

Lecture 17.

## Metropolis Algorithm

Note Title

5/7/2006

$$\pi(\underline{x}) = \frac{1}{Z} e^{-E(\underline{x})}$$

← size

Neighborhood  $N(\underline{x})$ , s.t. if  $y \in N(\underline{x})$  then  $\underline{x} \in N(y)$ .  $|N(\underline{x})|$  indep. of  $\underline{x}$

Transition kernel:

$$K(\underline{x}^{t+1} | \underline{x}^t) = \frac{1}{|N(\underline{x}^t)|} \min \left\{ 1, \frac{\pi(\underline{x}^{t+1})}{\pi(\underline{x}^t)} \right\}, \quad \begin{array}{l} \underline{x}^{t+1} \in N(\underline{x}^t) \\ \underline{x}^{t+1} \neq \underline{x}^t \end{array}$$

propose a move  $\underline{x}^t \rightarrow \underline{x}^{t+1} \in N(\underline{x}^t)$  with uniform probability

accept move with prob  $\min \left\{ 1, \frac{\pi(\underline{x}^{t+1})}{\pi(\underline{x}^t)} \right\}$

$$K(\underline{x}^t | \underline{x}^t) = 1 - \sum_{\underline{x}^{t+1} \in N(\underline{x}^t)} K(\underline{x}^{t+1} | \underline{x}^t) \rightarrow \text{prob you stay at } \underline{x}^t$$

Intuition for Metropolis:

$$\min \left\{ 1, \frac{\pi(\underline{x}^{t+1})}{\pi(\underline{x}^t)} \right\} = \min \left\{ 1, e^{E(\underline{x}^t) - E(\underline{x}^{t+1})} \right\}$$

If  $E(\underline{x}^{t+1}) \leq E(\underline{x}^t)$ , the move  $\underline{x}^t \rightarrow \underline{x}^{t+1}$  will always be accepted.

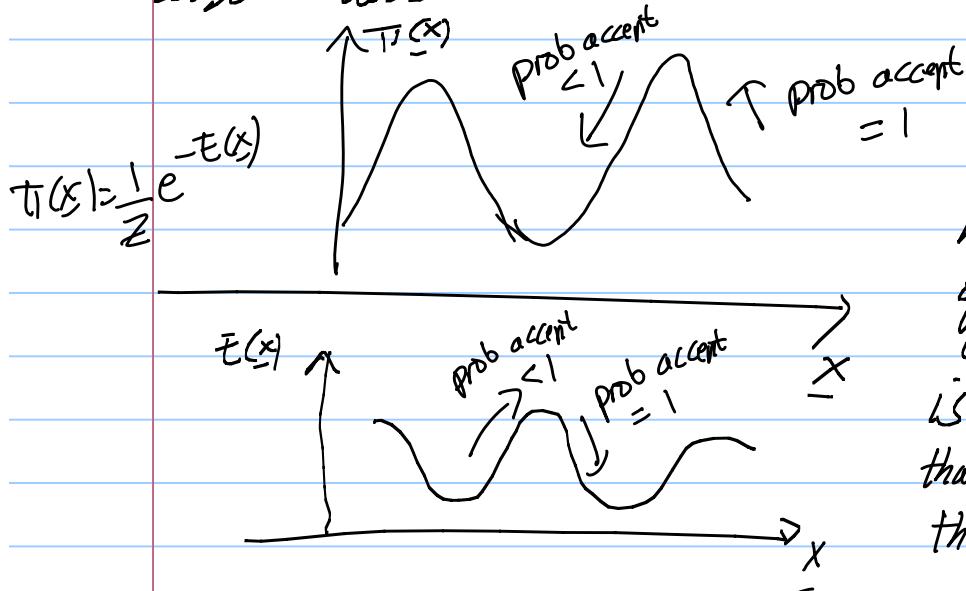
If  $E(\underline{x}^{t+1}) > E(\underline{x}^t)$ , the move  $\underline{x}^t \rightarrow \underline{x}^{t+1}$  will be accepted with probability  $e^{E(\underline{x}^t) - E(\underline{x}^{t+1})} < 1$ .

Hence proposed moves to lower energy / higher probability states will always be accepted.

But moves to higher energy / lower probability states have a probability of being accepted.

Page 2

Intuitively, Metropolis encourages movement to states of high probability, but also allows movement to low probability states.



