

# Markov Chain Monte Carlo

## Bayes methods

2023-02-11

Target distribution  $\underline{P}$ .

$x_0, x_1, x_2, \dots$

Markov chain. (Finite)

(State space)

① irreducible. ( $\forall i, j \in \Omega$ ,  $i \leftrightarrow j$ )  
( $\lambda$ -连続)  
( $\exists$  通)

② aperiodic ( $\forall i \in \Omega$ ,  $d_i = 1$ ,  $d = \gcd(n, p_{ii}^{(n)}) > 0$ ).  
( $\exists$  周期)

example,

$1 \xrightarrow{2} 2$        $d_1 = \gcd(3, 6, 9) = 3$   
 $\uparrow \downarrow$        $1 \leftrightarrow 2 \leftrightarrow 3$

$d_1 = d_2 = d_3 = 3$ , not aperiodic.

$1 \xleftrightarrow{2} 2$        $d_1 = \gcd(2, 3, \dots) = 1$   
 $\uparrow \downarrow$        $d_1 = d_2 = d_3 = 1$ , aperiodic

stationary distribution.

$P P = P$ .      ( $X_t \xrightarrow{\text{D}} p$ )  
↓  
Stationary distribution  
transition probability matrix

Retailed balance (细致平衡条件)

time reversible (时间可逆) ~~☆☆☆☆~~

$$\underbrace{P_i P_{i,j}}_{\downarrow} = \underbrace{P_j P_{j,i}}_{\downarrow} \quad \forall i, j \in \mathbb{N}.$$

$$\text{rate of } i \rightarrow j \quad \text{rate of } j \rightarrow i.$$

(\*)  $P_{ij}$  is a TPR of a irreducible chain if we can find a sequence  $\{\pi_i\}$ ,  $\sum_i \pi_i = 1$  and an another TPR  $Q_{ij}$  satisfying:

$$\underbrace{\pi_i P_{ij}}_{\pi_j Q_{ji}}$$

$Q_{ji}$  is the TPR of the reversible chain,

$\{\pi_i\}$  is the stationary probability.

Convergence and the Ergodic Theorem

(遍历)

$$\frac{\sum_{t=1}^T f(X_t)}{T} \xrightarrow{a.s.} E_{x \sim p}(f(x))$$

M-H Algorithm. (target distribution  $\varphi$ )

$$X^{(t)} \rightarrow X^{(t+1)}$$

(在每一步 转移, 若转移后能提高状态  $X_t$  在  
目标分布  $P$  中有密度值, 就接受转移结果;  
否则, 以一定的概率决定是转移还是不动).

Let  $q(j|i) = Q_{i,j}$ , proposal distribution.  
(simple)

$$q(j|i) > 0 \Leftrightarrow q(i|j) > 0$$

$X_t = i$  (discrete)

1. Draw  $j \sim q(\cdot|i)$  and  $u \sim U[0,1]$

2. If  $u \leq \alpha(j|i)$  where

(If  $i=j$ ,  $\alpha(j|i)=1$ )

$$\alpha(j|i) = \min \left\{ 1, \frac{p_j q(i|j)}{p_i q(j|i)} \right\}$$

(accept  $j$  with probability  $\alpha(j|i)$ )

then set  $X_{t+1} = j$ , otherwise set  $X_{t+1} = i$ .

(Start at  $X_0 = i_0$ )

Prove distribution of MCMC converges to  $p$

$$P_{i,j} = \frac{q(j|i)}{q(i|j)} \alpha(j|i) \quad (j \neq i).$$

$$P_i P_{i,j} = P_i q(j|i) \underline{\alpha(j|i)}$$

$$= \frac{P_i q(j|i)}{\min \left\{ 1, \frac{P_j q(i|j)}{P_i q(j|i)} \right\}}$$

$$= \min \left\{ P_i q(j|i), \frac{P_j q(i|j)}{\underline{P_j q(i|j)}} \right\}$$

$$= P_j q(i|j) \underline{\min \left\{ \frac{P_i q(j|i)}{P_j q(i|j)}, 1 \right\}}$$

$$= P_j q(i|j) \alpha(i|j)$$

$$= P_j P_{j,i}.$$

$\Rightarrow \{P_i\}$  is the stationary distribution.

Example. Hyper Geom ( $k; K, N, n$ )

$$\Omega = \{k \in \mathbb{Z} : \underbrace{\max\{0, n+K-N\}}_{B^-} \leq k \leq \underbrace{\min\{n, K\}}_{B^+}\}.$$

$$P_k = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}, \quad B^- \leq k \leq B^+$$

$$P_k = 0, \quad k \notin [B^-, B^+]$$

$$q(j|i) = \begin{cases} \frac{1}{2}, & \text{for } j = i \pm 1 \\ 0, & \text{otherwise} \end{cases}$$

① Draw  $j \sim \underline{U\{i-1, i+1\}}$ . and  $u \sim U(0, 1)$   
 (discrete uniform)

② If  $B^- \leq j \leq B^+$ ,

$$\begin{aligned} u &\leq \min\left\{1, \frac{P_j q(i|j)}{\sum_{j=1}^2 P_j q(j|i)}\right\} \\ &= \min\left\{1, \frac{P_j}{P_i}\right\} \\ &= \min\left\{1, \frac{\binom{K}{j} \binom{N-K}{n-j}}{\binom{K}{i} \binom{N-K}{n-i}}\right\}. \end{aligned}$$

then Set  $X_{t+1} = j$ , otherwise  $X_{t+1} = i$ .

