

Markov Chain Monte Carlo Sampling

Sargur Srihari
srihari@cedar.buffalo.edu

Topics

1. Limitations of simple Monte Carlo
2. Markov Chains
3. Basic Metropolis Algorithm
4. Metropolis-Hastings Algorithm

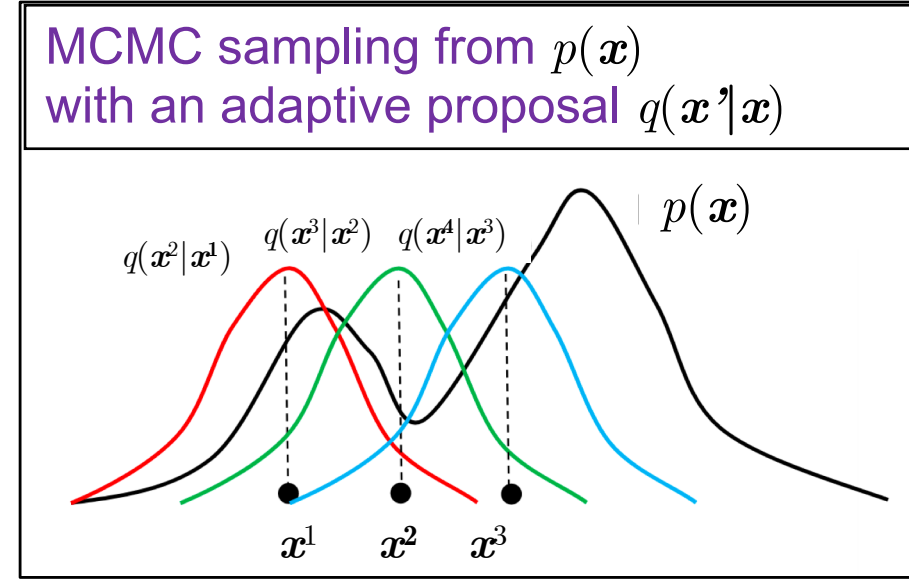
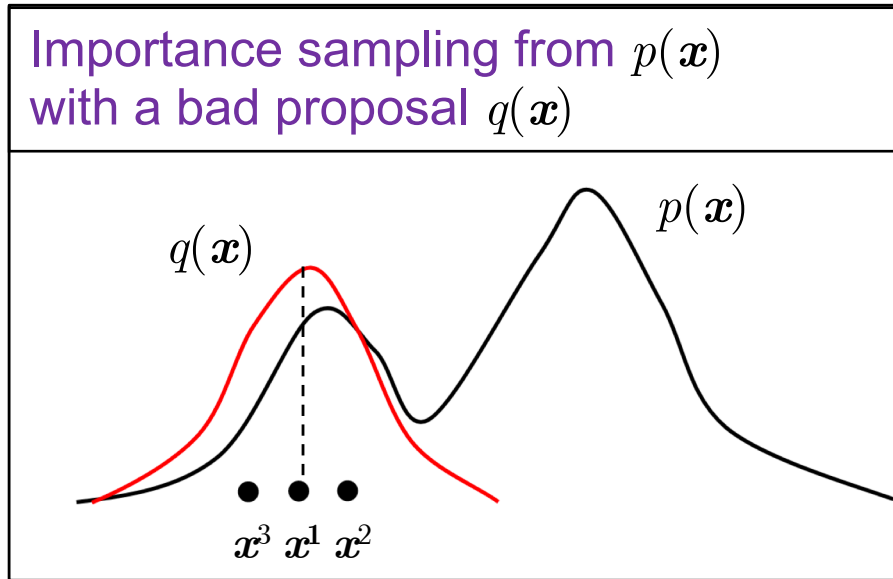
Limitations of plain Monte Carlo

- Direct (unconditional) sampling
 - Hard to get rare events in high-dimensional spaces
 - Infeasible for MRFs, unless we know partition function Z
- Rejection sampling, Importance sampling
 - Doesn't work well if proposal $q(\mathbf{x})$ is very different from $p(\mathbf{x})$
 - Yet constructing a $q(\mathbf{x})$ similar to $p(\mathbf{x})$ can be difficult
 - Making a good proposal usually requires knowledge of the analytic form of $p(\mathbf{x})$ – but if we had that, we wouldn't even need to sample!
- Intuition of MCMC
 - Instead of a fixed proposal $q(\mathbf{x})$, use an adaptive proposal