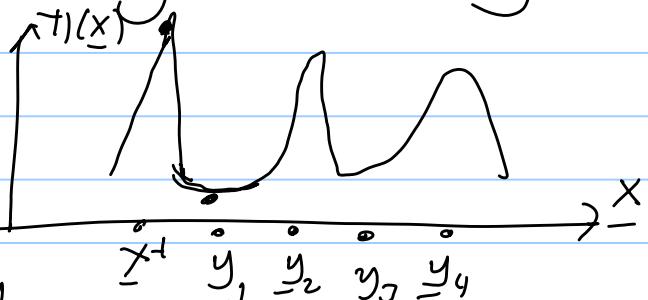
This ability to move in directions of lower probability (or higher energy) is needed to ensure that the MCMC explores the space of all  $x$ .

Metropolis obeys detailed balance (last lecture) and so is guaranteed to converge eventually to samples from  $\pi(x)$ .

But - WARNING - Convergence may take a long time.

Suppose after  $t$  iterations the chain is at  $x^t$

Then the probability to moving downhill in probability from  $\pi(x^t)$  to  $\pi(y_i)$  is small. It requires a sequence of moves which decrease the probability and which all have  $\text{prob} < 1$  being accepted.



Page 3.

## Metropolis-Example: Ising Model in 1-D.

Note: we already know how to sample a 1-D Ising model by sequential methods. This enables us to compare MCMC to sequential sampling.

$$P(x) = \frac{1}{Z} e^{-\beta \sum_{i=0}^{d-1} x_i x_{i+1} + \gamma \sum_{i=0}^d x_i c_i} \quad x_i \in \{-1\}$$

Choose neighborhood structure:

$$N(\underline{x}) = \{ (\underline{x}_0, x_1, \dots, \bar{x}_k, \dots, x_d) : k \in \{0, \dots, d\} \text{ with } \bar{x}_k = -x_k \}$$

i.e. flip the state  $x_k \rightarrow \bar{x}_k$ , of one node  $k$ .

Metropolis: At time  $t$ , state  $\underline{x}^t = (x_0^t, \dots, x_d^t)$

Select node  $k$  with prob  $1/d$

$$\text{Let } \underline{x}_{1,k}^t = (x_0^t, \dots, \bar{x}_k^t, \dots, x_d^t).$$

Calculate:  $E(\underline{x}_{1,k}^t) - E(\underline{x}^t)$

$$= 2\gamma x_k^t c_k + 2\beta x_k^t \langle x_{k+1}^t + x_{k-1}^t \rangle, \text{ for } k \neq 0, d \\ (\text{for } k=0, \text{ or } k=d \text{ we modify the second term}).$$

If  $E(\underline{x}_{1,k}^t) - E(\underline{x}^t) < 0$ , set  $\underline{x}^{t+1} = \underline{x}_{1,k}^t$

If  $E(\underline{x}_{1,k}^t) - E(\underline{x}^t) > 0$ , set  $\underline{x}^{t+1} = \underline{x}_{1,k}^t$  with prob.  $e^{E(\underline{x}^t) - E(\underline{x}_{1,k}^t)}$

Set  $\underline{x}^{t+1} = \underline{x}^t$  with prob  $1 - e^{E(\underline{x}^t) - E(\underline{x}_{1,k}^t)}$ .