If  $j < \beta^-$  or  $j > \beta^+$ , reject.

continuous

$$P(X_{t+1} \in A | X_t = i)$$

$$\underline{k(v, d\theta')} = p(X_{t+1} = d\theta' | X_t = \theta)$$

MH algorithm (continuous) (target distribution  $p$ )

$q(\theta'| \theta)$ . proposal probability density.

$$q(\theta'| \theta) > 0 \Leftrightarrow q(\theta | \theta') > 0.$$

1. Draw  $\theta' \sim q(\cdot | \theta)$  and  $u \sim \mathcal{U}[0, 1]$

2. If  $u \leq \alpha(\theta' | \theta)$ , where

$$\alpha(\theta' | \theta) = \min \left\{ 1, \frac{p(\theta', q(\theta | \theta'))}{p(\theta) q(\theta' | \theta)} \right\}$$

Set  $X_{t+1} = \theta'$ , otherwise  $X_{t+1} = \theta$ .

$$\theta = (\theta_1, \dots, \theta_p) \xrightarrow{\text{component}}$$

Updating different components using different proposal distribution. Suppose we have  $N$  different proposal densities  $\underline{q_k(\theta'|\theta)}$  and we select  $k^{\text{th}}$  one to make a proposal with probability  $\underline{\pi_k}$  and then accept or reject.

1. Draw  $k \sim \text{Multinom}(\underline{\pi_1, \dots, \pi_N})$

$\theta' \sim q_k(\cdot | \theta)$  and  $u \sim \mathcal{U}[0, 1]$

2. If  $u \leq \alpha_k(\theta' | \theta)$ ,

$$\alpha_k(\theta' | \theta) = \min \left\{ 1, \frac{p(\theta') q_k(\theta | \theta')}{p(\theta) q_k(\theta' | \theta)} \right\}$$

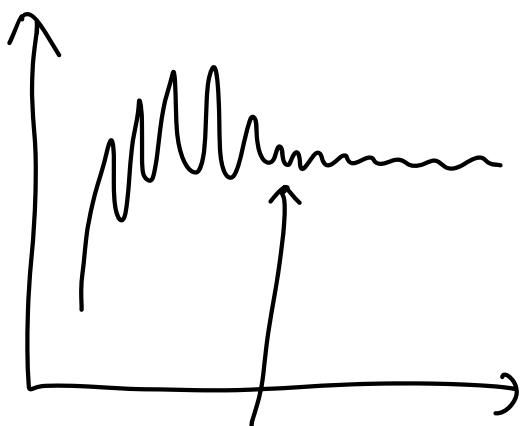
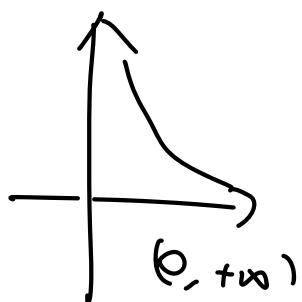
then set  $X_{t+1} = \theta'$ , otherwise set  $X_{t+1} = \theta$ .

If we have  $p$  components,  $\theta = (\theta_1, \dots, \theta_p)$ .

Take  $N=p$ , if we choose to update component  $k$ .

then propose  $\underline{\theta'_k} = \underline{\theta_k}$  and update  $\underline{\theta'_k} \sim q_{k|1}(\theta)$   
 $|\theta' - \theta|$  small,  $P(\theta') \approx P(\theta)$ .

Smaller changes will typically have a higher acceptance probability. The chain takes more steps to explore the support of the target density in  $\underline{\Omega}$ .



convergence time.

Example, mixture of two bivariate normals.

$$\pi(\theta) = (2\pi)^{-1} \left[ 0.5 e^{-\|\theta - \mu_1\|^2 / 2} + 0.5 e^{-\|\theta - \mu_2\|^2 / 2} \right]$$

$\theta = (\theta_1, \theta_2)$ .  $\mu_1, \mu_2, \Sigma_1, \Sigma_2$  known.

proposal distribution.  $[e_n]$

$$\begin{aligned} \theta'_i &\sim U(\theta_i - a, \theta_i + a) \\ q(\theta'_i | \theta_i) &= \frac{1}{2a} \end{aligned}$$

$$q(\theta' | \theta) = \left(\frac{1}{2a}\right)^2 = \frac{1}{4a^2}$$

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

1. for  $i = 1, 2$  simulate  $\theta'_i \sim U(\theta_i - a, \theta_i + a)$   
and  $u \sim \text{Unif}(0, 1)$

2. If  $u \leq \alpha(\theta' | \theta)$ ,

$$\alpha(\theta' | \theta) = \min \left\{ 1, \frac{\pi(\theta') q(\theta | \theta')}{\pi(\theta) q(\theta' | \theta)} \right\}$$

$$= \min \left\{ 1, \frac{\pi(\theta')}{\pi(\theta)} \right\}.$$

$$\theta^{(n+1)} = \theta', \text{ otherwise } \theta^{(n+1)} = \theta.$$

$$(\theta^{(n)} = [1])$$

# MH example: an equal mixture of bivariate normals

```
f<-function(x,mu1,mu2,S1i,S2i,p1=0.5) {  
  c1<-exp(-t(x-mu1)%*%S1i%*%(x-mu1))  
  c2<-exp(-t(x-mu2)%*%S2i%*%(x-mu2))  
  return(p1*c1+(1-p1)*c2)  
}  
a=3; n=2000; mu1=c(1,1); mu2=c(4,4);  
S=diag(2);  
S1i=S2i=solve(S);  
X=matrix(NA,2,n);  
X[,1]=x=mu1 #initial value  
for (t in 1:(n-1)) {  
  y<-x+(2*runif(2)-1)*a  
  MHR<-f(y,mu1,mu2,S1i,S2i)/f(x,mu1,mu2,S1i,S2i)  
  if (runif(1)<MHR) x<-y  
  X[,t+1]<-x  
}
```

