

Approximate inference: Sampling methods

Probabilistic Graphical Models

Tavassolipour

Approximate inference

- ▶ **Approximate inference techniques**
 - ▶ Deterministic approximation
 - ▶ Variational algorithms
 - ▶ Stochastic simulation / sampling methods

Sampling-based estimation

- ▶ Assume that $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$ shows the set of i.i.d. samples drawn from the desired distribution P
- ▶ For any distribution P , function f , we can estimate $E_P[f]$:

$$E_P[f] \approx \underbrace{\frac{1}{N} \sum_{n=1}^N f(x^{(n)})}_{\text{Empirical expectation}}$$

- ▶ Expectations reveal interesting properties about distribution P
 - ▶ Means and variance of P
 - ▶ Probability of events
 - ▶ E.g., we can find $\hat{P}(x = k)$ by estimating $E_P[f]$ where $f(x) = I(x = k)$
- ▶ We can use a stochastic representation of a complex distribution

The mean and variance of the estimator

- For samples drawn independently from the distribution P :

$$\hat{f} = \frac{1}{N} \sum_{n=1}^N f(x^{(n)})$$

$$E[\hat{f}] = E[f]$$

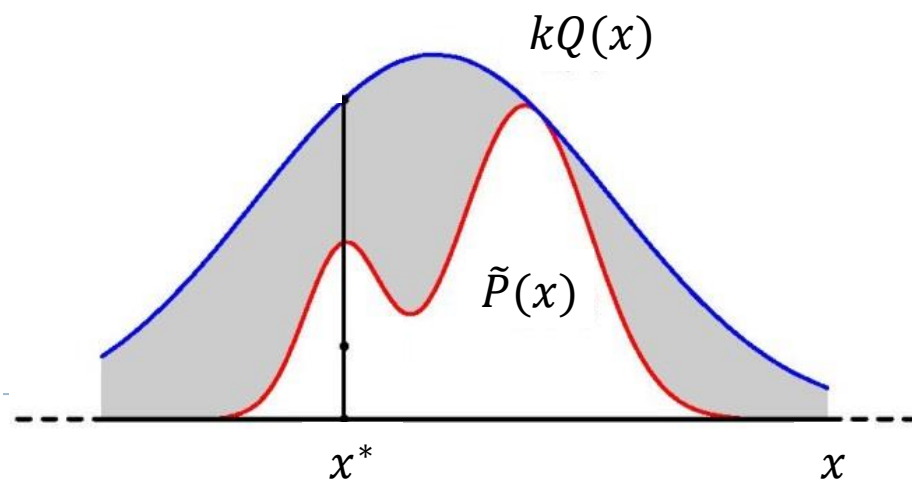
$$\text{var}[\hat{f}] = \frac{1}{N} E[(f - E[f])^2]$$

Monte Carlo methods

- ▶ Using a set of samples to find the answer of an inference query
 - ▶ expectations can be approximated using sample-based averages
- ▶ **Asymptotically** exact and easy to apply to arbitrary problems
- ▶ Challenges:
 - ▶ Drawing samples from many distributions is not trivial
 - ▶ Are the gathered samples enough?
 - ▶ Are all samples useful, or equally useful?

Rejection sampling

- ▶ Suppose we wish to sample from $P(\mathbf{x}) = \tilde{P}(\mathbf{x})/Z$.
 - ▶ $P(\mathbf{x})$ is difficult to sample, but $\tilde{P}(\mathbf{x})$ is easy to evaluate
 - ▶ We choose a simpler (proposal) distribution $Q(\mathbf{x})$ that we can sample from it more easily
 - ▶ Where $\exists k, kQ(\mathbf{x}) \geq \tilde{P}(\mathbf{x})$
 - ▶ Sample from $Q(\mathbf{x})$: $\mathbf{x}^* \sim Q(\mathbf{x})$
 - ▶ accept \mathbf{x}^* with probability $\frac{\tilde{P}(\mathbf{x}^*)}{kQ(\mathbf{x}^*)}$



