

Sampling in Machine Learning

Sargur Srihari

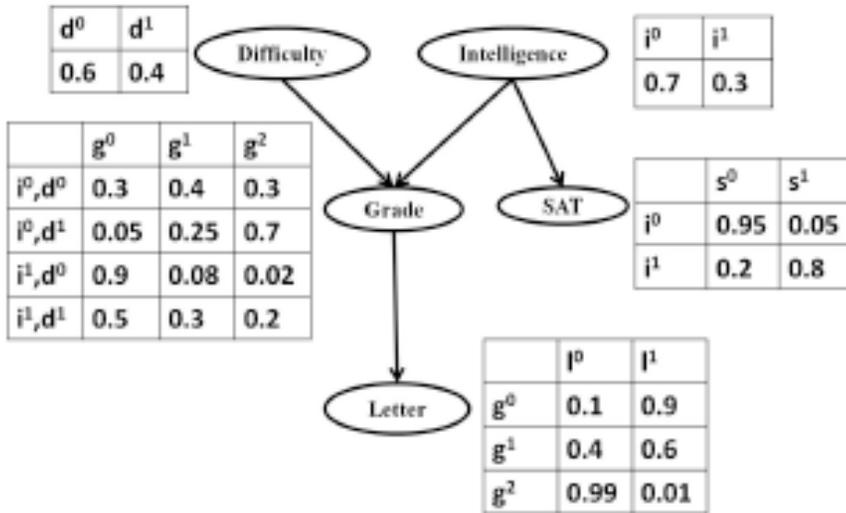
srihari@cedar.buffalo.edu

Topics

- Inference in Directed and Undirected models
- Variational vs Sampling methods
- Uses of Sampling Methods
- Sampling from a probability distribution
 - Rejection Sampling
 - Importance Sampling
 - MCMC Sampling
 - Metropolis-Hastings Sampling
 - Gibbs Sampling

Probability distributions as Graphical Models

Bayesian Networks

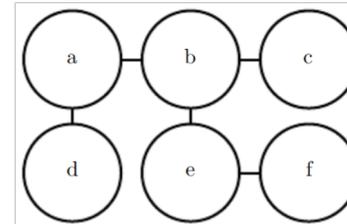


- CPDs: $p(x_i | pa(x_i))$
- Joint Distribution $\mathbf{x} = \{x_1, \dots, x_n\}$

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i | pa(x_i))$$

$$P(D, I, G, S, L) = P(D)P(I)P(G | D, I)P(S | I)P(L | G)$$

Markov Networks



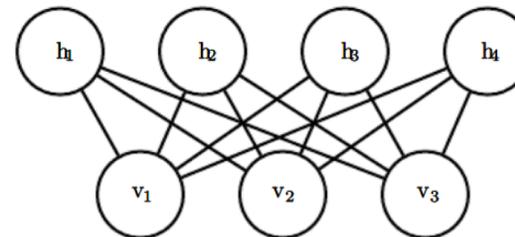
$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x}) \quad \tilde{p}(\mathbf{x}) = \prod_{C \in \mathcal{G}} \phi(C)$$

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

$$p(a, b, c, d, e, f) = \frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$$

Energy model $E(a, b, c, d, e, f) = E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$
 $\phi_{a,b}(a, b) = \exp(-E(a, b))$

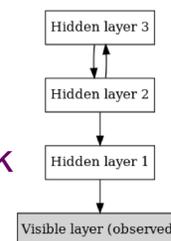
Restricted Boltzmann machine



$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$

$$p(\mathbf{h} | \mathbf{v}) = \prod_i p(h_i | \mathbf{v}) \quad \text{and} \quad p(\mathbf{v} | \mathbf{h}) = \prod_i p(v_i | \mathbf{h})$$

Deep Belief Network



Inference in a Bayesian Network

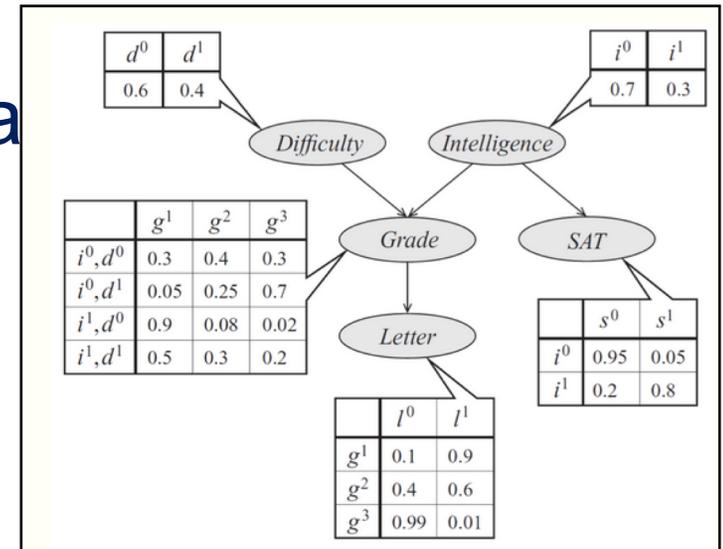
- Knowing *Letter* ($l=l^0$) we have a need to infer *Intelligence* (i)
 - Can be done using Bayes' rule

