

# Computer Intensive Statistics and Applications

## Chapter 1 and 2: Random Variable Generation

Shaobo Jin

Department of Mathematics

# Statistical Model

In statistical inference,

- the possible observations are from a **sample space**  $\mathcal{X}$ ,
- observation  $x$  is a realization of a random variable/vector  $X$ .

Statistical inference is concerned with methods of using the observations to obtain information about the distribution of  $X$ .

A **statistical model** is a class of possible probability measures  $\mathcal{P}$  on the sample space  $\mathcal{X}$ , parameterized by parameter(s)  $\theta$ :

$$\mathcal{P} = \{P_{\theta} : \theta \in \Theta\},$$

where  $\Theta$  is the **parameter space**.

# Simulate Data

A **sample**  $X = (X_1, \dots, X_n)$  is a collection of independent random variables where  $X_i$  is distributed according to a distribution  $P_{i,\theta}$ .

- We can view data as “simulated” from  $P_\theta$  and we would like to know how data are generated by  $P_\theta$ .
- This is similar to conducting many experiments in order to collect data.

Once we can simulate data from  $P_\theta$ , we will be able to perform various simulation-based inference:

- estimate unknown quantity (e.g., Monte Carlo methods),
- quantify uncertainty (e.g., bootstrap).

# Some Examples

Some examples that you should be able to solve by the end of the course:

- approximate expectations  $E[X] = \int x f(x) dx$ , if the closed form expression is difficult/impossible to find,
  - approximate integrals  $\int f(x, y) dx$ ,
  - simulate a stochastic process,
  - approximate distributions of functions of random variables  $h(X)$ ,
- among others.

The focus of our course is neither 100% programming nor 100% math. We will show you

- how simulations are performed,
- why simulations work or do not work.

# Uniform Distribution

We start with the fundamental component: **random number generation**.

## Definition (Definition 1.2)

A **uniform pseudo-random number generator** is an algorithm which, starting from an initial value  $u_0$  and a transformation  $D$ , produces a sequence with  $u_i = D^i(u_0)$  in  $[0, 1]$ . For all  $n$ ,  $(u_1, \dots, u_n)$  reproduce the behavior of an i.i.d. sample of uniform random variables when compared through a set of tests.

# Congruential Generators

To generate Uniform  $[0, 1]$ , we can use the classic [congruential generator](#)  $(a, b, M)$  on  $\{0, \dots, M - 1\}$ .

---

**Algorithm 1:** Generate pseudo random numbers from  $U[0, 1]$ 

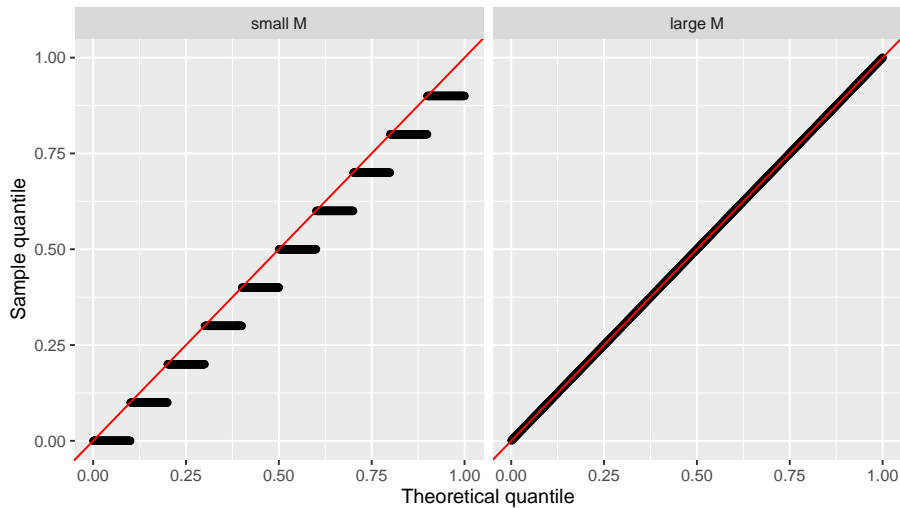
---

- 1 Choose a large  $M > 0$  and  $(a, b)$  from  $\{0, \dots, M - 1\}$ ;
  - 2 Choose a [seed](#)  $x_0$  ;
  - 3 **for**  $i = 1$  *in*  $1 : n$  **do**
  - 4      $x_i = (ax_{i-1} + b) \bmod M$ ;
  - 5 **end**
  - 6 Dividing the sequence by  $M$  gives a sequence that are regarded as random numbers from Uniform  $[0, 1]$  .
- 

The operator  $y = z \bmod M$  means that, for an integer  $z$ , there is an integer  $k$  such that  $y = kM + z$  but  $0 \leq y < M$ .

- If  $z$  is positive, then it is just the remainder.

# Illustration of QQ plot



## Expectation and Variance

Suppose that we generated  $(x_1, \dots, x_M)$  from a congruential generator  $(a, b, M)$ . Then, the ordered outcome must be

$$\left[0 \quad \frac{1}{M} \quad \dots \quad \frac{M-1}{M}\right].$$

The sample mean is

$$\frac{1}{M} \sum_{i=1}^M x_i = \frac{1}{M} \sum_{j=0}^{M-1} \frac{j}{M} = \frac{M-1}{2M} \rightarrow \frac{1}{2},$$

and the sample variance with denominator  $M$  is

$$\frac{1}{M} \sum_{i=1}^M x_i^2 - \left( \frac{1}{M} \sum_{i=1}^M x_i \right)^2 = \frac{M^2 - 1}{12M^2} \rightarrow \frac{1}{12},$$

as  $M$  increases.