Rejection sampling

► Correctness:

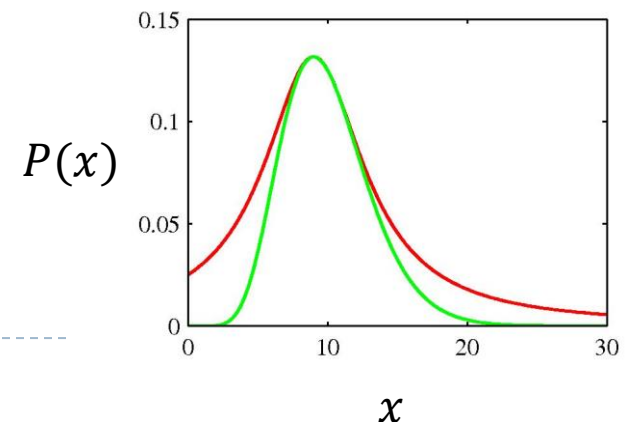
$$\frac{\frac{\tilde{P}(\mathbf{x})}{kQ(\mathbf{x})} Q(\mathbf{x})}{\int \frac{\tilde{P}(\mathbf{x})}{kQ(\mathbf{x})} Q(\mathbf{x}) d\mathbf{x}} = \frac{\tilde{P}(\mathbf{x})}{\int \tilde{P}(\mathbf{x}) d\mathbf{x}} = P(\mathbf{x})$$

Probability of acceptance:

$$P(\text{accept}) = \int \frac{\tilde{P}(\mathbf{x})}{kQ(\mathbf{x})} Q(\mathbf{x}) d\mathbf{x} = \frac{\int \tilde{P}(\mathbf{x}) d\mathbf{x}}{k}$$

High dimensional rejection sampling

- ▶ Problem: low acceptance rate of rejection sampling in high dimensional spaces
 - ▶ exponential decrease of acceptance rate with dimensionality
- ▶ Example:
 - ▶ Using $Q = N(\boldsymbol{\mu}, \sigma_q^2 \mathbf{I})$ to sample $P = N(\boldsymbol{\mu}, \sigma_p^2 \mathbf{I})$
 - ▶ If σ_q exceeds σ_p by 1%, and $d = 1000$
 - ▶ $\left(\frac{\sigma_q}{\sigma_p}\right)^d \approx 20,000$ and so the optimal acceptance rate is $1/20,000$ that is too small



Limitations of Monte Carlo

- ▶ Direct sampling: only when we can sample from $P(\mathbf{x})$
 - ▶ can be wasteful for rare events
- ▶ Rejection sampling uses a proposal distribution $Q(\mathbf{x})$ and can also be used when we can not sample $P(\mathbf{x})$ directly
 - ▶ In rejection sampling, when the proposal $Q(\mathbf{x})$ is very different from $P(\mathbf{x})$, most samples are rejected

Problem: Finding a good proposal $Q(\mathbf{x})$ that is similar to $P(\mathbf{x})$ usually requires knowledge of the analytic form of $P(\mathbf{x})$ that is not available

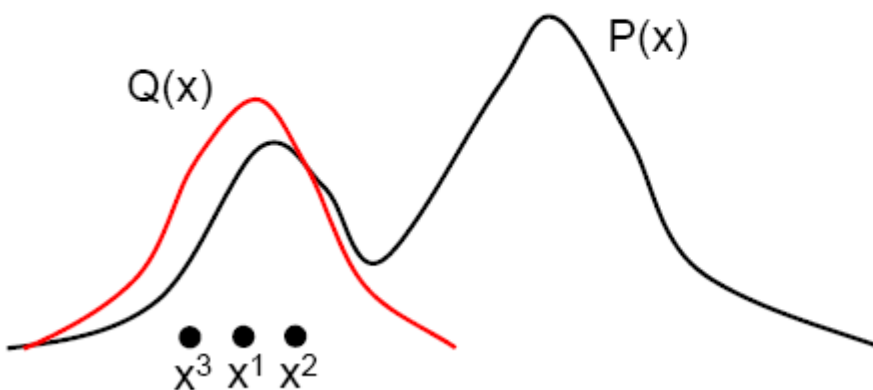
Markov Chain Monte Carlo (MCMC)