$$y = x_1 + \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} a.$$

$$= x_1 + [v_{(0,2)} - 1] a$$

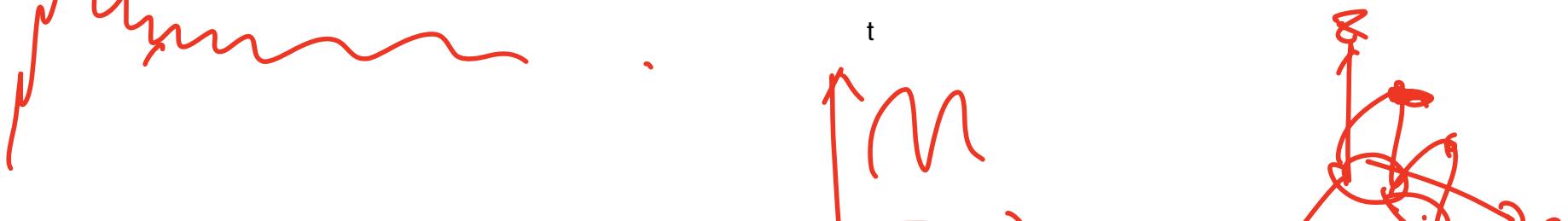
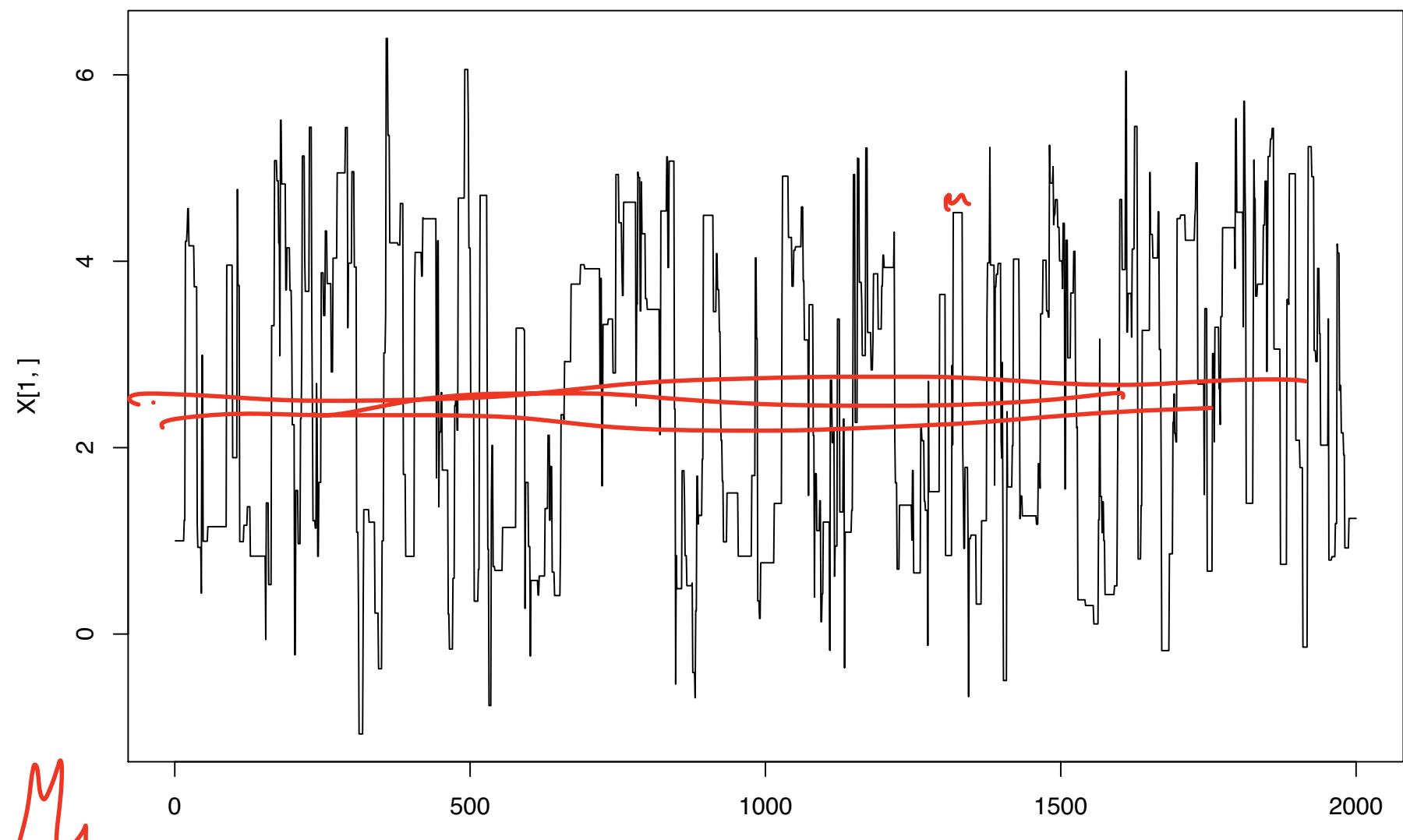
$$= x_1 + v_{(-1,1)} a$$