# Quantile

In this course, we will mostly assume that we can generate Uniform  $[0, 1]$  and study how we can generate random numbers of other distributions.

## Definition (Definition 1.1)

For a distribution function  $F$  on  $\mathbb{R}$ , the **generalized inverse**  $F^-$  of  $F$  is the function defined by

$$F^-(u) = \inf \{x; F(x) \geq u\} \quad \text{for } 0 \leq u \leq 1.$$

The generalized inverse is often called the **quantile function**, denoted by  $Q$ .

## Example

Let  $X \sim \text{Exp}(\lambda)$  with mean  $1/\lambda$ . Find  $F^-$ .

# Quantile Method

## Lemma (Lemma 1.1)

*If  $U \sim \text{Uniform}[0, 1]$ , then the random variable  $F^{-1}(U)$  has distribution  $F$ .*

This lemma means that as long as we can generate  $U \sim \text{Uniform}[0, 1]$ , we can obtain  $X \sim F$  by letting  $X = F^{-1}(U)$ .

## Example

For exponential distribution, we can generate data via

$$X = -\frac{\log(1 - U)}{\lambda}.$$

# Transformation Method

The quantile method is a special case of the [transformation method](#) that means that, if we are interested in  $Y$  and we know that  $Y = h(X)$ , then we can generate random numbers for  $X$  and  $h(X)$  will be the random numbers for  $Y$ .

## Example

[Box-Muller algorithm](#) for normal distribution:  $U_1 \sim \text{Uniform}[0, 1]$ ,  $U_2 \sim \text{Uniform}[0, 1]$ , and  $U_1$  is independent of  $U_2$ . Then,

$$\begin{aligned}X_1 &= \sqrt{-2 \log(U_1)} \cos(2\pi U_2) \sim N(0, 1), \\X_2 &= \sqrt{-2 \log(U_1)} \sin(2\pi U_2) \sim N(0, 1),\end{aligned}$$

and  $X_1$  is independent of  $X_2$ .

# Conditional Method

Suppose that we can partition  $X$  into  $X = (X_1, X_2, \dots, X_d)$  such that

$$p(x_1, \dots, x_d) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) \cdots p(x_d | x_1, \dots, x_{d-1}).$$

Then, we can generate  $X_1, X_2 | X_1, X_3 | (X_1, X_2), \dots, X_d | (X_1, \dots, X_{d-1})$  in order to generate  $X$ .

## Example

Suppose that  $Y \sim \text{Bernoulli}(\theta)$ , where  $\log\left(\frac{\theta}{1-\theta}\right) = a + bX$  with  $X \sim N(0, 1)$ . To generate  $(X, Y)$ , we generate

$$\begin{aligned} X &\sim N(0, 1), \\ Y | X = x &\sim \text{Bernoulli}\left(\frac{\exp\{a + bx\}}{1 + \exp\{a + bx\}}\right). \end{aligned}$$

# Dominance

Suppose that we want to generate a sample with **target density**  $p(x)$  but we know that it is easier to generate random numbers from another **instrumental density**  $g(x)$ . Let  $M$  be a constant such that

$$p(x) \leq M g(x), \quad \text{for all } x.$$

## Example

Let  $X \sim \text{Beta}(a, b)$  with density  $p(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$ , where  $x \in (0, 1)$ . We focus on  $a > 1$  and  $b > 1$ , and let  $g(x)$  be the uniform density.

# Rejection Algorithm

---

**Algorithm 2:** Rejection algorithm

---

```
1 Choose a seed (omitted in all following algorithms) ;
2 for  $i = 1$  in  $1 : n$  do
3   while new trial is needed do
4     Draw  $X = x$  from  $g(\cdot)$  ;
5     Compute the ratio  $r(x) = \frac{p(x)}{Mg(x)}$  ;
6     Draw  $U = u$  from Uniform  $[0, 1]$  ;
7     If  $u > r(x)$ , start a new trial. Otherwise, let  $x^{(t+1)} = x$  ;
8   end
9 end
```

---

The algorithm works even when  $x$  is multidimensional.

But why does the rejection algorithm work?

# Effect of $M$

- The probability of accepting a candidate is

$$P(U \leq r(X)) = \frac{1}{M}.$$

For large  $M$ , the proportion of useless draws that are discarded is large.

- But if  $M$  is not big enough, then we cannot guarantee  $p(x) \leq Mg(x)$  for all  $x$ .

## A Variant

Sometimes it is not easy to get the exact expression of  $p(x)$ , but it is easy to get  $f(x)$  where  $p(x) = cf(x)$ .

- That is, the normalizing constant  $c$  is often hard to obtain, such as in a posterior distribution.

Instead of using  $p(x)$  (hard to get) in the accept-reject algorithm, we can use  $f(x)$  (easy to get) instead.