Markov Chain Monte Carlo (MCMC)

- Simple Monte Carlo methods (Rejection sampling and importance sampling) are for evaluating expectations of functions
 - They suffer from severe limitations, particularly with high dimensionality
- MCMC is a very general and powerful framework
 - Markov refers to sequence of samples rather than the model being Markovian
 - Allows sampling from large class of distributions
 - Scales well with dimensionality
 - MCMC origin is in statistical physics (Metropolis 1949)

MCMC and Adaptive Proposals



- MCMC:
 - Instead of $q(\mathbf{x})$, use $q(\mathbf{x}'|\mathbf{x})$ where \mathbf{x}' is the new state being sampled, and \mathbf{x} is previous sample
 - As \mathbf{x} changes, $q(\mathbf{x}'|\mathbf{x})$ changes as a function of \mathbf{x}'

Markov chains

- First order Markov chain is a sequence of random variables $z^{(1)}, \dots, z^{(M)}$ such that

- conditional independence property holds:

$$p(z^{(m+1)} | z^{(1)}, \dots, z^{(m)}) = p(z^{(m+1)} | z^{(m)})$$

Each sample dependent
only on previous sample

- Represented in a directed graph as a chain
- Markov chain specified by

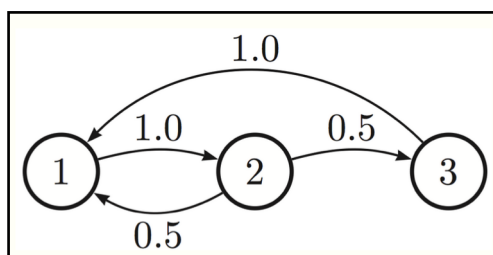
- Distribution of initial variable $p(z^{(0)})$
 - Conditional (transition) probabilities

$$T_m(z^{(m)}, z^{(m+1)}) = p(z^{(m+1)} | z^{(m)})$$

- Markov chain is homogeneous if all transition probabilities are the same for all m

Markov Chain

- A sequence of random variables S_0, S_1, S_2, \dots with each $S_i \in \{1, 2, \dots, d\}$ taking one of d possible values representing state of a system
 - Initial state distributed according to $p(S_0)$
 - Subsequent states generated from a conditional distribution that depends only on the previous state
 - i.e. S_i is distributed according to $p(S_i/S_{i-1})$



A Markov chain over three states.
The weighted directed edges indicate probabilities of transitioning to a different state

Idea of MCMC

- Idea of MCMC is
 - Construct a Markov chain whose states will be joint assignments to the variables in the model
 - And whose stationary distribution will equal the model probability p

Metropolis-Hastings

- User to specify a transition kernel $q(x'|x)$ and acceptance probability $A(x'|x)$
- M-H Algorithm
 - Draw sample x' from $q(x'|x)$, where x is previous sample
 - New sample x' accepted/rejected with probability $A(x'|x)$
 - This acceptance probability is

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

- It encourages us to move towards more likely points in the distribution

Acceptance probability

- $A(x'|x)$ is like ratio of importance sampling weights

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

$p(x') / q(x'|x)$ is the importance weight for x'