$$p(i/l) = p(i, l) / p(l)$$

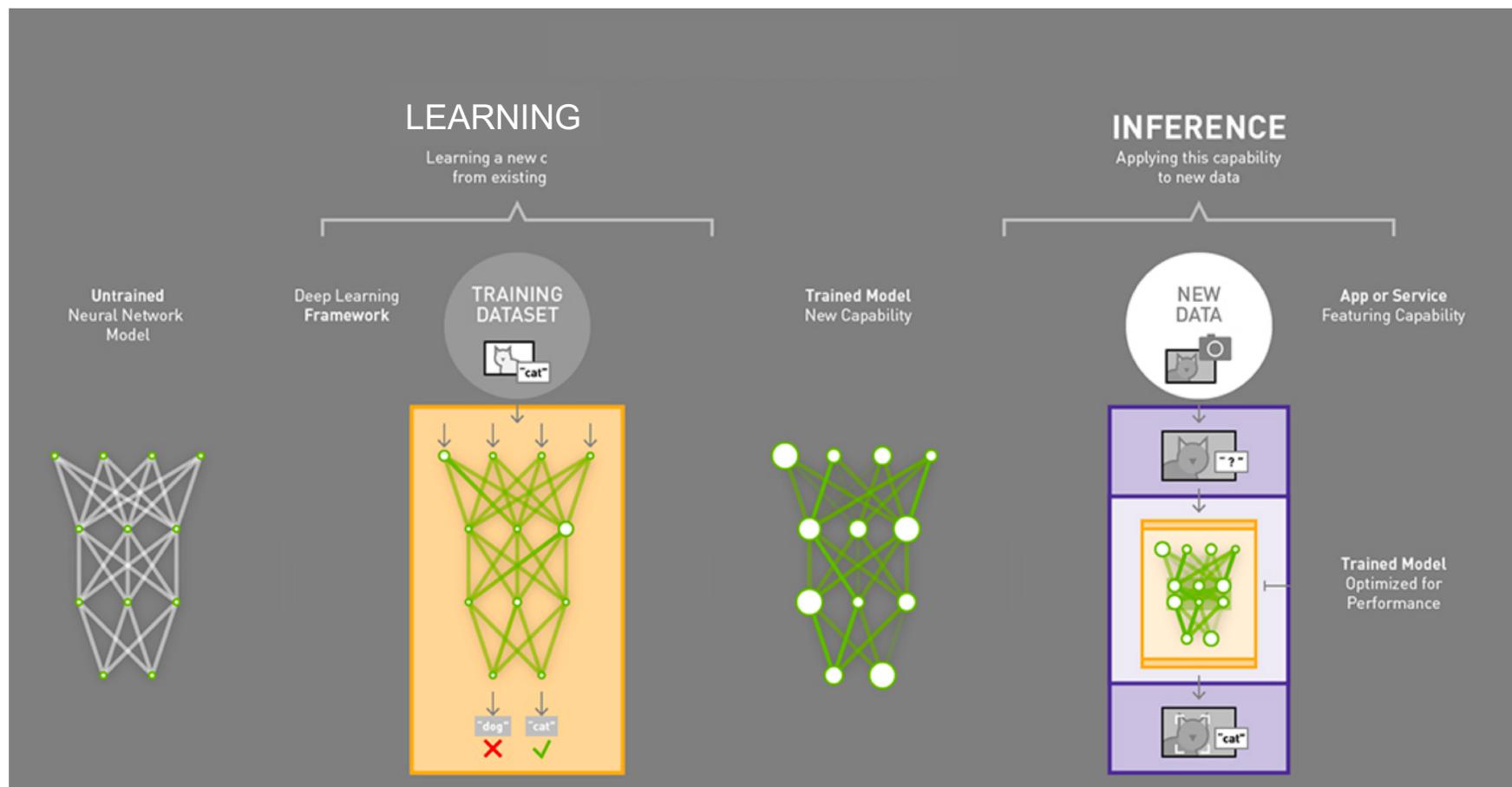
- Marginal $p(l)$ can be obtained using sum rule

$$p(l) = \sum_g \sum_s \sum_i \sum_d p(l | g) p(s | i) p(i) p(g | i, d) p(d)$$

- With n variables and k values each:
 - inference is $O(k^n)$
 - An intractable problem



Learning vs Inference



Parameter Learning:

For MRFs (including neural nets):
 performed using gradient descent to maximize log-likelihood
 For Bayes Nets: closed form solutions exist

Inference is performed
 using variational methods or sampling

Need for Approximate Inference

- Probabilistic models are quite complex
 - Simple algorithms (variable elimination) too slow
- Intractability of inference in a Bayesian network
 - Due to NP-hardness of inference
- Intractability of inference in a Neural network
 - In a generative model, need the partition function
- No exact polynomial time solutions
 - Thus much research effort spent on developing approximate solutions to inference problems

Variational vs Sampling methods

- Two families of approx. inference algorithms

1. Variational methods

- Formulate inference as an optimization problem
- Name comes from variational calculus:
 - Optimizing functions taking functions as arguments

2. Sampling methods

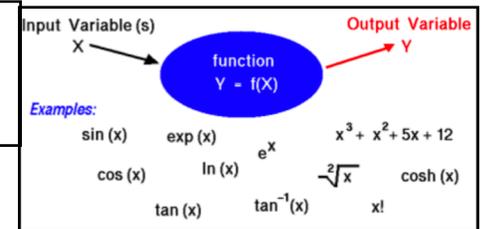
- Produce answers by repeatedly generating random numbers from a distribution of interest
- Sampling has been main method for inference
 - But variational methods emerged as often superior

Uses of Sampling methods

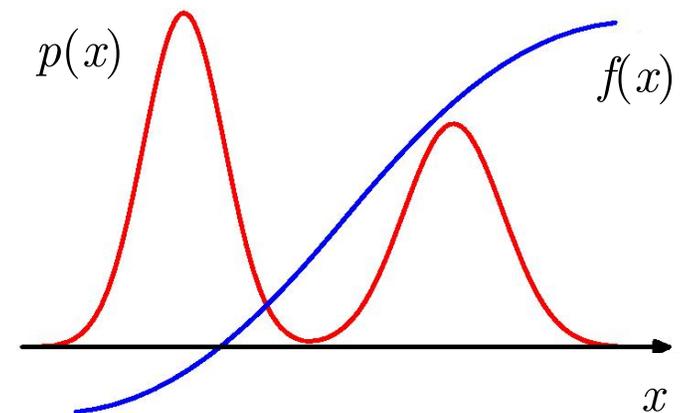
1. Used for both *marginal* and *MAP* inference
2. Compute interesting quantities, e.g., expectation
 - Mean of $f(\mathbf{x})$ where \mathbf{x} is distributed according to a given probabilistic model $p(\mathbf{x})$

$$E_{x \sim p}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})$$

Examples
of $f(\mathbf{x})$:



- Example of $p(x)$ and $f(x)$:
 - where $p(x)$ may be intractable



Computing mean from samples

- Mean of $f(\mathbf{x})$ where \mathbf{x} is distributed as $p(\mathbf{x})$

$$E_{\mathbf{x} \sim p}[f(\mathbf{x})] = \sum_{\mathbf{x}} f(\mathbf{x})p(\mathbf{x})$$

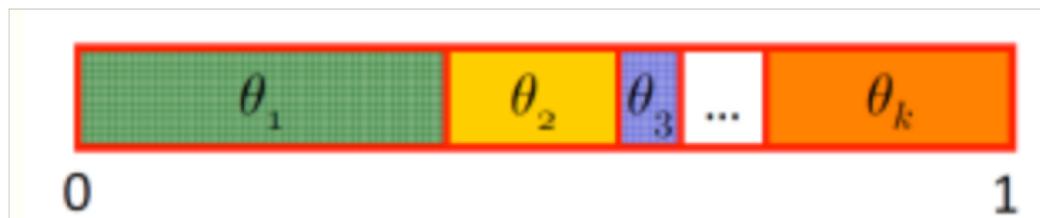
- Computed with samples using the approximation

$$E_{\mathbf{x} \sim p}[f(\mathbf{x})] \approx \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)}) \quad \text{where } \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)} \text{ are samples from } p$$

