

Markov Chain Monte Carlo Methods

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Monte Carlo Methods

1. Sampling and Monte Carlo Methods
2. Importance Sampling
3. Markov Chain Monte Carlo Methods
4. Gibbs Sampling
5. Mixing between separated modes

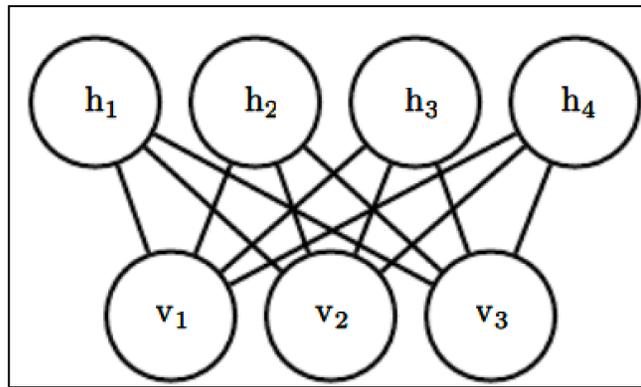
Topics in Markov Chain Monte Carlo

- Limitations of plain Monte Carlo methods
- Markov Chains
- MCMC and Energy-based models
- Metropolis-Hastings Algorithm
- Theoretical basis of MCMC

MCMC sampling from deep models

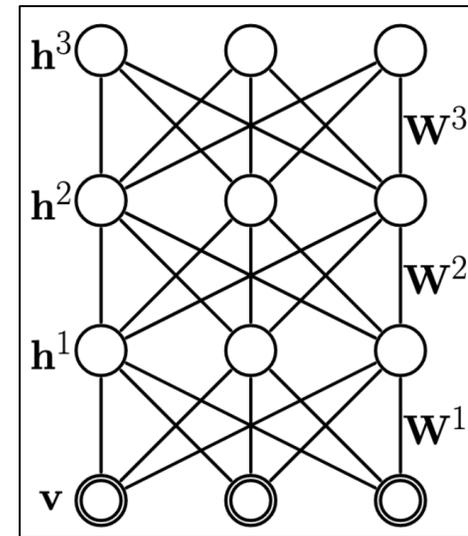
- Energy-based models and samples:

A Restricted Boltzmann machine



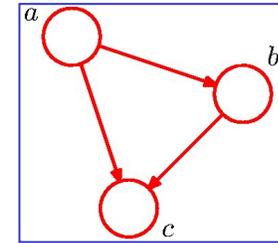
$$\mathbf{x} = (\mathbf{v}, \mathbf{h}), p_{\text{model}}(\mathbf{x})$$

A deep Boltzmann machine

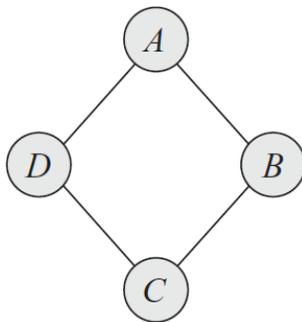


BN vs MN sampling

- For BNs: sample from $p(a,b,c)$
 - Start with lowest numbered node
 - Draw a sample from $p(x_1)$ which we call \hat{x}_1
 - Work through nodes in order
 - For node n we draw a sample from conditional distribution $p(x_n|pa_n)$
 - Defines a single pass algorithm to get sample from $p(\mathbf{x})$
- For MNs: chicken-egg problem to sample from $p(A,B,C,D)$
 - In order to sample from A we need to draw from $p(A|B,D)$
 - In order to sample from B we need to sample from $p(B|A,C)$



$$p(\mathbf{x}) = \prod_i p(x_i | Pa_G(x_i))$$



$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x}) \quad \tilde{p}(\mathbf{x}) = \prod_{C \in G} \phi(C) \quad Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