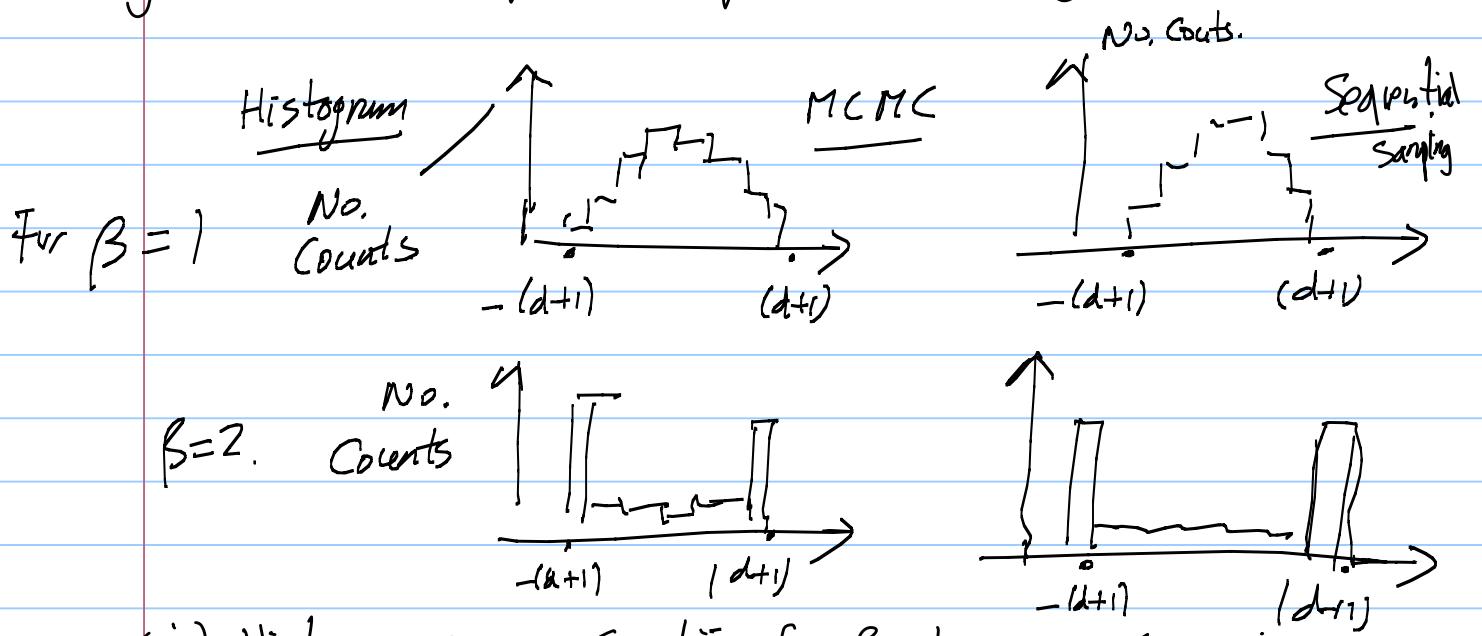
Page 4.

The book gives examples (Liu, p 109) with  $\delta=0$ . Sequential Sampling can be used — ie express  $\Pi(x) = \Pi(x_d) \Pi(x_{d-1}|x_d) \dots \Pi(x_0|x_1)$  — to get samples directly for comparison.

Suppose we want to estimate  $\frac{1}{\sum_x} \left( \sum_{i=0}^d x_i \right) \Pi(x)$ . This is the Magnetization  $M$ .

Consider  $M^t = \frac{1}{d} \sum_{i=0}^d x_i^t$ , where  $(x_0^t, \dots, x_d^t)$  note  $-(d+1) \leq M^t \leq (d+1)$ . is the  $t^{\text{th}}$  sample of the MCMC.

Run Metropolis for 1,000,000 iterations  
 Choose samples  $x^{50}, x^{100}, \dots, x^{50k}, \dots$   $k=1 \text{ to } 20,000$   
 to get 20,000 samples and plot their histograms.



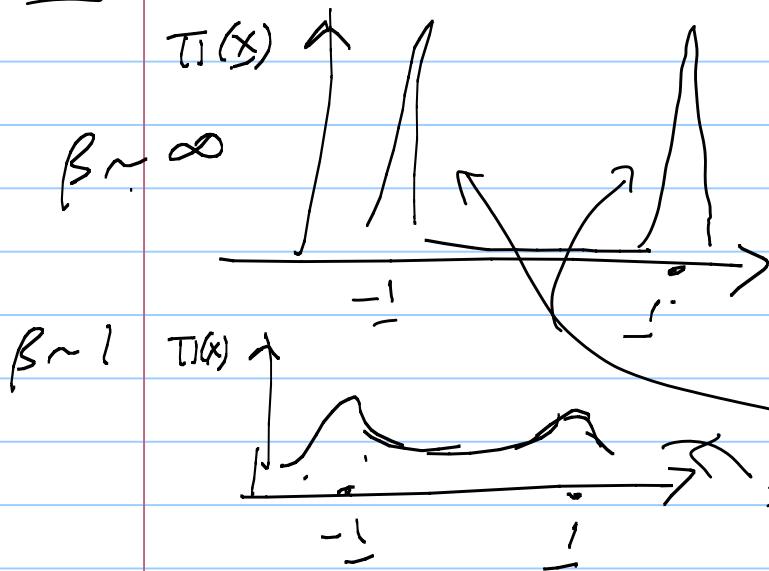
(i) Histograms are similar for  $\beta=1$ . MCMC gives similar results to sequential sampling (but slower computation)

Page 5

(ii) Histograms for MCMC and Sequential Sampling differ more for  $\beta=2$ . Also both histograms differ greatly from the histograms at  $\beta=1$ .  
Why?