- ▶ Instead of using a fixed proposal $Q(\boldsymbol{x})$, we can use an adaptive proposal $Q(\boldsymbol{x}|\boldsymbol{x}^{(t)})$ that depends on the last previous sample $\boldsymbol{x}^{(t)}$
 - ▶ The proposal distribution is adapted as a function of the last accepted sample
- ▶ MCMC methods
 - ▶ Metropolis-Hasting
 - ▶ Gibbs

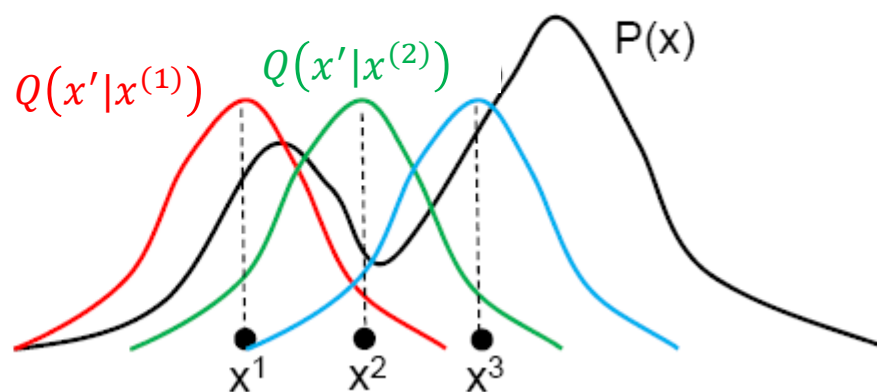
MCMC

- ▶ MCMC algorithms feature adaptive proposals
 - ▶ Instead of $Q(x')$, they use $Q(x'|x)$ where x' is the new state being sampled, and x is the previous sample
 - ▶ As x changes, $Q(x'|x)$ can also change (as a function of x')

Importance sampling with
a (bad) proposal $Q(x)$



MCMC with adaptive
proposal $Q(x'|x)$



Markov chains

- ▶ A Markov chain is a sequence of random variables $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ with the Markov property

$$P(\mathbf{x}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-1)}) = P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$$

- ▶ We focus on homogeneous Markov chains, in which $P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$ is fixed with time
 - ▶ Let \mathbf{x} be the previous state and \mathbf{x}' be the next state, we call $P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$ as $T(\mathbf{x}' | \mathbf{x})$
 - ▶
- ▶ Thus, at each time point, $\mathbf{x}^{(t)}$ is a state showing the configuration of all the variables in the model

Markov chains: invariant or stationary dist.

- ▶ $\pi^t(\mathbf{x})$: Probability distribution over state \mathbf{x} , at time t
 - ▶ Transition probability $T(\mathbf{x}'|\mathbf{x})$ redistributes the mass in state \mathbf{x} to other states \mathbf{x}' .

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

$$T(\mathbf{x}'|\mathbf{x}) = P(\mathbf{x}^{(t)} = \mathbf{x}' | \mathbf{x}^{(t-1)} = \mathbf{x})$$

- ▶ $\pi(\mathbf{x})$ is **invariant** or **stationary** if it does not change under the transitions:

$$\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) T(\mathbf{x}'|\mathbf{x}) \quad \forall \mathbf{x}'$$

Invariant distributions are of great importance in MCMC methods.

More specific than stationary distribution

- ▶ There is also no guarantees that the stationary distribution is unique
- ▶ In some chains, the stationary distribution reached depends on our starting distribution $\pi^0(x)$
- ▶ We want to restrict our attention to MCs that have a unique stationary distribution, which is reached from any starting distribution $\pi^0(x)$.
 - ▶ There are various conditions that suffice to guarantee this property.
 - ▶ The most commonly used condition is ergodicity.

