

Lecture Notes for STAT546: Computational Statistics

—Lecture 2: Monte Carlo

Faming Liang

Purdue University

August 21, 2024

Why do we need MCMC?

For some problems, the exact sampler is not available or very expensive, we have to resort to MCMC for simulating some approximate samples for inference.

Some convergence theorems (Tierney, 1994)

Theorem 1

Suppose $\{X_n\}$ is an irreducible, aperiodic Markov chain with transition kernel P and invariant distribution π . Then

$$\|P^n(x, \cdot) - \pi(\cdot)\| \rightarrow 0$$

for π -almost all x .

Theorem 2

Suppose $\{X_n\}$ is an irreducible Markov chain with transition kernel P and invariant distribution π , and let h be a real-valued function such that $E_\pi|h(x)| < \infty$. Then $P(\bar{h}(x) \rightarrow E_\pi h(x)) = 1$, for π -almost all x , where $\bar{h}(x) = \frac{1}{n} \sum_{i=1}^n h(x_i)$.

How to construct such a Markov chain with the desired invariant distribution $\pi(x)$?

General Metropolis-Hastings Recipe

- ▶ Start with any $\mathbf{x}^{(0)}$ and a proposal distribution $T(\cdot \rightarrow \cdot)$.
- ▶ Compute the ratio

$$r = \frac{\pi(\mathbf{y})}{\pi(\mathbf{x}^{(t)})} \frac{T(\mathbf{y} \rightarrow \mathbf{x}^{(t)})}{T(\mathbf{x}^{(t)} \rightarrow \mathbf{y})}$$

- ▶ Accept/Rejection decision:

$$\mathbf{x}^{(t+1)} = \begin{cases} \mathbf{y} & \text{with } p = \min(1, r) \\ \mathbf{x}^{(t)} & \text{otherwise} \end{cases}$$

Why Does It Work?

It satisfies the detailed balance condition:

$$\begin{aligned}\pi(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{y}) \min\left\{1, \frac{\pi(\mathbf{y})T(\mathbf{y} \rightarrow \mathbf{x})}{\pi(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{y})}\right\} \\ &= \min\{\pi(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{y}), \pi(\mathbf{y})T(\mathbf{y} \rightarrow \mathbf{x})\} \\ &= \pi(\mathbf{y})T(\mathbf{y} \rightarrow \mathbf{x}) \min\left\{1, \frac{\pi(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{y})}{\pi(\mathbf{y})T(\mathbf{y} \rightarrow \mathbf{x})}\right\}\end{aligned}$$

Write $P(x, y) = T(\mathbf{x} \rightarrow \mathbf{y}) \min\left\{1, \frac{\pi(\mathbf{y})T(\mathbf{y} \rightarrow \mathbf{x})}{\pi(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{y})}\right\}$, then we have

$$\pi(\mathbf{x})P(x, y) = \pi(\mathbf{y})P(y, x),$$

and

$$\int \pi(\mathbf{x})P(x, y)dx = \int \pi(\mathbf{y})P(y, x)dx = \pi(\mathbf{y}),$$

i.e., $\pi(\mathbf{x})$ is invariant for the Markov chain.

Summary:

In the standard Markov chain theory, if the chain is irreducible, aperiodic [this is almost sure for the Metropolis-Hastings algorithm, (Tierney, 1994)], and possesses an invariant distribution, then the chain will become stationary at its target distribution $\pi(x)$.

Therefore, if we run the chain long enough, the samples x_{t_0+1}, \dots, x_T (after a burn-in period of t_0 iterations), can be regarded as approximately following the target distribution $\pi(x)$.

Proposal distribution

- ▶ Independence sampler:

$$T(X, Y) = q(Y)$$

In this case, the Metropolis ratio is

$$r = \frac{\pi(y)q(x)}{\pi(x)q(y)} = \frac{w(y)}{w(x)},$$

where $w(y) = \pi(y)/q(y)$.

- ▶ Random-walk proposal:

$$T(X, Y) = q(Y - X).$$

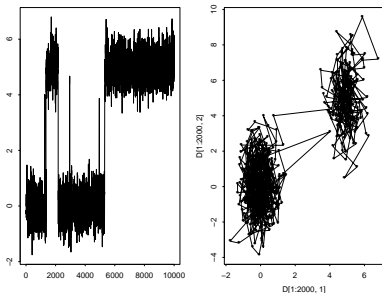
This is very popular, for example,

$$Y|X \sim N(X, \sigma^2).$$

Example: Metropolis-Hastings algorithm

- The moves are very “local”.
- Tend to be trapped in a local mode.

$$\pi(x, y) \propto e^{-\frac{1}{2}\left(\frac{x^2}{0.25} + \frac{y^2}{2}\right)} + 2e^{-\frac{1}{2}\left(\frac{(x-5)^2}{0.25} + \frac{(y-5)^2}{2}\right)}$$



Gibbs sampler

- Purpose: Draw from a joint distribution $\pi(x_1, \dots, x_k)$.
- Method: Iterative conditional sampling
 $\forall i$, Draw $x_i \sim \pi(x_i | \mathbf{x}_{[-i]})$.

