

Markov Chain Monte Carlo

Galin L. Jones

School of Statistics

University of Minnesota

Draft: March 26, 2023

Contents

1	Introduction	1
2	Markov Chains	2
2.1	Stability	4
3	Constructing MCMC Algorithms	7
3.1	Metropolis-Hastings	7
3.1.1	Common Strategies for Choosing a Proposal	9
	Exercises	9
	Appendix	10

1 Introduction

Later.

2 Markov Chains

Let $(\mathbf{X}, \mathcal{B})$ be a measurable space. A sequence of \mathbf{X} -valued random variables $\{X_1, X_2, X_3, \dots\}$ is a Markov chain if for all g

$$E[g(X_{n+1}, X_{n+2}, \dots) \mid X_n, X_{n-1}, \dots, X_1] = E[g(X_{n+1}, X_{n+2}, \dots) \mid X_n].$$

Then P is a *Markov kernel* if $P : \mathbf{X} \times \mathcal{B} \rightarrow \mathbb{R}$ satisfying (i) for each fixed $x \in \mathbf{X}$, $P(x, \cdot)$ is a probability measure and (ii) for each fixed $B \in \mathcal{B}$, $P(\cdot, B)$ is a measurable function.

When \mathbf{X} is a discrete set a Markov kernel can be represented as a square matrix whose entries are nonnegative and whose rows sum to 1.

Example 2.1. Suppose

$$P = \begin{pmatrix} 1/2 & 1/2 \\ 1/3 & 2/3 \end{pmatrix}$$

Then P is a Markov matrix on two states $\{0, 1\}$, say. The first row for example, is interpreted as the probability of moving in one step from state 0 to state 0 is $1/2$ which is the same as the probability of moving in one step from state 0 to state 1.

Example 2.2. Let $\mathbf{X} = \mathbb{Z}$ and let $0 < \theta < 1$. If $x \geq 1$, then a Markov kernel is defined by the matrix P with elements

$$P(x, x+1) = P(-x, -x-1) = \theta, \quad P(x, 0) = P(-x, 0) = 1 - \theta,$$

and $P(0, 1) = P(1, 0) = 1/2$.

Often \mathbf{X} will be uncountable and \mathcal{B} will be countably generated. If \mathbf{X} is topological, then \mathcal{B} will be the Borel σ -algebra generated by \mathbf{X} .

Example 2.3. Let $\mathsf{X} = (0, 1)$ and consider the Markov chain that evolves as follows. Draw $U \sim \text{Uniform}(0, 1)$. If $u \leq 0.5$, $X_{n+1} \sim \text{Uniform}(0, X_n)$, but if $u > 0.5$, $X_{n+1} \sim \text{Uniform}(X_n, 1)$. Then if $X_n = x$ and $B \in \mathcal{B}$

$$P(x, B) = \int_B \left[\frac{1}{2} \frac{1}{x} I_y((0, x)) + \frac{1}{2} \frac{1}{1-x} I_y(x, 1) \right] dy.$$

In Example 2.3, the integrand in the Markov kernel is a conditional density on X . This is a setting that will be encountered repeatedly throughout. If there is a conditional density $k(y \mid x)$, with respect to a measure λ , such that the Markov kernel satisfies for $B \in \mathcal{B}$

$$P(x, B) = \int_B k(y \mid x) \lambda(dy),$$

then say k is a *Markov transition density*.

Example 2.4. Suppose $f(x, y)$ is a joint density with support \mathbb{R}^2 and conditional densities $f_{X|Y}(x \mid y)$ and $f_{Y|X}(y \mid x)$. Then

$$k(x', y' \mid x, y) = f_{X|Y}(x' \mid y) f_{Y|X}(y' \mid x')$$

is a Markov transition density. The Markov chain evolves from $(X_k = x, Y_k = y)$ to (X_{k+1}, Y_{k+1}) by drawing $X_{k+1} \sim F_{X|Y}(\cdot \mid y)$ followed by $Y_{k+1} \sim F_{Y|X}(\cdot \mid X_{k+1})$. This is a special case of the so-called two-variable Gibbs sampler.

Suppose λ is a positive measure on $(\mathsf{X}, \mathcal{B})$, define

$$\lambda P(B) = \int_{\mathsf{X}} \lambda(dx) P(x, B). \tag{1}$$

When λ is a probability measure, the encouraged interpretation is that $X_{n+1} \mid X_n \sim P(X_n, \cdot)$ and $X_n \sim \lambda$, the product $\lambda(dx) P(x, \cdot)$ is the joint distribution of (X_n, X_{n+1}) and λP is the marginal distribution of X_{n+1} .