First: why are MCMC and sequential samples similar for  $\beta=1$  but less similar for  $\beta=2$ ?

Note: The distribution  $\pi(x)$  gets sharper as  $\beta$  gets larger.



$$\begin{aligned} \frac{1}{\beta} &= (1, 1, \dots, 1) \\ \frac{-1}{\beta} &= (-1, -1, \dots, -1) \end{aligned}$$

Harder for MCMC to sample for large  $\beta$ . It is difficult to move between the two peaks.  
Easier to sample from.

So, for  $\beta=1$ , the MCMC has converged to good samples from  $\pi(x)$  at  $N=50$  and so sampling at  $N=100, 150, 200, \dots, 50k, \dots$  are also good.

But, for  $\beta=2$ , the MCMC needs a longer time to converge because of the difficulty of sampling both of the sharp peaks.

You would get better results by using  $N = 1000, 2000, 3000, 4000, \dots, 1000k, \dots$ .

Page 6. Second, why are the magnetization histograms very different from  $\beta=1$  and  $\beta=2$ ?

Answer: Phase Transition.

$$Pr(M=z) = \sum_{\substack{x \\ \text{any value}}} \delta_{x_0 \dots x_d, z} \pi(x)$$

Kronecker delta function  
 Indicator Function

Phase factor - number of ways you can get  $M=z$ .

For  $M=d+1$ , only one way  $\underline{x} = (1, \dots, 1)$

$M=-d+1$ , only one way  $\underline{x} = (-1, \dots, -1)$

Hence  $P(M=d+1) = P(M=-d+1) = \pi(1) = \pi(-1)$

For  $M=d$ , impossible

$M=d-1, (d+1)$  ways  $(-1, 1, 1, \dots)$   $(1, -1, 1, \dots)$  etc

$$P(M=d-1) = \pi(-1, 1, \dots) + \pi(1, -1, 1, \dots) + \dots + \pi(1, 1, \dots, 1, -1)$$

For  $M=0$ , an exponentially large no. of ways provided  $d+1$  is an even number.

So small phase factor for  $M=\pm d+1$   
 enormous phase factor for  $M=0$ .

$P(M=0)$  is the sum of an exponential number of terms  
 Each term may be small, but they can add up to a big number.

Conflicting "forces"

Phase factor wants  $M=0$ .

Probability  $\pi(x)$  wants  $M=d+1$  or  $-d+1$

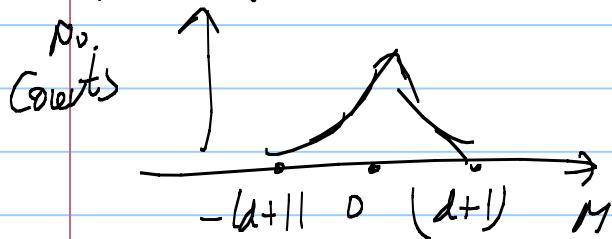
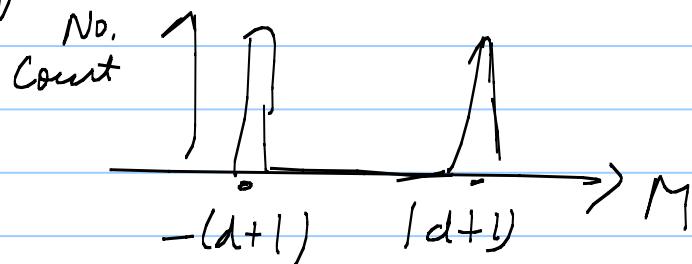
(Page 7) For large  $\beta$ , the probability  $\pi(x)$  is very sharp. The probabilities 'win' and so  $P(M)$  is peaked at  $M = \pm(d+1)$

For small  $\beta$ , the probability  $\pi(x)$  is smoother. The phase factor 'wins' and so  $P(M)$  is peaked at  $M=0$ .

This agrees with the histograms of  $M$  (our estimate of  $P(M)$  make from the samples)

So for large  $\beta$ , probabilities win

For small  $\beta$ , phase factors win.



Summary (1) the number of samples needed for the MCMC to converge will be vary depending on the form of the distribution (e.g. differ from  $\beta=1$  to  $\beta=2$ )

(2) MCMC can give similar results to sequential sampling, but can be extended to cases where sequential doesn't work.

(3) Probability distributions may have phase transitions. Properties - e.g.  $P(M)$  - can change greatly depending on parameters