---

**Algorithm 3:** Rejection algorithm

---

```
1 for  $i = 1$  in  $1 : n$  do
2   while new trial is needed do
3     Draw  $X = x$  from  $g(\cdot)$  ;
4     Compute the ratio  $r(x) = \frac{f(x)}{Mg(x)}$  ;
5     Draw  $U = u$  from Uniform  $[0, 1]$  ;
6     If  $u > r(x)$ , start a new trial. Otherwise, let  $x^{(t+1)} = x$  ;
7   end
8 end
```

---



## Another Variant: Envelope Accept-Reject Method

Suppose that  $p(x)$  can be computed by can be time consuming. We can find an extra lower bound  $g_l(x)$  such that

$$g_l(x) \leq p(x) \leq Mg(x), \quad \text{for all } x.$$

- In the accept-reject algorithm, we accept the candidate value if

$$u \leq \frac{p(x)}{Mg(x)}.$$

- In the [envelop accept-reject algorithm](#), we first check if

$$u \leq \frac{g_l(x)}{Mg(x)}.$$

# Envelope Accept-Reject Method

---

**Algorithm 4:** Envelope accept-reject algorithm

---

```
1 for  $i = 1$  in  $1 : n$  do
2   while new trial is needed do
3     Draw  $X = x$  from  $g(\cdot)$  ;
4     Draw  $U \sim \text{Uniform}[0, 1]$  ;
5     if  $U \leq \frac{g_I(x)}{Mg(x)}$  then
6       | Accept  $X$  ;
7     else
8       | if  $U \leq \frac{p(x)}{Mg(x)}$  then
9         | Accept  $X$ ;
10      | else
11      | Start a new trial ;
12      | end
13    end
14  end
15 end
```

---

# Markov Property and Markov Chain

Another idea of generating random number is through a [Markov chain](#).

## Definition

Let  $X_i \in \Omega$  for all  $i$ , where the set  $\Omega$  is the state space. The random variables satisfy the [Markov property](#) if

$$P(X_{i+1} \in A \mid X_j = x_j, 0 \leq j \leq i) = P(X_{i+1} \in A \mid X_i = x_i).$$

## Definition

A [Markov chain](#) is a sequence of random variables  $X_i \in \Omega$  that satisfy the Markov property.

## Transition Kernel

Consider a sequence of discrete random variables  $\{X_n\}$ . The Markov property means that

$$P(X_n = y \mid X_{n-1} = x, X_{n-2}, \dots, X_0) = P(X_n = y \mid X_{n-1} = x).$$

In this case, a **transition kernel** is a matrix  $\mathbf{K}$  with elements  $P(X_n = y \mid X_{n-1} = x)$ .

Consider a sequence of continuous random variables  $\{X_n\}$ . The Markov property means that

$$P(X_n \in A \mid X_{n-1} = x, X_{n-2}, \dots, X_0) = \int_{y \in A} K(x, y) dy,$$

$$p(X_n = y \mid X_{n-1} = x, X_{n-2}, \dots, X_0) = p(X_n = y \mid X_{n-1} = x) = K(x, y),$$

where the **transition kernel**  $K(x, y)$  is the conditional density of  $Y$  given  $X = x$ .

# Discrete Marginal Distribution: Time 1

Take the sequence of discrete random variables  $\{X_n\}$  as an example. Suppose that each  $X_n$  takes values from  $1, \dots, M$ . Then,

$$P(X_1 = j) = \sum_{i=1}^M P(X_0 = i) P(X_1 = j \mid X_0 = i).$$

Expressing it as a matrix, we have

$$\mathbf{p}_1 = \mathbf{p}_0 \mathbf{K},$$

where  $\mathbf{p}_k$  is a row vector with  $t$ th element  $P(X_k = t)$ , and the [transition kernel](#)  $\mathbf{K}$  is a  $M \times M$  matrix with  $(i, j)$ th element  $P(X_1 = j \mid X_0 = i)$ .

## Discrete Marginal Distribution: Time $n$

If we consider a [homogenous Markov chain](#), where

$$P(X_i = y \mid X_{i-1} = x) = P(X_1 = y \mid X_0 = x)$$

for all  $i$ , then

❶ At time  $t = 2$ ,

$$\begin{aligned} P(X_2 = j) &= \sum_{i=1}^M P(X_1 = i) P(X_2 = j \mid X_1 = i) \\ \Rightarrow \mathbf{p}_2 &= \mathbf{p}_1 \mathbf{K} = \mathbf{p}_0 \mathbf{K}^2. \end{aligned}$$

❷ In general,

$$\begin{aligned} P(X_n = j) &= \sum_{i=1}^M P(X_{n-1} = i) P(X_n = j \mid X_{n-1} = i) \\ \Rightarrow \mathbf{p}_n &= \mathbf{p}_{n-1} \mathbf{K} = \mathbf{p}_0 \mathbf{K}^n. \end{aligned}$$

# Stationary Distribution

## Definition

The distribution  $\pi$  on  $\Omega$  is a **stationary distribution** (or **invariant distribution**) of the Markov chain with the transition kernel  $K$ , if

$$\pi(y) = \sum_{x \in \Omega} \pi(x) K(x, y), \quad \text{discrete case,}$$

where  $\pi(x) = P(X = x)$ , or

$$\pi(y) = \int_{x \in \Omega} \pi(x) K(x, y) dx, \quad \text{continuous case,}$$

where  $\pi(x) = p(x)$ .

- The stationary distribution means that if the initial state  $X_0 \sim \pi$ , then  $X_n \sim \pi$  for all  $n \geq 0$ , the same distribution.

# Stationary Distribution: Existence

The stationary distribution does not always exist. Our goal in MCMC is to design the algorithm such that the stationary distribution is the distribution that we want to sample random numbers from.