- It is an unbiased estimate (expected value is true mean)
- Its variance $(1/M)\text{Var}_{\mathbf{x} \sim p}[f(\mathbf{x})]$ can be made arbitrarily small

Sampling from a multinomial

- Sampling from distribution with k possible outcomes with probabilities $\theta_1, \dots, \theta_k$
- Sampling in general is not easy
- Computers can generate samples from very simple distributions, such as Uniform over $[0,1]$
- Sampling from a multinomial can be done by dividing the unit interval into k regions
 - Return value of region in which sample falls



Methods discussed

1. Forward Sampling

- Simplest possible method
- Simply generates samples from original network
- Rejection sampling

2. Importance sampling

- Significantly improved method
- Generates samples closer to posterior distribution

3. Markov chain Monte Carlo

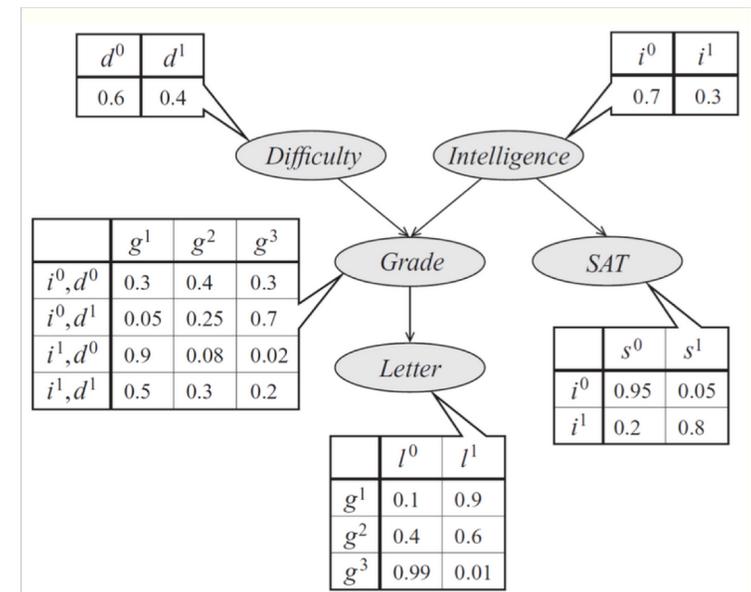
- Sampling process that generates, as it converges, samples arbitrarily close to posterior

Forward (Ancestral) sampling

- Given probability $p(x_1, \dots, x_n)$ specified by a BN, sample variables in topological order
 - Start with variables with no parents
 - Sample from the next generation by conditioning these variables CPDs to values samples at first step

We sample from $p(\mathbf{x})$ in $O(n)$ time by taking exactly one multinomial sample from each CPD

At each step we use multinomial sampling

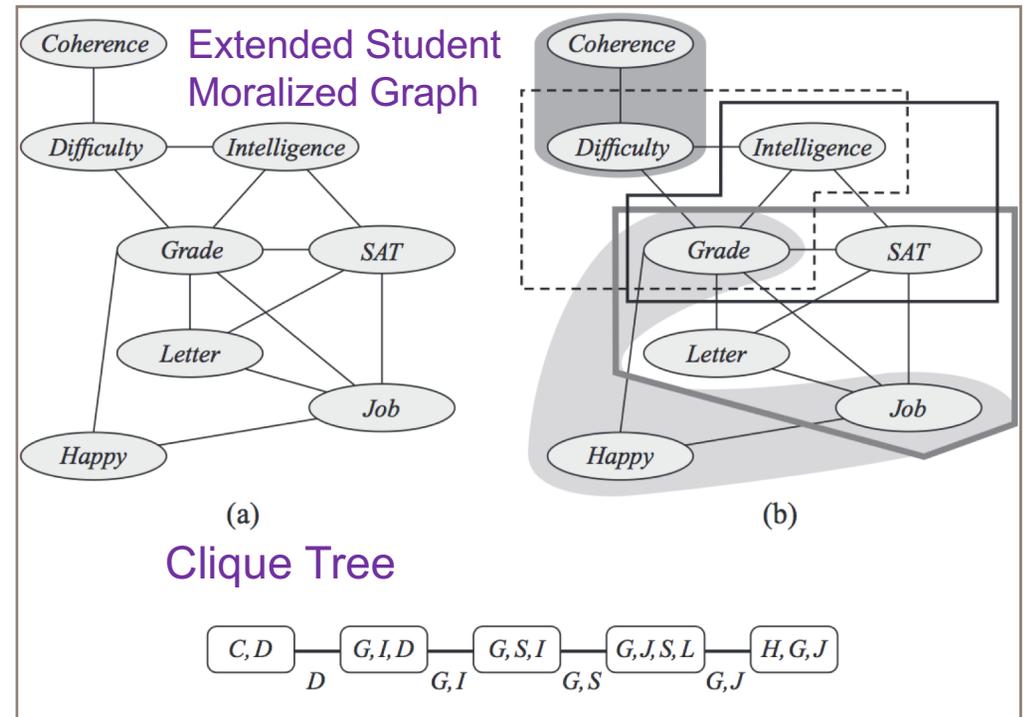


Forward sampling from MNs

- Can be performed efficiently if clique tree has small no of variables per node

Calibrate the clique tree
 Which gives marginal distribution over each node
 Choose a node to be the root
 Then marginalize variables in the root node

Once the marginal for a single variable $x_1 \sim p(X_1|E=e)$ has been sampled from the root node, the newly sampled value $X_1=x_1$ can be incorporated as evidence.



Finish sampling other variables from the same node,
 each time incorporating the newly sampled nodes as evidence,
 i.e. $x_2 \sim p(X_2=x_2|X_1=x_1, E=e)$ and $x_3 \sim p(X_3=x_3|X_1=x_1, X_2=x_2, E=e)$ and so on

Rejection Sampling

- If we want to compute $p(E=e)$ we can draw samples from p and reject those samples that are inconsistent with the evidence

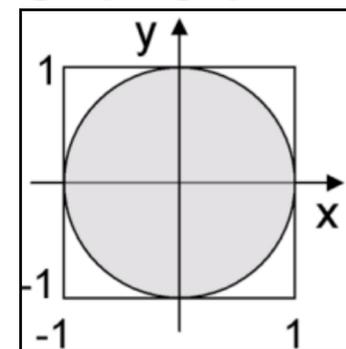
- To determine area of circle

- Reject samples outside of circle

- If $\chi = Z \cup E$ then to compute marginal $p(E=e)$

$$p(E = e) = \sum_z p(Z = z, E = e) = \sum_x p(x) \mathbb{I}(E = e) = \mathbb{E}_{x \sim p}[\mathbb{I}(E = e)]$$