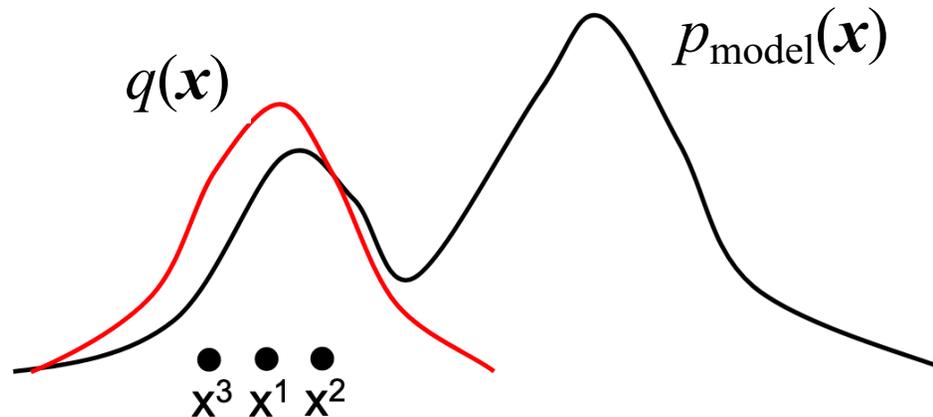
We avoid chicken-egg problem using a Markov chain

Difficulties in plain Monte Carlo

1. Direct (unconditional) sampling $p(\mathbf{x})$
 - Hard to get rare events in high-dimensional spaces
 - Infeasible for MRFs, unless we know partition function Z
2. Importance sampling
 - Doesn't work well if proposal $q(\mathbf{x})$ is very different from $p(\mathbf{x})$
 - Constructing a $q(\mathbf{x})$ similar to $p(\mathbf{x})$ can be difficult
 - Good proposal requires knowledge of 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

Shortcoming of Direct Sampling

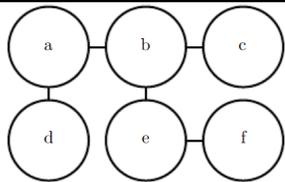
- No tractable method for either
 1. Sampling directly from $p_{\text{model}}(\mathbf{x})$
 - Chicken-and-egg problem described on next slide
 2. Importance sampling with low variance proposal distribution $q(\mathbf{x})$



Energy-Based Model

- MCMC requires no zero probability to any state
 - Guaranteed by an energy-based model (EBM)

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x}))$$



Graph, with five cliques, implies that

$$E(a,b,c,d,e,f) = E_{a,b}(a,b) + E_{b,c}(b,c) + E_{a,d}(a,d) + E_{b,e}(b,e) + E_{e,f}(e,f)$$

By setting, $\phi_{a,b}(a,b) = \exp(-E(a,b))$ etc we get

$$P(a,b,c,d,e,f) = (1/Z) \exp[-(E(a,b) + E(b,c) + E(a,d) + E(b,e) + E(e,f))]$$

$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x})$$

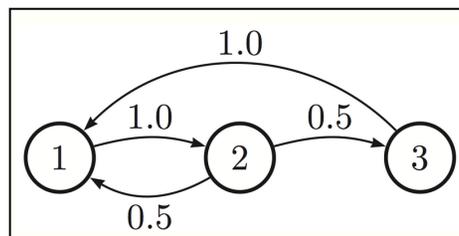
$$\hat{p}(\mathbf{x}) = \prod_{C \in \mathcal{G}} \phi(C)$$

$$Z = \int \hat{p}(\mathbf{x}) d\mathbf{x}$$

$$\hat{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$$

Markov Chain States

- 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

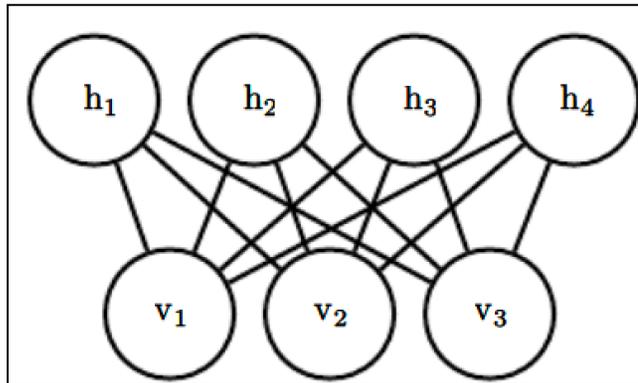
Transition Distribution

- Core idea of Markov chain
 - Have a state \mathbf{x} that begins with an arbitrary value
 - Over time we repeatedly update \mathbf{x}
 - Eventually \mathbf{x} becomes a fair sample from $p(\mathbf{x})$
- Markov chain is defined by random state \mathbf{x} and *transition distribution* $T(\mathbf{x}' | \mathbf{x})$
 - Specifying probability that a random update will go to state \mathbf{x}' if it starts in state \mathbf{x}
- Running the Markov chain means repeatedly updating the state \mathbf{x} to a value \mathbf{x}' sampled from $T(\mathbf{x}' | \mathbf{x})$