Example: Simulating $\pi(x_1, \dots, x_k)$.

Given the current sample $x^t = (x_1^t, \dots, x_k^t)$,

- ▶ Draw $x_1^{t+1} \sim \pi(x_1 | x_2^t, \dots, x_k^t)$;
- ▶ draw $x_2^{t+1} \sim \pi(x_2 | x_1^{t+1}, x_3^t, \dots, x_k^t)$;
- ▶ ...
- ▶ draw $x_k^{t+1} \sim \pi(x_k | x_1^{t+1}, \dots, x_{k-1}^{t+1})$.

Alternatively, x_1, \dots, x_k can be updated in a random order.

Illustrating the Gibbs sampler

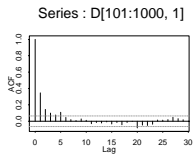
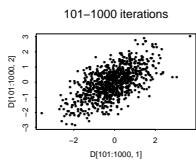
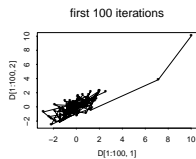
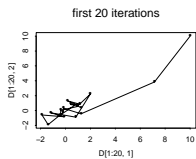
- Suppose the target distribution is

$$(x, y) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\mu} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, and $\boldsymbol{\Sigma} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$.

- Gibbs sampler:
 - ▶ $[X|Y = y] \sim N(\rho y, 1 - \rho^2)$
 - ▶ $[Y|X = x] \sim N(\rho x, 1 - \rho^2)$

Start from, say, $(X, Y) = (10, 10)$, we can take a look at the trajectories. We took $\rho = 0.6$.



Why does the Gibbs sampler work?

Compute the Metropolis ratio:

$$\begin{aligned} r &= \frac{\pi(x, y^*)\pi(y|x)}{\pi(x, y)\pi(y^*|x)} \\ &= \frac{\pi(x, y^*)\pi(x, y)/\pi(x)}{\pi(x, y)\pi(x, y^*)/\pi(x)} \\ &= 1 \end{aligned}$$

A special case of the Metropolis-Hastings algorithm!

Data Augmentation

The DA algorithm can be described in the context of Bayesian analysis of incomplete data.

- ▶ X_{obs} : the observed data
- ▶ X_{mis} : the missing-data
- ▶ $X_{\text{com}} = (X_{\text{obs}}, X_{\text{mis}})$: the complete data

Suppose that the complete-data model has density $g(X_{\text{obs}}, X_{\text{mis}}|\theta)$ with the parameter $\theta \in \Theta \subseteq \mathbb{R}^d$ for some positive integer d . The objective is to make Bayesian inference with a prior distribution $p(\theta)$ for the parameter θ . Let $f(X_{\text{obs}}|\theta)$ be the observed-data model, *i.e.*,

$$f(X_{\text{obs}}|\theta) = \int_{\mathbb{X}_{\text{mis}}} g(X_{\text{obs}}, X_{\text{mis}}|\theta) dX_{\text{mis}} \quad (\theta \in \Theta) \quad (1)$$

For Bayesian inference about θ using MCMC methods, it requires to sampling the true or observed-data posterior

$$p(\theta|X_{\text{obs}}) \propto f(X_{\text{obs}}|\theta)p(\theta) \quad (\theta \in \Theta) \quad (2)$$

or more generally the joint distribution of θ and X_{mis} ,

$$p(\theta, X_{\text{mis}}|X_{\text{obs}}) \propto g(X_{\text{obs}}, X_{\text{mis}}|\theta)p(\theta) \quad (\theta \in \Theta) \quad (3)$$

Let $h(X_{\text{mis}}|\theta, X_{\text{obs}})$ be the conditional distribution of X_{mis} given θ and X_{obs} . Suppose that both $h(X_{\text{mis}}|\theta, X_{\text{obs}})$ and $p(\theta|X_{\text{obs}}, X_{\text{mis}}) \propto g(X_{\text{obs}}, X_{\text{mis}}|\theta)p(\theta)$ are easy to draw samples from. The two-step Gibbs sampler based on these two conditionals becomes a natural choice and is known as the DA algorithm. More precisely, the DA algorithm can be summarized as follows.