Reject samples inconsistent with $E=e$



- Method breaks down

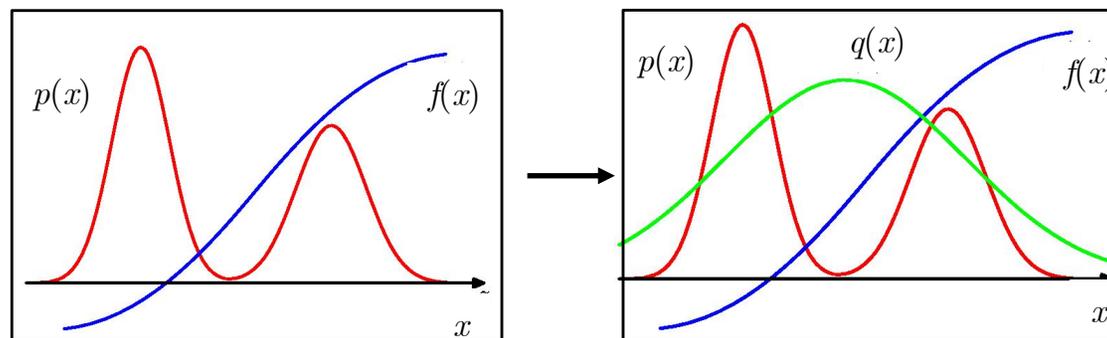
- If circle is small compared to square

Importance Sampling

- Better way to compute integrals than rejection
- Sample from a distribution $q(\mathbf{x})$ and reweigh the samples in a principled way

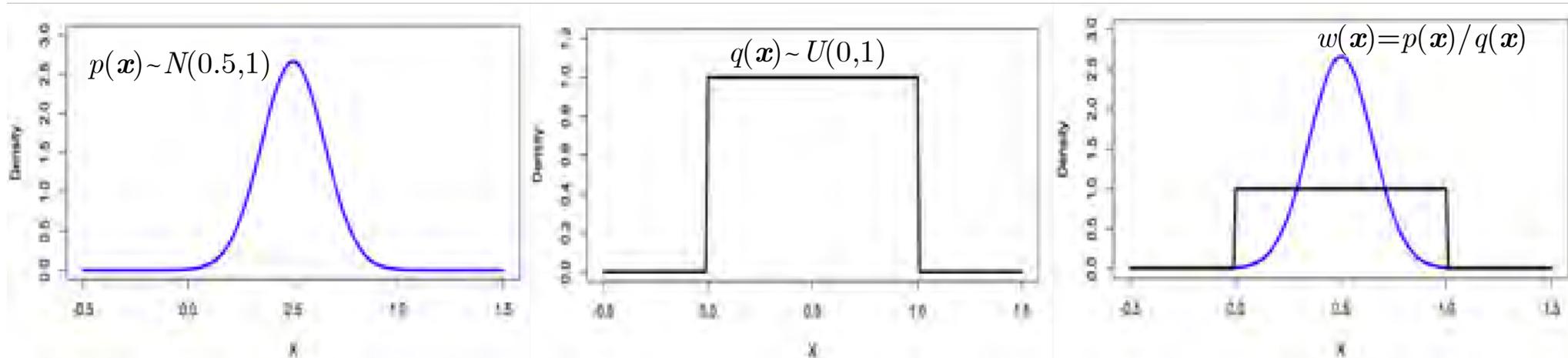
$$\begin{aligned}
 E_{\mathbf{x} \sim p}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\
 &= \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} \\
 &= E_{\mathbf{x} \sim q}[f(\mathbf{x})w(\mathbf{x})] \\
 &\approx \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)})w(\mathbf{x}^{(m)})
 \end{aligned}$$

This can be treated as computing expectation of $p(\mathbf{x})/q(\mathbf{x})$ wrt $q(\mathbf{x})$



- where $w(\mathbf{x}) = p(\mathbf{x})/q(\mathbf{x})$ and samples $\mathbf{x}^{(m)}$ drawn from q
- So take samples from $q(\mathbf{x})$ and reweigh them with $w(\mathbf{x})$
 - But $p(x)$ is intractable, but there is a solution

Case of Uniform $q(\mathbf{x})$



If $\chi = Z \cup E$ then to compute marginal $p(E=e)$

$$\begin{aligned}
 p(E=e) &= \mathbb{E}_{z \sim p}[p(e|z)] \\
 &= \mathbb{E}_{z \sim q} \left[p(e|z) \frac{p(z)}{q(z)} \right] \\
 &= \mathbb{E}_{z \sim q} \left[\frac{p(e, z)}{q(z)} \right] \\
 &= \mathbb{E}_{z \sim q}[w_e(z)] \\
 &\approx \frac{1}{M} \sum_{m=1}^M w_e(z^{(m)})
 \end{aligned}$$

where $w_e(z) = p(e, z)/q(z)$ is the sample weight

Unlike rejection sampling importance sampling will use all samples

Dealing with intractable $p(\mathbf{x})$

- Importance sampling seems to assume that we are able to compute the importance weight $w(\mathbf{x}) = p(\mathbf{x}) / q(\mathbf{x})$
- While $q(\mathbf{x})$ is left as our choice $p(\mathbf{x})$ may be intractable.
- We can use the unnormalized $\tilde{p}(\mathbf{x})$ as shown next

Unnormalized Importance sampling

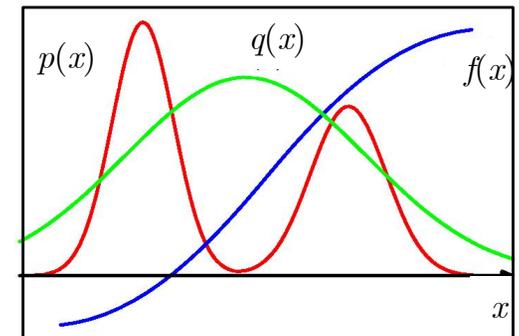
- Define weight using unnormalized distribution

$$w(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})}$$

- Then

$$E_{x \sim q} [w(\mathbf{x})] = \sum_x q(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} = \sum_x \tilde{p}(\mathbf{x}) = \alpha$$

- And



$$\begin{aligned} E_{x \sim p} [f(\mathbf{x})] &= \sum_x p(\mathbf{x}) f(\mathbf{x}) \\ &= \sum_x q(\mathbf{x}) f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} \\ &= \frac{1}{\alpha} \sum_x q(\mathbf{x}) f(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \\ &= \frac{1}{\alpha} E_{x \sim q} [f(\mathbf{x}) w(\mathbf{x})] \\ &= \frac{E_{x \sim q} [f(\mathbf{x}) w(\mathbf{x})]}{E_{x \sim q} [w(\mathbf{x})]} \end{aligned}$$

