

Markov Chain Monte Carlo. (MCMC)

Note Title

5/7/2006

MCMC is a way to sample from any distribution $\pi(\underline{x})$

It proceeds by designing a Markov process that converges to samples from $\pi(\underline{x})$.

The original MCMC is the Metropolis algorithm (Metropolis, Rosenbluth & Rosenbluth, Teller & Teller 1953).

Generalized to the Metropolis-Hastings algorithm

strictly speaking, Metropolis-Hastings is a design principle for algorithm.

What is a Markov Chain?

Define a function

$$P(\underline{x}^{t+1} | \underline{x}^t) = K(\underline{x}^{t+1} | \underline{x}^t)$$

obeys $\sum_{\underline{x}^{t+1}} K(\underline{x}^{t+1} | \underline{x}^t) = 1, \forall \underline{x}^t$

Sample $\underline{x}^0, \underline{x}^1, \dots, \underline{x}^t, \dots$

\underline{x}^1 from $P(\underline{x}^1 | \underline{x}^0), \dots, \underline{x}^t$ from $P(\underline{x}^t | \underline{x}^{t-1}) \dots$

This is a Markov Chain. E.g. A random walk.

(*) Aperiodic - almost always satisfied - (*) Irreducible, $\forall x, y \exists n \text{ s.t.}$
 $\sum_{x_1, \dots, x_n} K(x|x_1)K(x_1|x_2)\dots K(x_n|y) > 0.$

MCMC is a special Markov Chain
 where $K(\underline{x}^{t+1}|\underline{x}^t)$ - the transition probability
 - is designed to satisfy the condition:

$$\sum_{\underline{y}} K(\underline{x}|\underline{y})\pi(\underline{y}) = \pi(\underline{x}).$$

In practice, $K(\underline{x}^{t+1}|\underline{x}^t)$ is almost
 always required to satisfy a stronger
 condition called "detailed balance"

$$K(\underline{x}|\underline{y})\pi(\underline{y}) = K(\underline{y}|\underline{x})\pi(\underline{x})$$

Detailed Balance implies that

$$\sum_{\underline{y}} K(\underline{x}|\underline{y})\pi(\underline{y}) = \sum_{\underline{y}} K(\underline{y}|\underline{x})\pi(\underline{x}) = \pi(\underline{x}) //$$

Provided $K(\underline{x}|\underline{y})$ obeys additional technical
 conditions^(*), then samples from the Markov
 Chain will converge to samples from $\pi(\underline{x})$

$\underline{x}^0, \underline{x}^1, \underline{x}^2, \dots, \underline{x}^t, \dots$
 For sufficiently large t , \underline{x}^t is a sample
 from $\pi(\underline{x})$.

Metropolis

Suppose $\pi(\underline{x}) = \frac{1}{Z} e^{-E(\underline{x})}$
(Z hard/impossible to calculate.)

For each \underline{x} , define a neighbourhood structure $N(\underline{x})$ st. $\underline{y} \in N(\underline{x}) \Leftrightarrow \underline{x} \in N(\underline{y})$
and $|N(\underline{x})|$ is independent of \underline{x} .

Loop over t ,

At time t with state \underline{x}^t :

(1) propose a move $\underline{x}^t \rightarrow \underline{x}^{t+1} \in N(\underline{x}^t)$
with uniform probability.
(i.e. $p(\underline{x}^{t+1}) = \frac{1}{|N(\underline{x}^t)|}$)

(2) accept the move with probability
 $\min \left\{ 1, \frac{\pi(\underline{x}^{t+1})}{\pi(\underline{x}^t)} \right\}$. (No need to compute Z !)