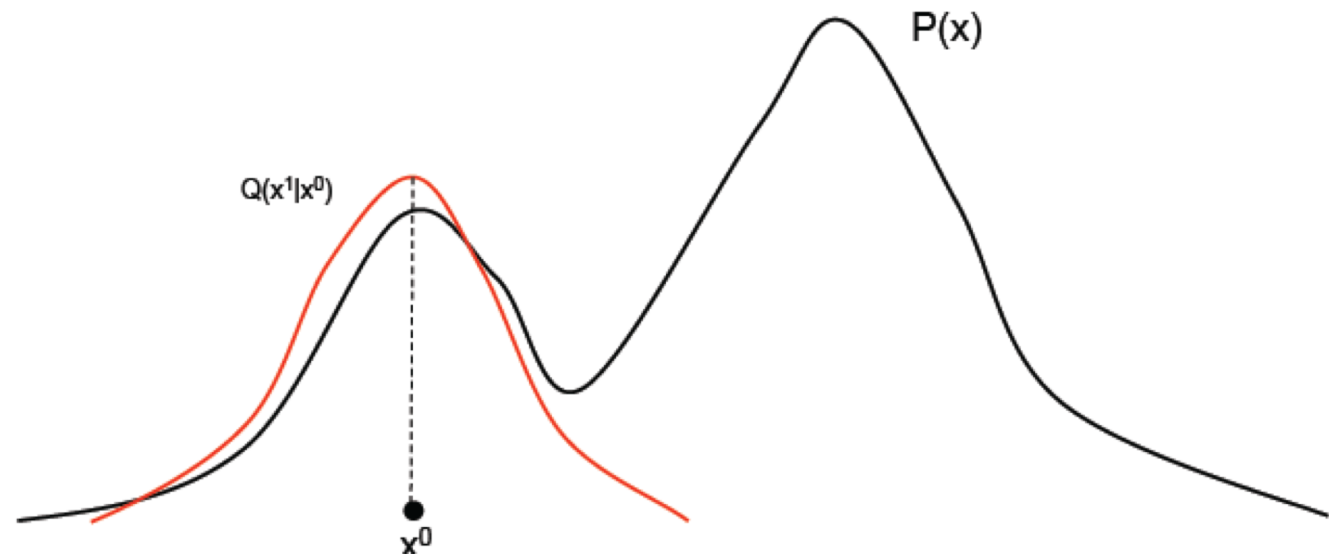
$p(x) / q(x|x')$ is the importance weight for x

- We divide the importance weight for x' by that of x
- Notice that we only need to compute $p(x') / p(x)$ rather than $p(x')$ or $p(x)$ separately
- $A(x'|x)$ ensures that, after sufficient draws, samples will come from true distribution $p(x)$

The M-H Algorithm (1)

- Let $q(x'|x)$ be a Gaussian centered on x
- We are trying to sample from bimodal $p(x)$

Initialize $x^{(0)}$

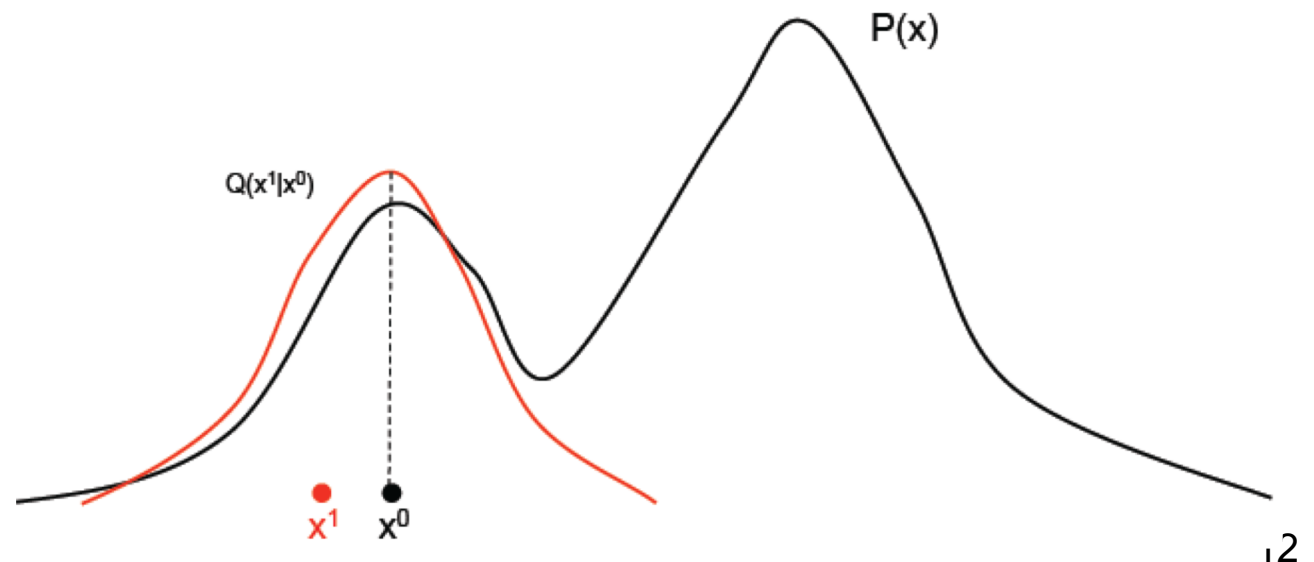


The M-H Algorithm (2)