## Example

- ① The stationary distribution of  $\mathbf{K} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$  is  $\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$ .
- ② But if we let  $\Omega = \mathbb{Z}$  and consider

$$\begin{aligned} P(X_n = i - 1 \mid X_{n-1} = i) &= \frac{1}{4}, \\ P(X_n = i + 1 \mid X_{n-1} = i) &= \frac{3}{4}, \end{aligned}$$

then the stationary distribution does not exist.



# Irreducible Markov Chain

From now on, we consider the finite-state Markov chain for simplicity.

- The states  $x$  and  $y$  **communicate** if there exists an integer  $n \geq 0$  such that  $P(X_n = x \mid X_0 = y) > 0$  and  $P(X_n = y \mid X_0 = x) > 0$ .
- The transition matrix  $\mathbf{K}$  is **irreducible** if  $x$  communicates with  $y$  whenever  $x, y \in \Omega$ .

## Example

$\mathbf{K} = \begin{bmatrix} 0.5 & 0.5 & 0 & 0 \\ 0.5 & 0.5 & 0 & 0 \\ 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 0.5 \end{bmatrix}$  is reducible and  $\mathbf{K} = \begin{bmatrix} 0.5 & 0.5 & 0 & 0 \\ 0.5 & 0.25 & 0.25 & 0 \\ 0 & 0 & 0.5 & 0.5 \\ 0.25 & 0 & 0.25 & 0.5 \end{bmatrix}$  is irreducible.

# Periodic

For a transition matrix, the state  $w$  is **periodic** with period  $t > 1$  if

- 1)  $P(X_n = w \mid X_0 = w) > 0$  only holds for  $n = tm$  with integer  $m > 0$  and
- 2)  $t$  is the largest number of  $\{2, 3, \dots\}$  for which 1) holds.

The state  $w$  is **aperiodic** if  $P(X_n = w \mid X_0 = w) > 0$  for all  $n$ .

## Example

$$K = \begin{bmatrix} 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 & 0 \\ 0.5 & 0.5 & 0 & 0 \end{bmatrix} \text{ has period } t = 2.$$

# Long-Run Property: Finite-State Markov Chain

## Theorem

*Consider a finite-state Markov chain. Suppose that the transition matrix  $\mathbf{K}$  is irreducible, aperiodic, and has stationary distribution  $\pi$ . Then, for all starting state  $w_0 \in \Omega$ ,*

$$\lim_{n \rightarrow \infty} P(X_n = w \mid X_0 = w_0) = \pi(w).$$

- The stationary distribution is reached after large enough iterations. When we simulate random numbers using a Markov chain, we need to discard the early iterations. The discarded ones are called the **burn-in** period.
- Samples obtained after the burn-in period can be considered as samples from  $\pi$ .

## Choose the Transition Kernel

Our goal is to simulate data from some distribution with density  $p(x)$ . Now the question becomes how to choose  $K$  such that the stationary distribution is  $\pi(x) = p(x)$ .

### Fact

If  $p(x)$  and  $K(x, y)$  satisfies the *detailed balance condition*, i.e.,

$$K(x, y)p(x) = K(y, x)p(y),$$

for any  $x, y \in \Omega$ , then  $p(x)$  is the stationary distribution of the Markov chain with the transition kernel  $K$ .

# Proposal Distribution

When we simulate random numbers from a Markov chain, we need a **proposal distribution**

$$T(x, y) = p(Y = y \mid X = x).$$

Find a proposal distribution  $T(x, y)$  that satisfies the detailed balance condition is difficult.

- So with probability  $A(x, y)$  we let  $X_{n+1} = y$  (accept), and probability  $1 - A(x, y)$  we let  $X_{n+1} = x$  (reject).
- For  $X_{n+1} \neq x$ , the transition is

$$K(x, y) = T(x, y) A(x, y).$$

Hence, we should seek  $A$  such that the detailed balance condition is fulfilled.

## Detailed Balance: Symmetric Proposal

Start with the case where the proposal distribution is **symmetric**, i.e.,  $T(x, y) = T(y, x)$ .

### Example

$Y \mid X = x \sim N(x, \sigma^2)$  is symmetric, since

$$T(x, y) = \frac{1}{\sqrt{2\sigma^2}} \exp \left\{ -\frac{(x - y)^2}{2\sigma^2} \right\}$$

The transition kernel is

$$K(x, y) = T(x, y) A(x, y) + [1 - A(x, y)] \delta_x(y),$$

where  $\delta_x(y) = 1$  if  $x = y$ , and 0 otherwise. For  $y \neq x$ , the detailed balance condition reduces to

$$\pi(x) A(x, y) = \pi(y) A(y, x).$$

## Deriving $A(x, y)$ : Symmetric Proposal

The detailed balance condition is fulfilled, if we choose

$$A(x, y) = \lambda(x, y) \pi(y), \quad A(y, x) = \lambda(x, y) \pi(x),$$

such that  $\lambda(x, y) = \lambda(y, x)$ . Since  $A(\cdot, \cdot)$  is a probability, we must have

$$\lambda(x, y) \leq \min \left\{ \frac{1}{\pi(y)}, \frac{1}{\pi(x)} \right\}.$$

The maximum of  $A(\cdot, \cdot)$  with respect to  $\lambda$  is attained at

$$\lambda(x, y) = \min \left\{ \frac{1}{\pi(y)}, \frac{1}{\pi(x)} \right\}.$$

Hence,

$$A(x, y) = \lambda(x, y) \pi(y) = \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\}.$$