Therefore, we can determine $E_{x \sim p} [f(\mathbf{x})]$ using samples drawn from q as follows:

$$E_{x \sim p} [f(\mathbf{x})] = \frac{\sum_m f(x[m]) w(x[m])}{\sum_m w(x[m])}$$

Normalized Importance Sampling

- Unnormalized importance sampling is unsuitable for conditional probabilities

$$P(X_i = x_i | E = e) = \frac{P(X_i = x_i, E = e)}{P(E = e)}$$

– With unnormalized we could estimate numerator as

$$\begin{aligned} P(X_i = x_i, E = e) &= \sum_z \delta(z) p(e, z) \\ &= \sum_z \delta(z) w_e(z) q(z) \\ &= \mathbb{E}_{z \sim q}[\delta(z) w_e(z)] \\ &\approx \frac{1}{M} \sum_{m=1}^M \delta(z^{(m)}) w_e(z^{(m)}) \end{aligned}$$

where

$$\delta(z) = \begin{cases} 1 & \text{if } z \text{ is consistent with } X_i = x_i \\ 0 & \text{otherwise} \end{cases}$$

and $w_e(z) = p(e, z) / q(z)$

– and denominator from samples as

$$P(E = e) \approx \frac{1}{M} \sum_{m=1}^M w_e(z^{(m)})$$

Which is shown on Slide 16

- But separate Numerator & Denominator compounds errors
 - » If Numerator is an underestimate, Denominator is an overestimate, overall result is a severe underestimate

Using same set of samples

- If we use the same set of M samples $z^1, \dots, z^M \sim q$
 - For both numerator and denominator, we avoid this issue of compounding errors
 - Final form of normalized importance sampling is

$$\hat{P}(X_i = x_i \mid E = e) = \frac{\frac{1}{T} \sum_{t=1}^T \delta(z^t) w_e(z^t)}{\frac{1}{T} \sum_{t=1}^T w_e(z^t)} \quad \text{where} \quad \begin{cases} \delta(z) = \begin{cases} 1 & \text{if } z \text{ is consistent with } X_i = x_i \\ 0 & \text{otherwise} \end{cases} \\ \text{and } w_e(z) = p(e, z) / q(z) \end{cases}$$

- However, it is biased, for $M=1$

$$\mathbb{E}_{z \sim q}[\hat{P}(X_i = x_i \mid E = e)] = \mathbb{E}_{z \sim q}[\delta(z)] \neq P(X_i = x_i \mid E = e)$$

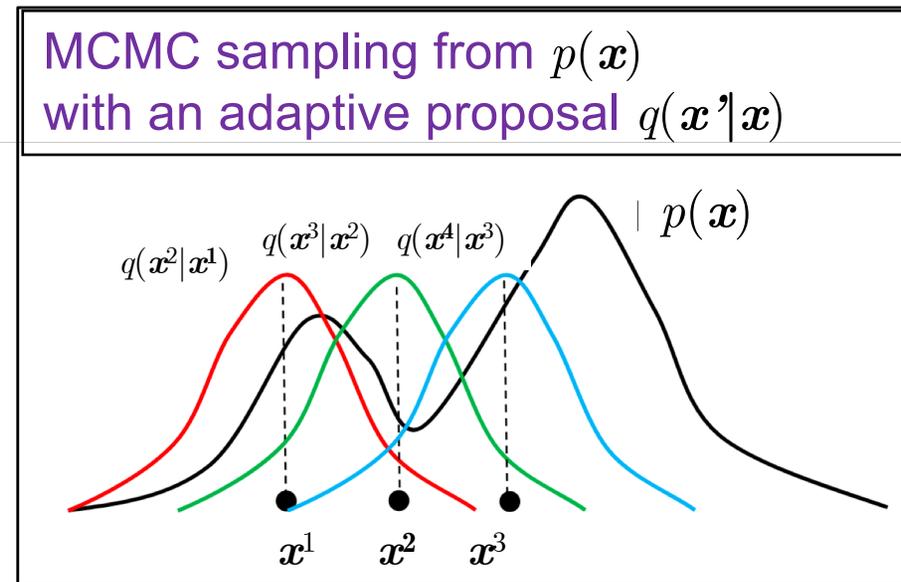
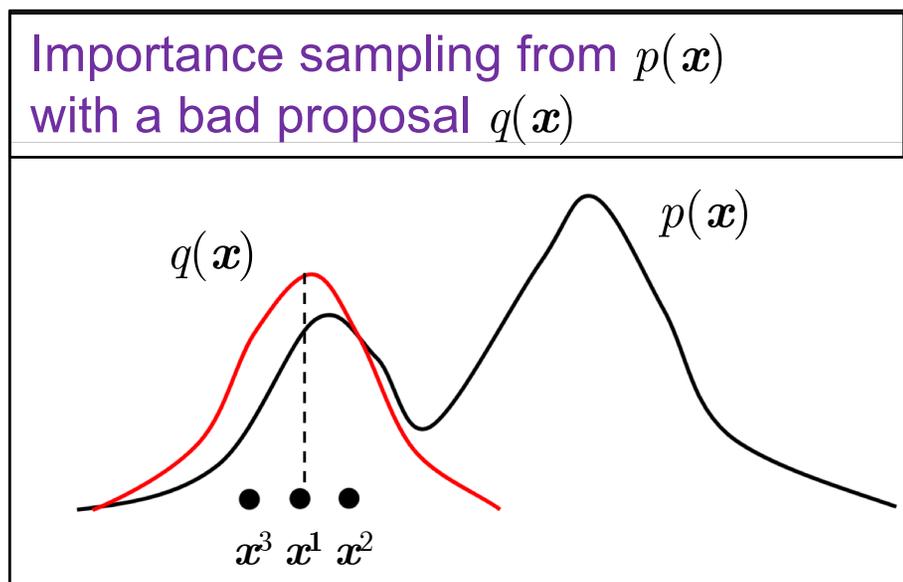
- Fortunately both numerator and denominator are unbiased so the result is asymptotically unbiased

$$\lim_{T \rightarrow \infty} \mathbb{E}_{z \sim q}[\hat{P}(X_i = x_i \mid E = e)] = P(X_i = x_i \mid E = e)$$

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

MCMC and Adaptive Proposals



- MCMC:

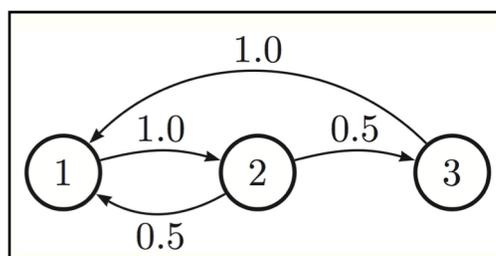
- Instead of $q(\mathbf{x})$, use $q(\mathbf{x}'|\mathbf{x})$ where \mathbf{x}' is the new state being sampled, and \mathbf{x} is previous sample
 - As \mathbf{x} changes, $q(\mathbf{x}'|\mathbf{x})$ changes as a function of \mathbf{x}'

Importance of MCMC

- We turn attention away from computing expectations to marginal and MAP inference
 - Solve these problems using a very powerful technique called Markov Chain Monte Carlo
 - MCMC is a top 10 algorithm of 20th century

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 $q(\mathbf{x}'|\mathbf{x})$ 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}

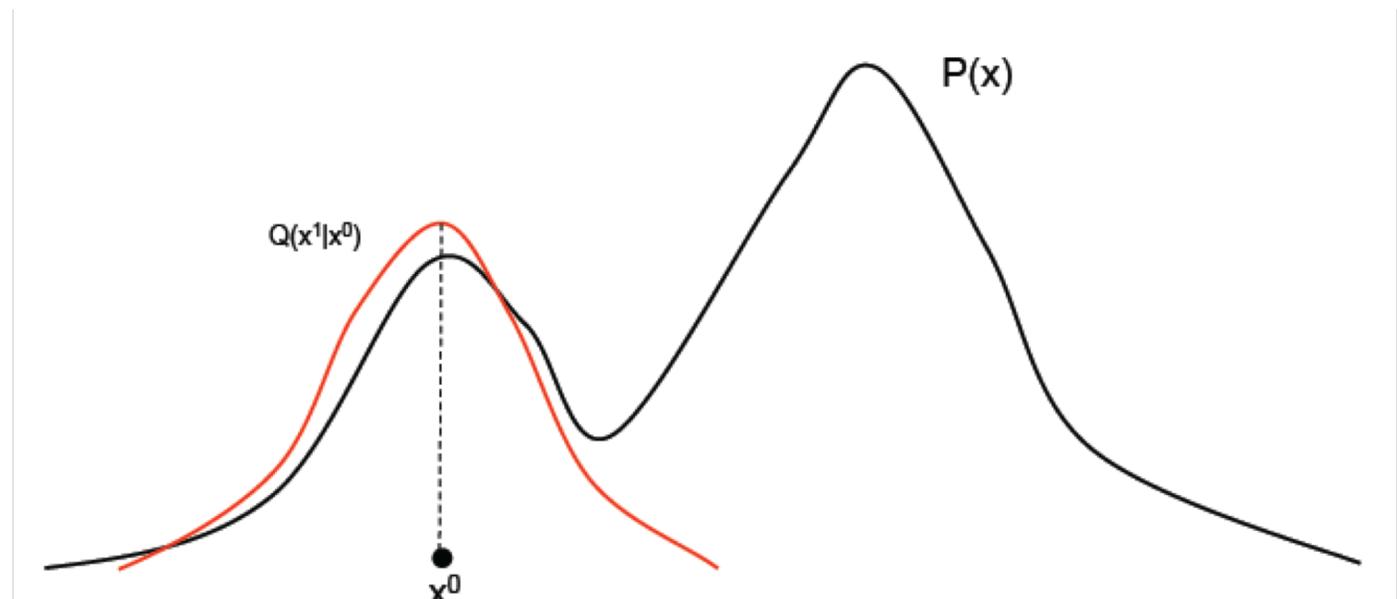
– Notice that we only need to compute $p(\mathbf{x}') / p(\mathbf{x})$ rather than $p(\mathbf{x}')$ or $p(\mathbf{x})$ separately

- $A(\mathbf{x}' | \mathbf{x})$ ensures that, after sufficient draws, samples will come from true distribution $p(\mathbf{x})$

The M-H Algorithm (1)

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

Initialize $x^{(0)}$

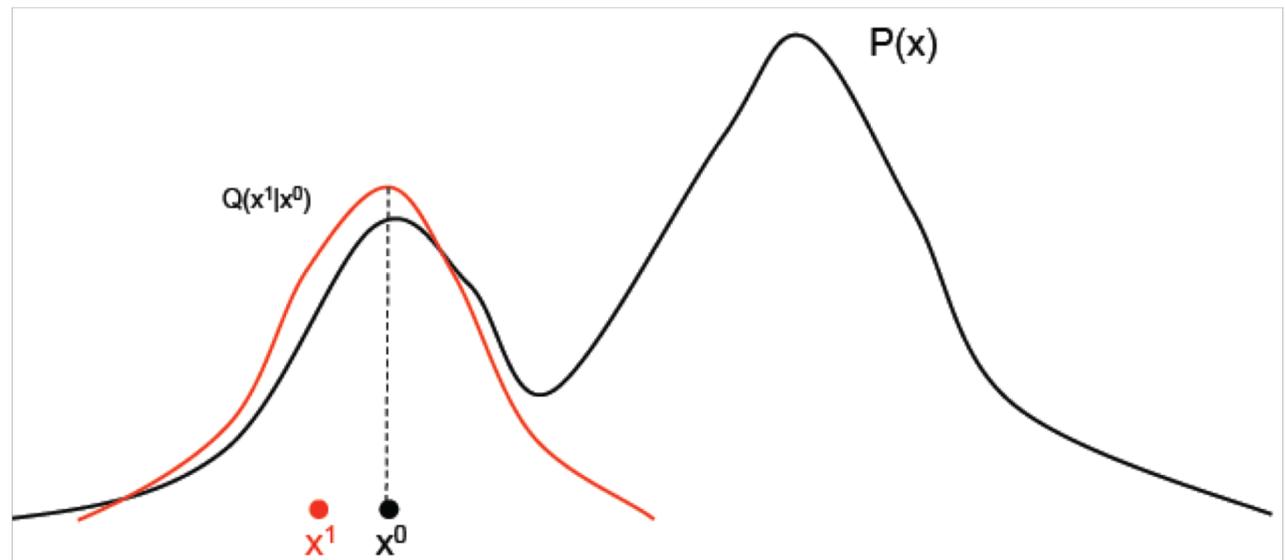


The M-H Algorithm (2)