$$= v(x_1 - a, x_2 + a)$$

$$X = \begin{pmatrix} & \end{pmatrix} \begin{pmatrix} 0, \\ 0, \end{pmatrix}$$

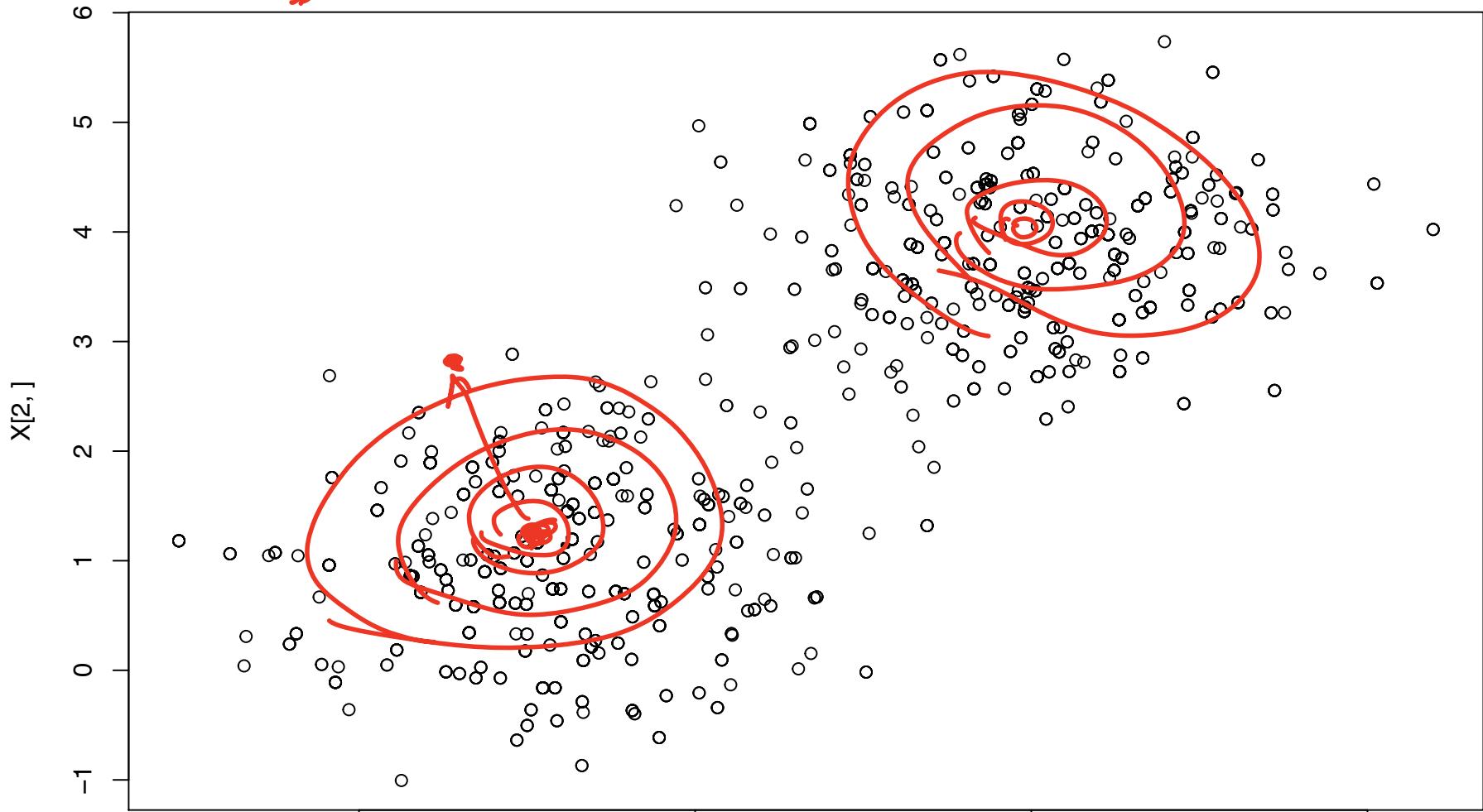
path of  $\theta_1$ ,  $a = 3$ ,  $n = 2000$

$n$ .



$$\theta_2 \sim \theta_1, a = 3, n = 2000$$

(best)