Since Markov kernels act to the left on measures (1),

$$P^2(x, B) = \int_{\mathsf{X}} P(x, dx_k) P(x_k, B).$$

Continuing in this fashion obtain for every $n \geq 2$

$$P^n(x, B) = \int_{\mathsf{X}} P(x, dx_k) P(x_k, dx_{k+1}) \cdots P(x_{k+n-2}, B).$$

More generally, the so-called Chapman-Kolmogorov equations hold for $n \geq m \geq 0$

$$P^n(x, B) = \int_{\mathbf{X}} P^m(x, dy) P^{n-m}(y, B).$$

If $\lambda = \lambda P$, then λ is *invariant* for P . Notice that if λ is invariant for P and $X_n \sim \lambda$, then $X_{n+1} \sim \lambda$. That is, the marginal distribution does not depend upon n in which case the Markov chain is *stationary*.

[Come back to these examples.](#)

Example 2.5.

Example 2.6. Recall the Markov chain defined in Example (2.4)

One common way of establishing invariance of MCMC Markov chains is to verify a *detailed balance condition*; see Exercise 3.1. Detailed balance holds if

$$\lambda(dx)P(x, dy) = \lambda(dy)P(y, dx). \quad (2)$$

When λ is a probability measure, one interpretation is that the joint distribution of (X_k, X_{k+1}) is the same as the distribution of (X_{k+1}, X_k) so that this is also often called the *reversibility condition*. Another name often encountered is that P is λ -*symmetric*.

2.1 Stability

MCMC applications typically are constructed so that a specific probability distribution F is invariant. However, in applications where MCMC is required it is typically difficult to simulate from the invariant distribution. The most that can be hoped for is that the simulation will eventually produce a representative sample from F . This long-run behavior is in no way guaranteed without additional assumptions. The following simple examples illustrate that the problems can arise due to either the way the kernel is specified or the properties of the state space \mathbf{X} .

Example 2.7. Suppose F lives on $\{1, 2\}$ with $F(1) = 1 - F(2) = 1/4$ and

$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

Since the Markov chain moves deterministically between the two states, it will over represent state 1 and underrepresent state 2 no matter how many iterations there are.

Example 2.8. Suppose F lives on $\{1, 2, 3\}$ with $F(1) = F(2) = F(3) = 1/3$ and

$$P = \begin{pmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 1/2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Thus starting at state $\{3\}$ the chain remains there forever while starting from $\{1, 2\}$ the chain will never visit $\{3\}$. Thus the chain cannot represent F in the long run.

Example 2.9. Suppose for $i = 1, 2$, f_i is a pf on $\mathbf{X}_i \subseteq \mathbb{R}$ and g_i is a pf on $\mathbf{Y}_i \subseteq \mathbb{R}$. Set

$$f(x, y) = \frac{1}{2}f_1(x)g_1(y) + \frac{1}{2}f_2(x)g_2(y).$$

Then

$$f_{X|Y}(x | y) = \frac{f_1(x)g_1(y) + f_2(x)g_2(y)}{g_1(y) + g_2(y)}$$

and

$$f_{Y|X}(y | x) = \frac{f_1(x)g_1(y) + f_2(x)g_2(y)}{f_1(y) + f_2(y)}$$

and the Gibbs sampler MTD is

$$k(x', y' | x, y) = f_{X|Y}(x' | y)f_{Y|X}(y' | x').$$

When $\mathbf{X}_i = \mathbf{Y}_i = \mathbb{R}$ this Gibbs sampler will produce a representative sample eventually. However, complications may arise if the spaces are constrained.

Suppose $\mathbf{X}_1 = \mathbf{X}_2 = (0, 1)$ and that $f_1 = f_2$ is the Uniform density. Let $Y_1 = (0, 1)$ and $Y_2 = (2, 3)$ and g_1 and g_2 be Uniform densities. Easy calculation yields that

$$f_{X|Y}(x | y) = I(0 < x < 1) [I(0 < y < 1) + I(2 < y < 3)]$$

and

$$f_{Y|X}(y | x) = \frac{1}{2}I(0 < x < 1) [I(0 < y < 1) + I(2 < y < 3)]$$

and that, no matter which square, $\mathbf{X}_1 \times \mathbf{Y}_1$ or $\mathbf{X}_2 \times \mathbf{Y}_2$, the current state is in there is a positive probability of the next state being in either square. This Gibbs sampler will eventually produce a representative sample.

