ_t denote the current draw of θ by the algorithm. Let $y_1^{(t)}, \dots, y_m^{(t)}$ denote the auxiliary samples simulated from the distribution $f(y|\theta_t)$. The MCMH algorithm works by iterating between the following steps:

1. Draw ϑ from some proposal distribution $Q(\theta_t, \vartheta)$.
2. Estimate the normalizing constant ratio $R(\theta_t, \vartheta) = \kappa(\vartheta)/\kappa(\theta_t)$ by

$$\widehat{R}_m(\theta_t, \mathbf{y}_t, \vartheta) = \frac{1}{m} \sum_{i=1}^m \frac{g(y_i^{(t)}, \vartheta)}{g(y_i^{(t)}, \theta_t)},$$

where $g(y, \theta) = \exp\{-U(y, \theta)\}$, and $\mathbf{y}_t = (y_1^{(t)}, \dots, y_m^{(t)})$ denotes the collection of the auxiliary samples.

3. Calculate the Monte Carlo MH ratio

$$\tilde{r}_m(\theta_t, \mathbf{y}_t, \vartheta) = \frac{1}{\widehat{R}_m(\theta_t, \mathbf{y}_t, \vartheta)} \frac{g(x, \vartheta) f(\vartheta)}{g(x, \theta_t) f(\theta_t)} \frac{Q(\vartheta, \theta_t)}{Q(\theta_t, \vartheta)},$$

where $f(\theta)$ denotes the prior distribution imposed on θ .

4. Set $\theta_{t+1} = \vartheta$ with probability

$\tilde{\alpha}(\theta_t, \mathbf{y}_t, \vartheta) = \min\{1, \tilde{r}_m(\theta_t, \mathbf{y}_t, \vartheta)\}$, and set $\theta_{t+1} = \theta_t$ with the remaining probability.

5. If the proposal is rejected in step 4, set $\mathbf{y}_{t+1} = \mathbf{y}_t$. Otherwise, draw samples $\mathbf{y}_{t+1} = (y_1^{(t+1)}, \dots, y_m^{(t+1)})$ from $f(y|\theta_{t+1})$ using either a MCMC algorithm or an automated rejection sampling algorithm.

Since the algorithm involves a Monte Carlo step to estimate the unknown normalizing constant ratio, it is termed as “Monte Carlo MH”. Clearly, the samples $\{(\theta_t, \mathbf{y}_t)\}$ forms a Markov chain, for which the transition kernel can be written as

$$\begin{aligned} \tilde{P}_m(\theta, \mathbf{y}; d\vartheta, d\mathbf{z}) &= \tilde{\alpha}(\theta, \mathbf{y}, \vartheta) Q(\theta, d\vartheta) f_{\vartheta}^m(d\mathbf{z}) \\ &+ \delta_{\theta, \mathbf{y}}(d\vartheta, d\mathbf{z}) \left[1 - \int_{\Theta \times \mathbb{Y}} \tilde{\alpha}(\theta, \mathbf{y}, \vartheta) Q(\theta, d\vartheta) f_{\vartheta}^m(d\mathbf{y}) \right], \end{aligned} \quad (11)$$

where $f_{\vartheta}^m(\mathbf{y}) = f(y_1, \dots, y_m | \vartheta)$ denotes the joint density of y_1, \dots, y_m .

However, since here we are mainly interested in the marginal law of θ , in what follows we will consider only the marginal transition kernel

$$\begin{aligned}\tilde{P}_m(\theta, d\vartheta) &= \int_{\mathbb{Y}} \tilde{\alpha}(\theta, \mathbf{y}, \vartheta) Q(\theta, d\vartheta) f_{\theta}^m(d\mathbf{y}) \\ &\quad + \delta_{\theta}(d\vartheta) \left[1 - \int_{\Theta \times \mathbb{Y}} \tilde{\alpha}_m(\theta, \mathbf{y}, \vartheta) Q(\theta, d\vartheta) f_{\theta}^m(d\mathbf{y}) \right],\end{aligned}\tag{12}$$

showing that $\tilde{P}_m(\theta, d\vartheta)$ will converge to the posterior distribution $f(\theta|x)$ when m is large and the number of iterations goes to infinity.