# Metropolis Algorithm

---

**Algorithm 5:** Metropolis Algorithm

---

```
1 Choose an initial state  $x^{(0)}$  ;  
2 for  $t = 1$  in  $1 : n$  do  
3   Sample a candidate  $y$  from  $T(x^{(t)}, y)$  ;  
4   Calculate the ratio  $R(x^{(t)}, y) = \frac{\pi(y)}{\pi(x^{(t)})}$  ;  
5   Draw  $U \sim \text{U}[0, 1]$  ;  
6   Update  

$$x^{(t+1)} = \begin{cases} y, & \text{if } U \leq R(x^{(t)}, y), \\ x^{(t)}, & \text{otherwise.} \end{cases}$$
  
7 end
```

---



## Metropolis Algorithm: Example

Since we have  $R(x^{(t)}, y) = \frac{\pi(y)}{\pi(x^{(t)})}$  in the Metropolis algorithm, we only need to know  $\pi()$  up to a normalizing constant.

### Example

Suppose that we want to sample random numbers from a distribution with density

$$p(x) = c \exp(-x^2) [2 + \sin(5x) + \sin(2x)],$$

where  $c$  is a normalizing constant.

# Metropolis-Hastings Algorithm

The Metropolis-Hastings algorithm allows **asymmetric** proposal distributions, as long as  $T(x, y) > 0$  if and only if  $T(y, x) > 0$ .

---

## Algorithm 6: Metropolis-Hastings Algorithm

---

1 Choose an initial state  $x^{(0)}$  ;

2 **for**  $t = 1$  *in*  $1 : n$  **do**

3     Sample a candidate  $y$  from  $T(x^{(t)}, y)$  ;

4     Calculate the ratio  $R(x^{(t)}, y) = \frac{\pi(y)T(y, x^{(t)})}{\pi(x^{(t)})T(x^{(t)}, y)}$  ;

5     Draw  $U \sim \text{U}[0, 1]$  ;

6     Update

$$x^{(t+1)} = \begin{cases} y, & \text{if } U \leq R(x^{(t)}, y), \\ x^{(t)}, & \text{otherwise.} \end{cases}$$

7 **end**

---

## Examples of MCMC Algorithms

Many different MCMC algorithms differ mainly in how the candidate  $y$  is sampled.

- In the **random-walk Metropolis algorithm**,  $y = x^{(t)} + \epsilon$ , where  $\epsilon$  is sampled from some distribution, e.g., Uniform  $[-a, a]$ , Normal, etc.
- In **multiplicative random walk**,  $y = x^{(t)}\epsilon$ .
- In **independence sampler**,  $y$  is sampled from  $g(y)$  that does not depend on  $x^{(t)}$ .
- In **hit and run**, we first sample a direction  $\rho^{(t)}$  and then a random length  $\lambda^{(t)}$ . Then  $y = x^{(t)} + \lambda^{(t)}\rho^{(t)}$ .
- The **Langevin Metropolis-Hastings algorithm** explores the shape of the target distribution by  $y = x^{(t)} + d^{(t)} + \sigma\epsilon$ , where  $\epsilon \sim N(0, I)$  and

$$d^{(t)} = \frac{\sigma^2}{2} \frac{\partial \log \pi(x^{(t)})}{\partial x}.$$

# Gibbs Sampler: Conditioning

It can be the case that it is not easy to sample from the joint distribution of  $X \in \mathbb{R}^d$ . But it is easy to sample from the conditional distributions.

- Suppose that  $X = (X_1, \dots, X_p)$ , where  $X_i \in \mathbb{R}^{d_i}$ .
- Let  $\pi_{i|-i}(X_i | X_{-i})$  be the conditional distribution of  $X_i$  given  $X_{-i}$ , where  $X_{-i} = (X_1 \ \cdots \ X_{i-1} \ X_{i+1} \ \cdots \ X_p)$ .

---

## Algorithm 7: Gibbs Sampler

---

```

1 Choose an initial state  $x^{(0)}$  ;
2 for  $t = 1$  in  $1 : n$  do
3   | for  $i = 1$  in  $1 : p$  do
4   |   | Draw  $X_i^{(t+1)} \sim \pi_{i|-i}(X_i | X_1^{(t+1)}, \dots, X_{i-1}^{(t+1)}, X_{i+1}^{(t)}, \dots, X_p^{(t)})$  ;
5   |   end
6 end
```

---

# Gibbs Sampler: Example

## Example

Suppose that we can generate data from  $N(\mu, \sigma^2)$ . For a bivariate normal distribution  $\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}\right)$ , we know that

$$X_2 \mid X_1 = x_1 \sim N(\mu_2 + \rho\sigma_1^{-1}\sigma_2(x_1 - \mu_1), \sigma_2^2(1 - \rho^2)).$$

Sample random numbers from the bivariate normal with Gibbs sampler.

# Why Does Gibbs Sampler Work?

The transition kernel  $K(x, y)$  is

$$\pi_{1|-1}(y_1 | x_2, \dots, x_p) \pi_{2|-2}(y_2 | y_1, x_3, \dots, x_p) \cdots \pi_{p|-p}(y_p | y_1, \dots, y_{p-1}).$$

If we consider the continuous random variables, the transition kernel satisfies

$$\pi(y) = \int K(x, y) \pi(x) dx,$$

where  $\pi(x)$  is the marginal distribution. Here  $\pi()$  is not treated as a generic symbol.

## Gibbs Sampler: Another Example

Gibbs sampler simulates random number from “correct” conditional distribution. But it can work poorly if the correlations between components are high, especially when the conditional distribution has a low variance so that the step size becomes small. The same issue can occur to Metropolis-Hastings as well.