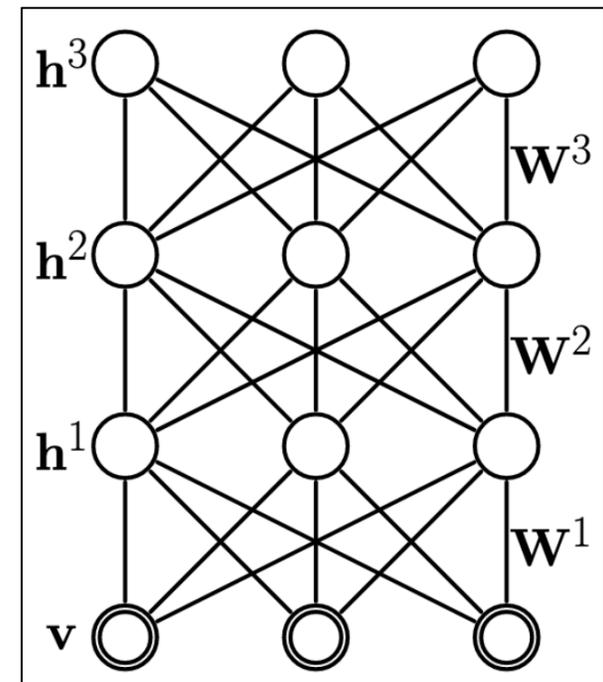
Undirected models in deep learning

- $p_{\text{model}}(\mathbf{x})$ is an undirected model

A Restricted Boltzmann machine



A deep Boltzmann machine



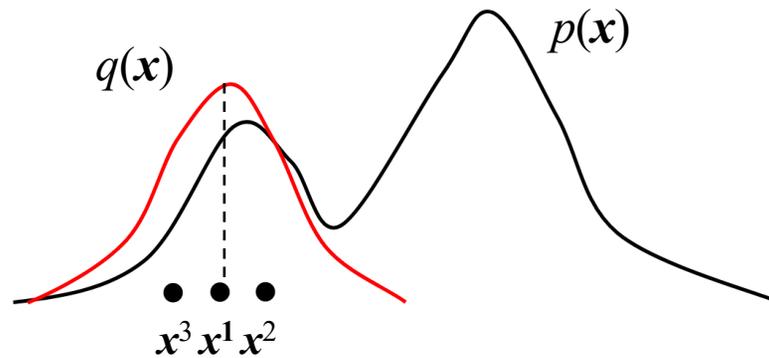
State: an instance of values assigned to all variables

Our goal is to obtain samples from the trained model

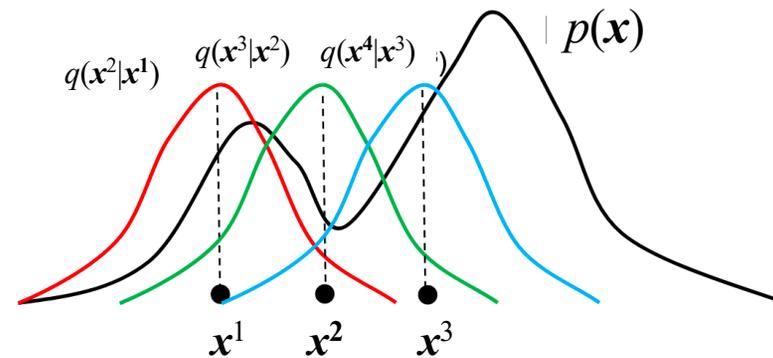
Assume that for any $\mathbf{x}=(\mathbf{v},\mathbf{h})$ we are able to compute $p(\mathbf{x})$

Importance Sampling and MCMC

Importance sampling from $p(\mathbf{x})$
with a bad proposal $q(\mathbf{x})$



MCMC sampling from $p(\mathbf{x})$
with an adaptive proposal $q(\mathbf{x}'|\mathbf{x})$



- **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}'
- The adaptive proposal $q(\mathbf{x}'|\mathbf{x})$ allows us to sample both modes of $p(\mathbf{x})$