Detailed balance

- ▶ A sufficient (but not necessary) condition for ensuring that $\pi(\mathbf{x})$ is **stationary distribution** of an MC is the **detailed balance** condition:

$$\pi(\mathbf{x})T(\mathbf{x}'|\mathbf{x}) = \pi(\mathbf{x}')T(\mathbf{x}|\mathbf{x}')$$

Detailed balance means the sequences \mathbf{x}', \mathbf{x} and \mathbf{x}, \mathbf{x}' are equally probable (although the probability of $\mathbf{x}' \rightarrow \mathbf{x}$ and $\mathbf{x} \rightarrow \mathbf{x}'$ can be different)

Reversible Chains

- ▶ Theorem: **Detailed balance** implies the stationary distribution

Proof:

$$\pi(\mathbf{x})T(\mathbf{x}'|\mathbf{x}) = \pi(\mathbf{x}')T(\mathbf{x}|\mathbf{x}')$$

$$\Rightarrow \sum_{\mathbf{x}} \pi(\mathbf{x})T(\mathbf{x}'|\mathbf{x}) = \sum_{\mathbf{x}} \pi(\mathbf{x}')T(\mathbf{x}|\mathbf{x}')$$

$$\Rightarrow \sum_{\mathbf{x}} \pi(\mathbf{x})T(\mathbf{x}'|\mathbf{x}) = \pi(\mathbf{x}') \sum_{\mathbf{x}} T(\mathbf{x}|\mathbf{x}') = \pi(\mathbf{x}')$$

- ▶ Theorem: If detailed balance holds and T is ergodic, then T has a unique stationary distribution

Stationary distribution: summary

- ▶ $\pi^t(\mathbf{x})$: Probability distribution over state \mathbf{x} , at time t
 - ▶ Transition probability $T(\mathbf{x}'|\mathbf{x})$ redistributes the mass in state \mathbf{x} to other states \mathbf{x}' .

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

- ▶ $\pi^*(\mathbf{x})$ is **invariant** or **stationary** if it does not change under the transitions:

$$\pi^*(\mathbf{x}') = \sum_{\mathbf{x}} \pi^*(\mathbf{x}) T(\mathbf{x}'|\mathbf{x}) \quad \forall \mathbf{x}'$$

- ▶ A sufficient condition for ensuring that $\pi^*(\mathbf{x})$ is **stationary distribution** of an MC is the **detailed balance** condition:

$$\pi^*(\mathbf{x}) T(\mathbf{x}'|\mathbf{x}) = \pi^*(\mathbf{x}') T(\mathbf{x}|\mathbf{x}')$$

How to use Markov chains for sampling from $P(\mathbf{x})$?

- ▶ Our goal is to use Markov chains to sample from a given distribution.
- ▶ We can achieve this if we set up a Markov chain whose unique stationary distribution is P .
- ▶ We design the transition distribution $T(\mathbf{x}'|\mathbf{x})$ so that the chain has a unique stationary distribution $P(\mathbf{x})$ (independent of P)
 - ▶ The ergodic condition is a sufficient condition

Sample $\mathbf{x}^{(0)}$ randomly

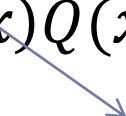
For $t = 0, 1, 2, \dots$

Sample $\mathbf{x}^{(t+1)}$ from $T(\mathbf{x}'|\mathbf{x}^{(t)})$

Metropolis-Hastings

- ▶ Draws a sample x' from $Q(x'|x)$, where x is the previous sample

- ▶ The new sample x' is accepted with probability $A(x'|x)$:

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


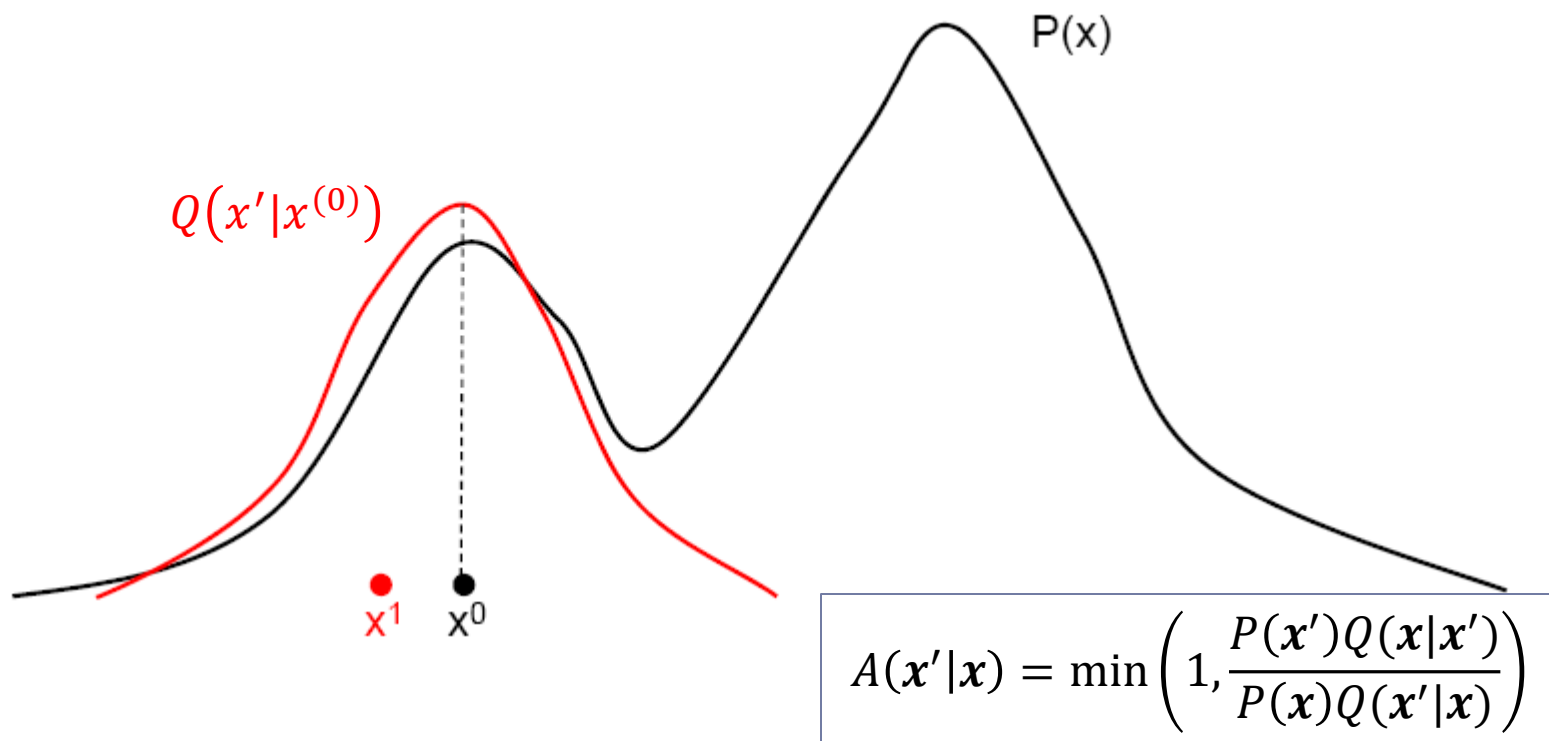
we only need to compute $\frac{P(x')}{P(x)}$ rather than $P(x')$ or $P(x)$ separately

- ▶ We use $A(x'|x)$ to ensure that after sufficiently many draws, our samples will come from the true distribution $P(x)$

Metropolis-Hastings algorithm: Example

- ▶ Let $Q(x'|x)$ be a Gaussian centered on x
- ▶ We're trying to sample from a bimodal distribution $P(x)$