Monte Carlo MH Algorithm II

1. Draw ϑ from some proposal distribution $Q(\theta_t, \vartheta)$.
2. Draw auxiliary samples $\mathbf{y}_t = (y_1^{(t)}, \dots, y_m^{(t)})$ from $f(y|\theta_t)$ using a MCMC algorithm.
3. Estimate the normalizing constant ratio $R(\theta_t, \vartheta) = \kappa(\vartheta)/\kappa(\theta_t)$ by

$$\widehat{R}_m(\theta_t, \mathbf{y}_t, \vartheta) = \frac{1}{m} \sum_{i=1}^m \frac{g(y_i^{(t)}, \vartheta)}{g(y_i^{(t)}, \theta_t)}.$$

4. Calculate the Monte Carlo MH ratio

$$\tilde{r}_m(\theta_t, \mathbf{y}_t, \vartheta) = \frac{1}{\widehat{R}_m(\theta_t, \mathbf{y}_t, \vartheta)} \frac{g(x, \vartheta) f(\vartheta)}{g(x, \theta_t) f(\theta_t)} \frac{Q(\vartheta, \theta_t)}{Q(\theta_t, \vartheta)}.$$

5. Set $\theta_{t+1} = \vartheta$ with probability

$\tilde{\alpha}(\theta_t, \mathbf{y}_t, \vartheta) = \min\{1, \tilde{r}_m(\theta_t, \mathbf{y}_t, \vartheta)\}$ and set $\theta_{t+1} = \theta_t$ with the remaining probability.

Convergence of the MCMH sampler

(A₁) Assume that P defines an irreducible and aperiodic Markov chain such that $\pi P = \pi$, and for any $\theta_0 \in \Theta$,
 $\lim_{k \rightarrow \infty} \|P^k(\theta_0, \cdot) - \pi(\cdot)\| = 0$.

(A₂) For any $(\theta, \vartheta) \in \Theta \times \Theta$,

$$\gamma_m(\theta, \mathbf{u}, \vartheta) > 0, \quad f_\theta^m(\cdot) - a.s.$$

(A₃) For any $\theta \in \Theta$ and any $\epsilon > 0$,

$$\lim_{m \rightarrow \infty} Q(\theta, f_\theta^m(\lambda_m(\theta, \mathbf{u}, \vartheta) > \epsilon)) = 0,$$

where

$$Q(\theta, f_\theta^m(\lambda_m(\theta, \mathbf{u}, \vartheta) > \epsilon)) = \int_{\{(\vartheta, \mathbf{u}): \lambda_m(\theta, \mathbf{u}, \vartheta) > \epsilon\}} f_\theta^m(d\mathbf{u}) Q(\theta, d\vartheta).$$

Theorem 1.

Assume (A_1) , (A_2) and (A_3) hold. For any $\varepsilon \in (0, 1]$, there exist $M \in \mathbb{N}$ and $K \in \mathbb{N}$ such that for any $m > M$ and $k > K$

$$\|\tilde{P}_m^k(\theta_0, \cdot) - \pi(\cdot)\| \leq \varepsilon,$$

where $\pi(\cdot)$ denotes the posterior density of θ .

Table 1: Computational results for the U.S. cancer mortality data (Liang and Jin, 2010).

Algorithm	Setting	$\hat{\alpha}$	$\hat{\beta}$	CPU
MCMH I	$m = 20$	$-0.3018 (1.2 \times 10^{-3})$	$0.1232 (6.0 \times 10^{-4})$	11
	$m = 50$	$-0.3018 (1.1 \times 10^{-3})$	$0.1230 (5.3 \times 10^{-4})$	24
	$m = 100$	$-0.3028 (6.7 \times 10^{-4})$	$0.1225 (3.8 \times 10^{-4})$	46
MCMH II	$m = 20$	$-0.3028 (1.2 \times 10^{-3})$	$0.1226 (5.9 \times 10^{-4})$	26
	$m = 50$	$-0.3019 (1.0 \times 10^{-3})$	$0.1228 (5.3 \times 10^{-4})$	63
	$m = 100$	$-0.3016 (8.2 \times 10^{-4})$	$0.1231 (3.8 \times 10^{-4})$	129
Exchange	—	$-0.3015 (4.3 \times 10^{-4})$	$0.1229 (2.3 \times 10^{-4})$	33