

# Chapter 0: A general introduction to the course

The example considers the study of the impact advertising expenditure may have on the interest of a product, which is measured by the percentage of  $\pi$  of people who remembered having watched the advertisement on TV. A common model for this example is logistic regression as the following:

$$\log \left( \frac{\pi}{1 - \pi} \right) = \alpha + \beta \mathbf{x} \quad (0.1)$$

where  $\mathbf{x}$  is the expenditure of the advertisement. The larger the value of  $\beta$ , the more effective the advertisement campaign is in boosting awareness.

Once the model is built, our task is to use methods in statistics to estimate  $\alpha, \beta$ . In the course 'Applied Statistics', we learnt linear regression method, i.e., given observed data  $\{(\pi_i, x_i)\}_{1 \leq i \leq n}$ , where  $i$  denotes some place. we use least square method as the following:

$$\min_{\alpha, \beta} \left\{ \frac{1}{n} \sum_{i=1}^n \left[ \log \frac{\pi_i}{1 - \pi_i} - \alpha - \beta x_i \right]^2 \right\}.$$

In this course, we shall study the inference problem of  $\alpha, \beta$  via Bayesian approach, in which we take the observed data  $\{(\pi_i, x_i)\}_{1 \leq i \leq n}, \alpha, \beta$  as random variables. We call the parameters  $\alpha, \beta$  (to be estimated) as unknowns, and the observed data  $\{(\pi_i, x_i)\}_{1 \leq i \leq n}$  as knowns. The distribution of  $(\alpha, \beta)$  given the known is called the **posterior distribution**. Obtaining posterior distribution is an important step but not the final one, one must draw meaningful information from it.

The model (0.1) is extremely simple but usually fails to describe the real situation, for instance, the parameters  $(\alpha, \beta)$  may change with the time  $t$ , in which the

model needs to be modified as

$$\log \left( \frac{\pi_t}{1 - \pi_t} \right) = \alpha_t + \beta_t \mathbf{x}_t.$$

Now, the number of unknown quantities has risen dramatically from 2 to  $2n$  if our observations last for  $n$ , consequently the inference problem will be much more complex to handle.

For complex problems in statistics, people often develop approximation methods to get a not exact but acceptable solution. In this course, we only study **stochastic approximation**, in which stochastic simulation plays a central role. We will learn techniques for sampling random variables and Monte Carlo (MC) method.

However, a simple MC is not enough to handle high dimensional statistical problems. So we will study a much more advanced stochastic simulation, Markov Chain Monte carlo (MCMC), to sample high dimensional random variables. We shall learn two typical sampling methods: Gibbs sampling, Metropolis-Hastings algorithms.

## Chapter 1. stochastic Simulation

We start from a random quantity  $U$  generating from uniform distribution on the interval  $[0, 1]$ , denoted by  $U \sim U([0, 1])$ .  $U$  can be easily sampled by randomly choosing a number from  $[0, 1]$ .

In the generic case of a random quantity  $x$  assigning values  $\{x_1, \dots, x_k\}$  with respective probabilities  $p_1, \dots, p_k$  subject to  $p_1 + \dots + p_k = 1$ , the interval  $[0, 1]$  is split

Chapter 1:

\* Section 1.2.2. is not required!

\* Be careful on accept-reject sampling in Chapter 1.

\* Section 2.2 is not required (Example 2.2 is required to do the calculation)

\* Section 2.3. is not required!

\* Theorem 3.1 in Section 3.4 is not required.

\* Section 4.2 (decomposition of state space) is not required.

### 1.1. Generation of discrete random quantities

\* Metropolis algorithm (describe it!)

it is a reversible. M.C.

into  $k$  intervals  $I_1, \dots, I_k$  with  $I_i = (F_{i-1}, F_i]$  where  $F_0 = 0, F_i = p_1 + \dots + p_i$  for  $1 \leq i \leq k$ . We take a random quantity  $u$  from  $U([0, 1])$ , if  $u \in I_i$ , the random quantity  $X$  takes the value  $x_i$ . We can see that  $P(X = x_i) = P(u \in I_i) = F_i - F_{i-1} = p_i$ .

- Bernoulli distribution

$$P(X = 1) = p \quad P(X = 0) = 1 - p.$$

$$k = 2, x_1 = 1, x_2 = 0, F_0 = 0, F_1 = p, F_2 = 1.$$

- Binomial distribution  $\text{Bin}(n, p)$

$$p_i = P(X = i) = \binom{n}{i} p^i (1-p)^{n-i}, i = 0, 1, \dots, n.$$

$$k = n + 1; x_1 = 0, x_2 = 1, \dots, x_k = n$$

$$F_0 = 0, F_i = p_1 + \dots + p_i \text{ for } i = 1, 2, \dots, n + 1$$

## 1.2. Generation of continuous random quantities

### 1.2.1. Probability integral transforms

Consider drawing a random quantity  $X$  from a continuous probability distribution with the distribution function  $F$ . We know  $F$  is a continuous nondecreasing function if  $F$  has an inverse  $F^{-1}$ , then  $X = F^{-1}(U)$ , where  $U$  is a random quantity drawn from  $U([0, 1])$ , is a random quantity as desired. Indeed,

$$P(X \leq z) = P(F^{-1}(U) \leq z) = P(U \leq F(z)) = F(z), \forall z \in \mathbb{R}$$

**Example:** Exponential distribution  $\text{Exp}(1)$ .

$\text{Exp}(1)$  has a probability density function:  $f(z) = \begin{cases} e^{-z}, & z \geq 0, \\ 0, & z < 0. \end{cases}$

Its distribution function is  $F(z) = \begin{cases} 1 - e^{-z}, & z \geq 0, \\ 0, & z < 0. \end{cases}$

We only need to concentrate on  $F(z)$  on  $[0, \infty)$ , and have

$$F^{-1}(z) = -\log(1 - z).$$

So  $F^{-1}(U) = -\log(1 - U)$  has a probability distribution Exp (1). Because  $1 - U \sim U([0, 1])$ , we have  $-\log U \sim \text{Exp}(1)$ .

For a distribution function which does not have an inverse, we define a generalized inverse as the following:

$$F^-(z) = \inf\{x \in \mathbb{R} : F(x) \geq z\}.$$

**Assignment:** Sample a random quantity  $Z \sim \text{Exp}(\lambda)$  for some  $\lambda > 0$ .

### 1.2.2. Bivariate techniques

If  $X = (X_1, X_2)$  has a joint density  $f_X(x_1, x_2)$  and  $g: (x_1, x_2) \rightarrow (y_1, y_2)$  is a differentiable transform with the inverse  $(x_1, x_2) = g^{-1}(y_1, y_2)$ . Then the random vector  $(Y_1, Y_2) = g(X_1, X_2)$  has a density

$$f_Y(y_1, y_2) = f_X(g^{-1}(y_1, y_2)) J(y_1, y_2)$$

where  $J(y_1, y_2)$  is the Jacobi matrix  $J(y_1, y_2) = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_2}{\partial y_1} \\ \frac{\partial x_1}{\partial y_2} & \frac{\partial x_2}{\partial y_2} \end{vmatrix}$ .

**Example:** Let  $U_1 \sim U([0, 1])$ ,  $U_2 \sim U([0, 1])$  be independent, define

$$x_1 = \sqrt{-2 \log u_1} \cos(2\pi u_2), x_2 = \sqrt{-2 \log u_1} \sin(2\pi u_2)$$

i.e.  $(x_1, x_2) = g(u_1, u_2) = (\sqrt{-2 \log u_1} \cos(2\pi u_2), \sqrt{-2 \log u_1} \sin(2\pi u_2))$ .

$$g^{-1}(x_1, x_2) = \left( \exp\left(-\frac{x_1^2 + x_2^2}{2}\right), \frac{1}{2\pi} \tan^{-1}\left(\frac{x_2}{x_1}\right) \right),$$

$$f_U(u_1, u_2) = 1 \text{ for } (u_1, u_2) \in [0, 1] \times [0, 1].$$

$$J(x_1, x_2) = \begin{vmatrix} \frac{\partial u_1}{\partial x_1} & \frac{\partial u_2}{\partial x_2} \\ \frac{\partial u_1}{\partial x_2} & \frac{\partial u_2}{\partial x_2} \end{vmatrix} = \begin{vmatrix} \frac{\partial g_1^{-1}(x_1, x_2)}{\partial x_1} & \frac{\partial g_2^{-1}(x_1, x_2)}{\partial x_1} \\ \frac{\partial g_1^{-1}(x_1, x_2)}{\partial x_2} & \frac{\partial g_2^{-1}(x_1, x_2)}{\partial x_2} \end{vmatrix} = \frac{1}{2\pi} e^{-\frac{x_1^2 + x_2^2}{2}}$$

$$\text{Hence } f_X(x_1, x_2) = \frac{1}{2\pi} \exp\left(-\frac{x_1^2 + x_2^2}{2}\right).$$

**Example:** Ratio of uniforms methods (Not required in theory, but required in python program)

In data science, we often come across a distribution whose density is known up to a constant, e.g. a random variable  $X$  with a density

$$f_X(x) = \frac{1}{C} e^{-x^4}$$

the constant  $C = \int_{-\infty}^{+\infty} e^{-x^4} dx$  is not known. Taking

$$f^*(x) = e^{-x^4},$$

define  $C_f = \{(x_1, x_2) : 0 \leq x_1 \leq \sqrt{f^*(x_2/x_1)}\}$ . Let  $(U_1, U_2)$  be the random vector uniformly distributed on  $C_f$ . Then  $Y = \frac{U_2}{U_1}$  has a density  $\frac{f^*(y)}{\int_{-\infty}^{+\infty} f^*(y) dy}$ .

Indeed, let  $x = u_1, y = \frac{u_2}{u_1}$ , i.e.

$$(x, y) = g(u_1, u_2) = \left(u_1, \frac{u_2}{u_1}\right).$$

$$(u_1, u_2) = g^{-1}(x, y) = (x, xy).$$

$$J(x, y) = \begin{vmatrix} \frac{\partial g_1^{-1}(x, y)}{\partial x} & \frac{\partial g_2^{-1}(x, y)}{\partial x} \\ \frac{\partial g_1^{-1}(x, y)}{\partial y} & \frac{\partial g_2^{-1}(x, y)}{\partial y} \end{vmatrix} = \begin{vmatrix} 1 & y \\ 0 & x \end{vmatrix} = x.$$

Let  $f_{(X,Y)}(x,y) = x \cdot k$ , where  $k = \frac{1}{\text{Area}(C_f)}$ , since  $f(U_1, U_2)(u_1, u_2) = k$  for  $(u_1, u_2) \in C_f$ . Now we consider the distribution of  $Y$ , which is given by

$$\begin{aligned} f_Y(y) &= \int_{-\infty}^{+\infty} f_{(X,Y)}(x,y)dx \\ &= \int_{-\infty}^{+\infty} x \cdot k 1_{C_f} dx \\ &= \int_0^{\sqrt{f^*(y)}} k x dx \\ &= k f^*(y) = \frac{f^*(y)}{\int_{-\infty}^{+\infty} f^*(y)dy}. \end{aligned}$$

**Assignment:** Make a python or  $R$  program to generate 1000 random quantities from the probability distribution with a density  $f_X(x) = \frac{1}{C} e^{-\frac{|x|^5}{2}}$  with  $C = \int_{-\infty}^{+\infty} e^{-\frac{|x|^5}{2}} dx$ .

### 1.2.3. Methods based on mixture

**Example 1:** Mixture of Gaussian distribution

Let  $X_1 \sim N(\mu_1, \sigma_1^2)$ ,  $X_2 \sim N(\mu_2, \sigma_2^2)$  with  $\sigma_1 > 0$  and  $\sigma_2 > 0$ , and let  $X_1$  and  $X_2$  be independent. Let  $\Theta \sim \text{Ber}(p)$  with  $p \in (0, 1)$  and let  $\Theta$  be independent of  $X_1$  and  $X_2$ . Then,  $(\Theta)X_1 + (1 - \Theta)X_2$  is a simplest Gaussian mixture model.

**Assignment:**  $\Theta X_1 + (1 - \Theta)X_2$  has a probability density

$$\frac{p}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{(x - \mu_1)^2}{2\sigma_1^2}\right) + \frac{1-p}{\sqrt{2\pi}\sigma_2} \exp\left(-\frac{(x - \mu_2)^2}{2\sigma_2^2}\right).$$

(Hint:  $P(\Theta X_1 + (1 - \Theta)X_2 \leq x) = P(X_1 \leq x | \Theta = 1)P(\Theta = 1) + P(X_2 \leq x | \Theta = 0)P(\Theta = 0)$  and the independence of  $\Theta$  of  $X_1$  and  $X_2$ .)

We now aim to sample a random quantity from the distribution of the above

Gaussian Mixture Model.

### Pseudo-code

- Sample a random quantity  $\Theta \sim \text{Ber}(p)$

If  $\Theta = 1$ , sample a random quantity  $X \sim N(\mu_1, \sigma_1^2)$ ,

If  $\Theta = 0$ , sample a random quantity  $X \sim N(\mu_2, \sigma_2^2)$ .

- Return  $Y = X$ .

**Assignment:** Realize the above pseudocode by programming to create 1000 random quantities, and draw a histogram. (Hint: To draw a random quantity  $X$  from  $N(\mu, \sigma^2)$ , first draw a random quantity  $Z \sim N(0, 1)$  and take  $X = \mu + \sigma Z$ .)

**Example 2:**  $t$ -distribution.

The  $t$ -distribution arises as the sampling distribution of the  $t$ -test. Let  $X_1, \dots, X_n$  be i.i.d with the distribution  $N(\mu, \sigma^2)$ , the sample mean and the sample variance are respectively as the below:

$$\bar{X} = \frac{1}{n} (X_1 + \dots + X_n), \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

The  $t$ -statistic is defined as

$$T = \frac{\bar{X} - \mu}{\sqrt{S^2/n}},$$

$T$  has a  $t$ -distribution with a degree  $n - 1$ , whose density function reads as

$$f(x) = \frac{\tau\left(\frac{n}{2}\right)}{\sqrt{(n-1)\pi\tau\left(\frac{n-1}{2}\right)}} \left(1 + \frac{x^2}{n-1}\right)^{-\frac{n}{2}}$$

where  $\tau(\cdot)$  is defined as  $\tau(z) = \int_0^\infty x^{z-1} e^{-x} dx$  for  $z > 0$ .

How to sample random quantity  $Z$  with  $t$ -distribution with  $(n - 1)$  degree?

A remarkable theorem about  $T$  is

$$T = \frac{Z}{\sqrt{\frac{Y}{n-1}}}$$

where  $Z \sim N(0, 1)$  and  $Y \sim \chi^2(n - 1)$ , chi-square distribution with a degree  $n - 1$ , moreover  $Z$  and  $Y$  are independent.

**Pseudo-Code:**

- Sample a random quantity  $Y$  from the distribution  $\chi^2(n - 1)$ ,
- Sample a random quantity  $Z \sim N(0, \frac{Y}{n-1})$
- Return  $X = Z$ .

**Assignment:** Realize the above pseudo-code to generate 1000 random quantities by programming, and draw a histogram.

### 1.3. Resampling methods

Resampling method consists of two steps, the first one provides a random quantity from an approximate distribution, while the second is a correction mechanisms. In what follows, we denote by  $f$  the target distribution and by  $g$  the auxiliary distribution.

We only study the classical accept-rejection method in this course.