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 $q(\mathbf{x}'|\mathbf{x})$ and acceptance probability $A(\mathbf{x}'|\mathbf{x})$
- M-H Algorithm
 - Draw sample \mathbf{x}' from $q(\mathbf{x}'|\mathbf{x})$, where \mathbf{x} is previous sample
 - New sample \mathbf{x}' accepted/rejected with probability $A(\mathbf{x}'|\mathbf{x})$
 - This acceptance probability is

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

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

Acceptance probability

- $A(\mathbf{x}'|\mathbf{x})$ is like ratio of importance sampling weights

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

$p(\mathbf{x}')/q(\mathbf{x}'|\mathbf{x})$ is the importance weight for \mathbf{x}'

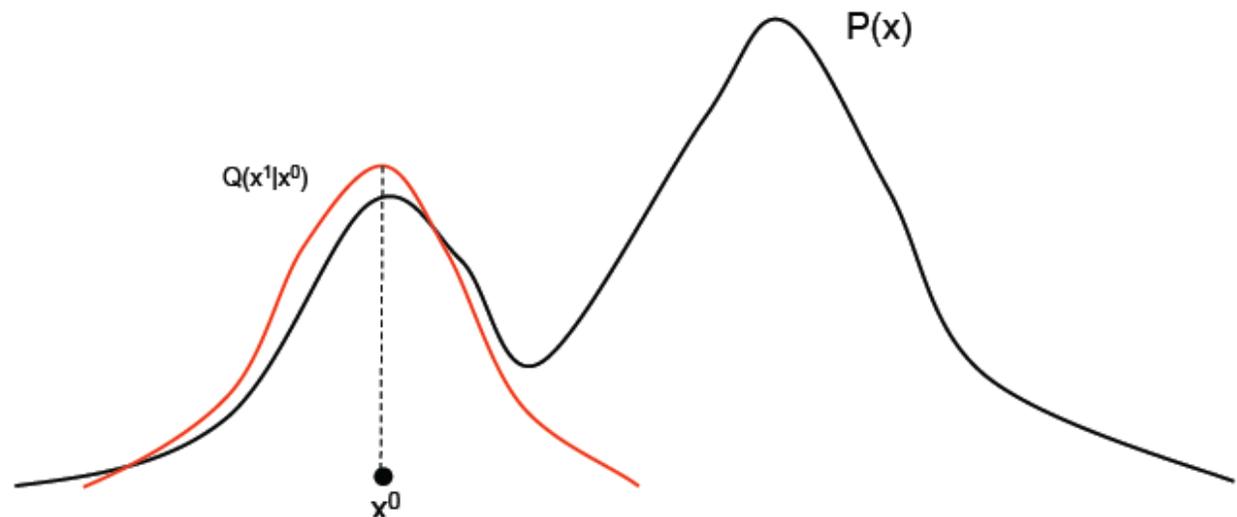
$p(\mathbf{x})/q(\mathbf{x}|\mathbf{x}')$ is the importance weight for \mathbf{x}

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

The M-H Algorithm (1)

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

Initialize $x^{(0)}$

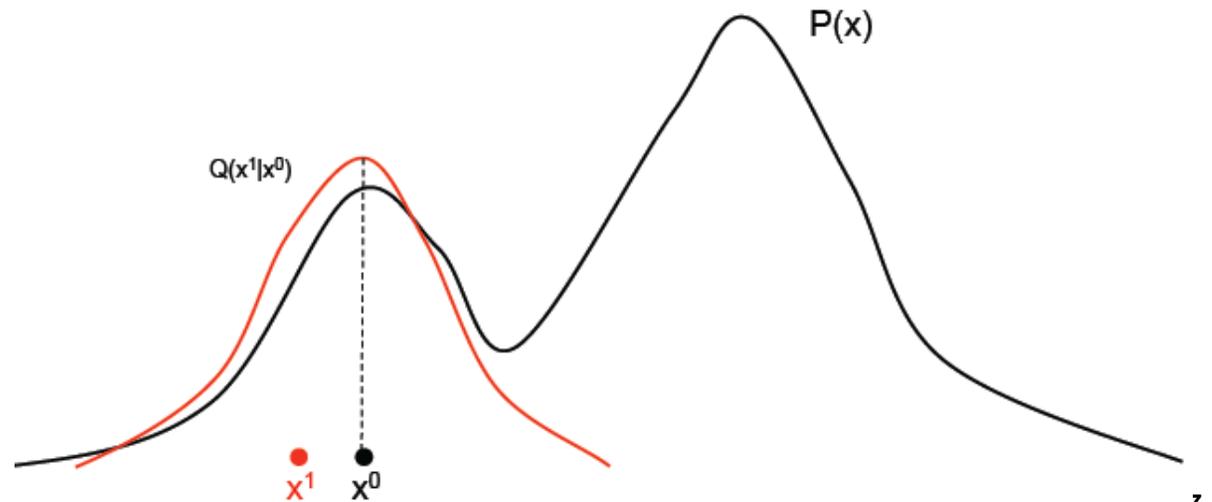


The M-H Algorithm (2)