- Let $q(x'|x)$ be a Gaussian centered on x
- We are trying to sample from bimodal $p(x)$

Initialize $x^{(0)}$

Draw, accept $x^{(1)}$



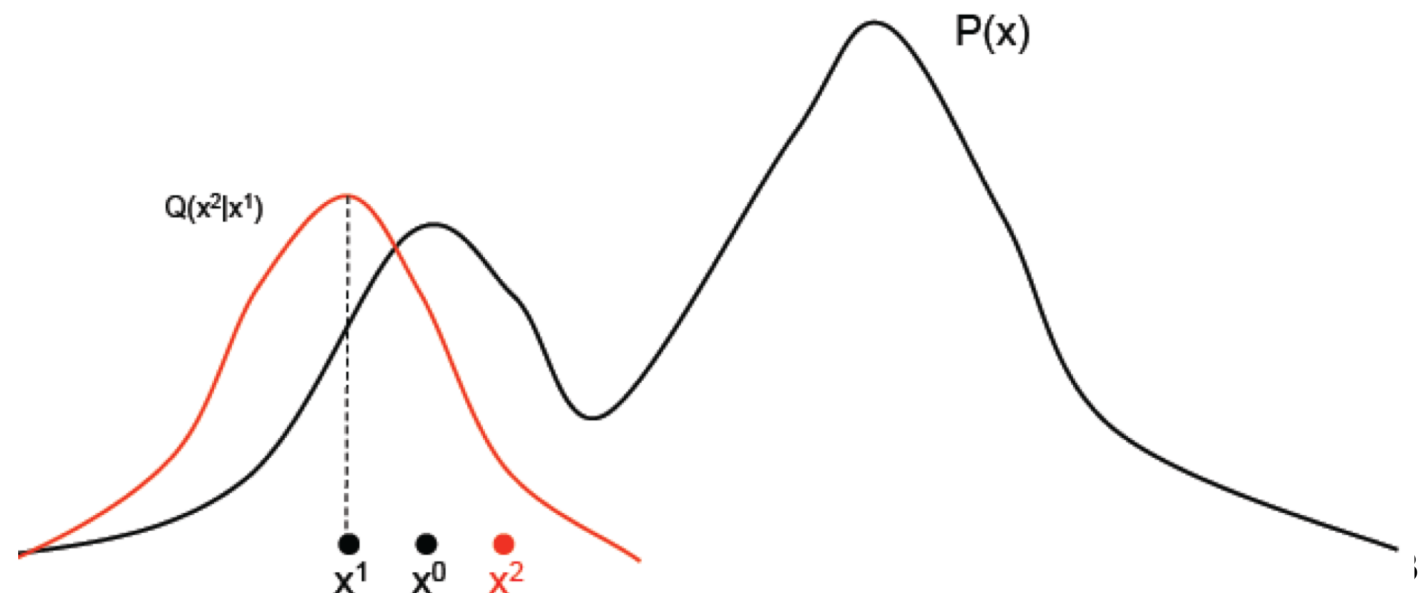
The M-H Algorithm (3)

- Let $q(x'|x)$ be a Gaussian centered on x
- We are trying to sample from bimodal $p(x)$

Initialize $x^{(0)}$

Draw, accept $x^{(1)}$

Draw, accept $x^{(2)}$



The M-H Algorithm (4)

- Let $q(x'|x)$ be a Gaussian centered on x
- We are trying to sample from bimodal $p(x)$