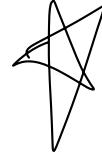
We aim to draw a random quantity  $X$  from a distribution whose probability density is  $f$ , which may be difficult to be drawn. We introduce an auxiliary probability

distribution  $g$  and take the following assumption:

(A) Assume that there exists some  $M > 0$  such that

$$\frac{f(x)}{g(x)} \leq M \quad \text{for all } x.$$

**Pseudo code for accept-rejection method:**



- Sample a random quantity  $Y$  from the distribution  $g$
- Sample a random quantity  $U$  from  $U([0, 1])$ .

If  $U \leq \frac{f(Y)}{Mg(Y)}$ , accept, ie.  $X = Y$

If  $U > \frac{f(Y)}{Mg(Y)}$ , reject and repeat the procedure.

Theoretical proof for this algorithm: Let us consider the one dimension case. For any  $x \in \mathbb{R}$ , consider

$$\begin{aligned} P(X \leq x) &= P\left(Y \leq x \mid U \leq \frac{f(Y)}{Mg(Y)}\right) \\ &= \frac{P\left(Y \leq x, U \leq \frac{f(Y)}{Mg(Y)}\right)}{P\left(U \leq \frac{f(Y)}{Mg(Y)}\right)} \\ &= \frac{\int_{-\infty}^x \left( \int_0^{\frac{f(y)}{Mg(y)}} du \right) g(y) dy}{\int_{-\infty}^{+\infty} \left( \int_0^{\frac{f(y)}{Mg(y)}} du \right) g(y) dy} \left( \because \frac{f(y)}{Mg(y)} \leq 1 \text{ for all } y \right) \\ &= \frac{\int_{-\infty}^x f(y) dy}{\int_{-\infty}^{+\infty} f(y) dy} \end{aligned}$$



So,  $X$  has a distribution with the density  $f$ .

**Assignment:** Repeat the proof of the algorithm.

**Remark 1.** Accept-rejection method can be used to sample a random quantity whose distribution is known up to a constant. An example is to sample a random quantity  $X$  from the distribution with the following density

$$\frac{1}{C} \cdot \frac{1}{1 + |x - 2|^3} \quad (1.3.1)$$

where  $C = \int_{-\infty}^{+\infty} \frac{1}{1 + |x - 2|^3} dx$ , which can not be explicitly computed out.

**Assignment:** Write the pseudo-code for sampling a random quantity from the distribution (1.3.1) via accept-rejection method (Hint: take  $M = 5$ ). Make a program to run the accept-rejection procedure 1000 times, count how many random quantities are created, and draw a histogram.

## Chapter 2. Bayesian Inference

### 2.1. Bayes' Theorem

In statistics, we often need to use the observed data  $x_1, \dots, x_n$  to estimate the parameter  $\theta$  in a model by some statistical methods such as maximum likelihood function.

**Example 2.1:** Consider a series of measurements about a physical quantity  $\mu$  with measurement error  $\varepsilon_i \sim N(0, \sigma^2)$ , where  $\sigma$  is not known. In this case, the measurements are

$$x_i = \mu + \varepsilon_i, \quad i = 1, 2, \dots, n.$$

The likelihood function of  $\theta = (\mu, \sigma)$  is

$$\ell(\theta) = \prod_{i=1}^n P_\theta(x_i) = \prod_{i=1}^n \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \right].$$

The maximum likelihood estimator is

$$\hat{\theta} = \operatorname{argmax}_{\theta} \ell(\theta)$$

In this estimation, we only consider it from the observed data  $x_1, \dots, x_n$ .

However, in practice, we may have some prior knowledge about  $\theta$ , which is natural to be taken into account in the estimate. We need to use Bayes' theorem.

### 2.1.1. Prior, posterior and predictive distributions

For the parameter  $\theta$  to be estimated, we assign it a **prior distribution  $P(\theta)$**  according to our prior knowledge about it.

Given a  $\theta$ , the observed data  $x$  has a conditional distribution  $f(x | \theta)$ . The distribution of  $\theta$  after observing  $x$  is denoted by  $p(\theta | x)$ . By Bayes' Theorem,

$$p(\theta | x) = \frac{f(x | \theta)p(\theta)}{f(x)}$$

where  $f(x) = \int f(x | \theta)p(\theta)d\theta$ . Indeed,

$$\begin{aligned} p(\theta | x) &= \frac{f(x, \theta)}{f(x)} = \frac{f(x | \theta)p(\theta)}{\int_{-\infty}^{+\infty} f(x, \theta)d\theta} \\ &= \frac{f(x | \theta)p(\theta)}{\int_{-\infty}^{+\infty} f(x | \theta)p(\theta)d\theta}. \end{aligned}$$

Because  $f(x)$  does not depend on  $\theta$ , we may write

$$p(\theta | x) \propto f(x | \theta)p(\theta)$$

In statistical models,  $f(x | \theta)$  is the likelihood function  $\ell(\theta)$  given the observed data  $x$ . (e.g.  $x = (x_1, \dots, x_n)$  in Example 2.1). We usually suppress the dependence on  $x$  and write

$$\ell(\theta) = f(x | \theta), \pi(\theta) = p(\theta | x).$$

Then

$$\pi(\theta) \propto \ell(\theta)p(\theta)$$

**Example 2.1** (continued) Suppose that  $\sigma^2$  is known, so  $\theta = \mu$  is the parameter to be estimated. Let us consider this model in Bayesian inference setting. Suppose that a prior distribution  $p(\mu) = \frac{1}{\sqrt{2\pi}\tau_0} e^{-\frac{(\mu-\mu_0)^2}{2\tau_0^2}}$ , where  $\mu_0$  and  $\tau_0$  are both known.

The likelihood function is

$$\begin{aligned}\ell(\mu) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) \\ &\propto \exp\left(-\frac{n}{2\sigma^2} (\bar{x} - \mu)^2\right)\end{aligned}$$

where  $\bar{x} = \frac{1}{n} (x_1 + \dots + x_n)$ , and the second  $\propto$  is due to

$$\begin{aligned}\sum_{i=1}^n (x_i - \mu)^2 &= n\mu^2 - 2\sum_{i=1}^n x_i \cdot \mu + \sum_{i=1}^n x_i^2 \\ &= n(\mu^2 - 2\bar{x}\mu) + \sum_{i=1}^n x_i^2 \\ &= n(\mu - \bar{x})^2 + \sum_{i=1}^n x_i^2 - n\bar{x}^2 \\ &= n(\mu - \bar{x})^2 + n \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2\right)}_{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}.\end{aligned}$$

Hence,

$$\begin{aligned}\pi(\mu) &= \ell(\mu)p(\mu) \\ &\propto \exp\left(-\frac{1}{2} \cdot \frac{(\mu - \bar{x})^2}{\sigma^2/n}\right) \exp\left(-\frac{1}{2} \cdot \frac{(\mu - \mu_0)^2}{\tau_0^2}\right) \\ &= \exp\left(-\frac{1}{2} \cdot \frac{(\mu - \mu_1)^2}{\tau_1^2} - \frac{1}{2} \cdot \frac{(\bar{x} - \mu_0)^2}{n^{-1}\sigma^2 + \tau_0^2}\right)\end{aligned}$$

where  $\tau_1^{-2} = n\sigma^{-2} + \tau_0^{-2}$ ,  $\mu_1 = \tau_1^2 (n\sigma^{-2}\bar{x} + \tau_0^{-2}\mu_0)$ . So

$$\pi(\mu) \propto \exp\left(-\frac{1}{2} \cdot \frac{(\mu - \mu_1)^2}{\tau_1^2}\right),$$

this means that the posterior distribution of  $\mu$  is  $N(\mu_1, \tau_1^2)$  with

$$\mu_1 = \tau_1^2 (n\sigma^{-2}\bar{x} + \tau_0^{-2}\mu_0), \tau_1^2 = \frac{1}{n\sigma^{-2} + \tau_0^{-2}}.$$

So, we can see that the information of the data has effect on the parameter  $\mu$ .

If we increase  $\tau_0$ , we get less information from the prior distribution  $p(\mu)$ .

### 2.1.2. Summarizing the information from the posterior distribution.

Once the posterior distribution is available, one may seek to summarize its information through a few elements such as the location and dispersion measures. The main location measures are mean, mode and median, while the main dispersion measures are variance, standard deviation, precision, interquantile range and curvature at the mode.

Let  $\theta$  be a scalar, we give the definitions of the above concepts:

- Posterior mean =  $\int \theta \pi(\theta) d\theta$
- Posterior variance =  $\int (\theta - \int \theta \pi(\theta) d\theta)^2 \pi(\theta) d\theta$
- Posterior median is the median of the distribution  $\pi$ , ie. a point  $\theta_0$  such that

$$\int_{-\infty}^{\theta_0} \pi(\theta) d\theta = \frac{1}{2}.$$

- Credibility interval  $C$  : given an  $\alpha > 0$  (e.g.  $\alpha = 0.05$ ), an interval  $C$  is called  $100(1 - \alpha)\%$  credibility interval for  $\theta$  if  $\int_C \pi(\theta) d\theta = 1 - \alpha$ . We usually look for a credibility interval with the shortest length. This interval is called highest posterior density (HPD) interval.

**Example 2.1** (continued) We have shown that the posterior distribution of  $\mu \sim N(\mu_1, \tau_1^2)$ . So the posterior mean, median and mode are all  $\mu_1$ , the posterior variance is  $\tau_1^2$  and the HPD interval of  $(1 - \alpha)$  is  $(\mu_1 - \tau_1 z_{\alpha/2}, \mu_1 + \tau_1 z_{\alpha/2})$  where  $z_{\alpha/2}$  satisfies

$$\int_{z_{\alpha/2}}^{\infty} \varphi(x) dx = \frac{\alpha}{2}$$

## 2.2. A brief discussion on conjugate distributions

From the above Example 2.1, we choose the prior distribution  $p(\mu) \propto \exp\left(-\frac{(\mu-\mu_0)^2}{2\tau_0^2}\right)$ , which makes the derivation of the posterior distribution very simple as a change of parameters:

$$\tau_1^{-2} = n\sigma^{-2} + \tau_0^{-2}, \quad \mu_1 = \tau_1^2 (n\sigma^{-2}\bar{x} + \tau_0^{-2}\mu_0)$$

**Example 2.1** (continued). After observing  $n_1$  samples  $x_1^{(1)}, x_2^{(1)}, \dots, x_{n_1}^{(1)}$  from the population with the distribution  $N(\mu, \sigma^2)$ , we obtain a posterior distribution about  $\mu$  as  $N(\mu_1, \tau_1^2)$  with  $\tau_1^{-2} = n_1\sigma^{-2} + \tau_0^{-2}$ ,  $\mu_1 = \tau_1^2 (n_1\sigma^{-2}\bar{x}^{(1)} + \tau_0^{-2}\mu_0)$ , with  $\bar{x}^{(1)} = \frac{1}{n_1} (x_1^{(1)} + \dots + x_{n_1}^{(1)})$ . Now a new sequence of samples  $x_1^{(2)}, \dots, x_{n_2}^{(2)}$  are observed, we take  $N(\mu_1, \tau_1^2)$  as the prior distribution, and get the new posterior distribution as  $N(\mu_2, \tau_2^2)$  with

$$\tau_2^{-2} = n_2\sigma^{-2} + \tau_1^{-2}, \quad \mu_2 = \tau_2^2 (n_2\sigma^{-2}\bar{x}^{(2)} + \tau_1^{-2}\mu_1),$$

where  $\bar{x}^{(2)} = \frac{1}{n_2} (x_1^{(2)} + \dots + x_{n_2}^{(2)})$ . In other words, the new posterior distribution is obtained based on the observation  $x^{(1)}$  and  $x^{(2)}$ .

Because we take the prior distribution  $p(\theta) \sim N(\mu_0, \tau_0^2)$ , we can obtain the posterior distribution  $\pi(\theta) \sim N(\mu_1, \tau_1^2)$  and continue the update of posterior distributions for the data flow  $x^{(1)}, x^{(2)}, \dots, x^{(n)}, \dots$  in the following way:

$$\tau_{i+1}^{-2} = n_{i+1}\sigma^{-2} + \tau_i^{-2}, \mu_{i+1} = \tau_{i+1}^2 (n_{i+1}\sigma^{-2}\bar{x}^{(i+1)} + \tau_i^{-2}\mu_i)$$

This  $p(\theta)$  is a conjugate distribution w.r.t.  $\ell(\theta)$ . The conjugate distributions in Bayesian inference is a deep theory, and we will not enter it but give several examples.

*Required*

**Example 2.2** Consider the Poisson model: Let  $x_1, x_2, \dots, x_n$  be a sequence of observation from Poisson distribution  $\text{Poi}(\lambda)$ . Given a  $\lambda$ ,

$$\ell(\lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}.$$

Given a prior distribution  $p(\lambda) \sim G(\alpha, \beta)$ , gamma distr. with the parameter  $\alpha, \beta$ , ie.  $p(\lambda) = \lambda^{\alpha-1} e^{-\beta\lambda}$ .

$$\begin{aligned} \pi(\lambda) &\propto \ell(\lambda)p(\lambda) \\ &\propto \prod_{i=1}^n (\lambda^{x_i} e^{-\lambda}) \lambda^{\alpha-1} e^{-\beta\lambda} \\ &= \lambda^{\sum_{i=1}^n x_i + \alpha - 1} e^{-n\lambda - \beta\lambda} \end{aligned}$$

i.e.  $\pi(\lambda) \propto \lambda^{\sum_{i=1}^n x_i + \alpha - 1} e^{-n\lambda - \beta\lambda}$ . So  $\pi(\lambda) \sim G(\sum_{i=1}^n x_i + \alpha, n\lambda + \beta)$ . We say  $G(\alpha, \beta)$  is a conjugate distribution to Poisson model.

**Example 2.3:** Consider Bernoulli model: Let  $x_1, \dots, x_n$  be a sequence of observations from the Bernoulli distr.  $\text{Ber}(\theta)$ . Then,

$$\ell(\theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}$$

(recall  $x_i \in \{0, 1\}$ ). Take the prior distribution  $P(\theta) \sim \text{Beta}(\alpha, \beta)$  with  $\alpha > 0, \beta > 0$ , i.e.

$$p(\theta) = \frac{\tau(\alpha + \beta)}{\tau(\alpha)\tau(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}, \quad \theta \in [0, 1].$$

So

$$\begin{aligned} \pi(\theta) &\propto \ell(\theta)p(\theta) \\ &\propto \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &= \theta^{\sum_{i=1}^n x_i + \alpha - 1} (1-\theta)^{n - \sum_{i=1}^n x_i + \beta - 1}. \end{aligned}$$

i.e.  $\pi(\theta) \sim \text{Beta}(\sum_{i=1}^n x_i + \alpha, n - \sum_{i=1}^n x_i + \beta)$ .

### 2.3. Dynamic models (Time series)

Dynamic linear models are defined by a pair of equations, called the observation equation and the system equation, respectively given by

$$y_t = F_t' \beta_t + \varepsilon_t, \quad \varepsilon_t \sim N(0, \sigma_t^2), \quad (2.3.1)$$

$$\beta_t = G_t \beta_{t-1} + \omega_t, \quad \omega_t \sim N(0, W_t), \quad (2.3.2)$$

where  $y_t \in \mathbb{R}$  is an observation at the time  $t$ ,  $F_t \in \mathbb{R}^d$ ,  $\beta_t \in \mathbb{R}^d$ ,  $G_t \in \mathbb{R}^{d \times d}$ , and  $W_t \in \mathbb{R}^{d \times d}$  is the covariance matrix. Moreover, the errors  $\varepsilon_t$  and  $\omega_t$  are independent.