The DA algorithm: a two-step Gibbs sampler

Take $\theta^{(0)} \in \Theta$, and iterate for $t = 1, 2, \dots$

I-Step. Generate $X_{\text{mis}}^{(t)} \sim f_{\text{mis}}(X_{\text{mis}} | \theta^{(t-1)}, X_{\text{obs}})$.

P-Step. Generate $\theta^{(t)} \sim p(\theta | X_{\text{obs}}, X_{\text{mis}})$.

As a two-step Gibbs sampler, DA creates two interleaving Markov chains $\{\theta^{(t)} : t = 1, 2, \dots\}$ and $\{X_{\text{mis}}^{(t)} : t = 1, 2, \dots\}$.

Consider a complete-data set consisting of a sample of n , Y_1, \dots, Y_n , from the p -dimensional multivariate normal distribution $N_p(\mu, \Sigma)$ with unknown mean vector $\mu \in \mathbb{R}^p$ and $p \times p$ (positive definite) covariance matrix Σ . Each component of Y_i is either fully observed or missing for each $i = 1, \dots, n$. Let $Y_{\text{obs}}^{(i)}$ denote the observed components and let $Y_{\text{mis}}^{(i)}$ the missing components of Y_i . The conditional distribution of $Y_{\text{mis}}^{(i)}$ given $Y_{\text{obs}}^{(i)}$ and (μ, Σ) is

$$Y_{\text{mis}}^{(i)} | (Y_{\text{obs}}^{(i)}, \mu, \Sigma) \sim N_{\dim(Y_{\text{mis}}^{(i)})}(\mu_{\text{mis}}^{(i)} + \Sigma_{\text{mis,obs}}^{(i)} [\Sigma_{\text{obs,obs}}^{(i)}]^{-1} (Y_{i,\text{obs}} - \mu_{\text{obs}}^{(i)})) \quad (4)$$

Suppose that for Bayesian analysis, we use the prior distribution

$$p(\mu, \Sigma) \propto |\Sigma|^{-(q+1)/2}, \quad (5)$$

where q is a known integer. With $q = p$, this prior becomes the the Jeffreys prior for Σ .

Let $\bar{Y} = n^{-1} \sum_{i=1}^n Y_i$ and let $S = \sum_{i=1}^n (Y_i - \bar{Y})(Y_i - \bar{Y})'$. The complete-data posterior distribution $p(\mu, \Sigma | Y_1, \dots, Y_n)$ can be characterized by

$$\Sigma | Y_1, \dots, Y_n \sim \frac{1}{|\Sigma|^{(n+q)/2}} \exp \left\{ -\frac{1}{2} \text{trace} (\Sigma^{-1} S) \right\}, \quad (6)$$

that is, the inverse Wishart distribution, and

$$\mu | \Sigma, Y_1, \dots, Y_n \sim N_p(\bar{Y}, \Sigma/n). \quad (7)$$

Thus, the DA algorithm has the following I and P steps.

- I-step.* For $i = 1, \dots, n$, draw $Y_{\text{mis}}^{(i)}$ from (4).
- P-step.* First draw Σ from (6) given Y_1, \dots, Y_n and then draw μ from (7) given Y_1, \dots, Y_n and Σ .

Convergence Diagnostic: Trace plot

- ▶ One intuitive and easily implemented diagnostic tool is a trace plot, which plots the sample at time t against the iteration number.
- ▶ A clear sign of non-convergence with a trace plot occurs when we observe some trending in sample space.
- ▶ For the single modal distribution, the traceplot will move snake around the mode after convergence.
- ▶ The problem with the trace plot is that it may appear we have converged, however, the chain is actually trapped in a local region rather exploring the full distribution.

Convergence Diagnostic: Gelman-Rubin Statistic

- ▶ Gelman argues that the best way to identify non-convergence is to simulate multiple sequences with over-dispersed starting points.
- ▶ The intuition is that the behavior of all of the chains should be basically the same. Or, the variance within the chain should be the same as the variance across the chains.
- ▶ Formulas:
 - ▶ Within chain variance $W = \frac{1}{m(n-1)} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2$.
 - ▶ Between chain variance $B = \frac{n}{m-1} \sum_{i=1}^m (\bar{x}_i - \bar{x})^2$.
 - ▶ Estimated variance $\hat{V} = (1 - \frac{1}{n})W + \frac{1}{n}B$.
 - ▶ Gelman-Rubin Statistic: $R = \sqrt{\hat{V}/W}$.
- ▶ The chain is considered as converged when $R < 1.1$ or 1.2 .