Now consider the setting with $\mathbf{X}_1 = \mathbf{Y}_1 = (0, 1)$ and $\mathbf{X}_2 = \mathbf{Y}_2 = (2, 3)$ so that

$$f_{X|Y}(x | y) = \frac{I(0 < x < 1)I(0 < y < 1) + I(2 < x < 3)I(2 < y < 3)}{I(0 < y < 1) + I(2 < y < 3)}$$

and

$$f_{Y|X}(y | x) = \frac{I(0 < x < 1)I(0 < y < 1) + I(2 < x < 3)I(2 < y < 3)}{I(0 < x < 1) + I(2 < x < 3)}.$$

If $y \in (0, 1)$, then $f_{X|Y}(x | y) = I(0 < x < 1)$ and, similarly, if $X \in (0, 1)$, then $f_{Y|X}(y | x) = I(0 < y < 1)$. Thus if the current state is in the square $\mathbf{X}_1 \times \mathbf{Y}_1$, then the next step will be in $\mathbf{X}_1 \times \mathbf{Y}_1$. That is, there is no chance for the chain to visit $\mathbf{X}_2 \times \mathbf{Y}_2$. This Gibbs sampler will not produce a representative sample from the target distribution.

The above examples demonstrate that one way problems arise is when the Markov chain cannot access all of the space eventually and hence properties which avoid this are required.

Let ϕ be a non-trivial positive measure on \mathcal{B} . Then $A \in \mathcal{B}$ is ϕ -communicating if for all $B \subseteq A$ such that $\phi(B) > 0$ and for all $x \in A$, there exists n such that $P^n(x, B) > 0$.

This is a weak property that does not alone ensure desirable long run properties. Consider the Gibbs sampler from Example 2.9 with $\mathbf{X}_1 = \mathbf{Y}_1 = (0, 1)$ and $\mathbf{X}_2 = \mathbf{Y}_2 = (2, 3)$. If ϕ denotes Lebesgue measure, then this Markov chain is ϕ -communicating but does not have desirable long run properties.

The Markov kernel, P , is ϕ -irreducible if for all $x \in \mathbf{X}$ and for all $A \in \mathcal{B}$ such that $\phi(A) > 0$ there exists n such that $P^n(x, A) > 0$. This is a key property in many MCMC settings, but is not enough to ensure desirable long run properties; consider the Markov chain in Example 2.7 which is irreducible.

There is some arbitrariness in the definition of ϕ -irreducibility, but if for some ϕ , P is ϕ -irreducible, then there exists a maximal irreducibility measure ψ (meaning that $\psi(A) = 0$ implies $\phi(A) = 0$ for all irreducibility measure ϕ).

Proposition 2.1. *If λ is an invariant measure for the Markov kernel P and, for some ϕ , P is ϕ -irreducible, then P is λ -irreducible.*

Corollary 2.1. *If P is λ -symmetric and ϕ -irreducible, then P is λ -irreducible.*

Often inspection of the kernel P and the underlying space X is enough to establish ϕ -irreducibility. For example, many kernels obviously satisfy a positivity condition so that $P(x, A) > 0$ for all $x \in \mathsf{X}$ and $A \in \mathcal{B}$ with $\phi(A) > 0$. The following is a special case of a more general result [1, Theorem 3], but is often useful.

Proposition 2.2. *Suppose X is a connected, separable metric space. If every nonempty, open set A satisfies $\phi(A) > 0$ and every point has a ϕ -communicating neighborhood, the Markov kernel is ϕ -irreducible.*

3 Constructing MCMC Algorithms

3.1 Metropolis-Hastings

[Give credit to R?](#)

The fundamental MCMC algorithm is Metropolis-Hastings [2, 3]. It serves as a building block for many, if not most, applications of MCMC. The following will be generalized later, but, for now, suppose $\mathsf{X} \subseteq \mathbb{R}^d$ and that $q(x, y)$ is a proposal conditional density that is easy to sample. This notation can take an adjustment for those accustomed to the alternative $q(y \mid x)$, which is more common in statistics. Define the *Hastings ratio*

$$r(x, y) = \frac{f(y)q(y, x)}{f(x)q(x, y)}.$$

Algorithm 1 Metropolis-Hastings

- 1: *Input:* Current value $X_n = x$.
 - 2: Draw $Y \sim Q(x, \cdot)$
 - 3: Draw $U \sim \text{Uniform}(0, 1)$
 - 4: If $u \leq r(x, y) \wedge 1$, accept y and set $X_{n+1} = y$, else set $X_{n+1} = x$.
-

Algorithm 2 is the formal definition, but is not at all how the algorithm should be implemented in practice. The following implementation will help avoid overflow issues.