**Example 2.4** The simplest time series model is the first order model, which is given by

$$\begin{aligned} y_t &= \beta_t + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma_t^2) \\ \beta_t &= \beta_{t-1} + \omega_t \quad \omega_t \sim N(0, W_t) \end{aligned}$$

where  $\beta_t \in \mathbb{R}$  and  $d = 1$ .

A typical linear growth model is as the following:

$$y_t = \beta_{1,t} + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma_t^2)$$

$$\beta_{1,t} = \beta_{1,t-1} + \beta_{2,t-1} + \omega_{1,t}$$

$$\beta_{2,t} = \beta_{2,t-1} + \omega_{2,t}$$

$$\omega_t = (\omega_{1,t}, \omega_{2,t}) \sim N(0, \omega_t)$$

this model is a special case of (2.3.1)-(2.3.2) with  $F_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  and  $G_t = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$   
for all  $t$  therein.

**Sequential inference:** Denote by  $y^0$  the initial information available, and  $y^i = \{y^0, y_1, \dots, y_i\}$ . For dynamic model, the inference is based on the updated distribution  $\beta_t | y^t$ . There are three basic operations involved here: evolution, prediction and updating. Let us demonstrate these three operations as below.

Consider at the time  $t - 1$ , the updated distribution is

$$\beta_{t-1} | y^{t-1} \sim N(m_{t-1}, C_{t-1})$$

Recall the system equation (2.3.2):  $\beta_t = G_t \beta_{t-1} + \omega_t$ , from which we know

$$\beta_t | \beta_{t-1} \sim N(G_t \beta_{t-1}, W_t)$$

By a standard calculation (not required in this course), we know

$$\beta_t | y^{t-1} \sim N(a_t, R_t)$$

where  $a_t = G_t m_{t-1}$  and  $R_t = G_t C_{t-1} G_t' + W_t$ .

**Assignment:** Let

$$y | x \sim N(\mu_1, \sigma_1^2)$$

$$z | y \sim N(\mu_2 y, \sigma_2^2)$$

where  $\mu_1 \in \mathbb{R}, \mu_2 \in \mathbb{R}$  and  $\sigma_1^2 > 0, \sigma_2^2 > 0$ . Show that

$$z \mid x \sim N(\mu_1\mu_2, \sigma_1^2\mu_2^2 + \sigma_2^2).$$

Now we would like to compute the updating distribution  $\beta_t \mid y^t$ .

$$\begin{aligned} p(\beta_t \mid y^t) &= p(\beta_t \mid y_t, y^{t-1}) \\ &= \frac{p(\beta_t, y_t, y^{t-1})}{p(y_t, y^{t-1})} = \frac{f(y_t \mid \beta_t) p(\beta_t \mid y^{t-1}) p(y^{t-1})}{p(y_t, y^{t-1})} \\ &= \frac{p(y^{t-1})}{p(y_t, y^{t-1})} f(y_t \mid \beta_t) p(\beta_t \mid y^{t-1}) \\ &\propto f(y_t \mid \beta_t) p(\beta_t \mid y^{t-1}) \end{aligned}$$

## Chapter 3. Approximate method of inference

Recall  $\pi(\theta)$  is the posterior distribution, we would like to draw the useful information from  $\pi(\theta)$ , such as posterior mean, posterior median and credibility region, which can be generally summarized as

$$J = \int t(\theta) \pi(\theta) d\theta$$

e.g.  $t(\theta) = \theta$ , one obtains the posterior mean,  $t(\theta) = 1(\theta + C)$  and  $J = 1 - \alpha$  give the credibility regime. However, it is often difficult to compute this integral, so people have to develop approximation methods. There are two type of approximation methods, one deterministic approximation and the other stochastic approximation.

### 3.1. Asymptotic approximation

#### 3.1.2. Normal approximation

Let  $m$  be the mode of the distribution  $\pi(\theta)$ , i.e.

$$m = \operatorname{argmax}_{\theta} \pi(\theta),$$

note that there may be more then one modes. Consider the Taylor expansion of  $\log \pi(\theta)$  at  $\theta = m$ , we get

$$\begin{aligned} \log \pi(\theta) &= \log \pi(m) + (\nabla \log \pi(m))'(\theta - m) \\ &\quad + \frac{1}{2}(\theta - m)' (\nabla^2 \log \pi(m)) (\theta - m) + R(\theta). \\ &= \log \pi(m) + \frac{1}{2}(\theta - m)' (\nabla^2 \log \pi(m)) (\theta - m) + R(\theta) \end{aligned}$$

Since  $\nabla \log \pi(m) = 0$ . Note that  $\nabla^2 \log \pi(m)$  is the Hessian matrix of  $\log \pi(\theta)$  at  $\theta = m$ . More precisely, for  $\theta = (\theta_1, \dots, \theta_d)$ ,

$$\begin{aligned} \nabla \log \pi(\theta) &= \begin{bmatrix} \frac{\partial}{\partial \theta_1} \log \pi(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_d} \log \pi(\theta) \end{bmatrix} \in \mathbb{R}^d, \\ \nabla^2 \log \pi(\theta) &= \begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} \log \pi(\theta) & \frac{\partial^2}{\partial \theta_1 \partial \theta_2} \log \pi(\theta) & \cdots & \frac{\partial^2}{\partial \theta_1 \partial \theta_d} \log \pi(\theta) \\ \frac{\partial^2}{\partial \theta_1 \partial \theta_2} \log \pi(\theta) & \frac{\partial^2}{\partial \theta_2^2} \log \pi(\theta) & \cdots & \frac{\partial^2}{\partial \theta_2 \partial \theta_d} \log \pi(\theta) \\ \cdots & \cdots & & \\ \frac{\partial}{\partial \theta_1 \partial \theta_d} \log \pi(\theta) & \frac{\partial^2}{\partial \theta_2 \partial \theta_d} \log \pi(\theta) & \cdots & \frac{\partial^2}{\partial \theta_d^2} \log \pi(\theta) \end{bmatrix} \in \mathbb{R}^{d \times d}. \end{aligned}$$

$\nabla^2 \log \pi(\theta)$  is a symmetric matrix and  $\nabla^2 \log \pi(m) = (\nabla^2 \log \pi(\theta))|_{\theta=m}$ . We drop the remainder  $R(\theta)$  and write

$$\begin{aligned}\log \pi(\theta) &\approx \log \pi(m) + \frac{1}{2}(\theta - m)' (\nabla^2 \log \pi(m)) (\theta - m) \\ &= \log \pi(m) - \frac{1}{2}(\theta - m)' (-\nabla^2 \log \pi(m)) (\theta - m),\end{aligned}$$

which implies

$$\pi(\theta) \propto \exp \left\{ -\frac{1}{2}(\theta - m)' (-\nabla^2 \log \pi(m)) (\theta - m) \right\}$$

Denote  $\Sigma = (-\nabla^2 \log \pi(m))^{-1}$ , then

$$\pi(\theta) \propto \exp \left( -\frac{1}{2}(\theta - m)' \Sigma^{-1} (\theta - m) \right).$$

**Example 3.1:** Let  $\pi(\theta) = \frac{1}{C} \exp \left( -\frac{1}{2}(\theta + \sin \theta)^2 \right)$  with  $C = \int \exp \left( -\frac{1}{2}(\theta + \sin \theta)^2 \right) d\theta$

$$\log \pi(\theta) = \log C - \frac{1}{2}(\theta + \sin \theta)^2. \quad m = 0$$

$$(\log \pi(\theta))' = -(\theta + \sin \theta)(1 + \cos \theta),$$

$$(\log \pi(\theta))'' = (\theta + \sin \theta) \sin \theta - (1 + \cos \theta)^2$$

$$(\log \pi(0))'' = -(1 + \cos 0)^2 = -4.$$

So we have the following normal approximation for  $\pi(\theta)$  :

$$\tilde{\pi}(\theta) \propto \exp(-2\theta^2)$$

**Assignment:** Let  $\pi(\theta) = \frac{1}{C} \exp \left( -\frac{1}{2}(\theta + \cos \theta)^2 \right)$  where  $C = \int e^{-\frac{1}{2}(\theta + \cos \theta)^2} d\theta$ . Derive its normal approximation.

**Example 3.2:** Consider the observation  $x \sim poi(\lambda)$  and a conjugate prior distribution  $\lambda \sim G(\alpha, \beta)$ , i.e.  $p(\lambda) = \frac{1}{C} \lambda^{\alpha-1} e^{-\beta\lambda}$  with  $C = \int_0^\infty \lambda^{\alpha-1} e^{-\beta\lambda} d\lambda$ . According to

Chapter 2, we have

$$\begin{aligned}\pi(\lambda) &\propto p(\lambda)\ell(x) \\ &\propto \lambda^{\alpha-1}e^{-\beta\lambda}\lambda^x e^{-\lambda} \\ &= \lambda^{\alpha+x-1}e^{-(\beta+1)\lambda}\end{aligned}$$

Hence  $\pi(\lambda)$  is  $G(\alpha_1, \beta_1)$  with  $\alpha_1 = \alpha + x$  and  $\beta_1 = \beta + 1$ . Let us consider the normal approximation of  $\pi(\lambda)$ . The mode is

$$\begin{aligned}m &= \frac{\alpha + x - 1}{\beta + 1} \\ \log \pi(\lambda) &= C + (\alpha + x - 1) \log \lambda - (\beta + 1)\lambda \\ (\log \pi(\lambda))'' &= -\frac{(\alpha + x - 1)}{\lambda^2} \\ (\log \pi(m))'' &= -(\alpha + x - 1) \cdot \left(\frac{\beta + 1}{\alpha + x - 1}\right)^2 = -\frac{(\beta + 1)^2}{\alpha + x - 1}.\end{aligned}$$

Hence the normal approximation is

$$\tilde{\pi}(\lambda) \propto \exp\left(-\frac{1}{2} \cdot \frac{(\beta + 1)^2}{\alpha + x - 1}(\lambda - m)^2\right).$$

So the approximate posterior variance is  $\frac{\alpha+\alpha-1}{(\beta+1)^2}$ .

### 3.2. Monte Carlo integration

Our aim is to compute the integral:

$$J = \int t(\theta)\pi(\theta)d\theta.$$

If we can sample a sequence of data  $\theta_1, \dots, \theta_n$  from  $\pi$ , then we have the approximation  $\hat{J}$  for  $J$  as the following:

$$\hat{J} = \frac{1}{n} \sum_{i=1}^n t(\theta_i).$$

By Law of Large Number (LLN), we know  $\hat{J} \rightarrow J$  a.s.s. as  $n \rightarrow \infty$ .

Quite often, sampling from  $\pi(\theta)$  is difficult, one considers first sampling from an auxiliary distribution  $q(\theta)$  and compute

$$\hat{J} = \frac{1}{n} \sum_{i=1}^n t(\theta_i) \frac{\pi(\theta_i)}{q(\theta_i)}, \theta_1, \dots, \theta_n \sim q(\theta),$$

because  $\int t(\theta)\pi(\theta)d\theta = \int \left[ t(\theta) \frac{\pi(\theta)}{q(\theta)} \right] q(\theta)d\theta$ .

Another problem in Bayesian inference is that we only know

$$\pi(\theta) \propto \ell(\theta)p(\theta).$$

Denote  $\pi^*(\theta) = \ell(\theta)p(\theta)$ , then  $\pi(\theta) = \frac{\ell(\theta)p(\theta)}{\int \pi^*(\theta)d\theta}$ . and

$$J = \frac{\int t(\theta)\pi^*(\theta)d\theta}{\int \pi^*(\theta)d\theta}.$$

Rewrite it as  $J = \frac{\int \frac{t(\theta)\pi^*(\theta)}{q(\theta)}q(\theta)d\theta}{\int \frac{\pi^*(\theta)}{q(\theta)}q(\theta)d\theta}$ , it is natural to approximates  $J$  as the following:  
draw  $\theta_1, \dots, \theta_n$  from  $q(\theta)$  and compute

$$\hat{J} = \frac{\frac{1}{n} \sum_{i=1}^n \frac{t(\theta_i)\pi^*(\theta_i)}{q(\theta_i)}}{\frac{1}{n} \sum_{i=1}^n \frac{\pi^*(\theta_i)}{q(\theta_i)}}$$

### 3.3. Bayes' theorem via the rejection method

Recall the classical accept-reject method:

We aim to draw a random quantity  $Z$  from a distribution whose probability density is  $f$ , which may be difficult to be drown. We introduce an auxiliary probability distribution  $g$  and take the following assumption:

(A) Assume that there exists some  $M > 0$  such that

$$\frac{f(x)}{g(x)} \leq M \quad \text{for all } x.$$

Pseudo code for accept-rejection method:

- Sample a random quantity  $Y$  from the distribution  $g$
- Sample a random quantity  $U$  from  $U([0, 1])$ .

If  $U \leq \frac{f(Y)}{Mg(Y)}$ , accept, i.e.  $X = Y$

If  $U > \frac{f(Y)}{Mg(Y)}$ , reject and repeat the procedure.

In the Bayes' theorem setting, we recall that the posterior distribution  $\pi(\theta)$  satisfies

$$\pi(\theta) \propto p(\theta)\ell(\theta).$$

where  $p(\theta)$  is the prior distribution and  $\ell(\theta)$  is the likelihood function. Denote

$$\pi^*(\theta) = p(\theta)\ell(\theta),$$

We aim to sample random quantities from  $\pi(\theta)$  by the accept-reject method.

- We choose  $p(\theta)$  as the auxiliary distribution, i.e.  $g(\theta)$  above,
- We choose  $M = \ell_{\max} = \max_{\theta} \ell(\theta)$ .

The corresponding accept-reject sampling is as follows:

- (1) Sample a random quantity  $Y$  from the prior distribution  $p(\theta)$ ,

- (2) Sample a random quantity  $U$  from the uniform distribution  $U([0, 1])$ 
  - (2.1) If  $U \leq \frac{\pi^*(Y)}{\ell_{\max} p(Y)} = \frac{\ell(Y)}{\ell_{\max}}$ . Accept, i.e.  $X = Y$
  - (2.2) If  $U > \frac{\ell(Y)}{\ell_{\max}}$ . Reject, i.e. do not assign a quantity to  $X$ , return to (1).
- (3) Return  $X$ .

In the practice, it is often not easy to get  $\ell_{\max}$ .

### 3.4. Bayes' theorem via weighted resampling.

Resampling is a very important and popular method for drawing random quantity from a complicated probability distribution. The procedure for drawn a random quantity from the posterior distribution is as follows:

- (1) Draw  $n$  random quantities  $\theta_1, \theta_2, \dots, \theta_n$  from a distribution  $q(\theta)$ , which can be easily sampled, and create

$$\omega_i = \frac{\pi(\theta_i)/q(\theta_i)}{\sum_{j=1}^n \pi(\theta_j)/q(\theta_j)}.$$

- (2) For these sampled  $\theta_1, \theta_2, \dots, \theta_n$ , resample them according to the probability  $P(\theta = \theta_i) = \omega_i$  for  $i = 1, 2, \dots, n$ . (It is sampled in the following way: We split the interval  $[0, 1]$  into the sub-intervals  $I_1 = [0, \omega_1], I_2 = [\omega_1, \omega_1 + \omega_2], \dots, I_k = [\omega_1 + \dots + \omega_{k-1}, \omega_1 + \dots + \omega_k], \dots, I_n = [\omega_1 + \dots + \omega_{n-1}, 1]$ . One randomly draws a quantity  $u$  from  $[0, 1]$ , if  $u \in I_i$ , then we take  $\theta = \theta_i$ ).