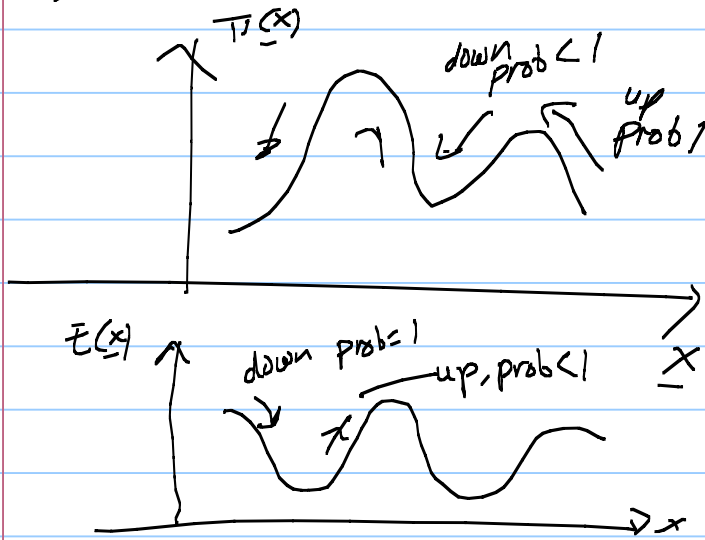
Note: you accept the move for certain

if $\pi(\underline{x}^{t+1}) > \pi(\underline{x}^t)$ (i.e. $E(\underline{x}^{t+1}) < E(\underline{x}^t)$)

but you also accept it with prob. $\frac{\pi(\underline{x}^{t+1})}{\pi(\underline{x}^t)} = e^{E(\underline{x}^t) - E(\underline{x}^{t+1})}$

if $\pi(\underline{x}^{t+1}) < \pi(\underline{x}^t)$ (i.e. $E(\underline{x}^{t+1}) > E(\underline{x}^t)$).

Intuitively, Metropolis's 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 MC explores the full space.

Check that Metropolis obeys Detailed Balance.

$$\begin{aligned}
 K(\underline{x}|\underline{y})\pi(\underline{y}) &= 0 \quad \text{if } \underline{x} \notin N(\underline{y}) \\
 \text{otherwise} &= \frac{1}{|N(\underline{y})|} \pi(\underline{y}) \min\left\{1, \frac{\pi(\underline{x})}{\pi(\underline{y})}\right\} \\
 &= \frac{1}{|N(\underline{y})|} \min\{\pi(\underline{y}), \pi(\underline{x})\} \quad \text{Symmetry} \\
 &= K(\underline{y}|\underline{x})\pi(\underline{x})
 \end{aligned}$$

Recall $\underline{y} \in N(\underline{y}) \Leftrightarrow \underline{y} \in N(\underline{x})$
& $N(\underline{x})$ indep of \underline{x}

Also the MC must be able to explore every part of the space (Irreducible)

Metropolis Example

Ising Model:

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

$x_i \in \{\pm 1\}$

Define neighbourhood structure.

$$N(\underline{x}) = \{ (x_0, x_1, \dots, \bar{x}_k, \dots, x_d) : k = 0, \dots, d \}$$

where $\bar{x}_k = -x_k$.

At time t : \underline{x}^t
select node k at random
(prob $1/d+1$)

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

$$E(\underline{x}_{,k}^t) - E(\underline{x}^t) = 2\gamma x_k^t C_k \quad k \neq 0, d \\ + 2\beta x_k^t (x_{k+1}^t + x_{k-1}^t).$$

Set $\underline{x}^{t+1} = \underline{x}_{,k}^t$ with prob $\frac{1}{2}$
if $E(\underline{x}_{,k}^t) - E(\underline{x}^t) < 0$

$\underline{x}^{t+1} = \underline{x}_{,k}^t$ with prob $e^{\frac{E(\underline{x}^t) - E(\underline{x}_{,k}^t)}{2}}$
if $E(\underline{x}^t) - E(\underline{x}_{,k}^t) < 0$

otherwise $\underline{x}^{t+1} = \underline{x}^t$.