Algorithm 2 Metropolis-Hastings Implementation

- 1: *Input:* Current value $X_n = x$.
 - 2: Draw $Y \sim Q(x, \cdot)$
 - 3: Draw $U \sim \text{Uniform}(0, 1)$
 - 4: If $\log r(x, y) \geq 0$ set $X_{n+1} = y$, else if $u \leq r(x, y)$, set $X_{n+1} = y$, else set $X_{n+1} = x$.
-

Metropolis-Hastings defines a Markov kernel. Specifically, if $\alpha(x, y) = 1 \wedge r(x, y)$, then

$$P(x, dy) = Q(x, dy) + \delta_x(dy) \int [1 - \alpha(x, u)] Q(x, du). \quad (3)$$

Proposition 3.1. *The Metropolis-Hastings kernel (3) is F -symmetric.*

Proof. It suffices to consider $x \neq y$. Then

$$\begin{aligned} F(dx)P(x, dy) &= f(x)q(x, y) [1 \wedge r(x, y)] \mu(dx)\mu(dy) \\ &= [f(x)q(x, y) \wedge f(y)q(y, x)] \mu(dx)\mu(dy) \\ &= f(y)q(y, x) [1 \wedge r(y, x)] \mu(dx)\mu(dy) \\ &= F(dy)P(y, dx). \end{aligned}$$

□

The proof of the following result is easy and is left as an exercise.

Proposition 3.2. *If $q(x, y) > 0$ for all $x, y \in \mathbf{X}$, then the Metropolis-Hastings kernel P is F -irreducible.*

Example 3.1. Suppose F is the standard normal distribution. If the proposal distribution is the Student's t distribution with $d \geq 1$ degrees of freedom, then by Proposition 3.2 the resulting Metropolis-Hastings kernel is F -irreducible.

If the proposal distribution is $\text{Uniform}(0, 1)$, then the resulting Metropolis-Hastings kernel is obviously reducible since no value outside of $(0, 1)$ will be proposed.

Now if $X_n = x$ and the proposal is $\text{Uniform}(x - 1, x + 1)$, then Proposition 2.2 can be used to establish F -irreducibility. The first two conditions of the proposition are easy consequences of the properties of the real line. Let $0 < \delta < 1$ so that $N_{x,\delta} = (x - \delta, x + \delta)$ is a neighborhood of $\{x\}$ and

$$P(x, N_{x,\delta}) = \int_{x-\delta}^{x+\delta} \frac{1}{2} \left(1 \wedge e^{-0.5(x^2+y^2)} \right) dy > 0$$

from which it is easy to see that $N_{x,\delta}$ is an F -communicating neighborhood.

3.1.1 Common Strategies for Choosing a Proposal

Exercises

Exercise 3.1. Prove that if Equation 2 holds, then λ is invariant for P .

Exercise 3.2. What is the invariant distribution of the Markov chain defined in Example 2.7?

Exercise 3.3. Consider the Gibbs samplers in Example 2.9. establish that the Gibbs sampler on $\mathbf{X}_1 = \mathbf{Y}_1 = (0, 1)$ and $\mathbf{X}_2 = \mathbf{Y}_2 = (2, 3)$ is not irreducible while the one on $\mathbf{X}_1 = \mathbf{Y}_1 = \mathbf{X}_2 = (0, 1)$ and $\mathbf{Y}_2 = (2, 3)$ is irreducible.

Exercise 3.4. Establish Proposition 2.1.

Exercise 3.5. Establish Proposition 2.2.

Exercise 3.6. Prove that Metropolis-Hastings has Markov kernel given in (3). Hint: Note that if $A \in \mathcal{B}$ and $j \geq 1$ then

$$\Pr(X_{j+1} \in A \mid X_j) = \Pr(X_{j+1} \in A, U \leq \alpha(x, y) \mid X_j) + \Pr(X_{j+1} \in A, U > \alpha(x, y) \mid X_j).$$

Appendix

References

- [1] Geyer, C. J. (2014). Stat 8501 lecture notes.
- [2] Hastings, W. K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1).
- [3] Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). Equations of state calculations by fast computing machine. *Journal of Chemical Physics*, 21.