**Theorem 3.1** The random quantity  $\theta$  which is drown according to (1) and (2) has a distribution  $\pi$ . as  $n \rightarrow \infty$ .

**Proof:** For any  $x \in \mathbb{R}$ , consider  $P(\theta \leq x)$ . It can be decomposed into

$$\begin{aligned}
P(\theta \leq x) &= \sum_{i=1}^n P(\theta \leq x, \theta = \theta_i) \\
&= \sum_{i=1}^n P(\theta_i \leq x, \theta = \theta_i) \\
&= \sum_{i=1}^n \int_{-\infty}^x P(\theta = \theta_i | \theta_i = y) P_{\theta_i}(y) dy \\
&= \sum_{i=1}^n \int_{-\infty}^x P(\theta = y | \theta_i = y) q(y) dy \\
&\approx \sum_{i=1}^n \int_{-\infty}^x \frac{\pi(y)}{\sum_{j=1}^n \frac{\pi(\theta_j)}{q(\theta_j)}} q(y) dy \\
&= \frac{n \int_{-\infty}^x \pi(y) dy}{\sum_{j=1}^n \frac{\pi(\theta_j)}{q(\theta_j)}} \\
&= \frac{\int_{-\infty}^x \pi(y) dy}{\frac{1}{n} \sum_{j=1}^n \frac{\pi(\theta_j)}{q(\theta_j)}} = \frac{\int_{-\infty}^x \pi(y) dy}{\int_{-\infty}^{+\infty} \pi(\theta) d\theta}.
\end{aligned}$$

As  $n \rightarrow \infty$ , by law of large number, we know

$$\frac{1}{n} \sum_{j=1}^n \frac{\pi(\theta_j)}{q(\theta_j)} = \int \frac{\pi(\theta)}{q(\theta)} \cdot q(\theta) d\theta = \int_{-\infty}^{+\infty} \pi(\theta) d\theta.$$

In the Bayesian inference, we let  $q(\theta)$  be the prior distribution  $p(\theta)$ . The procedure of sampling random quantity with distribution can be done as the following:

- (1) Sampling a sequence of i.i.d. random quantities  $\theta_1, \dots, \theta_n$ , according to the probability  $p(\theta)$  and define

$$\omega_i = \frac{\pi^*(\theta_i) / p(\theta_i)}{\sum_{j=1}^n \frac{\pi^*(\theta_j)}{p(\theta_j)}}, \quad i = 1, 2, \dots, n.$$

- (2) Resample  $\theta_1, \dots, \theta_n$  obtained from (1) according to the probability  $P(\theta = \theta_i) = \omega_i$  for  $i = 1, 2, \dots, n$ .

By the same method, we can show that

$$P(\theta \leq a) = \frac{\int_{-\infty}^a \pi^*(\theta) d\theta}{\int_{-\infty}^{+\infty} \pi^*(\theta) d\theta} = \int_{-\infty}^a \pi(\theta) d\theta.$$

**Example 3.3.** A total of animals are categorized into three types with  $y = (y_1, y_2, y_3)$ , the three types are assigned probability  $(\frac{2+\theta}{4}, \frac{2-2\theta}{4}, \frac{\theta}{4})$  with  $\theta \in [0, 1]$ . Given a  $\theta$ , the conditional probability is

$$l(\theta) \quad p(y | \theta) \propto (2 + \theta)^{y_1} (2 - 2\theta)^{y_2} \theta^{y_3} \propto (2 + \theta)^{y_1} (1 - \theta)^{y_2} \theta^{y_3}$$

Suppose that the observed data are  $y = (125, 38, 34)$ , take  $\theta \sim U([0, 1])$ , then the posterior distribution is

$$\pi(\theta) \propto (2 + \theta)^{125} (1 - \theta)^{38} \theta^{34}$$

However, we do not know the number  $C = \int_0^1 (2 + \theta)^{125} (1 - \theta)^{38} \theta^{34} d\theta$ . We shall use resampling algorithm, take  $q(\theta) = p(\theta) \sim U([0, 1])$ .

- Sample  $n$  (a large number)  $\theta_1, \theta_2, \dots, \theta_n$  from  $U([0, 1])$  and compute the values

$$\omega_i = \frac{\ell(\theta_i)}{\sum_{j=1}^n \ell(\theta_j)} = \frac{(2 + \theta_i)^{125} (1 - \theta_i)^{38} \theta_i^{34}}{\sum_{j=1}^n (2 + \theta_j)^{125} (1 - \theta_j)^{38} \theta_j^{34}}$$

- Taking the obtained  $\{\theta_1, \dots, \theta_n\}$  as the sample space and resample them according to the probability  $\{\omega_i\}_{1 \leq i \leq n}$ , i.e.

$$P(\theta = \theta_i) = \omega_i.$$

**Assignment:** Make a python program to realize the resampling in Example 3.3, choose  $n = 10000$ .

# Chapter 4. Markov chains

In this chapter, we consider a special stochastic process, Markov chain, denoted by  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  with  $\mathbb{N}_0 = \{0\} \cup \mathbb{N}$ .

## 4.1. Definition and transition probabilities

A Markov chain  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  is a stochastic process such that given the present state, the past and future states are independent, more precisely,

$$\begin{aligned} & P(\theta^{(n+1)} \in A \mid \theta^{(n)} = x, \theta^{(n-1)} \in A_{n-1}, \dots, \theta^{(0)} \in A_0) \\ &= P(\theta^{(n+1)} \in A \mid \theta^{(n)} = x) \end{aligned}$$

where  $A_0, A_1, \dots, A_{n-1} \subset S$  and  $x \in S$ ,  $S$  is the state space. Generally, the conditioned probability  $P(\theta^{(n+1)} \in A \mid \theta^{(n)} = x)$  depends on  $x, A$  and  $n$ . If it does not depend on  $n$ , then the Markov chain is time homogeneous.

In this lecture, without specific statement, we only consider time homogeneous Markov chain. In this case,  $P(\theta^{(1)} \in A \mid \theta^{(0)} = x) = P(\theta^{(n+1)} \in A \mid \theta^{(n)} = x)$ , we denote  $P(x, A) = P(\theta^{(1)} \in A \mid \theta^{(0)} = x)$  for  $x \in S$  and  $A \subset S$ . We call  $P(x, A)$  a transition probability:

- For  $x \in S, P(x, \cdot)$  is a probability measure over  $S$ .
- For  $A \subset S, P(\cdot, A)$  is a function on  $S$ .

If  $S$  is a discrete space, we write

$$P(x, y) = P(x, \{y\}) \text{ for } x, y \in S.$$

It is easy to see that

$$P(x, y) \geq 0 \quad \text{for all } x, y \in S.$$

$$\sum_{y \in S} P(x, y) = 1 \quad \text{for all } x \in S$$

**Example 4.1: (Random walk on  $\mathbb{Z}$ )** Consider a particle starting from the origin and moving on  $\mathbb{Z}$ , more precisely, the positions of the particle is as follows:

$$\theta^{(0)} = 0,$$

$$\theta^{(n)} = \theta^{(n-1)} + \omega_n, n \geq 1,$$

where  $\omega_1, \omega_2, \dots$ , are i.i.d. integer valued random variables. Denote by  $f$  the probability function of  $\omega_1$ .

If  $f(1) = p, f(-1) = q, f(0) = r$  with  $p+q+r = 1$ , then the transition probability are given by

$$P(x, y) = \begin{cases} p, & y = x + 1 \\ q, & y = x - 1 \\ r, & y = x \\ 0, & \text{otherwise} \end{cases}$$

In the case of discrete space  $S = \{x_1, x_2, \dots\}$ , we can put all transition probabilities  $P(x_i, x_j)$  together to define transition probability matrix

$$\{P(x_i, x_j)\}_{x_i \in S, x_j \in S}$$

For instance, the transition probability matrix of above random walk is

$$\begin{bmatrix} \ddots & \ddots & \ddots & & \\ & q & r & p & \\ & q & r & p & \\ & q & r & p & \\ & \ddots & \ddots & \ddots & \end{bmatrix}$$

### Chapman - Kolmogorov relation

Denote  $P^m(x, y) = P(\theta^{(m)} = y \mid \theta^{(0)} = x)$ . When  $m = 2$ , we see

$$\begin{aligned} P^2(x, y) &= P(\theta^{(2)} = y \mid \theta^{(0)} = x) \\ &= \sum_{x_1} P(\theta^{(2)} = y, \theta^{(1)} = x_1 \mid \theta^{(0)} = x) \\ &= \sum_{x_1} \frac{P(\theta^{(2)} = y, \theta^{(1)} = x_1, \theta^{(0)} = x)}{P(\theta^{(0)} = x)} \\ &= \sum_{x_1} \frac{P(\theta^{(2)} = y \mid \theta^{(1)} = x_1, \theta^{(1)} = x) P(\theta^{(1)} = x_1, \theta^{(1)} = x)}{P(\theta^{(0)} = x)} \end{aligned}$$

$$\begin{aligned} \text{Markov Property} &= \sum_{x_1} P(\theta^{(2)} = y \mid \theta^{(1)} = x_1) P(\theta^{(1)} = x_1 \mid \theta^{(0)} = x) \\ &= \sum_{x_1} P(x_1, y) P(x, x_1) \end{aligned}$$

Hence.  $P^2 = P \cdot P$ .

**Assignment:** Using a similar argument to show that  $P^{m+1} = P^m P$ . and  $P^m = \underbrace{P \dots P}_m$  by an induction. Moreover,  $P^{m+n} = P^m P^n$  for all  $m \in \mathbb{N}$  and  $n \in \mathbb{N}$ .

**Remark:** It is easy to see that  $P^{m+n} = P^m P^n$  can be presented as

$$P^{m+n}(x, y) = \sum_{z \in S} P^m(x, z) P^n(z, y), \quad \forall x, y \in S.$$

For all  $n \in \mathbb{N}_0$ , denote by  $\pi^{(n)}$  the distribution of  $\theta^{(n)}$ . For  $y \in S$ ,

$$\begin{aligned}\pi^{(n)}(y) &= P(\theta^{(n)} = y) \\ &= \sum_{x \in S} P(\theta^{(n)} = y, \theta^{(0)} = x) \\ &= \sum_{x \in S} P(\theta^{(n)} = y \mid \theta^{(0)} = x) P(\theta^{(0)} = x) \\ &= \sum_{x \in S} P^n(x, y) \pi^{(0)}(x)\end{aligned}$$

for all  $x, y \in S$ . The above relation can be written as

$$\pi^{(n)} = \pi^{(0)} P^n. \quad (*)$$

When  $S = \{x_1, \dots, x_r\}$  is finite,  $\pi^{(n)} = (\pi^{(n)}(x_1), \dots, \pi^{(n)}(x_r))$ , then  $(*)$  can be presented as

$$(\pi^{(n)}(x_1), \dots, \pi^{(n)}(x_r)) = (\pi^{(0)}(x_1), \dots, \pi^{(0)}(x_r)) \begin{bmatrix} P^n(x_1, x_1) & \cdots & P^n(x_1, x_r) \\ P^n(x_2, x_1) & \cdots & P^n(x_2, x_r) \\ \vdots & \ddots & \vdots \\ P^n(x_r, x_1) & \cdots & P^n(x_r, x_r) \end{bmatrix}$$

**Example 4.2.** Consider a Markov chain  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  on the state space  $S = \{0, 1\}$  with initial distribution  $\pi^{(0)} = (\pi^{(0)}(0), \pi^{(0)}(1))$  and transition probability matrix  $P = \begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix}$ . Let us compute  $P^n$ , the n-step transition probability matrix. To this end, we compute the eigenvalue of  $P$ .

$$\begin{aligned}\text{Consider } \det(\lambda I - P) &= \det \begin{bmatrix} \lambda - 1 + p & -p \\ -q & \lambda - 1 + q \end{bmatrix} = (\lambda - 1 + p)(\lambda - 1 + q) - pq \\ &= \lambda^2 - (2 - p - q)\lambda + (1 - p)(1 - q) - pq = \lambda^2 - (2 - p - q)\lambda + 1 - p - q.\end{aligned}$$

Let  $\det(\lambda I - P) = 0$ , we obtain

$$\lambda_1 = 1 \quad \lambda_2 = 1 - p - q$$

It is easy to see that

$$\begin{aligned} P \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix} &= \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1-p-q \end{bmatrix} \\ \Rightarrow P &= \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1-p-q \end{bmatrix} \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix}^{-1} \\ P^n &= \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (1-p-q)^n \end{bmatrix} \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix}^{-1} \\ &= \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (1-p-q)^n \end{bmatrix} \begin{bmatrix} \frac{q}{p+q} & \frac{p}{p+q} \\ \frac{-1}{p+q} & \frac{1}{p+q} \end{bmatrix} \end{aligned}$$

If  $0 < p + q < 2$ . As  $n \rightarrow \infty$ ,

$$P^n \rightarrow \begin{bmatrix} 1 & -p \\ 1 & q \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{q}{p+q} & \frac{p}{p+q} \\ \frac{-1}{p+q} & \frac{1}{p+q} \end{bmatrix} = \begin{bmatrix} \frac{q}{p+q} & \frac{p}{p+q} \\ \frac{q}{p+q} & \frac{p}{p+q} \end{bmatrix}.$$

Given  $\pi^{(0)} = (\pi^{(0)}(0), \pi^{(0)}(1))$ , we see that

$$\pi^{(n)} = \pi^{(0)} P^n$$

$$\text{As } n \rightarrow \infty, \pi^{(n)} = \pi^{(0)} \begin{bmatrix} \frac{q}{p+q} & \frac{p}{p+q} \\ \frac{q}{p+q} & \frac{p}{p+q} \end{bmatrix} = \begin{bmatrix} \frac{q}{p+q} & \frac{p}{p+q} \end{bmatrix}.$$

If a Markov chain starts from  $x \in S$ , i.e.  $\theta^{(0)} = x$ , We denote

$$P_x(\cdot) = P(\cdot | \theta^{(0)} = x).$$

For instance, let  $B \subset S$ , we have

$$P_x(\theta^{(n)} \in B) = P(\theta^{(n)} \in B \mid \theta^{(0)} = x).$$

**Definition (First hitting time)** Let  $A \subset S$ , define

$$T_A = \min \{n \geq 1 : \theta^{(n)} \in A\}$$

$T_A$  is the moment that the Markov chain  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  first visits the set  $A$ . It is called the first hitting time of  $A$ .  $T_A$  is a random variable. If the Markov chain can never hit  $A$ , we define  $T_A = \infty$ .

**Example 4.2 (continued):** Let  $\theta^{(0)} = 0$ , we now compute  $P_0(T_0 = n)$ .

Observe

$$\begin{aligned} P_0(T_0 = n) &= P_0(\theta^{(n)} = 0, \theta^{(n-1)} = 1, \theta^{(n-2)} = 1, \dots, \theta^{(1)} = 1) \\ &= P(0, 1) \underbrace{p(1, 1) \cdots p(1, 1)}_{n-2} p(1, 0) \\ &= p(1 - q)^{n-2} \cdot q \quad P_0(T_0 = 1) = P(0, 0) = 1 - p. \end{aligned}$$

$$\begin{aligned} P_0(T_0 < \infty) &= \sum_{n=1}^{\infty} P(T_0 = n) \\ &= (1 - p) + \sum_{n=2}^{\infty} p(1 - q)^{n-2} q = 1 - p + pq \frac{1}{q} = 1 \end{aligned}$$

This means that the Markov chain starting 0 and will come back to 0 in finite time with probability 1.