Modes :  $-E_2^r$  path.  $\theta^{(t)}$

$X[1,]$

$\tau_1(\theta^{(t+1)}) \gg \tau_1(\theta^{(t)})$

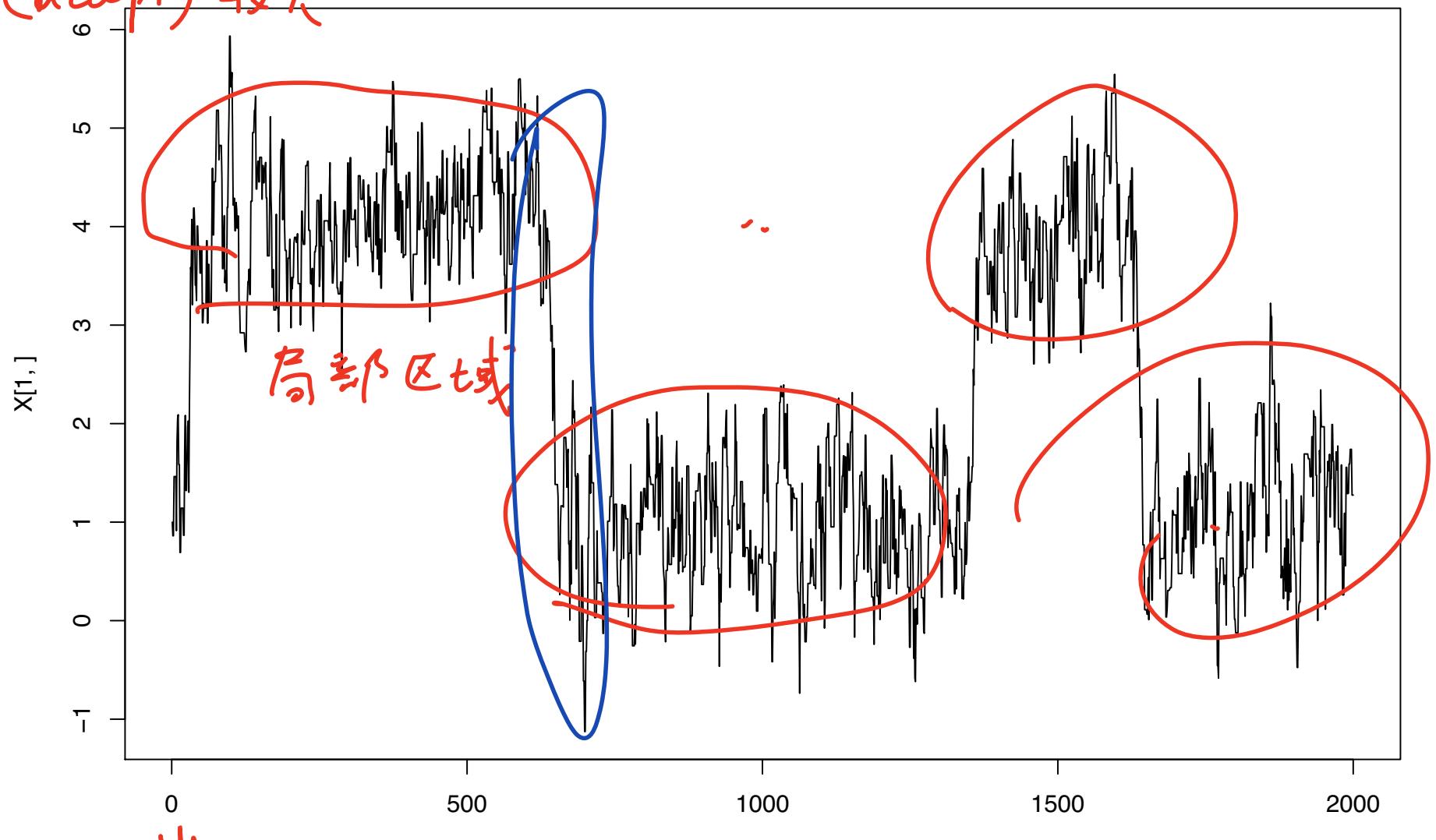
path of  $\theta_1$ ,  $a = 1, n = 2000$

$\theta_{t+1}$

(多峰分布容易出现长时间)

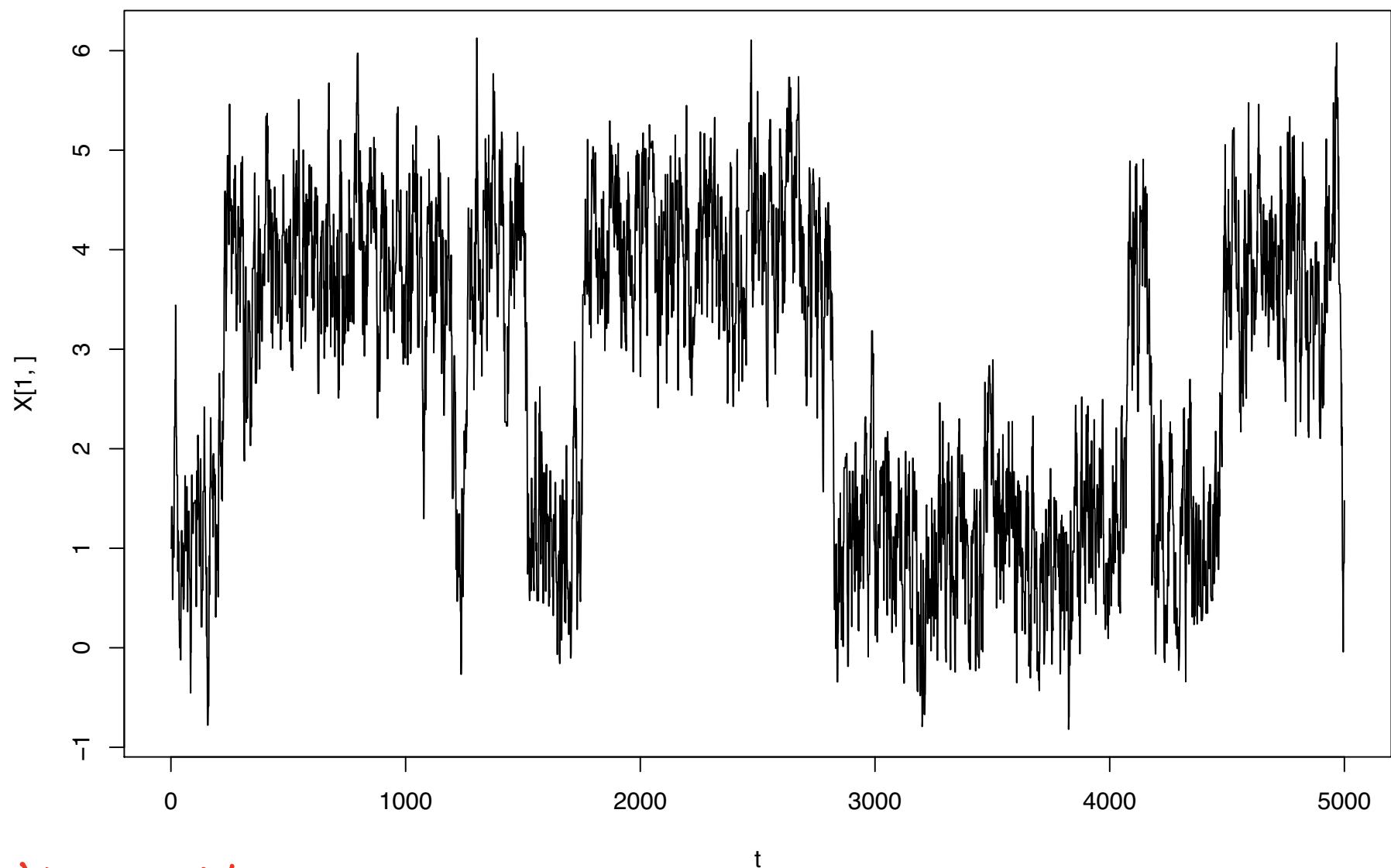
停留在一个  $\theta$  的局部区域内)