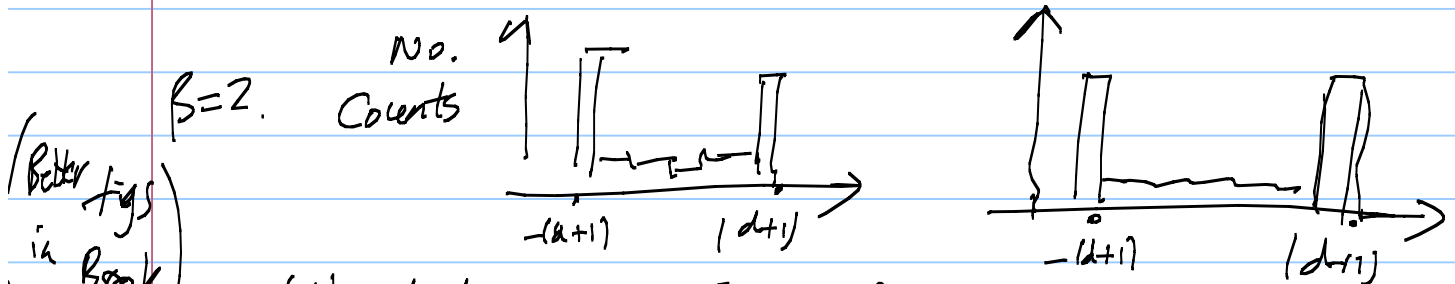
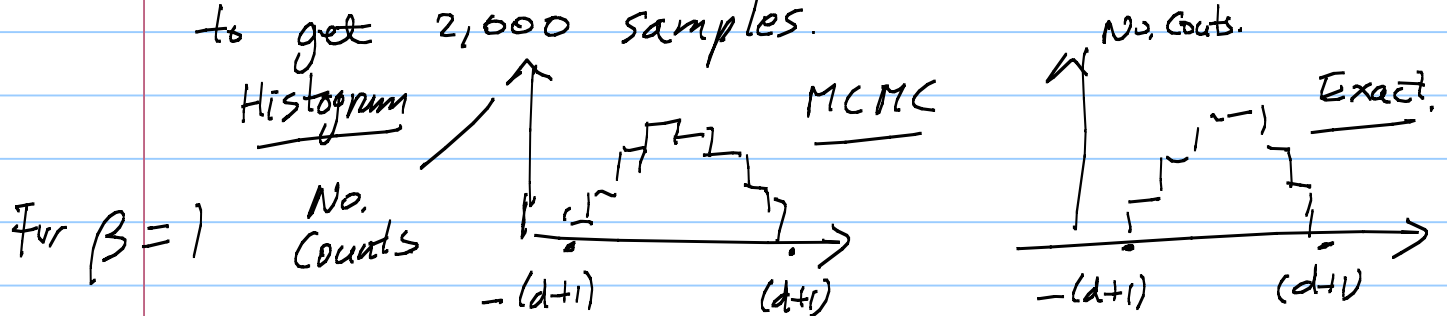
Metropolis Example (cont)

The book gives examples (with $\gamma=0$), p109.

For this problem, you can do exact sampling - express $\pi(x) = \pi(x_d) \pi(x_{d-1}|x_d) \dots \pi(x_0|x_1)$ (using DP), then sample x_d from $\pi(x_d)$, $x_{d-1} \sim \pi(x_{d-1}|x_d)$ etc.

Summary Statistic "magnetization" $M^{(t)} = \frac{d}{Z} \sum_{i=1}^d X_i^{(t)}$

Run Metropolis for 1,000,000 iterations
Choose $x^{50}, x^{100}, \dots, x^{50k}, \dots$ $k=1, \dots, 20,000$
to get 2,000 samples.

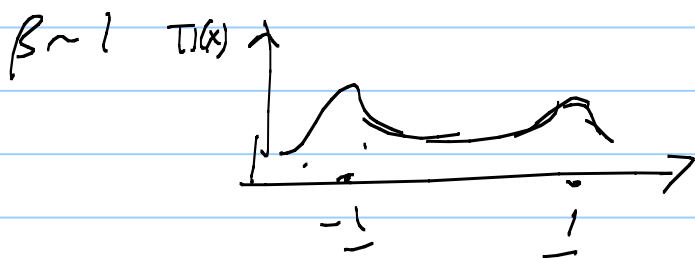
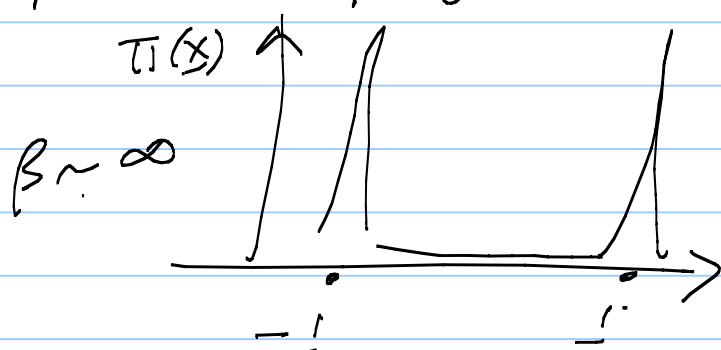