## 4.2. Decomposition of state space

Let  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  be a M.C. on the state space  $S$ . Given  $y \in S$ , recall the hitting time  $T_y$  of  $y$  is defined by

$$T_y = \inf \{n \geq 1 : \theta^{(n)} = y\}$$

Let us define two quantities:

- (1)  $\rho_{xy} = P_x(T_y < \infty)$ , it is the probability that the Markov chain starting from  $x$  hits the state  $y$  in finite time.
- (2) The number of the visits of the M.C. to  $y$ , ie.

$$N(y) = \#\{n \geq 1 : \theta^{(n)} = y\} = \sum_{k=1}^{\infty} 1(\theta^{(k)} = y)$$

There are two useful relations:

$$\begin{aligned}\mathbb{E}[T_y | \theta^{(0)} = x] &= \sum_{n=1}^{\infty} P_x(T_y > n), \\ \mathbb{E}[N(y) | \theta^{(0)} = x] &= \sum_{n=1}^{\infty} P^n(x, y).\end{aligned}$$

$$\begin{aligned}\text{Let us show the second relation. } \mathbb{E}[N(y) | \theta^{(0)} = x] &= \mathbb{E}\left[\sum_{n=1}^{\infty} 1(\theta^{(n)} = y) | \theta^{(0)} = x\right] \\ &= \sum_{n=1}^{\infty} \mathbb{E}[1(\theta^{(n)} = y) | \theta^{(n)} = x] = \sum_{n=1}^{\infty} P(\theta^{(n)} = y | \theta^{(0)} = x) = \sum_{n=1}^{\infty} P^n(x, y).\end{aligned}$$

A state  $y \in S$  is said to be **recurrent** if

$$\rho_{yy} = P_y(T_y < \infty) = 1.$$

This means that the M.C. starting from  $y$  will return to  $y$  in finite time with probability 1.

A state  $y \in S$  is said to be **transient** if

$$\rho_{yy} = P_y(T_y < \infty) < 1.$$

This means that  $P_y(T_y = \infty) > 0$ , ie. the Markov chain starting from  $y$  will never return to  $y$  in finite time with positive probability.

A recurrent state  $y$  can be further classified as

- positive recurrent  $\mathbb{E} [T_y | \theta^{(0)} = y] < \infty$  &  $\rho_{yy} = 1$ .
- null recurrent:  $\mathbb{E} [T_y | \theta^{(0)} = y] = \infty$  &  $\rho_{yy} = 1$ .

For a recurrent state  $y \in S$ , the Markov chain will revisit it infinitely many times, more precisely,  $P_y(N(y) = \infty) = 1$ .

For two states  $x, y \in S$ , if  $\rho_{xy} = P_x(T_y < \infty) > 0$ , we denote  $x \rightarrow y$  (' $x$  communicates to  $y$ ).  $C \subset S$  said to be **irreducible** if  $x \rightarrow y$  for every pair  $x, y \in C$ . A M.C. is irreducible if  $S$  is irreducible. The set  $C \subseteq S$  is said to be closed if  $\rho_{xy} = 0$  for all  $x \in C$  and  $y \notin C$ .

**An important claim:**  $\rho_{xy} > 0 \Leftrightarrow$  There exists some  $n \in \mathbb{N}$  s.t.  $P^n(x, y) > 0$ .

*Proof.* "  $\Rightarrow$  since  $\rho_{xy} = P_x(T_y < \infty) = \sum_{n=1}^{\infty} P_x(T_y = n) > 0$  There must be some  $n_0 \in \mathbb{N}$  so that  $P_x(T_y = n_0) > 0$ , i.e.

$$P_x(\theta^{(1)} \neq y, \dots, \theta^{(n_0-1)} \neq y, \theta^{(n_0)} = y) > 0,$$

this immediately implies  $P_x(\theta^{(n_0)} = y) > 0$ , i.e.  $P^{n_0}(x, y) > 0$ .

"  $\Leftarrow$  Observe  $\{\theta^{(n)} = y\} \subset \{T_y < \infty\}$  for some  $n \in \mathbb{N}$ .

$$P_x(T_y < \infty) > P_x(\theta^{(n)} = y) = P^n(x, y) > 0.$$

□

**Example 4.3 (a)**  $S = \{0, 1, 2\}$ , the (one-step) transition probability matrix

$$P = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & \frac{1}{3} & \frac{2}{3} \end{bmatrix}$$

The Markov chain is irreducible. To this end, from the form of  $P$ , we only need to show there exist some  $n \in \mathbb{N}, m \in \mathbb{N}$  so that  $P^n(0, 2) > 0$  and  $P^m(2, 0) > 0$ .  
Obviously

$$\begin{aligned} P^2(0, 2) &= P(\theta^{(2)} = 2 | \theta^{(0)} = 0) \\ &\geq P(\theta^{(2)} = 2, \theta^{(1)} = 1 | \theta^{(0)} = 0) \\ &= \frac{P(\theta^{(2)} = 2 | \theta^{(1)} = 1, \theta^{(0)} = 0) P(\theta^{(1)} = 1 | \theta^{(0)} = 0)}{P(\theta^{(0)} = 0)} \\ &= P(\theta^{(2)} = 2 | \theta^{(1)} = 1) P(\theta^{(1)} = 1 | \theta^{(0)} = 0) \\ &= P(1, 2)P(0, 1) = \frac{1}{4} \times \frac{1}{2} = \frac{1}{8} \end{aligned}$$

Similarly,  $P^2(2, 0) \geq \frac{1}{3} \cdot \frac{1}{4} = \frac{1}{12}$ .

**Assignment:** Show  $P^2(2, 0) \geq \frac{1}{12}$ .



### 4.3. Stationary distributions

A fundamental problem for M.C. in the context of simulation is the study of its stationary distributions (or invariant measures). Let  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  be a time homogeneous M.C. on the state space  $S$ , with transition probability matrix  $\{P(x, y)\}_{x \in S, y \in S}$ . A distribution  $\pi = \{\pi(x)\}_{x \in S}$  is stationary if  $\pi(y) = \sum_{x \in S} \pi(x)P(x, y), \forall y \in S$ . This relation can be represented as

$$\pi = \pi P.$$

If  $S$  has finite elements, e.g.  $S = \{x_1, \dots, x_m\}$ , then

$$[\pi(x_1) \dots \pi(x_m)] = [\pi(x_1) \dots \pi(x_m)] \begin{bmatrix} P(x_1, x_1) & \dots & P(x_1, x_m) \\ P(x_2, x_1) & \dots & P(x_2, x_m) \\ \vdots & \ddots & \vdots \\ P(x_m, x_1) & \dots & P(x_m, x_m) \end{bmatrix}.$$

So we can start from  $\theta^{(0)} = x$

and take  $\theta^{(n)}$  from  $P(\theta^{(n)}, \cdot)$

We could like design a M.C. which converges to  $\pi$  very fast, then the  $N$  can be

Small and thus we save computation resources.

Let  $\theta_n = \frac{1}{n}$  and  $b_n = e^{-n}$ .

$\theta_n \rightarrow 0$ ,  $b_n \rightarrow 0$

$$\theta_n = 0.005 \Rightarrow n = \frac{1}{0.005} = 200$$

$$b_n = 0.005 \Rightarrow n = \log \frac{1}{0.005} = 3 \log 5 \approx 6$$

If  $S = \{x_1, \dots, x_n\}$ .

$$[T(x_1) \dots T(x_n)] = [T(x_1) \dots T(x_n)] \begin{bmatrix} P(x_1, x_1) & P(x_1, x_2) & \dots & P(x_1, x_n) \\ P(x_2, x_1) & P(x_2, x_2) & \dots & P(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ P(x_n, x_1) & P(x_n, x_2) & \dots & P(x_n, x_n) \end{bmatrix}$$

A theorem tells us: Under certain condition,

If  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$  has a stationary dist.  $\pi$ , then as  $n \rightarrow \infty$

$$(*) \quad P(x, y) \rightarrow \pi(y) \text{ for all } x, y.$$

$$P^n(x, y) = P(\theta^{(n)}=y | \theta^{(n)}=x)$$

Given a target dist.  $\pi$  (i.e. we would like to sample random quantities from it), as long as we can design a M.C.  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$  whose stationary dist. is  $\pi$ , then with  $P(x, y)$

by (\*), we can get a M.C. in the following way.

$$\theta^{(n)} = x.$$

$$\theta^{(n)} \sim P(\theta^{(n)}).$$

$$\theta^{(n)} \sim P(\theta^{(n)}).$$

$$\vdots$$

$$\theta^{(n)} \sim P(\theta^{(n)}).$$

Theoretically,  $\theta^{(n)}$  has dist.  $\pi$ . In practice, we stop the above procedure at some large  $N$ .

i.e. We take  $\theta^{(n)}$  as a random quantity from  $\pi$ .

$$\vdots$$

$$\theta^{(n)} \sim P(\theta^{(n)}).$$

One of most speedy M.C. is the reversible M.C. Given a  $\pi$ , we would like

to design such a M.C. What is the reversible M.C.? Given a M.C.

part

below

$\{\theta^{(n)}\}_{n \in \mathbb{N}}$ , we consider  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(k)}, \dots, \theta^{(n)}$ .

the reverse

$\{\theta^{(n)}, \theta^{(n-1)}, \dots, \theta^{(1)}\}$  is also Markov. It is

transition probability

$$P(\theta^{(k)}=y | \theta^{(k+1)}=x) = \frac{P(y, x) \pi^{(k)}(y)}{\pi^{(k+1)}(x)}$$

where  $\pi^{(k)}$  is the dist. of  $\theta^{(k)}$  and  $\pi^{(k+1)}$  is the dist. of  $\theta^{(k+1)}$ .

$\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(k)}, \dots, \theta^{(n)}$

$\theta^{(n)}, \theta^{(n-1)}, \dots, \theta^{(1)}$

**Theorem 1.** Under some appropriate conditions, we have

$$P^n(x, y) \rightarrow \pi(y)$$

for all  $x, y \in S$ . Recall  $P^n(x, y)$  is the  $n$  step transition probability from  $x$  to  $y$ .

**Example 4.2 (continued):**  $S = \{0, 1\}$ . The transition probability is  $\begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix}$

$$\pi = \begin{bmatrix} \frac{q}{p+q} & \frac{p}{p+q} \end{bmatrix}$$

**Example 4.3 (A prototype of Gibbs sampler)** Let the state space be  $S = \{0, 1\}^2$ , i.e.  $S = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$ . For  $\theta \in S$ ,  $\theta = (\theta_1, \theta_2)$  with  $\theta_i = 0$  or 1 for  $i = 1, 2$ . For a prob. distribution  $\pi$  on  $S$ , we represented as

	0	1
0	$\pi_{00}$	$\pi_{01}$
1	$\pi_{10}$	$\pi_{11}$

$$\pi(i, j) = p(\theta_1 = i, \theta_2 = j), i \in \{0, 1\} \text{ and } j \in \{0, 1\}.$$

- For the first component  $\theta_1$ , given  $\theta_2 = j$ , the conditional probability  $\theta_1$  is

$$\pi_1(0 | j) = \frac{\pi_{0j}}{\pi_{+j}} \quad \pi_1(1 | j) = \frac{\pi_{1j}}{\pi_{+j}}$$

$$\text{where } \pi_{+j} = \pi_{0j} + \pi_{1j}.$$

- For the second component  $\theta_2$ , given  $\theta_1 = i$ , the conditional probability  $\theta_2$  is

$$\pi_2(0 | i) = \frac{\pi_{i0}}{\pi_{i+}} \quad \pi_2(1 | i) = \frac{\pi_{i1}}{\pi_{i+}}$$

$$\text{where } \pi_{i+} = \pi_{i0} + \pi_{i1}.$$

- Design a transition probability

$$P((i, j), (k, l)) = P(\theta^{(n)} = (k, l) \mid \theta^{(n-1)} = (i, j)) = \frac{\pi_{kl}}{\pi_{k+}} \cdot \frac{\pi_{kj}}{\pi_{+j}}$$

$\{P((i, j), (k, l))\}_{i,j,k,l \in \{0,1\}}$  is a  $4 \times 4$  matrix.

$$P = [P_1 \ P_2 \ P_3 \ P_4] = \begin{bmatrix} \frac{\pi_{00}\pi_{00}}{\pi_{0+}\pi_{+0}} & \times & \times & \times \\ \frac{\pi_{00}\pi_{01}}{\pi_{0+}\pi_{+1}} & \times & \times & \times \\ \frac{\pi_{00}\pi_{00}}{\pi_{0+}\pi_{+0}} & \times & \times & \times \\ \frac{\pi_{00}\pi_{01}}{\pi_{0+}\pi_{+1}} & \times & \times & \times \end{bmatrix}$$

we have  $\pi P = \pi$  with  $[\pi_{00} \ \pi_{01} \ \pi_{10} \ \pi_{11}]$ .

Let us here only verify that  $\pi P_1 = \pi_{00}$

$$\begin{bmatrix} \pi_{00} & \pi_{01} & \pi_{10} & \pi_{11} \end{bmatrix} \begin{bmatrix} \frac{\pi_{00}\pi_{00}}{\pi_{0+}\pi_{+0}} \\ \frac{\pi_{00}\pi_{01}}{\pi_{0+}\pi_{+1}} \\ \frac{\pi_{00}\pi_{00}}{\pi_{0+}\pi_{+0}} \\ \frac{\pi_{00}\pi_{01}}{\pi_{0+}\pi_{+1}} \end{bmatrix} = \frac{\pi_{00}^3 + \pi_{10}\pi_{00}^2}{\pi_{0+}\pi_{+0}} + \frac{\pi_{00}\pi_{01}^2 + \pi_{00}\pi_{01}\pi_{11}^2}{\pi_{0+}\pi_{+1}} \\ = \frac{\pi_{00}^2}{\pi_{0+}} + \frac{\pi_{00}\pi_{01}}{\pi_{0+}} = \pi_{00}$$

**Example 4.4.** Consider

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

Solving  $\pi P = \pi$ , we obtain  $\pi(0) = \pi(1) = \frac{1}{2}$ . On the other hand

$$P \cdot P = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad PPP = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \\ P^{2m} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad P^{2m+1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \text{for all } m \in \mathbb{N}.$$

$\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}, \theta^{(n+1)}, \dots, \theta^{(k)}$   
 ↓  
 past       $\theta^{(n)} = x$       future  
 post and future are mixed.

We can claim that  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}, \theta^{(n+1)}, \dots, \theta^{(k)}$  is also a M.C.  
 Let's show this claim rigorously.  
 $P(\theta^{(k)} = y | \theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k}) = P(\theta^{(k)} = y | \theta^{(k+1)} = x)$

$$\begin{aligned} ① &= \frac{P(\theta^{(k)} = y, \theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k})}{P(\theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k})} \\ &= \frac{(P(\theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k}) / P(\theta^{(k+1)} = x, \theta^{(k)} = y)) P(\theta^{(k+1)} = x, \theta^{(k)} = y)}{P(\theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k})} \\ &= \frac{P(\theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k} | \theta^{(k+1)} = x) P(\theta^{(k+1)} = x, \theta^{(k)} = y)}{P(\theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k} | \theta^{(k+1)} = x) P(\theta^{(k+1)} = x)} \\ &= \frac{P(\theta^{(k+1)} = x | \theta^{(k)} = y) P(\theta^{(k+1)} = y)}{P(\theta^{(k+1)} = x)} = \frac{P(y, x) \pi^{(k)}(y)}{\pi^{(k+1)}(x)} \end{aligned}$$

## 4.4. Reversible chains

Can I travel to the past? Can I conquer the time?