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

Initialize $x^{(0)}$

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



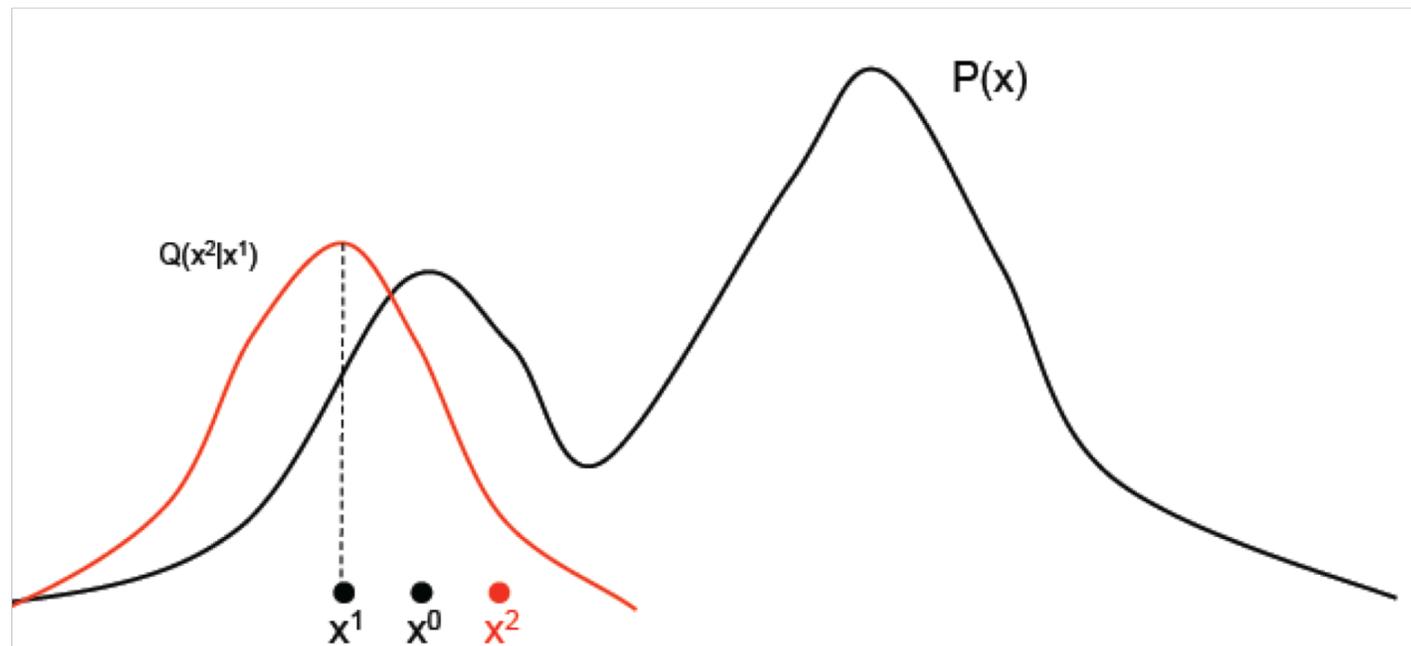
The M-H Algorithm (3)