(1) Histograms similar for $\beta=1$, MCMC gives similar results to Exact. Similar computation time

Metropolis Example: (cont, cont)

(ii) Histograms differ for $\beta=2$, also histogram for $\beta=2$ differ greatly to those for $\beta=1$. What is going on?

The distribution for $\pi(x)$ gets sharpened as β gets larger



$\underline{1} = (1, 1, \dots, 1)$
 $-\underline{1} = (-1, -1, \dots, -1)$

Hard for our algorithm to move from one peak to the other

← Easier...

Note: The Metropolis sampler we used is better at sampling from $\pi(\underline{x})$ with $\beta=1$.

For $\pi(\underline{x})$ with $\beta \geq 2$, it takes a very long time for the algorithm to move from peak at $-\underline{1}$ to peak at $\underline{1}$.

Why do histograms differ from $\beta=1$ & $\beta=2$?
 Answer: Phase Transition.

$$Pr(M=z) = \sum_{x_0 \dots x_d} \delta_{x_0 + \dots + x_d, z} \Pi(x)$$

x
 any value

Kronecker delta
 (Indicator Function)

There is a phase factor, no. of different ways that you can get $M=z$ from \underline{x} .

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

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

$$Pr(M = d+1) = Pr(M = -(d+1)) = \Pi(\underline{1}) = \Pi(\underline{-1})$$

For $M = d$, impossible

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

For $M = 0$, an exponential large no. of ways (provided $d+1$ is even)

In short, small phase factor for $M = \pm d+1$
 enormous phase factor for $M = 0$.

Two Conflicting "forces"

Phase factor wants $M = 0$.

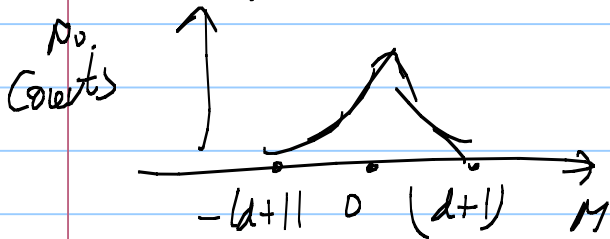
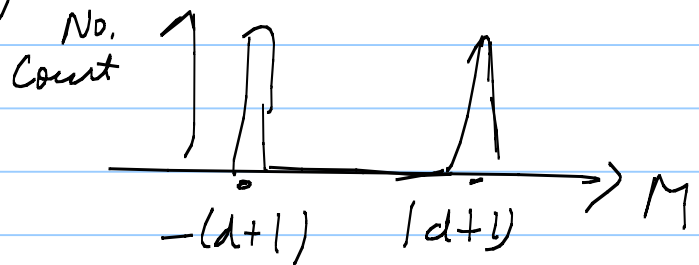
Probability wants $M = \pm (d+1)$.

For large β , probability biases strongly to ± 1 (i.e. $M = \pm (d+1)$).

For small β , probability biases weakly to ± 1 .

So for large β , probabilities win

For small β , phase factors win.



Comments, better MCMC can be designed for case with large β .