Let  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  be a time homogeneous M.C. with transition probability  $P(x, y)$  and stationary distribution  $\pi$ . Let us consider the M.C. from 0 to  $n$  in the reverse order,

i.e.  $\theta^{(n)}, \theta^{(n-1)}, \dots, \theta^{(0)}$ . (Diffusion model is a reverse M.C.) 

$$\begin{aligned} &P(\theta^{(k)} = y | \theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k}) = P(\theta^{(k)} = y | \theta^{(k+1)} = x) \\ &= \frac{P(\theta^{(k)} = y, \theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k})}{P(\theta^{(k+1)} = x, \theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k})} \\ &= \frac{P(\theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k} | \theta^{(k+1)} = x, \theta^{(k)} = y) P(\theta^{(k+1)} = x, \theta^{(k)} = y)}{P(\theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k} | \theta^{(k+1)} = x) P(\theta^{(k+1)} = x)} \\ &= \frac{P(\theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k} | \theta^{(k+1)} = x) P(\theta^{(k+1)} = x | \theta^{(k)} = y) P(\theta^{(k)} = y)}{P(\theta^{(k+2)} = x_2, \dots, \theta^{(n)} = x_{n-k} | \theta^{(k+1)} = x) / P(\theta^{(k+1)} = x)} \\ &= \frac{P(\theta^{(k+1)} = x | \theta^{(k)} = y) P(\theta^{(k)} = y)}{P(\theta^{(k+1)} = x)} = \frac{P(\theta^{(k+1)} = y | \theta^{(k+2)} = x) P(\theta^{(k+2)} = x)}{P(\theta^{(k+1)} = x)} \\ &= P(\theta^{(k)} = y | \theta^{(k+1)} = x) = \frac{P(y, x) \pi^{(k)}(y)}{\pi^{(k+1)}(x)}. \end{aligned}$$

Hence,  $\theta^{(n)}, \theta^{(n-1)}, \dots, \theta^{(1)}, \theta^{(0)}$  is also a M.C. If  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  is a time homogeneous M.C. with stationary distribution  $\pi$ , then the reverse M.C. is also time homogeneous with transition probability

for  $k = 0, 1, \dots, n-1$

we denote:  $P^*(x, y) = \frac{\pi(y) P(y, x)}{\pi(x)}$

Let's consider a spectral: reverse M.C. and forward M.C. have the same If  $P^*(x, y) = p(x, y)$  for all  $x \in S$  and  $y \in S$ , then we have transition probability.

case:

(\*) the M.C.  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  is reversible:  $\pi(x)p(x, y) = \pi(y)p(y, x)$  for all  $x, y \in S$

We call (\*) as the detailed balance condition. The corresponding M.C. is called time reversible M.C. Convergent very fast.

Given a posterior distribution  $\pi^*$ , we would design a M.C. whose (one-step) transition prob. satisfies (\*).

Example 4.5 Metropolis algorithm (1953. Metropolis) as old as my dear grandma

Given: We would like to sample random quantities from posterior dist.  $\pi^*$ .

From  $p^*(x, y) = p(x, y)$

we see that

$\pi^*(y) P(y, x) = \pi(x) P(x, y)$

Recall  $\pi^{(k)} = \pi$ , i.e.  $\pi^{(k)}(y) P(y, x) = \pi(x) P(x, y)$

i.e.  $P(\theta^{(k)} = y | \theta^{(k)} = x) = P(\theta^{(k)} = x | \theta^{(k)} = y) \pi(x) / \pi(y)$

This M.C. is in equilibrium (i.e. find  $P(x, y)$ )

We would like to design M.C. which has a stationary distribution

38

If  $\theta^{(k)} = \pi$ , where  $\pi$  is the stationary dist., then we know

$\pi^{(k)} = \pi^{(k)} P = \pi P = \pi$ .

$\pi^{(k)} = \pi^{(k)} P^2 = \pi P^2 = (\pi P) P = \pi P = \pi$ .

$\pi^{(k)} = \pi$

$\pi^{(k)} = \pi$ .

Hence  $P(\theta^{(k)} = y | \theta^{(k)} = x) = \frac{P(y, x) \pi(x)}{\pi(y)}$  does not

depend on  $k$ . This means that, if the forward M.C. starts from  $\pi$ , i.e.  $\theta^{(0)} = \pi$ , then, the reverse M.C. is time homogeneous.

Example 4.5: Metropolis algorithm (1953, Metropolis)

If  $\theta^n$   
Given a posterior distr.  $\{\pi(x) : x \in S\}$  (so we have  $\pi(x) > 0$ )  
and  $\sum_{x \in S} \pi(x) = 1$ . We would like to design a reversible M.C.  
satisfying (A\*). The key is to find a  $P(x, y)$  satisfying (\*\*) .

Accept-reject: Sample a random quantity from the distr.  $f$ .  
but it is difficult to do it.

- 1 Sample  $y$  from the distr.  $g$  (easy to be sampled)
- 2 Sample  $u \sim U([0, 1])$  if  $\frac{f(y)}{f(\theta^n)} > u$  accept  $y$   
if  $\frac{f(y)}{f(\theta^n)} \leq u$  reject  $y$  return to 1.

### posterior distribution

Consider a given distribution  $\{p_x : x \in S\}$  with  $p_x \geq 0$  for all  $x \in S$  and  $\sum_{x \in S} p_x = 1$ . We would like to design a M.C. reversible satisfies. The key is to find a  $P(x, y)$  satisfying (\*\*):

1. Let  $\{\theta^{(n)}\}$  be an irreducible M.C. with transition probability matrix  $\{Q(x, y)\}$ :  
 $x \in S, y \in S\}$  such that  $Q(x, y) = Q(y, x)$  for all  $x \in S, y \in S$ . Now let us modify the M.C.  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  by injecting a censorship:

(1) Given  $\theta^{(n)} = x$ , the M.C. transits from  $x$  to  $y$  according to  $Q(x, y)$ , i.e.  $\theta^{(n+1)} = y$  with probability  $Q(x, y)$ .

(2) The proposed value for  $\theta^{(n+1)}$  is accepted with probability  $\min\{1, p_y/p_x\}$ , and reject otherwise, leaving the chain staging at the state  $x$ .

Combining the two steps (1) & (2), we can figure out a transition probability  $P(x, y)$ .

- for  $y \neq x$ ,

$$\begin{aligned} P(x, y) &= P(\theta^{(n+1)} = y, \text{ accept } y \mid \theta^{(n)} = x) \\ &= P(\theta^{(n+1)} = y \mid \theta^{(n)} = x) P(\text{accept } y) \\ &= Q(x, y) \min \left\{ 1, \frac{p_y}{p_x} \right\}. \end{aligned}$$

the censorship is independent of  $\theta^{(n)} = x$ .

- for  $y = x$

$$\begin{aligned} P(x, x) &= P(\theta^{(n+1)} = x, \text{ accept } x \mid \theta^{(n)} = x) \\ &\quad + P(\theta^{(n+1)} \neq x, \text{ reject } y \mid \theta^{(n)} = x) \\ &= P(\theta^{(n+1)} = x \mid \theta^{(n)} = x) P(\text{accept } x) + \sum_{y \neq x} P(\theta^{(n+1)} = y, \text{ reject } y \mid \theta^{(n)} = x) \\ &= Q(x, x) + \sum_{y \neq x} P(\theta^{(n+1)} = y \mid \theta^{(n)} = x) P(\text{reject } y) \\ &= Q(x, x) + \sum_{y \neq x} Q(x, y) \left( 1 - \min \left\{ 1, \frac{p_y}{p_x} \right\} \right) \end{aligned}$$

We obtain a sequence  $\theta^{(0)}, \theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}, \dots$

We would like to show that  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  is a reversible M.C. with transition prob.  $P(x, y)$  s.t.

$$\pi(x) P(x, y) = \pi(y) P(y, x) \quad \forall x, y \in S$$

Let us compute one step transition prob. of  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  and verify (\*\*).

$$P(x, y)$$

We design a prob. transition prob.  $\theta(x, y)$

$$\text{s.t. } \theta(x, y) = Q(y, x)$$

this means a M.C.  $P(\tilde{\theta}^{(n+1)} = y \mid \tilde{\theta}^{(n)} = x) = Q(y, x)$

(1) Given  $\tilde{\theta}^{(n)} = x$ , we use  $\theta(x, y)$  to sample  $\tilde{\theta}^{(n+1)}$ .

(2) We will accept-reject, we accept  
sample a  $u \sim U[0, 1]$ .  
this  $y$  with prob.  $\min\{1, \frac{\pi(y)}{\pi(x)}\}$  accept, i.e.  $\tilde{\theta}^{(n)} = \tilde{\theta}^{(n+1)}$ .

and reject otherwise, this means that  $\tilde{\theta}^{(n)} = \tilde{\theta}^{(n+1)}$ .

(3) We use  $\theta(\tilde{\theta}^{(n)})$  to sample  $\tilde{\theta}^{(n+2)}$ ,  
accept  $\tilde{\theta}^{(n+1)}$  with prob.  $\min\{1, \frac{\pi(\tilde{\theta}^{(n+1)})}{\pi(\tilde{\theta}^{(n)})}\}$ , i.e.  $\tilde{\theta}^{(n)} = \tilde{\theta}^{(n+1)}$ .

reject with prob.  $\min\{1, \frac{\pi(\tilde{\theta}^{(n)})}{\pi(\tilde{\theta}^{(n+1)})}\}$ , i.e.  $\tilde{\theta}^{(n+1)} = \tilde{\theta}^{(n)}$ .

Couldn't directly design the M.C.  $P(x, y)$

(1) Initialize  $\theta^{(0)} = x \in S$ .

(2) Sample  $\tilde{\theta}^{(1)}$  from  $\theta(x, \cdot)$  ②

We shall use a censorship (i.e. accept or reject) to get  $\theta^{(1)}$ .

accept  $\tilde{\theta}^{(1)}$  with a prob.  $\min\{1, \frac{\pi(\tilde{\theta}^{(1)})}{\pi(x)}\}$ , i.e.  $\theta^{(1)} = \tilde{\theta}^{(1)}$ .

reject  $\tilde{\theta}^{(1)}$  with  $1 - \min\{1, \frac{\pi(\tilde{\theta}^{(1)})}{\pi(x)}\}$ , i.e.  $\theta^{(1)} = x$ .

Sample a  $u \sim U[0, 1]$ , if  $u \leq \min\{1, \frac{\pi(\tilde{\theta}^{(1)})}{\pi(x)}\}$ , accept

(3) Sample  $\tilde{\theta}^{(2)}$  from  $\theta(\theta^{(1)}, \cdot)$  if  $u > \min\{1, \frac{\pi(\tilde{\theta}^{(1)})}{\pi(x)}\}$ , reject

We shall use a censorship to get  $\theta^{(2)}$ .

accept  $\tilde{\theta}^{(2)}$  with a prob.  $\min\{1, \frac{\pi(\tilde{\theta}^{(2)})}{\pi(\theta^{(1)})}\}$ , i.e.  $\theta^{(2)} = \tilde{\theta}^{(2)}$ .

reject  $\tilde{\theta}^{(2)}$  with a prob.  $1 - \min\{1, \frac{\pi(\tilde{\theta}^{(2)})}{\pi(\theta^{(1)})}\}$ , i.e.  $\theta^{(2)} = \theta^{(1)}$ .

$$P(x,y) = P(\theta^n = y | \theta^0 = x).$$

If  $y \neq x$ :

$$\begin{aligned} P(x,y) &= P(\hat{\theta}^n = y, \text{accept } y | \theta^0 = x) \\ &= \theta(x,y) P(\text{accept } y) \\ &= \theta(x,y) \min\left\{1, \frac{\pi(x)}{\pi(y)}\right\} \end{aligned}$$

If  $y = x$ :

$$\begin{aligned} P(x,x) &= P(\hat{\theta}^n = x, \text{accept } x | \theta^0 = x) + P(\hat{\theta}^n = x, \text{reject } x | \theta^0 = x) \\ &+ \sum_{y \neq x} P(\hat{\theta}^n = y, \text{reject } y | \theta^0 = x) = \theta(x,x) \min\left\{1, \frac{\pi(x)}{\pi(y)}\right\} + \theta(x,x) (1 - \min\left\{1, \frac{\pi(x)}{\pi(y)}\right\}) + \sum_{y \neq x} \theta(x,y) (1 - \min\left\{1, \frac{\pi(y)}{\pi(x)}\right\}) = \theta(x,x) + \sum_{y \neq x} \theta(x,y) (1 - \min\left\{1, \frac{\pi(y)}{\pi(x)}\right\}) \end{aligned}$$

- (1) Initialize  $\hat{\theta}^0 = x \in S$ .
- (2) Sample  $\hat{\theta}^{(n)}$  from  $\theta(x, \cdot)$ .  
We shall use a censorship (i.e. accept or reject) to get  $\hat{\theta}^{(n)}$   
accept  $\hat{\theta}^{(n)}$  with a prob. min  $\left\{1, \frac{\pi(\hat{\theta}^{(n)})}{\pi(x)}\right\}$ , i.e.  $\hat{\theta}^{(n)} = \hat{\theta}^n$   
reject  $\hat{\theta}^{(n)}$  with  $1 - \min\left\{1, \frac{\pi(\hat{\theta}^{(n)})}{\pi(x)}\right\}$ , i.e.  $\hat{\theta}^{(n)} = x$ .

Let's show (\*):

If  $x = y$ ,  $\pi(x) p(x,y) = \pi(y) p(y,x)$  holds trivially.

If  $x \neq y$ . Suppose  $\pi(x) \leq \pi(y)$  without loss of generality.

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

Hence  $\pi(x) p(x,y) = \pi(y) p(y,x)$

If (\*) holds then

$\pi$  is a stationary distr. of M.C.  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$   
i.e.  $\pi(y) = \sum_{x \in S} \pi(x) p(x,y)$  for all  $x, y$

We can define a new Markov chain  $\{\hat{\theta}^{(n)}\}_{n \in \mathbb{N}_0}$  who transition prob. is  $P(x,y)$ . We claim that  $\{\theta^{(n)}\}_{n \in \mathbb{N}_0}$  is an reversible M.C with stationary distribution  $\{p_x\}$ . To show this, we only need to prove that

$$p_x P(x,y) = p_y P(y,x) \quad \forall x \in S, y \in S.$$

Indeed, for  $x \neq y$ , suppose  $p_y \geq p_x$ , then

$$P(x,y) = Q(x,y)$$

and thus

$$\begin{aligned} p_x p(x,y) &= p_x Q(x,y) = Q(y,x) \min\left\{1, \frac{p_x}{p_y}\right\} p_y \\ &= p(y,x) p_y. \end{aligned}$$

As  $p_x \geq p_y$ , by a similar argument, we have

$$p_x P(x,y) = p_y P(y,x).$$

A remark about Markov Chain on continuous state space  $S$ , e.g.  $S = \mathbb{R}$  or  $S = \mathbb{R}^d$ .