### Example

Suppose that we can generate data from  $N(\mu, \sigma^2)$ . For a bivariate normal distribution  $\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}\right)$ , we know that

$$X_2 \mid X_1 = x_1 \sim N(\mu_2 + \rho\sigma_1^{-1}\sigma_2(x_1 - \mu_1), \sigma_2^2(1 - \rho^2)).$$

Sample random numbers from the bivariate normal with Gibbs sampler when  $\rho = 0.99$ .

# Burn-In Period

The stationary distribution is reached after large enough iterations.

- If the iterations have not proceeded long enough, the simulated numbers may be unrepresentative of the target distribution.

To diminish the influence of the starting values, we can discard the early simulations, known as the **burn-in**.

- There is no golden standard on how long the burn-in period should be.



## Serial Correlation

It is obvious that  $X_{n+1}$  and  $X_n$  are not independent draws. Inference from autocorrelated draws is generally less precise than from the same number of independent draws.

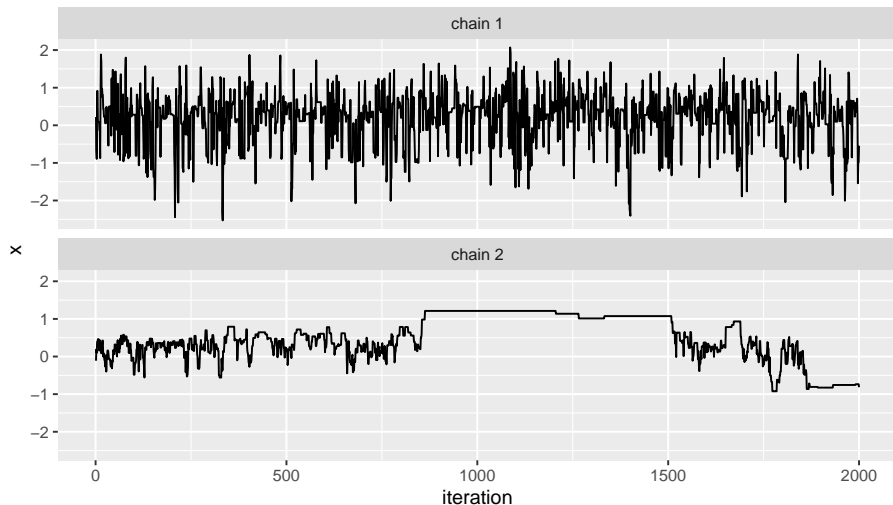
- However, such serial correlation is not necessarily a problem.
- Remember that, at convergence, we reach the stationary distribution.

Some prefer [thinning](#) the sequence by only keeping every  $k$ th draw from a sequence in order to reduce serial correlation.

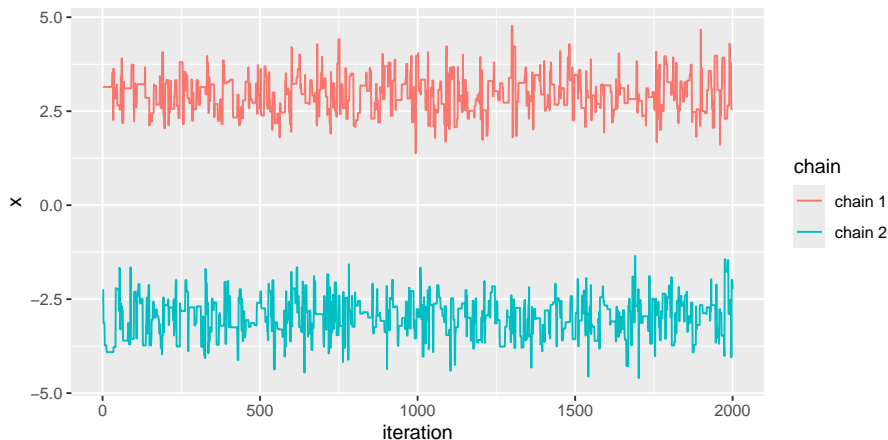
- But whether or not the Markov chain is thinned, it can be used for inferences, provided that it has reached convergence.

# Mixing

We want the Markov chain to show good [mixing](#).



# Several Markov Chains From Same Model



One suggestion is to generate several independent Markov chains, starting from widely separated places.

## Gelman-Rubin $\hat{R}$ Statistic: Variation

One way to assess convergence is the **Gelman-Rubin  $\hat{R}$  statistic**.

Suppose that we have simulated  $m$  chains each with  $n$  iterations after the burn-in period. Say we have a univariate statistic  $y_{ij} = f(x_{ij})$ , where  $x_{ij}$  is the  $i$ th value in the  $j$ th chain.

- The variation within the chains is measured by

$$W = \frac{1}{m} \sum_{j=1}^m \left[ \frac{1}{n-1} \sum_{i=1}^n (y_{ij} - \bar{y}_{\cdot j})^2 \right],$$

where  $\bar{y}_{\cdot j}$  is the average of  $\{y_{ij}\}$ .

- The variation between the chains is measured by

$$B = \frac{n}{m-1} \sum_{j=1}^m (\bar{y}_{\cdot j} - \bar{y}_{\cdot \cdot})^2,$$

where  $\bar{y}_{\cdot \cdot}$  is the average of all  $\bar{y}_{\cdot j}$ .