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

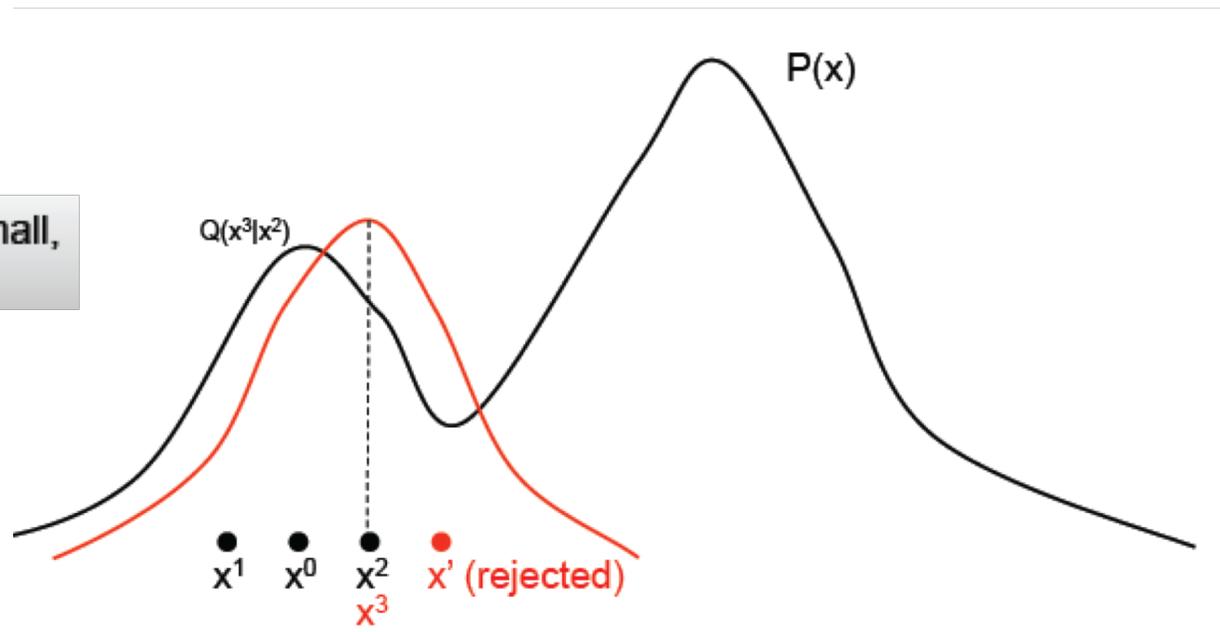
Initialize $x^{(0)}$

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 $p(x)$

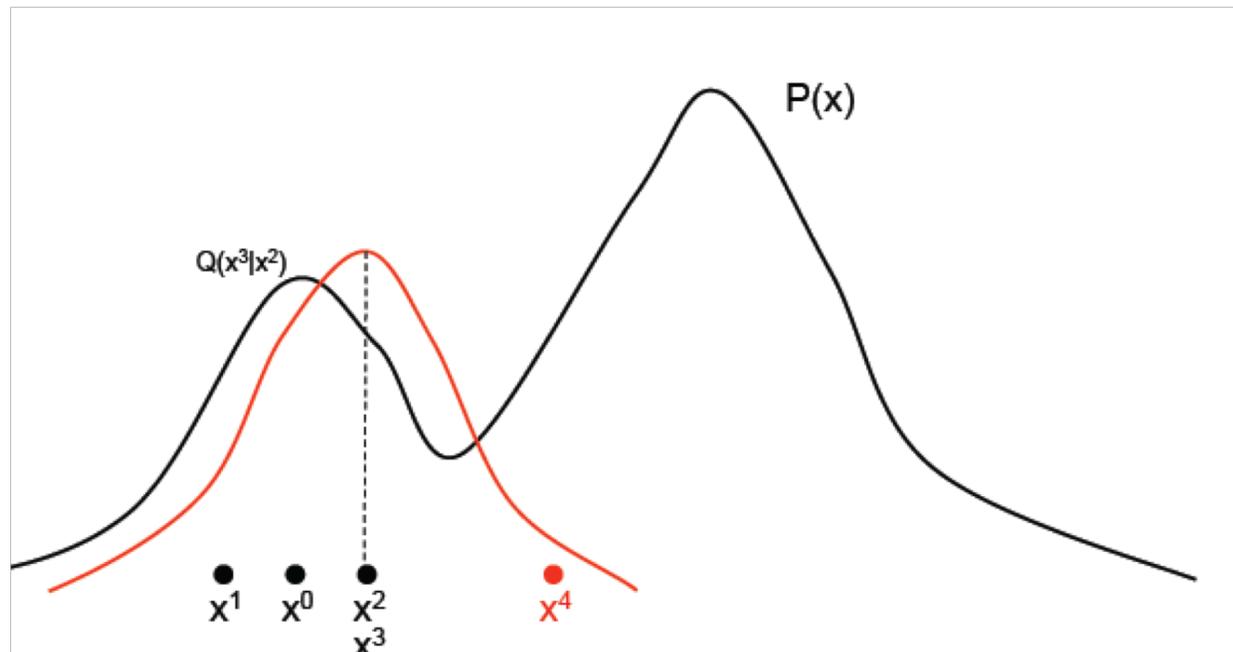
Initialize $x^{(0)}$

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 $x^{(1)}$

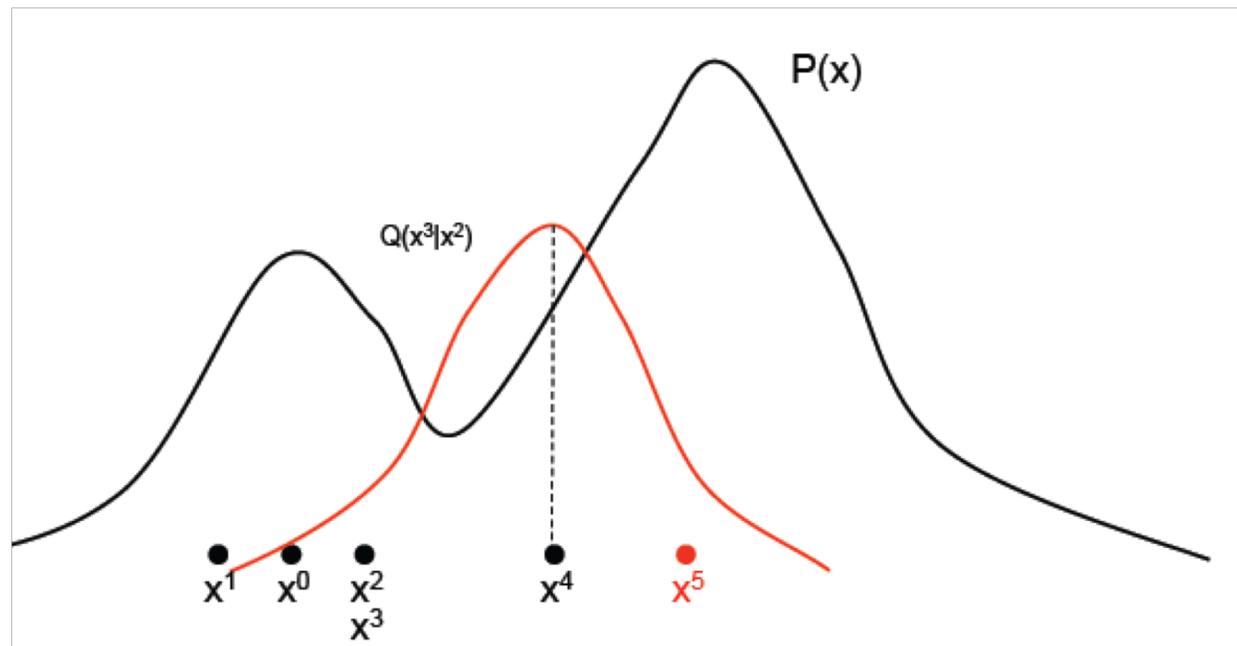
Draw, accept $x^{(2)}$

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 $P(x)$!



Gibbs Sampling

- Gibbs Sampling is an MCMC that samples each random variable of a PGM, one at a time
 - GS is a special case of the MH algorithm
- GS algorithms...
 - Are fairly easy to derive for many graphical models
 - e.g. mixture models, Latent Dirichlet allocation
 - Have reasonable computation and memory requirements, since they sample one r.v. at a time
 - Can be Rao-Blackwellized (integrate out some r.v.s)
 - to decrease sampling variance – called collapsed Gibbs sampling

The Gibbs Sampling Algorithm

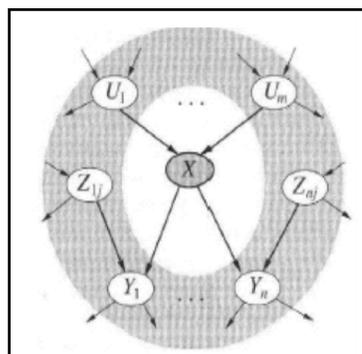
- 1. Suppose the PGM contains variables x_1, \dots, x_n
- 2. Initialize starting values for x_1, \dots, x_n
- 3. Do until convergence:
 - 1. Pick an ordering of the n variables (can be fixed or random)
 - 2. For each variable x_i in order:
 - 1. Sample $x \sim p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$,
 - i.e. the conditional distribution of x_i given the current values of all other variables
 - 2. Update $x_i \leftarrow x$
- When we update x_i , we immediately use its new value for sampling other variables x_j

Markov Blanket to the rescue

- The conditional $p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ looks intimidating, but recall Markov Blankets
- Let $MB(x_i)$ be the Markov Blanket of x_i , then

$$p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = p(x_i \mid MB(x_i))$$

- For a BN, the Markov Blanket of x_i is the set containing its parents, children, and co-parents



- For an MRF, the Markov Blanket of x_i is its immediate neighbors

Gibbs as a case of M-H

- Gibbs sampling can be seen as a special case of MH with proposal $q(x'_i, x_{-i}/x_i, x_{-i}) = p(x'_i/x_{-i})$
 - It is easy check that the acceptance probability simplifies to one
- Assuming the right transition operator, both Gibbs sampling and MH will eventually produce samples from their stationary distribution, which by construction is p

References

- D. Koller and N. Friedman, Probabilistic Graphical Models, MIT Press, 2009
- C. Bishop, Pattern Recognition and Machine Learning, Elsevier, 2006
- <https://ermongroup.github.io/cs228-notes/inference/sampling/>
- <http://people.csail.mit.edu/dsontag/courses/pgm13/slides/lecture9.pdf>