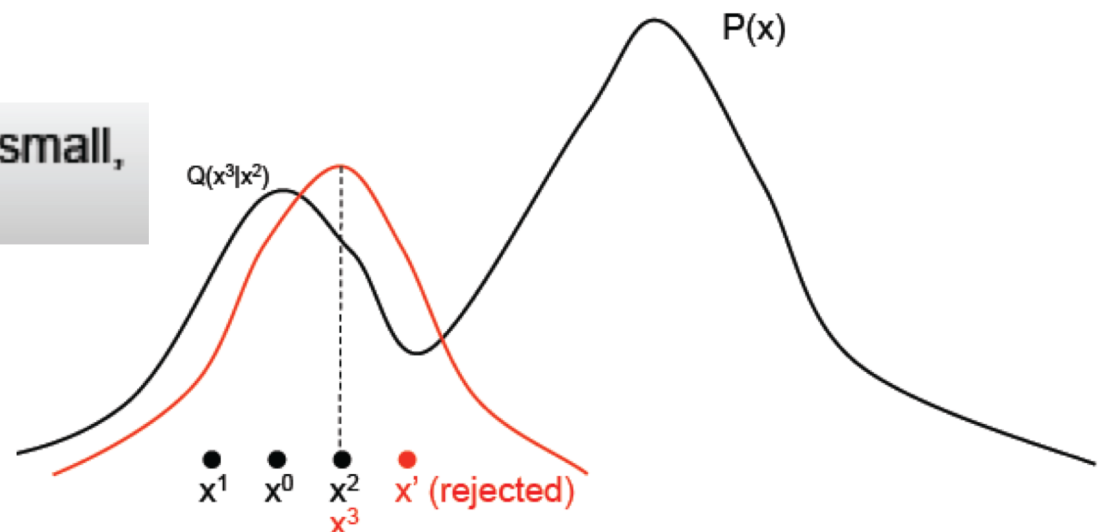
Initialize $x^{(0)}$

Draw, accept $x^{(1)}$

Draw, accept $x^{(2)}$

Draw, but reject, set $x^{(3)} = x^{(2)}$

We reject because $P(x')/P(x^2)$ is very small,
hence $A(x'|x^2)$ is close to zero!



The M-H Algorithm (5)

- Let $q(x'|x)$ be a Gaussian centered on x
- We are trying to sample from bimodal $p(x)$

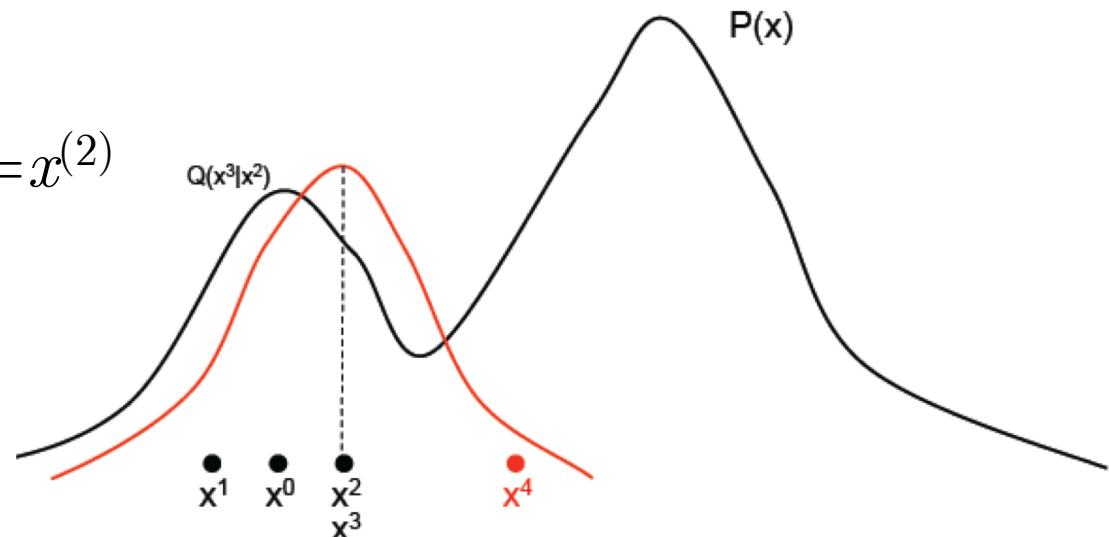
Initialize $x^{(0)}$

Draw, accept $x^{(1)}$

Draw, accept $x^{(2)}$

Draw, but reject, set $x^{(3)} = x^{(2)}$

Draw, accept $x^{(4)}$



The M-H Algorithm (6)

- Let $q(x'|x)$ be a Gaussian centered on x
- We are trying to sample from bimodal $p(x)$