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

Initialize $x^{(0)}$

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



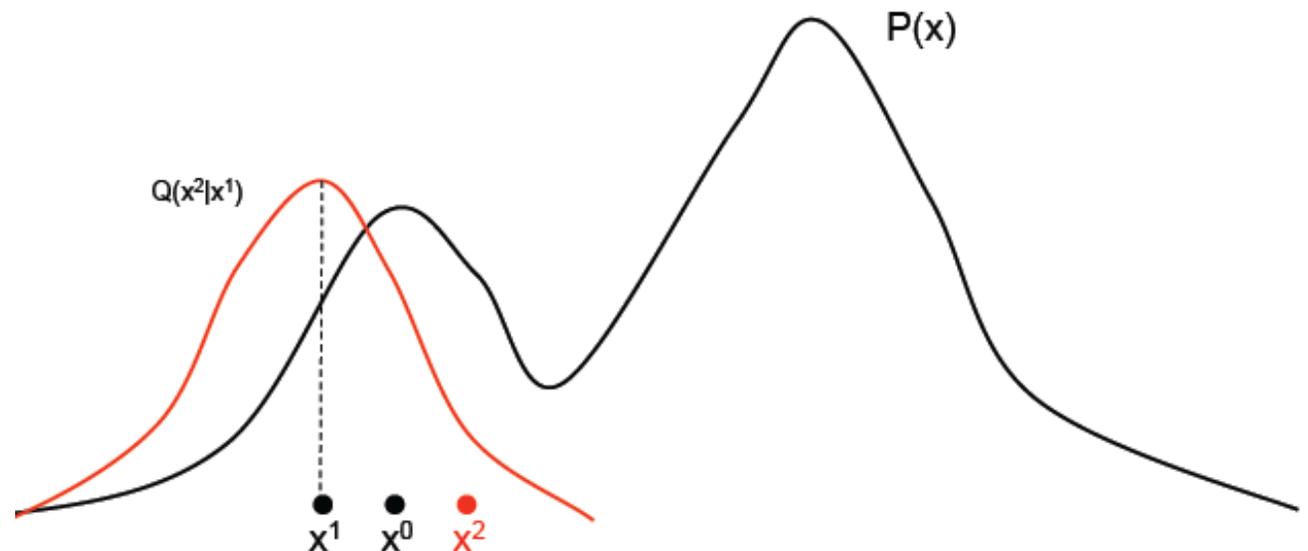
The M-H Algorithm (3)

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

Initialize $x^{(0)}$

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

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



The M-H Algorithm (4)