We shall replace the prob. mass function  $\{\pi(x) : x \in S\}$  with  $\sum_{x \in S} \pi(x) = 1$  by the prob. density function  $\pi(x)$  with  $\int_S \pi(x) dx = 1$ . The transition prob. matrix  $\{P(x,y)\}_{x \in S, y \in S}$  with  $\sum_{y \in S} P(x,y) = 1$  with the transition prob. kernel  $p(x,y)$  with  $\int_S p(x,y) dy = 1$ . The stationary measure  $\pi$  has the property  $\pi(y) = \int \pi(x)p(x,y) dx$ .

The detailed balance condition is

$$\pi(x)p(x,y) = \pi(y)p(y,x) \quad \forall x, y \in S.$$

## Chapter 5 Gibbs Sampling

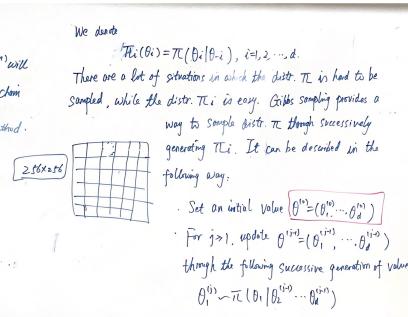


Gibbs sampling is a powerful tool for sampling posterior distribution and image processing. It was first proposed by German & German (1984) in image processing,

## Chapter 5: Gibbs Sampling

We obtain  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$ , it is a M.C. As  $n \rightarrow \infty$ , the dist. of  $\theta^{(n)}$  will tends to  $\pi$ . Under a certain (quite general) condition, the Markov Chain  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$  is reversible. So it is a very powerful sampling method.

Exa

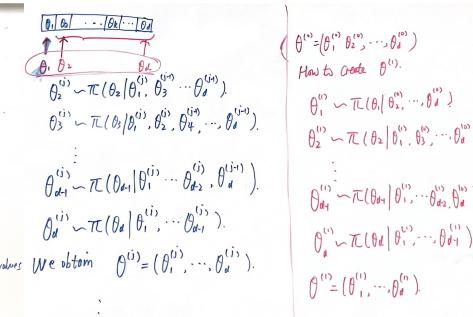


We denote

$\pi(\theta_i) = \pi(\theta_i | \theta_{-i}), i=1, 2, \dots, d$ .  
There are a lot of situations in which the dist.  $\pi$  is hard to be sampled, while the dist.  $\pi_i$  is easy. Gibbs sampling provides a way to sample dist.  $\pi$  through successively generating  $\pi_i$ . It can be described in the following way:

- Set an initial value  $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_d^{(0)})$
- For  $j \geq 1$ , update  $\theta^{(j+1)} = (\theta_1^{(j+1)}, \dots, \theta_d^{(j+1)})$  through the following successive generation of values

We obtain  $\theta^{(j)} = (\theta_1^{(j)}, \dots, \theta_d^{(j)})$ .



and later further developed by Gelfand & Smith (1990), Robert (2001).

## 5.1 Definition and properties

Assume that the distribution of interest is  $\pi(\theta)$  with  $\theta = (\theta_1, \dots, \theta_d)'$ , each  $\theta_i$  can be a scalar, a vector or a matrix, e.g.  $\theta_i$  is the color of the  $i$ -th pixel of an image. We denote  $\pi_i(\theta_i) = \pi(\theta_i | \theta_{-i})$  where  $\theta_{-i} = (\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_d)$  for  $i = 1, 2, \dots, d$ .

There are a lot of situations in which the distribution  $\pi$  is hard to be sampled, while the distribution  $\pi_i$  is easy. Gibbs sampling provides a way to sample distribution  $\pi$  through successive generating  $\pi_i$ . It can be described in the following way:

- Set an initial value  $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_d^{(0)})$
- For  $j \geq 1$ , update  $\theta^{(j+1)} = (\theta_1^{(j+1)}, \dots, \theta_d^{(j+1)})$  through the following successive generation of Values:

$$\begin{aligned}\theta_1^{(j)} &\sim \pi(\theta_1 | \theta_2^{(j-1)}, \dots, \theta_d^{(j-1)}) \\ \theta_2^{(j)} &\sim \pi(\theta_2 | \theta_1^{(j)}, \theta_3^{(j-1)}, \dots, \theta_d^{(j-1)}) \\ \theta_3^{(j)} &\sim \pi(\theta_3 | \theta_1^{(j)}, \theta_2^{(j)}, \theta_4^{(j-1)}, \dots, \theta_d^{(j-1)}) \\ &\dots \\ \theta_d^{(j)} &\sim \pi(\theta_d | \theta_1^{(j)}, \dots, \theta_{d-1}^{(j)})\end{aligned}$$

Denote  $\theta^{(j)} = (\theta_1^{(j)}, \dots, \theta_d^{(j)})$  for  $j \geq 0$ , it is a M.C. One needs to show the convergence of this M.C.  $\{\theta^{(i)}\}_{i \in \mathbb{N}_0}$ .

Chapter 5: Gibbs sampling

$$\text{Hence } P((i,j), (k,l)) = P(\theta_i^{(i)}, \theta_j^{(i)}) P(\theta_k^{(k)}, \theta_l^{(k)}) = P(\theta_i^{(i)}, \theta_j^{(i)}) P(\theta_k^{(k)}, \theta_l^{(k)})$$

$$= \frac{T_{ik}}{T_{il}} \cdot \frac{T_{jl}}{T_{kl}}$$

### Chapter 5: Gibbs sampling

We obtain  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$ , it is a M.C. As  $n \rightarrow \infty$ , the dist. of  $\theta^{(n)}$  will tends to  $\pi$ . Under a certain (quite general) condition, the Markov Chain  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$  is reversible. So it is a very powerful sampling method.

Example 5.1 (An example in Chapter 4 revisited, i.e. Example 4.3.)

Let the state space be  $S = \{0, 1\}^2 = \{0, 1\} \times \{0, 1\}$  as  $\omega$ .  
 $S = \{(0,0), (0,1), (1,0), (1,1)\}$  For  $\theta \in S$ ,  $\theta = (\theta_1, \theta_2)$  with  $\theta_1 = 0, 0$  or  $1$ , and  $\theta_2 = 0$  or  $1$ . For a prob. distr.  $\pi$  on  $S$ , we denote

$$\pi_{11} = P(\theta_1=1, \theta_2=1), \quad \pi_{10} = P(\theta_1=1, \theta_2=0),$$

$$\pi_{01} = P(\theta_1=0, \theta_2=1), \quad \pi_{00} = P(\theta_1=0, \theta_2=0).$$

To use Gibbs sampling, we write  $\pi(\cdot | \cdot)$  as  $\pi_{ij}(i|j)$ .

$$\text{Initialization: } \theta^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}) \text{ (e.g. } \theta^{(0)} = (0, 0))$$

Sample the first coordinate  $\theta_1^{(i)}$ , sample from the

Condition prob.  $\pi_{11}(\theta_1 | \theta^{(i)})$ , which is

$$\pi_{11}(\theta_1 | \theta^{(i)}) = \frac{\pi_{11}}{\pi_{11} + \pi_{10}}$$

$$\text{where } \pi_{11} = \pi_{11} + \pi_{01}.$$

Sample the second coordinate  $\theta_2^{(i)}$  from  $\pi_{11}(\theta_2 | \theta_1^{(i)})$

41

which is

$$\pi_{11}(\theta_2 | \theta_1^{(i)}) = \frac{\pi_{11}}{\pi_{11} + \pi_{10}}$$

where  $\pi_{11} = \pi_{11} + \pi_{01}$ , and  $i \in \mathbb{N}_0$ .

Hence  $\theta^{(i)} = (\theta_1^{(i)}, \theta_2^{(i)})$

Continue this procedure and obtain a M.C.  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$ . As  $n \rightarrow \infty$ , ideally

$\theta^{(n)}$  is a random quantity from  $\pi$ .

Let us figure out the transition prob. of this M.C.  $\{\theta^{(n)}\}_{n \in \mathbb{N}}$ .

$$P((i,j), (k,l)) = P(\theta_1^{(i)}=k, \theta_2^{(i)}=l | \theta_1^{(j)}=j, \theta_2^{(j)}=l)$$

from the procedure of Gibbs sampling. Fix  $\theta_2^{(i)}=l$ ,  $\theta_1^{(i)}=j$

$$\theta_1^{(i)}=j, \theta_2^{(i)}=l | \theta_1^{(i)}=j, \theta_2^{(i)}=l$$

$$\frac{\pi_{11}}{\pi_{11} + \pi_{10}}$$

## Example 5.1. (Example 4.3).

Hence  $P((i,j), (k,l)) = P((\theta_1^{(n)}, \theta_2^{(n)}) = (k,l) | (\theta_1^{(n)}, \theta_2^{(n)}) = (i,j))$

$$\{T_{ij}\} = \frac{T_{kj}}{T_{kj+}} \cdot \frac{T_{kl}}{T_{kl+}}$$

$$\text{Initialization } \theta^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}) = (0,0)$$

$$\text{Compute } T_{ij}(\theta_1^{(0)})$$

$$\text{Sample } \theta_1^{(1)} \text{ from the dist'n } T_{ij}(\theta_1^{(0)})$$

$$\text{Compute } T_{ij}(\theta_1^{(1)})$$

$$\text{Sample } \theta_2^{(1)} \text{ from } T_{ij}(\theta_1^{(1)})$$

## Example 5.1. (Example 4.3).

To us: Let the state space be  $S = \{0,1\}^2 = \{(0,0), (0,1), (1,0), (1,1)\}$

Initial Let  $\pi = (T_{00}, T_{01}, T_{10}, T_{11})$  be a distribution on  $S$ .

$$\frac{P(\theta_1=0, \theta_2=0)}{T_{00} + T_{01}}$$

$$\text{i.e. } T_{ij} = P(\theta_1=i, \theta_2=j) \quad i=0 \text{ or } 1, j=0 \text{ or } 1.$$

Condition We would like demonstrate how to use Gibbs sampling to sample random quantities from  $\pi$ .

Let's compute the transition prob. of Gibbs sampling i.e.

$$P(\theta^{(1)} = (k, l) | \theta^{(0)} = (i, j))$$

Where According the successive procedure, if  $\theta_2^{(0)} = j$ , sample  $\theta_1^{(1)}$  from  $T_{ij}(\theta_1 | \theta_2^{(0)=j}) = \begin{cases} \frac{\pi(\theta_1=k, \theta_2=j)}{\pi(\theta_1^{(0)=j})} & \text{if } \theta_1=k \\ \frac{\pi(\theta_1=l, \theta_2=j)}{\pi(\theta_1^{(0)=j})} & \text{if } \theta_1=l \end{cases}$ . Sample  $\theta_2^{(1)} = l$  given  $\theta_1^{(1)} = k$ , with  $\frac{T_{kj}}{T_{kj+}}$ .

Assume that we aim to sample a distribution  $\pi(\theta)$  with  $\theta = (\theta_1, \dots, \theta_d)$ . We denote  $T_i(\theta_i) = \pi(\theta_i | \theta_{-i})$  where  $\theta_{-i} = (\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_d)$ .

We know  $\theta_1^{(0)} = k$ , give  $\theta_1^{(0)} = k$ , sample  $\theta_1^{(1)}$  from  $T_{ij}(\theta_1 | \theta_2^{(0)=k})$ .

$$T_{ij}(\theta_1 | \theta_2^{(0)=k}) = \begin{cases} \frac{\pi(\theta_1=k, \theta_2=j)}{\pi(\theta_2^{(0)=j})} & \text{if } \theta_1=k \\ \frac{\pi(\theta_1=l, \theta_2=j)}{\pi(\theta_2^{(0)=j})} & \text{if } \theta_1=l \end{cases}$$

We know  $(\theta_1^{(0)}, \theta_2^{(0)}) = (k, l)$ , sample  $\theta_2^{(1)}$  with prob.  $\frac{T_{kj}}{T_{kj+}}$ .

## Example 5.1: (An example in chapter 4 revisited) Example 4.3 A proto-

**type of Gibbs sampler:** Let the state space be  $S = \{0, 1\}^2$ , ie.  $S = \{(0,0), (0,1), (1,0), (1,1)\}$

For  $\theta \in S$ ,  $\theta = (\theta_1, \theta_2)$  with  $\theta_i = 0$  or 1 for  $i = 1, 2$ . For a prob. distribution  $\pi$  on  $S$ ,

we represented as

	0	1
0	$\pi_{00}$	$\pi_{01}$
1	$\pi_{10}$	$\pi_{11}$

$$\pi(i, j) = p(\theta_1 = i, \theta_2 = j), \quad i \in \{0, 1\} \text{ and } j \in \{0, 1\}.$$

- For the first component  $\theta_1$ , given  $\theta_2 = j$ , the conditional probability  $\theta_1$  is

$$\pi_1(0 | j) = \frac{\pi_{0j}}{\pi_{+j}} \quad \pi_1(1 | j) = \frac{\pi_{1j}}{\pi_{+j}}$$

$$\text{where } \pi_{+j} = \pi_{0j} + \pi_{1j}.$$

- For the second component  $\theta_2$ , given  $\theta_1 = i$ , the conditional probability  $\theta_2$  is

$$\pi_2(0 | i) = \frac{\pi_{i0}}{\pi_{i+}} \quad \pi_2(1 | i) = \frac{\pi_{i1}}{\pi_{i+}}$$

$$\text{where } \pi_{i+} = \pi_{i0} + \pi_{i1}.$$

- Design a transition probability

$$P((i, j), (k, l)) = P(\theta^{(n)} = (k, l) | \theta^{(n-1)} = (i, j)) = \frac{\pi_{kl}}{\pi_{+j}} \cdot \frac{\pi_{kj}}{\pi_{+j}}$$

$$\theta^{(n-1)} = (i, j) : \text{First fix } \theta_2^{(n-1)} = j \rightarrow \theta_1^{(n)} = k \mid \theta_2^{(n-1)} = j$$

$$\text{Fix } \theta_1^{(n)} = k \rightarrow \theta_2^{(n)} = l \mid \theta_1^{(n)} = k$$

$$\theta_1^{(j)} \sim \pi(\theta_1 | \theta_2^{(j-1)}, \dots, \theta_d^{(j-1)})$$

$$\theta_2^{(j)} \sim \pi(\theta_2 | \theta_1^{(j)}, \theta_3^{(j-1)}, \dots, \theta_d^{(j-1)})$$

$$\theta_3^{(j)} \sim \pi(\theta_3 | \theta_1^{(j)}, \theta_2^{(j)}, \theta_4^{(j-1)}, \dots, \theta_d^{(j-1)})$$

$$\dots$$

$$\theta_d^{(j)} \sim \pi(\theta_d | \theta_1^{(j)}, \dots, \theta_{d-1}^{(j)})$$

**Example 5.2 (Carlin, Gelfand, Smith, 1992)** Let  $y_1, \dots, y_n$  be a sample from a poisson distribution for which there is a suspicion of a change point  $m$  among

So

$$P(\theta^{(n)} = (k, l) | \theta^{(n)} = (i, j)) = \frac{T_{kj}}{\frac{T_{0j} + T_{1j}}{T_{+j}}} \cdot \frac{T_{kl}}{\frac{T_{k0} + T_{k1}}{T_{+j}}} = \frac{T_{kj}}{T_{+j}} \cdot \frac{T_{kl}}{T_{+j}}$$