Initialize $x^{(0)}$

Draw, accept $x^{(1)}$

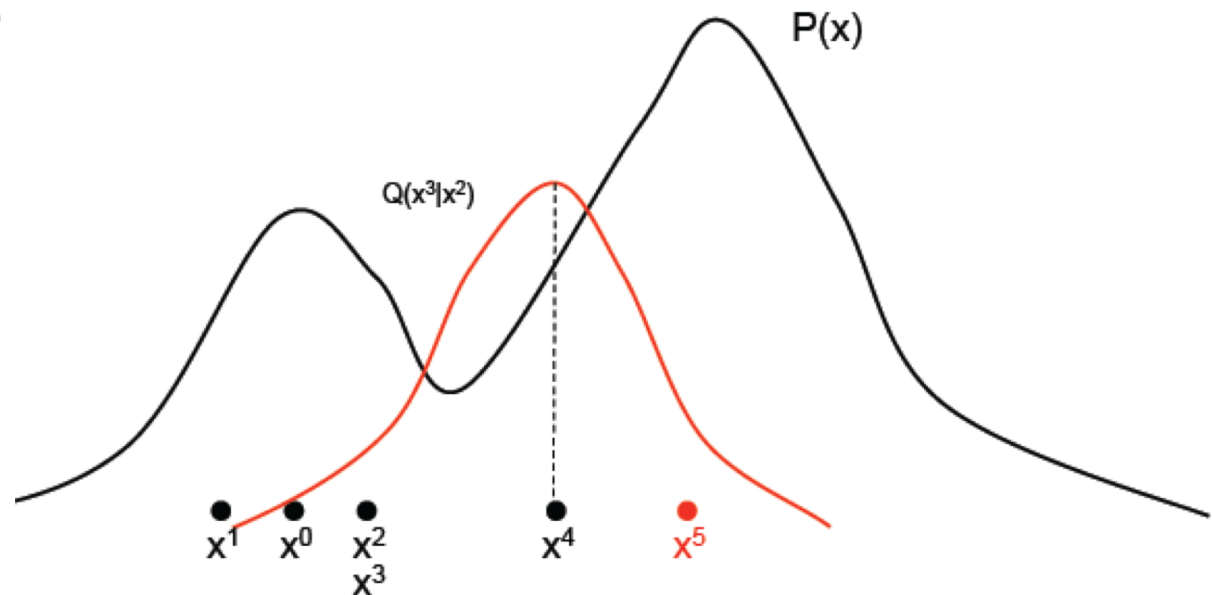
Draw, accept $x^{(2)}$

Draw, but reject, set $x^{(3)} = x^{(2)}$

Draw, accept $x^{(4)}$

Draw, accept $x^{(5)}$

The adaptive proposal $Q(x'|x)$ allows us to sample both modes of $P(x)$!