$p(\text{accept})$  较大



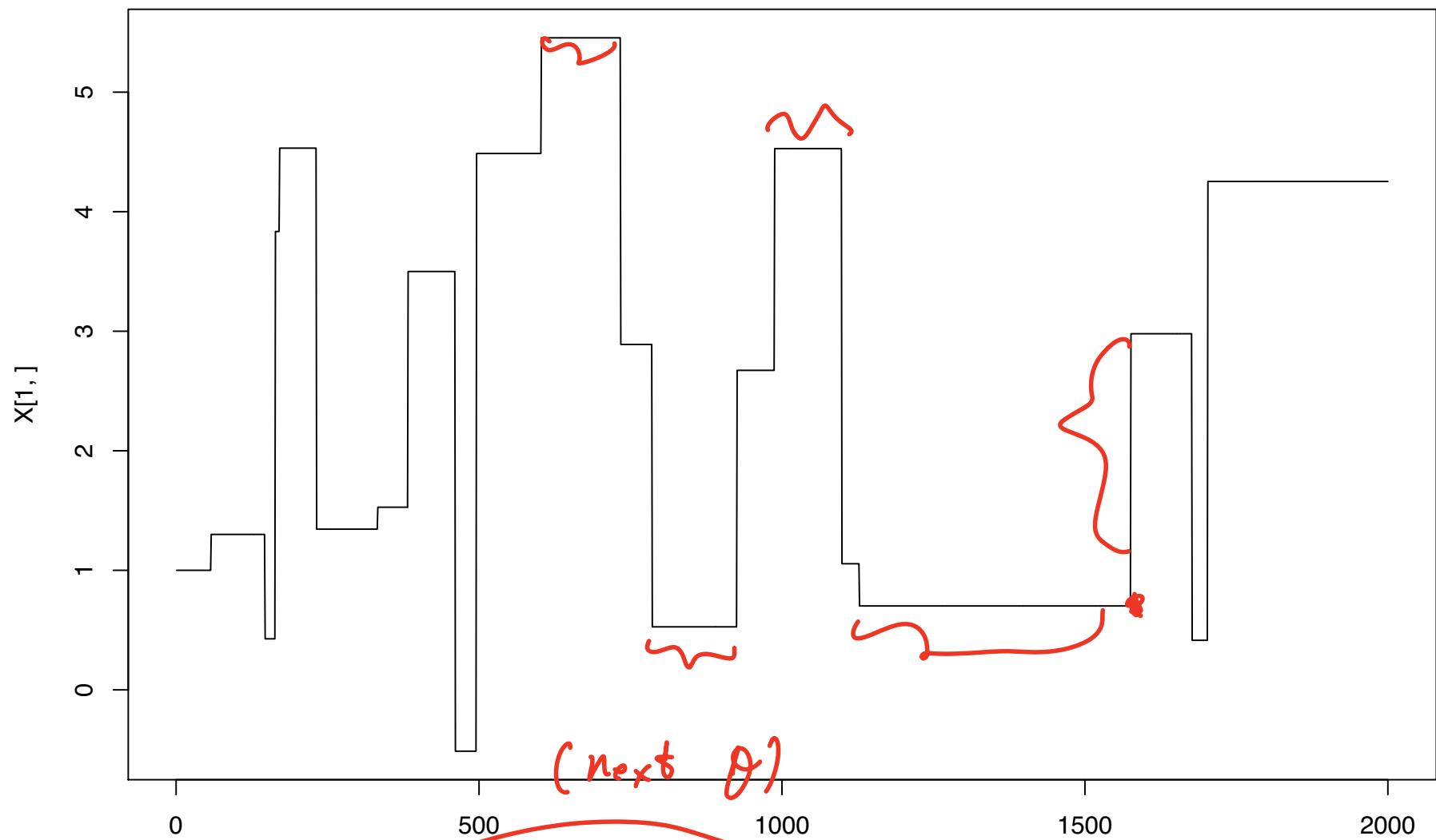
W<sub>11</sub>

path of  $\theta_1$ ,  $a = 1$ ,  $n = 5000$



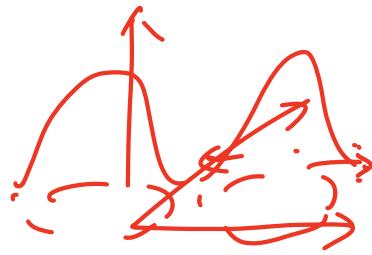
$$X_t \neq X_{t+1}$$

path of  $\theta_1$ ,  $a = 20$ ,  $n = 2000$   $P(\text{accept}) \neq 1$



It will make many proposals into very low density  
States in the tail of the density, which will

be rejected.