We can verify that

$$\pi_{ij} P((i, j), (k, l)) = \pi_{kl} P((k, l), (i, j))$$

the observation process where the means change,  $m = 1, \dots, n$ . Given the change point  $m$ , the observation distributions are  $y_i | \lambda \sim \text{Poi}(\lambda), i = 1, 2, \dots, m$  and  $y_i | \phi \sim \text{Poi}(\phi), i = m + 1, \dots, n$ . The model is completed with independent prior distributions  $\lambda \sim G(\alpha, \beta), \phi \sim G(\nu, \delta)$ , and the change point  $m$  uniformly distributed over  $\{1, 2, \dots, n\}$ , where  $\alpha, \beta, \nu$  and  $\delta$  are known constants. The posterior density is

$$\begin{aligned}\pi(\lambda, \phi, m) &\propto f(y_1, \dots, y_n | \lambda, \phi, m) p(\lambda, \phi, m) \\ &= \left( \prod_{i=1}^m \frac{\lambda^{y_i}}{y_i!} e^{-\lambda} \right) \left( \prod_{j=m+1}^n \frac{\phi^{y_j}}{y_j!} e^{-\phi} \right) \underbrace{\frac{1}{c_{\alpha, \beta}} \lambda^{\alpha-1} e^{-\beta \lambda}}_{\text{Constant}} \underbrace{\frac{1}{C_{\nu, \delta}} \phi^{\nu-1} e^{-\delta \phi}}_{\text{Constant}} \underbrace{\frac{1}{n}}_{\text{Constant}} \\ &\propto \left( \prod_{i=1}^m \lambda^{y_i} e^{-\lambda} \right) \left( \prod_{j=m+1}^n \phi^{y_j} e^{-\phi} \right) \lambda^{\alpha-1} e^{-\beta \lambda} \phi^{\nu-1} e^{-\delta \phi} \\ &= \lambda^{\alpha+S_m-1} e^{-(\beta+m)\lambda} \phi^{\nu+S_n-S_m-1} e^{-(\delta+n-m)\phi}\end{aligned}$$

where  $S_k = \underbrace{y_1 + \dots + y_k}$ . It is easy to see that given  $\phi, m, \lambda$  satisfies a Gamma distribution  $G(\alpha + S_m, \beta + m)$ , i.e.  $\pi(\lambda | \phi, m) = G(\alpha + S_m, \beta + m)$ . Similarly,  $\pi(\phi | \lambda, m) = G(\nu + S_n - S_m, \delta + n - m)$

$$\pi(m | \phi, \lambda) = \frac{\lambda^{\alpha+S_m-1} e^{-(\beta+m)\lambda} \phi^{\nu+S_n-S_m-1} e^{-(\delta+n-m)\phi}}{\sum_{l=1}^n \lambda^{\alpha+S_l-1} e^{-(\beta+l)\lambda} \phi^{\nu+S_n-S_l-1} e^{-(\delta+n-l)\phi}}$$

All these conditional distributions are easily sampled. So we can use Gibbs sampling to sample distribution  $\pi(\lambda, \phi, m)$ . The algorithm is as the following

1. Set the initial value  $(\lambda^{(0)}, \phi^{(0)}, m^{(0)})$
2. Obtain a new value  $(\lambda^{(j)}, \phi^{(j)}, m^{(j)})$  through the following successive generation of values

$$\lambda^{(j)} \sim G(\alpha + S_{m^{(j-1)}}, \beta + m^{(j-1)})$$

$$\phi^{(j)} \sim G(\nu + S_n - S_{m^{(j-1)}}, \delta + n - m^{(j-1)})$$

$$m^{(j)} \sim \pi(m | \lambda^{(j)}, \phi^{(j)})$$

Example 2 (Carlin, Gelfand, Smith, 1992).

Let  $y_1, \dots, y_n$  be a sample from a Poisson distribution for which there is a suspicion of a change point  $m$  among the observation process where the means change,  $m \in \{1, 2, \dots, n\}$ . Given the change point  $m$ , the observation distributions are  $y_i | \lambda \sim \text{Poi}(\lambda)$ ,  $i=1, \dots, m$  and  $y_i | \phi \sim \text{Poi}(\phi)$ ,  $i=m+1, \dots, n$ . The goal is to estimate  $\lambda, \phi, m$  by Bayesian inference.

Let us assume the prior distri. of  $\lambda$  and  $\phi$  are  $\lambda \sim G(d, \beta)$ ,  $\phi \sim G(y, \gamma)$ .

And  $m \sim \text{Uniform distribution in } \{1, 2, \dots, n\}$ .  
 Given  $\lambda, \phi, m$ , the likelihood of  $y_1, \dots, y_n$  is

$$f(y_1, \dots, y_n | \lambda, \phi, m) = \prod_{i=1}^m \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} \prod_{i=m+1}^n \frac{\phi^{y_i} e^{-\phi}}{y_i!}$$

The posterior distri. is

$$\pi(\lambda, \phi, m) \propto p(\lambda)p(\phi)p(m)f(y_1, \dots, y_n | \lambda, \phi, m)$$

*Derive*

$$\begin{aligned} & \propto \lambda^{d+1} e^{-\beta\lambda} \phi^{y_m} e^{-\delta\phi} \frac{m!}{\prod_{i=1}^m (y_i!)^{\lambda}} \frac{\prod_{i=1}^m (\frac{\lambda^{y_i}}{y_i!} e^{-\lambda})}{\prod_{i=m+1}^n (\frac{\phi^{y_i}}{y_i!} e^{-\phi})} \\ & \propto \lambda^{d+1} e^{-\beta\lambda} \phi^{y_m} e^{-\delta\phi} \lambda^{y_m} e^{-m\lambda} \phi^{y_{m+1}+...+y_n} e^{-(n-m)\phi} \\ & = \lambda^{d+1} e^{-\beta\lambda} \phi^{y_m} \underbrace{\lambda^{y_m} e^{-m\lambda} \phi^{y_{m+1}+...+y_n} e^{-(n-m)\phi}}_{S_n - S_m + y_{m+1} - (n-m)} \\ & = \lambda^{d+1} e^{-(\beta+m)\lambda} \phi^{S_n - S_m + y_{m+1} - (n-m) + \delta} \end{aligned}$$

We need to sample  $(\lambda, \phi, m)$

$$\begin{aligned} \pi(\lambda | \phi, m) & \sim G(d+S_m, \beta+m) \\ \pi(\phi | \lambda, m) & \sim G(y+y_{m+1} - (n-m)) \\ \pi(m | \phi, \lambda) & \sim \frac{\lambda^{d+1} e^{-(\beta+m)\lambda} \phi^{S_n - S_m + y_{m+1} - (n-m) + \delta}}{\sum_{k=1}^n \lambda^{d+1} e^{-(\beta+k)\lambda} \phi^{S_n - S_k + y_{k+1} - (n-k) + \delta}} \quad m=1, \dots, n \end{aligned}$$

So we obtain  $(\lambda^{(i)}, \phi^{(i)}, m^{(i)})$ .

3. Change counter  $j$  to  $j + 1$  and return to step 2 until convergence is reached.

## Chapter 6 Metropolis-Hastings algorithms

This algorithm stems from the papers by Metropolis et al. (1953) and Hastings (1970). The original paper by Metropolis et al. (1953) was published on Journal of Chemical Physics.

### 6.1 Definition and properties

Let us recall Metropolis algorithm (1953). Let  $\{\pi(x)\}_{x \in S}$  be a prob. distri. which is hard to be sampled directly. Metropolis proposed a Markov chain  $\{\tilde{\theta}_n\}_{n \in \mathbb{N}_0}$ , which is reversible and has  $\pi$  as its limiting distribution. Metropolis' procedure is as follows:

- (1) Introduce a Markov chain  $\{\theta_n\}_{n \in \mathbb{N}_0}$ , whose transition probability matrix is  $\{Q(x, y)\}_{x \in S, y \in S}$  such that  $\theta(x, y) = \theta(y, x), x \in S, y \in S$ .
- (2) Suppose that  $\tilde{\theta}_n = x$ , it jumps to  $\tilde{\theta}_{n+1}$  according to the following rule:

- (2.1)  $x$  first transits to  $y$  with a prob.  $Q(x, y)$ .
- (2.2) this transition is accepted with a prob.  $\min\left\{1, \frac{\pi(y)}{\pi(x)}\right\}$ , and rejected with a prob.  $1 - \min\left\{1, \frac{\pi(y)}{\pi(x)}\right\}$ , ie. the chain staying at  $x$ .

- (3) Replace  $n$  by  $n + 1$ , continuing step (2).

We have learnt that the transition probability  $\tilde{Q}(x, y)$  for  $\{\tilde{\theta}_n\}_{n \in \mathbb{N}_0}$  is computed

as the following:

$$\begin{aligned}
\text{As } y \neq x, \tilde{Q}(x, y) &= P(\theta_{n+1} = y, y \text{ accepted} \mid \theta_n = x) \\
&= \frac{P(\theta_{n+1} = y, y \text{ accepted}, \theta_n = x)}{P(\theta_n = x)} \\
&= \frac{P(y \text{ accepted} \mid \theta_{n+1} = y, \theta_n = x) P(\theta_{n+1} = y, \theta_n = x)}{P(\theta_n = x)} \\
&= \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\} \frac{P(\theta_{n+1} = y, \theta_n = x)}{P(\theta_n = x)} \\
&= \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\} Q(x, y).
\end{aligned}$$

$$\begin{aligned}
\text{As } y = x, \tilde{Q}(x, x) &= Q(x, x) + \sum_{y \neq x} P(\theta_{n+1} = y, y \text{ rejected} \mid \theta_n = x) \\
&= Q(x, x) + \sum_{y \neq x} \frac{P(\theta_{n+1} = y, y \text{ rejected}, \theta_n = x)}{P(\theta_n = x)} \\
&= Q(x, x) + \sum_{y \neq x} \frac{P(y \text{ rejected} \mid \theta_{n+1} = y, \theta_n = x) P(\theta_{n+1} = y, \theta_n = x)}{P(\theta_n = x)} \\
&= Q(x, x) + \sum_{y \neq x} \left( 1 - \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\} \right) Q(x, y)
\end{aligned}$$

It is very easy to verify that the following detailed balance condition holds:

$$\pi(x)\tilde{Q}(x, y) = \pi(y)\tilde{Q}(y, x) \quad \forall x, y \in S.$$

Indeed, if  $x = y$ , the relation is trivially true.

$$\begin{aligned}
\text{if } x \neq y, \pi(x)Q(x, y) \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\} &= \min\{\pi(x)Q(x, y), \pi(x)Q(x, y)\} \\
(*) &= \min\{\pi(x)Q(x, y), \pi(y)Q(y, x)\} = \pi(y)Q(y, x) \min \left\{ 1, \frac{\pi(x)Q(x, y)}{\pi(y)Q(y, x)} \right\} \\
(*) &= \pi(y)Q(y, x) \min \left\{ 1, \frac{\pi(x)}{\pi(y)} \right\}
\end{aligned}$$

where \* used  $Q(x, y) = Q(y, x)$ .

Metropdis-Hastings algorithm is a generalization of Metropolis algorithm, first proposed by Hastings, which remove the symmetry condition  $Q(x, y) = Q(y, x)$  and modify the accept rate  $\min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\}$  with the new one  $\min \left\{ 1, \frac{\pi(y)Q(y,x)}{\pi(x)Q(x,y)} \right\}$ .

More precisely, the procedure of MH algorithm is

(1) The auxiliary Markov chain  $\{\theta_n\}_{n \in \mathbb{N}}$  has a transition prob.  $Q(x, y)$ , and can be easily runned in practice. The condition  $Q(x, y) = Q(y, x)$  it NOT necessary.

(2) Suppose  $\tilde{\theta}_n = x$ , it jumps to  $\tilde{\theta}_{n+1}$  according to the following rule:

- (2.1) it first transits to a state  $y$  according to the prob.  $Q(x, y)$ .
  - (2.2) this transition is accepted according to the accept rate  $\min \left\{ 1, \frac{\pi(y)Q(y,x)}{\pi(x)Q(x,y)} \right\}$ , and is rejected with  $1 - \min \left\{ 1, \frac{\pi(y)Q(y,x)}{\pi(x)Q(x,y)} \right\}$  i.e. the chain staying at  $x$ .
- (3) Replace  $n$  with  $n + 1$ , continue the step 2.

**Example 6.1.** Consider the following mixture of Gaussian model

$$\pi(\theta, p) = p f_N(\theta; \mu_1, \Sigma_1) + (1 - p) f_N(\theta; \mu_2, \Sigma_2)$$

where  $p \in (0, 1)$ ,  $f_N(\theta; \mu_1, \Sigma_1) = \frac{1}{2\pi\sqrt{\det(\Sigma_1)}} \exp \left( -\frac{1}{2} (\theta - \mu_1)^\top \Sigma_1^{-1} (\theta - \mu_1) \right)$

$$f_N(\theta, \mu_2, \Sigma_2) = \frac{1}{2\pi\sqrt{\det(\Sigma_2)}} \exp \left( -\frac{1}{2} (\theta - \mu_2)^\top \Sigma_2^{-1} (\theta - \mu_2) \right)$$

$$\mu_1 = \begin{pmatrix} 4 \\ 5 \end{pmatrix}, \quad \mu_2 = \begin{pmatrix} 0.7 \\ 3.5 \end{pmatrix}, \quad \Sigma_1 = \begin{pmatrix} 1.0 & 0.7 \\ 0.7 & 1.0 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 1.0 & 0.7 \\ -0.7 & 1.0 \end{pmatrix},$$

$p$  is usually an unknown number. For simplicity, we assume  $p = 0.7$ . Usually we can sample a random variable  $\chi \sim \text{Ber}(p)$ ,  $X_1 \sim f_N(\theta; \mu_1, \Sigma_1)$ ,  $X_2 \sim f_N(\theta; \mu_2, \Sigma_2)$ ,

we take the combination

$$X = \chi X_1 + (1 - \chi) X_2$$

$$X \sim p f_N(\theta; \mu_1, \Sigma_1) + (1 - p) f_N(\theta; \mu_2, \Sigma_2).$$

We now use MH algorithm to sample  $\pi(\theta)$ . We choose the Markov chain whose transition prob. is  $q(x, y) = f_N(y; x, \gamma I_2) = \frac{1}{2\pi\sqrt{\gamma}} \exp\left(-\frac{1}{2\gamma}|y - x|^2\right)$ . The algorithm reads as

- 1. Initial  $\theta^{(0)} \in \mathbb{R}^2$ ;
- 2. For  $j \geq 1$ , sample  $\tilde{\theta}^{(j)}$  according to the transition prob.  $q\left(\theta^{(j-1)}, \tilde{\theta}^{(j)}\right)$ . Compute  $r_j = \min\left\{1, \frac{\pi(\tilde{\theta}^{(j)}) q(\tilde{\theta}^{(j)}, \theta^{(j-1)})}{\pi(\theta^{(j-1)}) q(\theta^{(j-1)}, \tilde{\theta}^{(j)})}\right\} = \min\left\{1, \frac{\pi(\tilde{\theta}^{(j)})}{\pi(\theta^{(j-1)})}\right\}$ . Sample a  $U \sim U([0, 1])$ , if  $U \leq r_j$ , accept  $\tilde{\theta}^{(j)}$ , i.e.  $\theta^{(j)} = \tilde{\theta}^{(j)}$ , if  $U > r_j$  reject  $\tilde{\theta}^{(j)}$ , i.e.  $\theta^{(j)} = \theta^{(j-1)}$ .
- 3. Replace  $j$  by  $j + 1$ , repeat step 2.

The parameter  $\gamma$  plays an important role in the algorithm;

- (1)  $\gamma$  is too small, the chain has difficulty moving across the parameter space.
- (2)  $\gamma$  is too large, the acceptance rate is small, leading to a computationally inefficient sampling scheme.