Basic Metropolis Algorithm

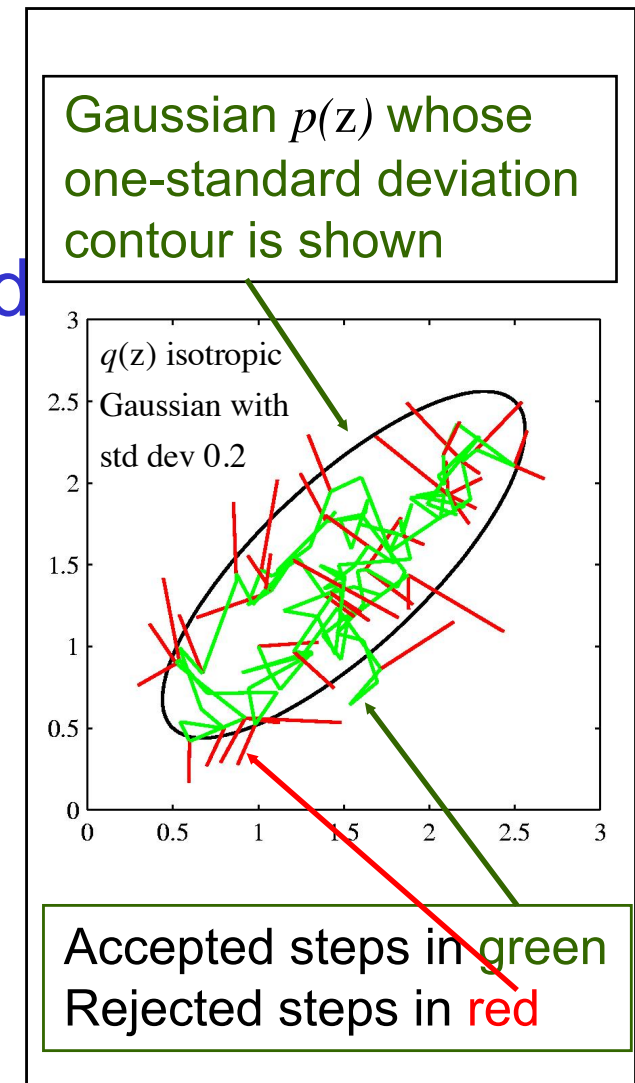
- As with rejection and importance sampling use a *proposal* distribution (simpler distribution)
- Maintain a record of current state $z^{(t)}$
- Proposal distribution $q(z|z^{(t)})$ depends on current state (next sample depends on previous one)
 - E.g., $q(z|z^{(t)})$ is a symmetric Gaussian with mean $z^{(t)}$ and a small variance
- Thus sequence of samples $z^{(1)}, z^{(2)} \dots$ forms a Markov chain
- Write $p(z) = \frac{1}{Z_p} \tilde{p}(z)$ where $\tilde{p}(z)$ is readily evaluated
- At each cycle generate candidate z^* and test for acceptance

Metropolis Algorithm

- Assumes simple proposal distribution, that is symmetric
 - e.g., an isotropic Gaussian $q(z)$,
 $q(z_A|z_B) = q(z_B|z_A)$ for all z_A, z_B
- New sample z^* from $q(z)$ is accepted with probability

$$A(z^*, z^{(\tau)}) = \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^{(\tau)})}\right)$$

- Done by choosing $u \sim U(0,1)$ and accepting if $A(z^*, z^{(t)}) > u$
- If accepted then $z^{(t+1)} = z^*$
- Otherwise:
 - z^* is discarded,
 - $z^{(t+1)}$ is set to $z^{(t)}$ and
 - another candidate drawn from $q(z|z^{(t+1)})$



Inefficiency of Random Walk

- Consider simple random walk
- State space z consisting of integers with probabilities

$$p(z^{(\tau+1)} = z^{(\tau)}) = 0.5$$

Stay in same state

$$p(z^{(\tau+1)} = z^{(\tau)} + 1) = 0.25$$

Increase state by 1

$$p(z^{(\tau+1)} = z^{(\tau)} - 1) = 0.25$$

Decrease state by 1

- If initial state is $z^{(1)}=0$
 - the expected state at time τ will be zero, $E[z^{(\tau)}] = 0$
 - Similarly $E[(z^{(\tau)})^2] = \tau / 2$
 - Thus after τ steps distance traveled is proportional to $\sqrt{\tau}$
- Random walks are inefficient in exploring state-space
 - MCMC algorithms try to avoid random walk behavior

Metropolis-Hastings Algorithm

- Generalizes Metropolis algorithm
 - Proposal distribution is no longer a symmetric function of arguments, i.e., $q(z_A|z_B) \neq q(z_B|z_A)$ for all z_A, z_B

1. At step t , in which current state is $z^{(t)}$ we draw a sample z^* from distribution $q_k(z|z^{(t)})$

2. Accept with probability $A_k(z^*, z^{(t)})$ where

$$A_k(z^*, z^{(\tau)}) = \min \left(1, \frac{\tilde{p}(z^*) q_k(z^{(\tau)} | z^*)}{\tilde{p}(z^{(\tau)}) q_k(z^* | z^{(\tau)})} \right)$$

k labels set of possible transitions considered

- Can show that $p(z)$ is an invariant distribution of the Markov chain defined by Metropolis-Hastings Algorithm
 - since detailed balance is satisfied

Choice of Proposal Distribution for Metropolis-Hastings

- Gaussian centered on current state
- Keeping rejection rate low
Scale ρ of proposal distribution should be of order σ_{\min}
- Independent sample
No. of steps needed to get independent sample is of order $(\sigma_{\max} / \sigma_{\min})$