## Gelman-Rubin $\hat{R}$ Statistic: Expression

Consider

$$\hat{V} = \frac{n-1}{n}W + \frac{1}{n}B.$$

- If the Markov chains have reached stationary as  $n \rightarrow \infty$ , then  $\hat{V}$  is unbiased for  $\text{Var}(Y)$ .
- $E[W]$  approaches  $\text{Var}(Y)$  when  $n \rightarrow \infty$ .

The Gelman-Rubin  $\hat{R}$  statistic is then

$$\hat{R} = \sqrt{\frac{\hat{V}}{W}},$$

which declines to 1 as  $n \rightarrow \infty$ .

- It is suggested that we keep simulating the Markov chain until  $\hat{R} < 1.1$ , or even 1.01.

## Variants of Gelman-Rubin $\hat{R}$

Several different versions of  $\hat{R}$  have been proposed.

- 1 One suggestion is to change  $\hat{V}$  to

$$\hat{V} = \frac{n-1}{n}W + \frac{1}{n} \left(1 + \frac{1}{m}\right) B.$$

to account for the possibility that  $\hat{V}$  is too low.

- 2 Another suggestion is to split each chain into two parts, yields  $2m$  chains of length  $n/2$  each. Then compute the  $\hat{R}$ , pretending that we have simulated  $2m$  chains of length  $n/2$ .
  - This can be useful to detect the case where each chain does not reach stationary but the chains cover a common distribution, e.g, two chains exhibit an X-shape.

## Geweke Diagnostic

Suppose that we have simulated a Markov chain  $\{x_i\}$ . The idea is that, the average of the first  $100p_1\%$  observations should be close to the last  $100p_2\%$  observations, where  $p_1 + p_2 < 1$  to ensure that these two segments do not overlap.

Let  $\bar{X}_i$  be the average of the  $i$ th segment. Then, the [Geweke](#) method computes

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\text{Var}(\bar{X}_1) + \text{Var}(\bar{X}_2)}},$$

which converges to  $N(0, 1)$ , as  $n \rightarrow \infty$ .

# Hamiltonian Monte Carlo

- The Metropolis algorithm and the Gibbs sampler often move too slowly through the target distribution when the dimension of the target distribution is high.
- **Hamiltonian Monte Carlo (HMC)** moves much quicker through the target distribution.
  - For each component in the target distribution, HMC adds a **momentum** variable and the proposal distribution largely depends on the momentum variable.
  - Both the component in the target distribution and the momentum are updated in the MCMC algorithm.



# Hamiltonian Dynamics

The idea of HMC originates from the [Hamiltonian dynamics](#) in physics.

- The state of a system consists of the [position](#)  $x \in \mathbb{R}^d$  and the [momentum](#)  $\phi \in \mathbb{R}^d$  of same dimension.
- The Hamiltonian is a function of  $x$  and  $\phi$ , denoted by  $H(x, \phi)$ .
- The position and the momentum can change over time  $t$ . The change is described by the [Hamilton's equations](#):

$$\frac{dx_i}{dt} = \frac{\partial H}{\partial \phi_i}, \quad \text{and} \quad \frac{d\phi_i}{dt} = -\frac{\partial H}{\partial x_i},$$

for  $i = 1, \dots, d$ .

# Potential and Kinetic Energy

For HMC, the Hamiltonian is usually

$$H(x, \phi) = U(x) + V(\phi),$$

where  $U(x) = -\log \pi(x)$  is called the **potential energy** and  $V(\phi)$  is called the **kinetic energy**.

- We want to sample from  $\pi(x)$ . Hence,  $\phi$  is artificial.
- We often let  $\phi \sim N(0, M)$ , independent of  $X$ , for a prespecified covariance matrix  $M$ , and  $V(\phi)$  the negative log density of  $\phi$ .
- The Hamilton's equations become

$$\frac{dx}{dt} = M^{-1}\phi, \quad \text{and} \quad \frac{d\phi}{dt} = \frac{\partial \log \pi(x)}{\partial x},$$

arranged as column vectors.

# Augmentation

Since  $X$  and  $\phi$  are independent, their joint density is

$$\begin{aligned} p(x, \phi) = \pi(x) p(\phi) &= \exp \{-U(x) - V(\phi)\} \\ &= \exp \{-H(x, \phi)\}. \end{aligned}$$

We have augmented the problem from sampling  $X$  from  $\pi(x)$  to sampling  $(X, \phi)$  from  $\exp \{-H(x, \phi)\}$ .

- ① We first sample  $\phi$  from  $N(0, M)$ , independent of current  $x$ .
  - Since  $\phi \sim N(0, M)$ , we already sample  $\phi$  from the desired distribution.
- ② We then sample  $X$ , where the new state is proposed by Hamiltonian dynamics by solving the differential equations.

# Solve Differential Equation

To solve the differential equations, we consider an approximation known as the **leapfrog method**. For some stepsize  $\epsilon > 0$ , we perform **half-step updates** as

$$\begin{aligned}\phi\left(t + \frac{\epsilon}{2}\right) &= \phi(t) + \frac{\epsilon}{2} \frac{\partial \phi(t)}{\partial t} = \phi(t) + \frac{\epsilon}{2} \frac{\partial \log \pi(x(t))}{\partial x}, \\ x(t + \epsilon) &= x(t) + \epsilon \frac{\partial x\left(t + \frac{\epsilon}{2}\right)}{\partial t} = x(t) + \epsilon M^{-1} \phi\left(t + \frac{\epsilon}{2}\right), \\ \phi(t + \epsilon) &= \phi\left(t + \frac{\epsilon}{2}\right) + \frac{\epsilon}{2} \frac{\partial \phi(t + \epsilon)}{\partial t} = \phi\left(t + \frac{\epsilon}{2}\right) + \frac{\epsilon}{2} \frac{\partial \log \pi(x(t + \epsilon))}{\partial x}.\end{aligned}$$

Starting from  $t = 0$ , we can get a trajectory at times  $\{\epsilon, 2\epsilon, \dots, L\epsilon\}$ , and approximate the values for  $x(L\epsilon)$  and  $\phi(L\epsilon)$ .