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

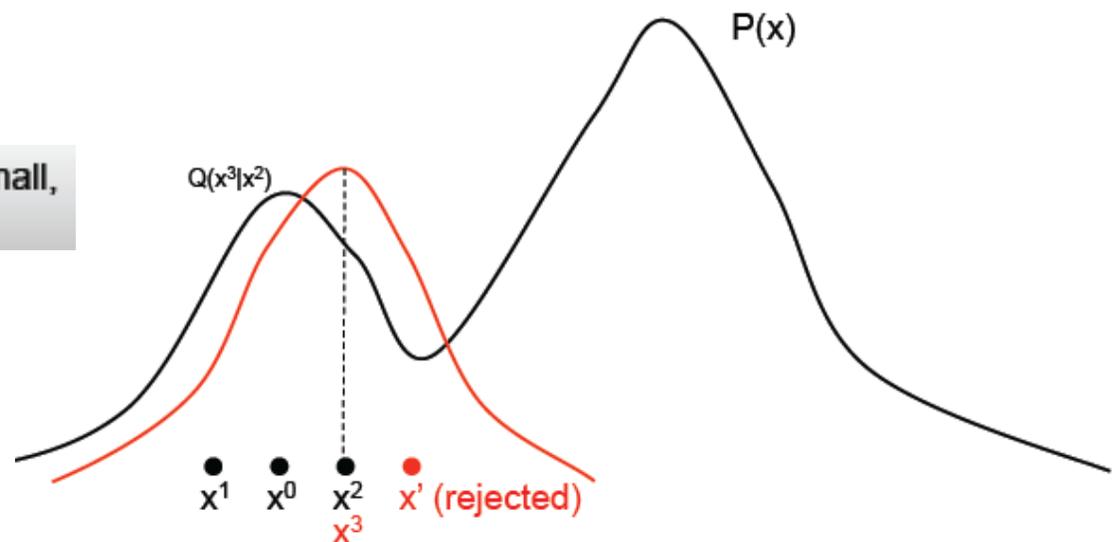
Initialize $x^{(0)}$

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

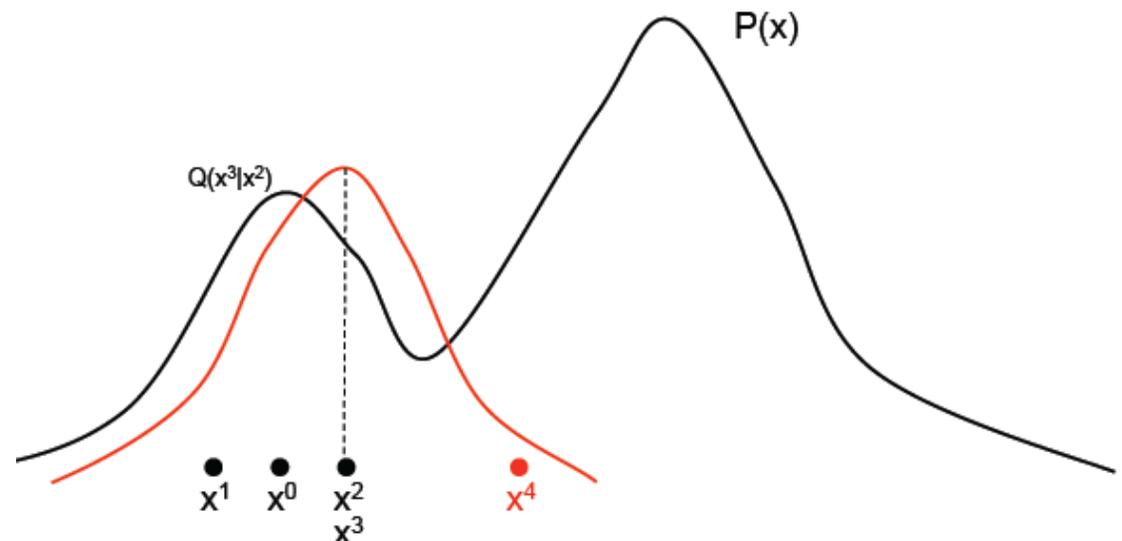
Initialize $x^{(0)}$

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

Initialize $x^{(0)}$

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

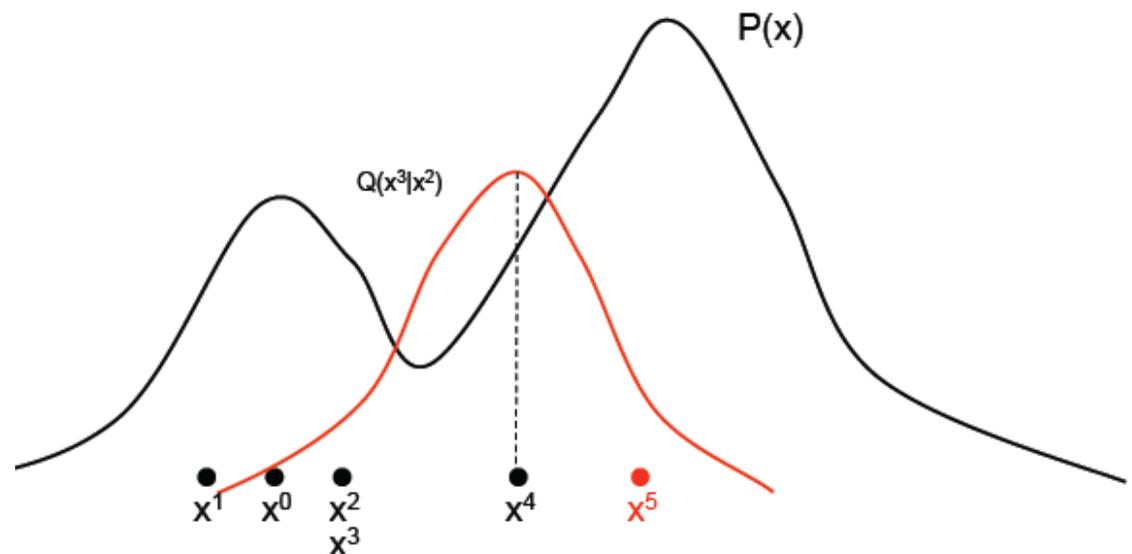
Draw, accept $x^{(2)}$

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

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

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

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



Theoretical understanding of MCMC

- Reparameterize the problem
 - Assume \mathbf{x} has countably many states
 - Represent state as an integer x
 - Different integer values of x map back to different states x in the original problem
 - We run infinitely many Markov chains in parallel
 - All states of Markov chains are drawn from some $q^{(t)}(x)$
 - where t is no of time steps
 - $q^{(0)}$ is some distribution that we initialize
- Our goal is for $q^{(t)}(x)$ to converge to $p(x)$
 - Because we have reparameterized in terms of positive integer \mathbf{x} , we can describe q using a vector \mathbf{v} with $q(x=i)=v_i$

Updating a Markov chain state

- When we update a single Markov chain's state \mathbf{x} to a new state \mathbf{x}'

- The probability of a single state landing in \mathbf{x}' is

$$q^{(t+1)}(\mathbf{x}') = \sum_{\mathbf{x}} q^{(t)}(\mathbf{x}) T(\mathbf{x}' | \mathbf{x})$$

- Using our integer parametrization, we represent the effect of the transition operator T using a matrix A

- We define A so that

$$A_{i,j} = T(\mathbf{x}' = i | \mathbf{x} = j)$$

- Using this definition, we can rewrite above equation for $q^{(t+1)}(\mathbf{x}')$

Equation state update

- Instead of q and T , we use \mathbf{v} and A
 - To describe how the entire distribution over all Markov chains (running in parallel) shifts as we apply an update: $\mathbf{v}^{(t)} = A\mathbf{v}^{(t-1)}$
 - Applying repeated Markov chain updates corresponds to multiplying by A repeatedly