Initial state. (sample) from target distribution  
 $pP = p$

(check 收敛的必要条件, 充分条件不行)

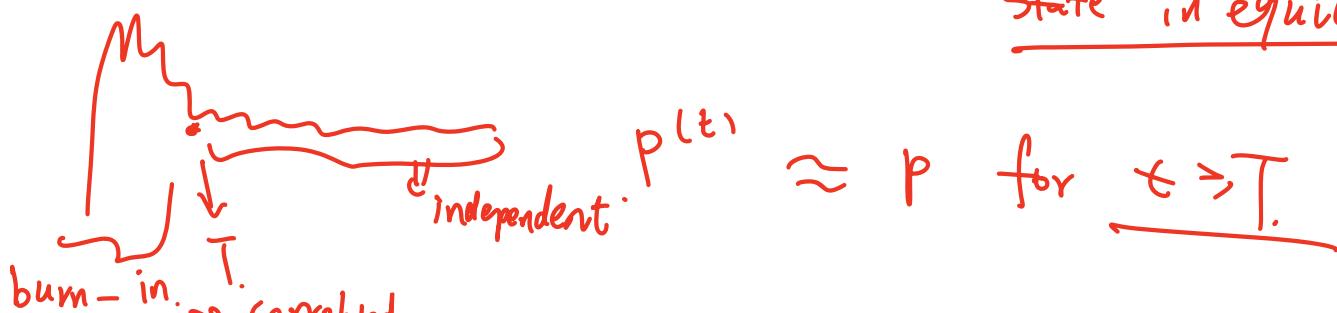
Convergence and mixing.

$$\text{MCMC samples } X_0, \dots, X_n \sim p(x) \\ \bar{f}_n = \frac{\sum_t f(X_t)}{n} \xrightarrow{\text{a.s.}} E_p[f(x)]$$

bias, variance.

Initial state. (大概率不是 Stationary State

State in equilibrium)



$$P^{(t)} \approx P \text{ for } t \geq T.$$

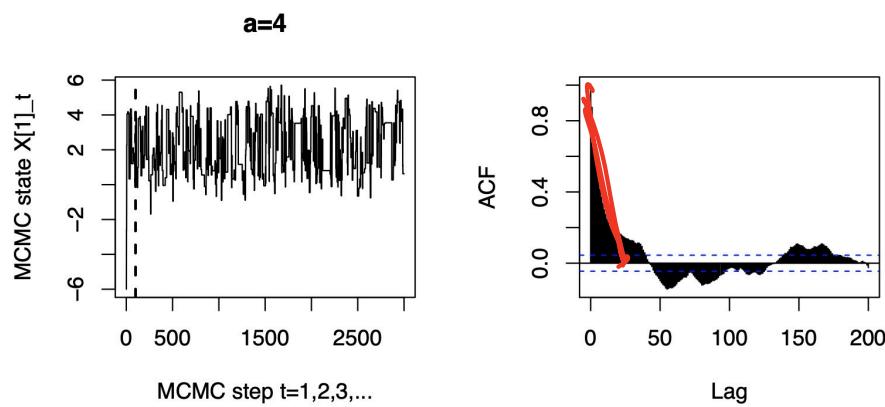
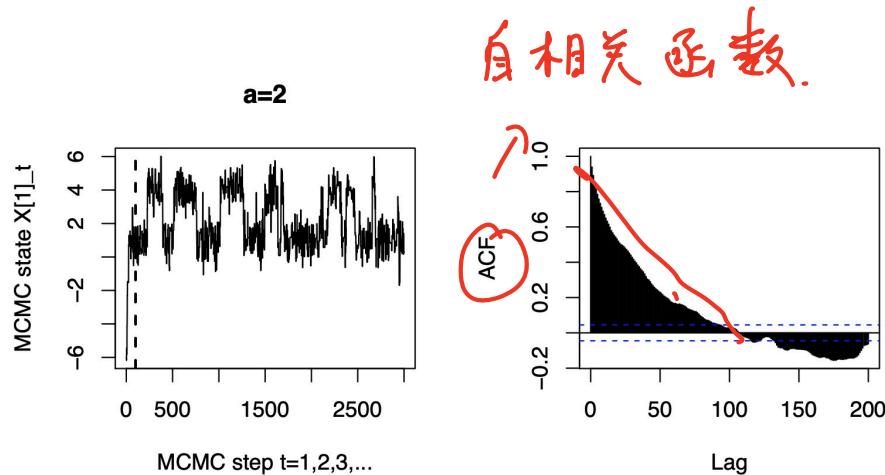
If  $n \gg T$ , the retained samples are representative  
of  $P$ .



Second, assume start the chain in equilibrium,

$\text{Var}(\bar{f}_n) \downarrow$  as  $n \uparrow$ ,  $\bar{f}_n$  has useful precision.

$\text{Var}(\bar{f}_n)$  calculation is complex as the  
MCMC samples are correlated.



— (no serial correlation)

$$\text{Var}(\bar{f}_n) = \frac{\text{Var}(f(x))}{\text{ESS}}$$

$x \sim p$

effective sample size.

(The number of independent samples which would give the same variance reduction as our n correlated samples.)

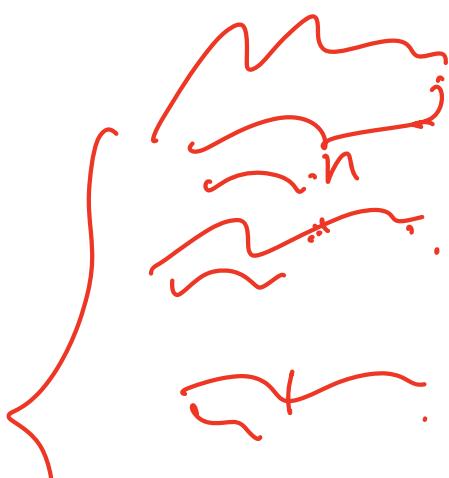
$\text{ESS} \ll n$ .

$\text{ESS} = n$  (MCMC Samples were independent)

$\hat{\sigma}_f^2$  is the estimate of  $\text{Var}(f(x)) = (\sigma_f^2)$

Let  $\hat{\sigma}_{f,n}^2 = \text{Var}(\bar{f}_n)$ ,  $\text{ESS} = \frac{\hat{\sigma}_f^2}{\hat{\sigma}_{f,n}^2}$ .