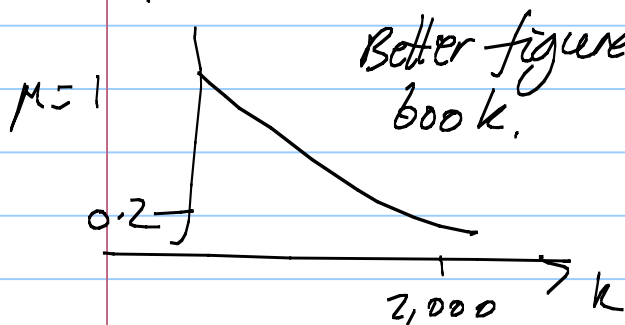
MCMC can be as effective as Exact Sampling. But we can use MCMC in cases where exact sampling is impossible

Burn-in and Stickyness

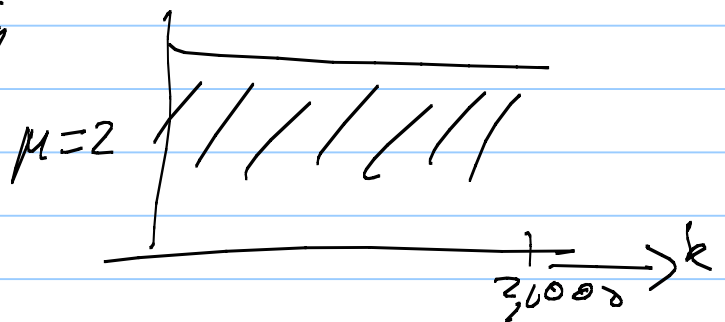
Burn-in is the amount of time for MCMC to produce samples from $\pi(x)$
Stickyness is a measure of how MCMC explores $\pi(x)$.

Practical Measure: Autocorrelation

$$\rho_k = \text{Correlation} (M^{(1)}, M^{(k+1)})$$



Better figures in book.



In short, this MCMC is sticky for $\mu=2$ since the autocorrelation decreases very slowly with time.

Samples after burn-in will be correlated. Must wait to get i.i.d samples

L-E. instead of samples $x^{50}, x^{100}, x^{150}, \dots$

need $x^{160}, x^{600}, x^{1100}, \dots$

Metropolis-Hastings

(*) To get Metropolis,
set $T(y|x)$ to be
uniform.

This improves the basic Metropolis algorithm by adding a proposal probability.

$$T(y|x) \text{ for } y \in N(x).$$

At time t : state \underline{x}^t

(i) propose $\underline{y} \in N(\underline{x})$ with
probability $T(\underline{y} | \underline{x}^t)$.

(ii) accept this proposal with
probability $r(\underline{x}^t, \underline{y}) = \min \left\{ 1, \frac{\pi(\underline{y}) T(\underline{x}^t | \underline{y})}{\pi(\underline{x}) T(\underline{y} | \underline{x}^t)} \right\}$

If accepted, set $\underline{x}^{t+1} = \underline{y}$.

Note: This is equivalent to setting

$$K(\underline{y} | \underline{x}) = T(\underline{y} | \underline{x}) \cdot \min \left\{ 1, \frac{\pi(\underline{y}) \cdot T(\underline{x} | \underline{y})}{\pi(\underline{x}) \cdot T(\underline{y} | \underline{x})} \right\}$$

Detailed Balance:

$$K(\underline{y} | \underline{x}) \pi(\underline{x}) = \min \left\{ \pi(\underline{x}) \pi(\underline{y} | \underline{x}), \pi(\underline{y}) \pi(\underline{x} | \underline{y}) \right\}$$

symmetric, so Detailed Balance holds

Good choice of proposals makes MCMC very fast.

Metropolis Hastings is the most popular way of designing an MCMC. (There are others, see later).

Difficulty is selecting the neighbourhood and the proposal probabilities.

E.g. Proposals for Ising model.

$$N(\underline{x}) = \{ (x_0, \dots, \bar{x}_k, \dots, x_d) ; k=0 \dots d \}$$

$$T(\underline{x}_{i,k} | \underline{x}) = \frac{e^{C_k}}{\sum_{i=0}^d e^{C_i}}, \text{ i.e. propose moves which involve flipping states where } C_k \text{ is largest.}$$

Designing a good MCMC is an art more than a science.

Study of the specific distribution $\pi(\underline{x})$, can suggest good proposal probabilities, good neighbourhood structures, etc.