# Leapfrog Method to Sample $X$

Suppose that the current state is  $(x, \phi)$ .

- ① Update  $\phi$  with a half-step update by

$$\phi \leftarrow \phi + \frac{\epsilon}{2} \frac{\partial \log \pi(x)}{\partial x}.$$

- ② For  $\ell = 1, \dots, L - 1$ ,

- ① Update the position:  $x \leftarrow x + \epsilon M^{-1} \phi$ .
- ② Update the momentum:

$$\phi \leftarrow \phi + \epsilon \frac{\partial \log \pi(x)}{\partial x}.$$

- ③ Make one last update on the position:  $x \leftarrow x + \epsilon M^{-1} \phi$ .
- ④ Make one last half-step update of the momentum

$$\phi \leftarrow \phi + \frac{\epsilon}{2} \frac{\partial \log \pi(x)}{\partial x}.$$

## Metropolis Step

Suppose that the state after such  $L$  updates is  $(x^*, \phi^*)$ . We **negate** the momentum and the new proposal state is  $(x^*, -\phi^*)$ .

- We determine whether to accept the proposal using the Metropolis algorithm, where the acceptance probability is

$$A((x, \phi), (x^*, -\phi^*)) = \min \left\{ 1, \frac{\exp \{-H(x^*, -\phi^*)\}}{\exp \{-H(x, \phi)\}} \right\}.$$

- If the proposed state is accepted, then we accept  $x^*$  as a new state for  $x$ , but don't care about  $\phi^*$ .
- No matter we accept or reject the proposal, we will draw a new momentum in the next iteration, independent of previous momentum.

# Properties of HMC

Some crucial properties of the Hamiltonian dynamics for MCMC updates include

- 1 **deterministic** updates. The Hamiltonian dynamics is deterministic. After running the leapfrog loop  $L$  times, we always move the initial state  $(x_0, \phi_0)$  to the same proposal  $(x^*, \phi^*)$ .
- 2 **reversible**. The mapping from the state at time  $t$ , denoted by  $(x(t), \phi(t))$ , to the state at time  $t + s$ , denoted by  $(x(t + s), \phi(t + s))$ , is one-to-one and has an inverse mapping. If we negate the momentum, we will come back from  $(x(t + s), -\phi(t + s))$  to  $(x(t), -\phi(t))$ .

# Tuning Parameters

Tuning of HMC can occur in several places such as

- ① the distribution for the momentum (often independent normal),
- ② the scaling factor  $\epsilon$ ,
- ③ the number of leapfrog steps  $L$  per iteration.

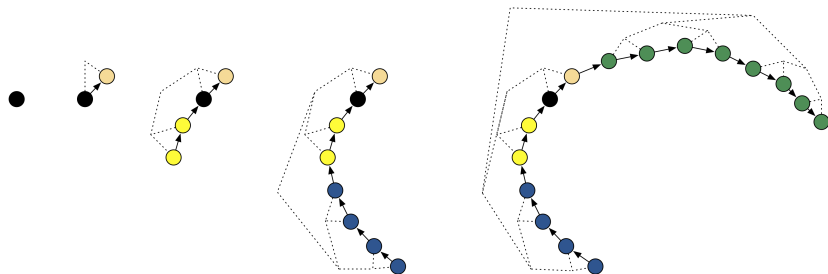
Some theory suggest that we can tune HMC such that the acceptance probability is around 65%.



# No-U-Turn Sampler

The **no-U-turn sampler** (NUTS) allows us to automatically tune the number of steps  $L$ : we increase  $L$  until the simulated dynamics is long enough such that the proposed position  $\theta^*$  starts to move back towards the initial position  $\theta$  if we run more steps.

- This is measured by the angle between  $\theta^* - \theta$  and current momentum  $\phi^*$ .



# Basic Algorithm of NUTS

A basic NUTS works as follows. Given the initial status,

- 1 Sample  $u \mid \theta, \phi \sim \text{Uniform} [0, \exp \{-H(\theta, \phi)\}]$ .
- 2 Apply the leapfrog method (with some modification) until a U-turn occurs.
- 3 Sample from the points in

$$\{(\theta, \phi) : \exp \{-H(\theta, \phi)\} \geq u\}$$

that the leapfrog step has visited and the detailed balance condition is fulfilled.

## Adaptively Tune $\epsilon$

A too small  $\epsilon$  will waste computation by taking needlessly tiny steps, and a too large will cause high rejection rates.

- In HMC, we tune  $\epsilon$  in the warm-up stage of MCMC such that the average acceptance probability  $\delta$  is the user specified value.
- In NUTS, there is no Metropolis accept/reject step. But we can still compute the ratio as if we were using the accept/reject step and set  $\epsilon$  such that the pseudo acceptance probability is the user specified value.

In [stan](#), the default is  $\delta = 0.8$ .