Simulate  $K \uparrow$  Markov chain, length =  $n$ .



$$\bar{f}_{1,n}, \dots, \bar{f}_{k,n}$$

Sample  $\underline{\theta}^{(k,t)}$ ,  $k = 1, \dots, K$ ,  
 $t = 1, \dots, n$ .

$$\bar{f}_{k,n} = \frac{\sum_t f(\underline{\theta}^{(k,t)})}{n}$$

$$\begin{aligned} \text{Var}(\bar{f}_n) &= \hat{\sigma}_{\bar{f}_n}^2 \\ &= \frac{1}{K-1} \sum_{k=1}^K \left( \bar{f}_{k,n} - \frac{\sum_{j=1}^K \bar{f}_{j,n}}{K} \right)^2. \end{aligned}$$

$$ESS = \frac{\hat{\sigma}_f^2}{\hat{\sigma}_{\bar{f}_n}^2} \quad (\text{precision gain!})$$

fix  $n$  the same, max  $ESS \Rightarrow$  algorithm.



equilibrium

分母 .  $n = 1000$



~~drop~~ burn-in

$$P_s = \frac{\text{Cov}(f(x_t), f(x_{t+s}))}{\sqrt{\text{Var}(f(x_t))} \sqrt{\text{Var}(f(x_{t+s}))}}$$

$$= \frac{\text{Cov}(f(\underline{x_t}), f(\underline{x_{t+s}}))}{\sqrt{\text{Var}(f(\underline{x_t}))}}$$

$$\text{Var}(\bar{f}_n) = \text{Var}\left(\frac{\sum_{t=1}^n f(x_t)}{n}\right)$$

$$= \text{Var}\left(\underbrace{f(x_1) + \dots + f(x_n)}_{\text{Sum of } n \text{ terms}}\right)$$

$$= \frac{\sum_{i=1}^n \sum_{j=1}^n \text{Cov}(f(x_i), f(x_j))}{n^2}$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P_{|i-j|} \cdot \sigma^2.$$

$$= \sum_{i=1}^n \sum_{j=i}^n + \sum_{i=1}^n \sum_{j>i}^n + \sum_{i=n}^1 \sum_{j=i}^n$$

$$= \sum_{i=1}^n \sum_{j=i}^n + \sum_{i=1}^n \sum_{j>i}^n + \sum_{i=n}^1 \sum_{j=i}^n$$

$$= \frac{\sigma^2}{n^2} \left( n + 2 \sum_{i=1}^n \sum_{j>i}^n p_{i-j} \right)$$

$$= \frac{\sigma^2}{n^2} \cdot \left( n + 2 \sum_{s=1}^{n-1} (n-s) p_s \right)$$

$$= \frac{\sigma^2}{n} \left( 1 + 2 \sum_{s=1}^{n-1} \left( 1 - \frac{s}{n} \right) p_s \right)$$

$$= \frac{\sigma^2}{n} \left( 1 + 2 \sum_{s=1}^{n-1} p_s - \frac{\sum_{s=1}^{n-1} s p_s}{n} \right)$$

$$\approx \frac{\sigma^2}{n} \left( 1 + 2 \sum_{s=1}^{n-1} p_s \right)$$

$\circlearrowleft J_f$

If  $p_s \rightarrow 0$  rapidly at large  $s$ .

$$\sum_{i=1}^n \sum_{j>i}^n p_{i-j}$$

$$\frac{p_1 + p_2 + p_3 + p_4 + p_5 + p_6}{+}$$

$$\sum_{s=1}^{n-1} (n-s) p_s$$

$\circlearrowleft (n-1)p_1 + (n-2)p_2 + \dots + p_{n-1}$

$$ESS = \frac{\text{Var}(f(x))}{\text{Var}(f_h)} = \frac{\sigma^2}{\frac{\sigma^2}{n} \cdot J_f} = \frac{n}{J_f}$$

Sample correlation estimate  $\hat{P}_S = \frac{\hat{r}_S}{\hat{r}_0}$

$$\hat{T}_f = 1 + 2 \sum_{S=1}^M \hat{P}_S$$

$$1 + 2 \sum_{S=1}^{n-1} \hat{P}_S$$

① the  $S \hat{r}_S$ ,  $\hat{P}_S$  本来就很小

②  $M \hat{r}_S$  noise

Geyer, for a markov chain:

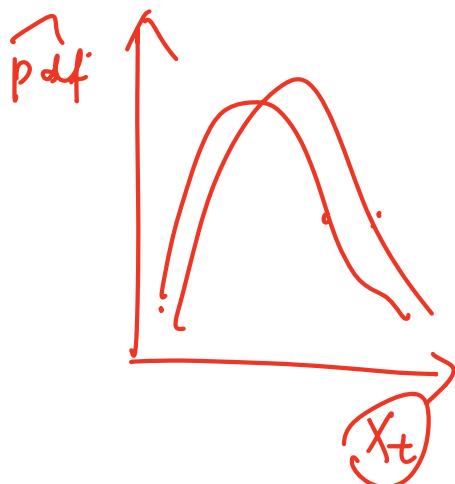
$P_t + P_{t+1}$  is positive, Monotone, Convex

choose  $M$  equal to the least  $s$  satisfying

$$\hat{P}_S + \hat{P}_{S+1} > 0 \text{ and } \hat{P}_S + \hat{P}_{S+1} > \hat{P}_{S-1} + \hat{P}_S$$

① start at different initial states.

check marginal distributions agree?



② ACT plot. 看 correlation falls rapidly?

Calculate ESS.

③ plot MCMC trace.