 - i.e., exponentiating the matrix A :

$$\mathbf{v}^{(t)} = A^t \mathbf{v}^{(0)}$$

- Matrix A has special structure because each of its columns represents a distribution
- Such matrices are called stochastic matrices

Stationary Distribution

- If there is a non-zero probability of transitioning from state \mathbf{x} to state \mathbf{x}' for some power t , then the Perron-Frobenius theorem guarantees that the largest eigenvalue is real and equal to 1
- Over time, all eigenvalues are exponentiated:
$$\mathbf{v}^{(t)} = (V \text{diag}(\lambda) V^{-1})^t \mathbf{v}^{(0)} = V \text{diag}(\lambda)^t V^{-1} \mathbf{v}^{(0)}$$
 - All eigenvalues not equal to 1 to decay to zero
 - Under mild conditions, A is guaranteed to have only one eigenvector with eigenvalue 1
- The process thus converges to a stationary distribution

Equilibrium Distribution

- The *stationary distribution*, also called the *equilibrium distribution*, is given by an eigenvector equation $\mathbf{v}' = A\mathbf{v} = \mathbf{v}$ and this same condition holds for every additional step.
 - This is an eigenvector equation
 - To be a stationary point, \mathbf{v} must be an eigenvector with eigenvalue 1
 - Condition guarantees that once we reach the stationary distribution, repeated applications of transition sampling procedure do not change the distribution over the states of all the various Markov chains (although transition operator does change each individual state, of course).

Continuous case: Harris Chain

- Markov Chains with countable states can be generalized to continuous variables
 - Also called a Harris chain
- In general, a Markov chain with transition operator T will converge, under mild conditions, to a fixed point described by the equation

$$q'(\mathbf{x}') = E_{\mathbf{x} \sim q} T(\mathbf{x}' | \mathbf{x})$$

- which in the discrete case is just rewriting equation $\mathbf{v}' = A\mathbf{v} = \mathbf{v}$

Choice of Transition Distribution

- If we have chosen T correctly then the stationary distribution q will be equal to the distribution p we wish to sample from
- We describe how to choose T next, when we discuss Gibbs sampling