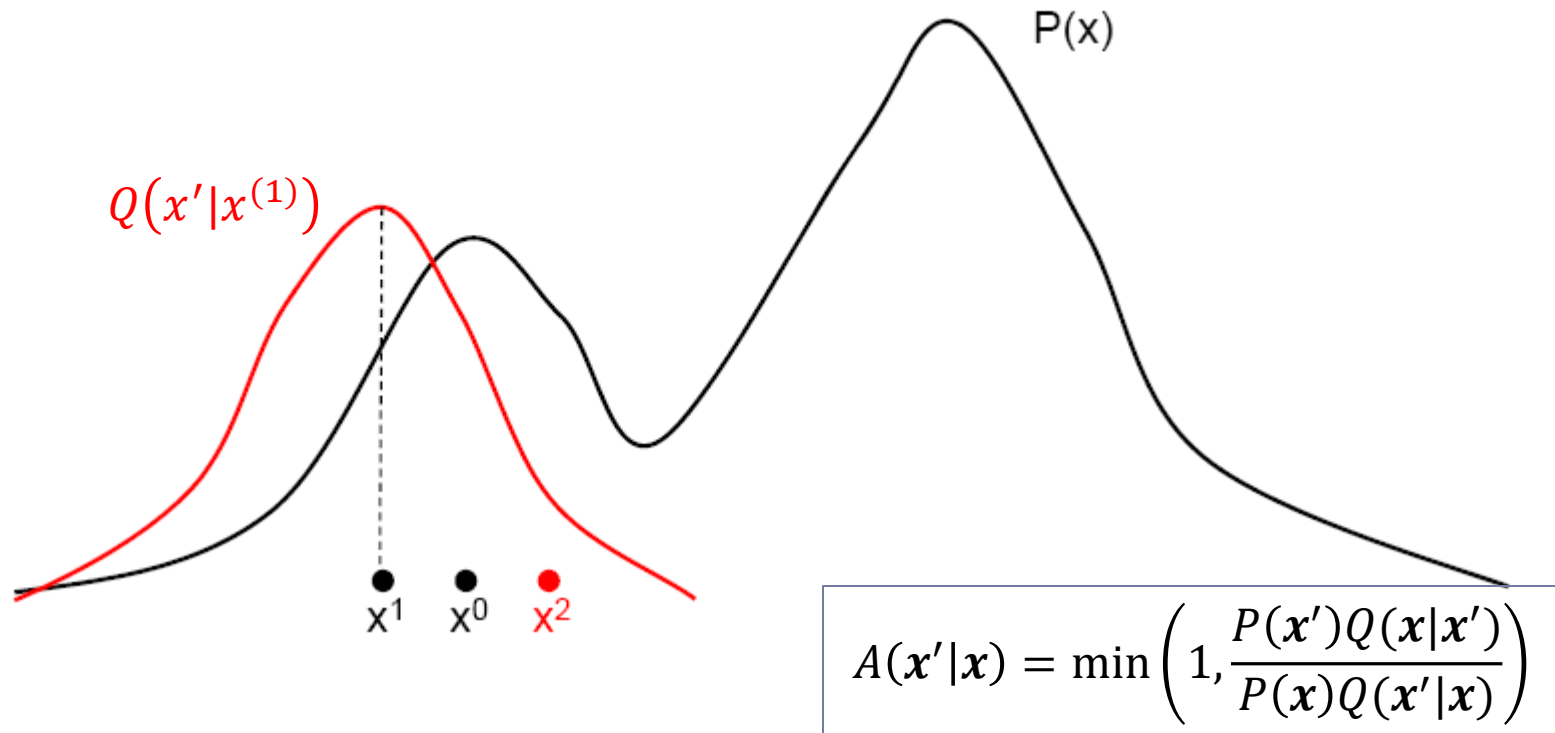
Initialize $x^{(0)}$
Draw, accept x^1



Metropolis-Hastings algorithm: Example

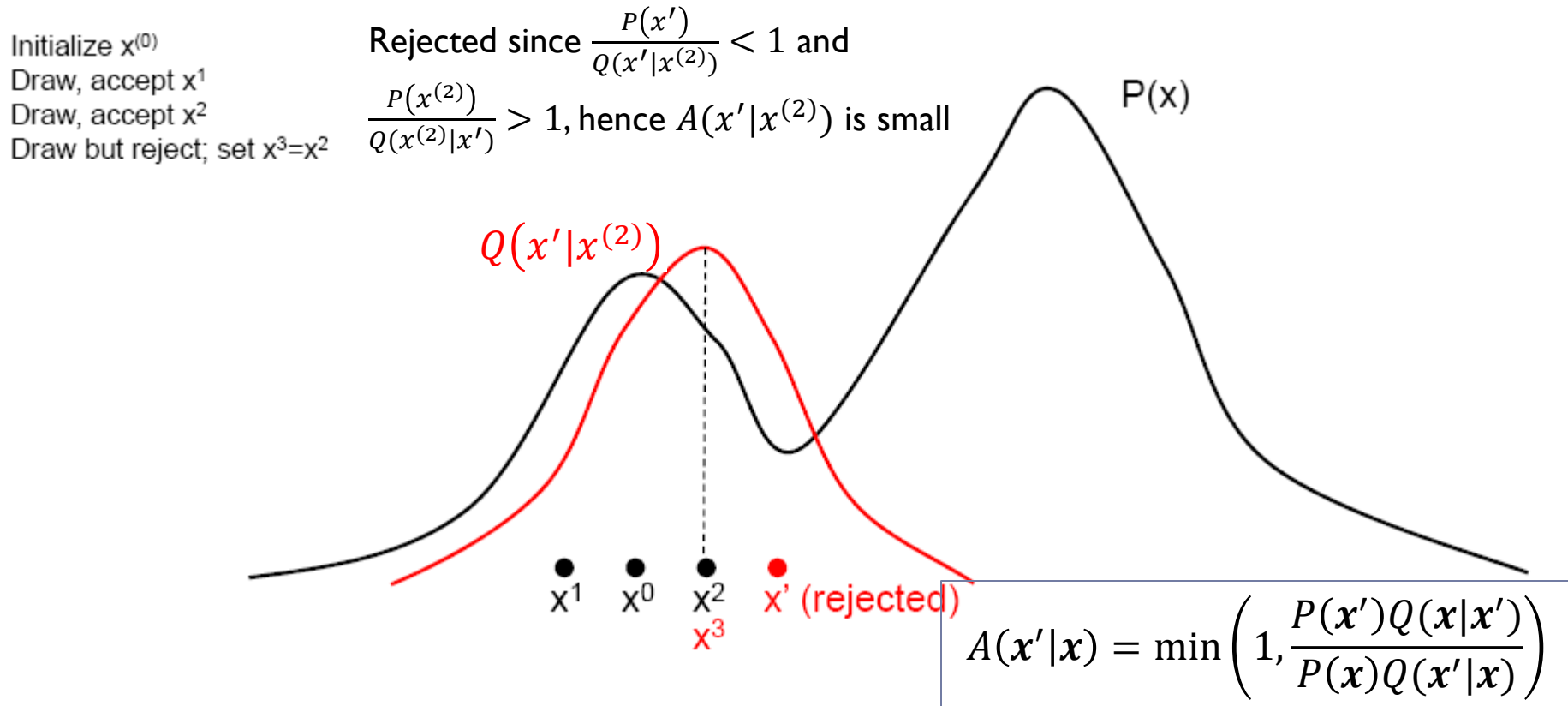
- ▶ Let $Q(x'|x)$ be a Gaussian centered on x
- ▶ We're trying to sample from a bimodal distribution $P(x)$

Initialize $x^{(0)}$
Draw, accept x^1
Draw, accept x^2



Metropolis-Hastings algorithm: Example

- ▶ Let $Q(x'|x)$ be a Gaussian centered on x
- ▶ We're trying to sample from a bimodal distribution $P(x)$



Metropolis-Hastings algorithm: Example

- ▶ Let $Q(x'|x)$ be a Gaussian centered on x
- ▶ We're trying to sample from a bimodal distribution $P(x)$

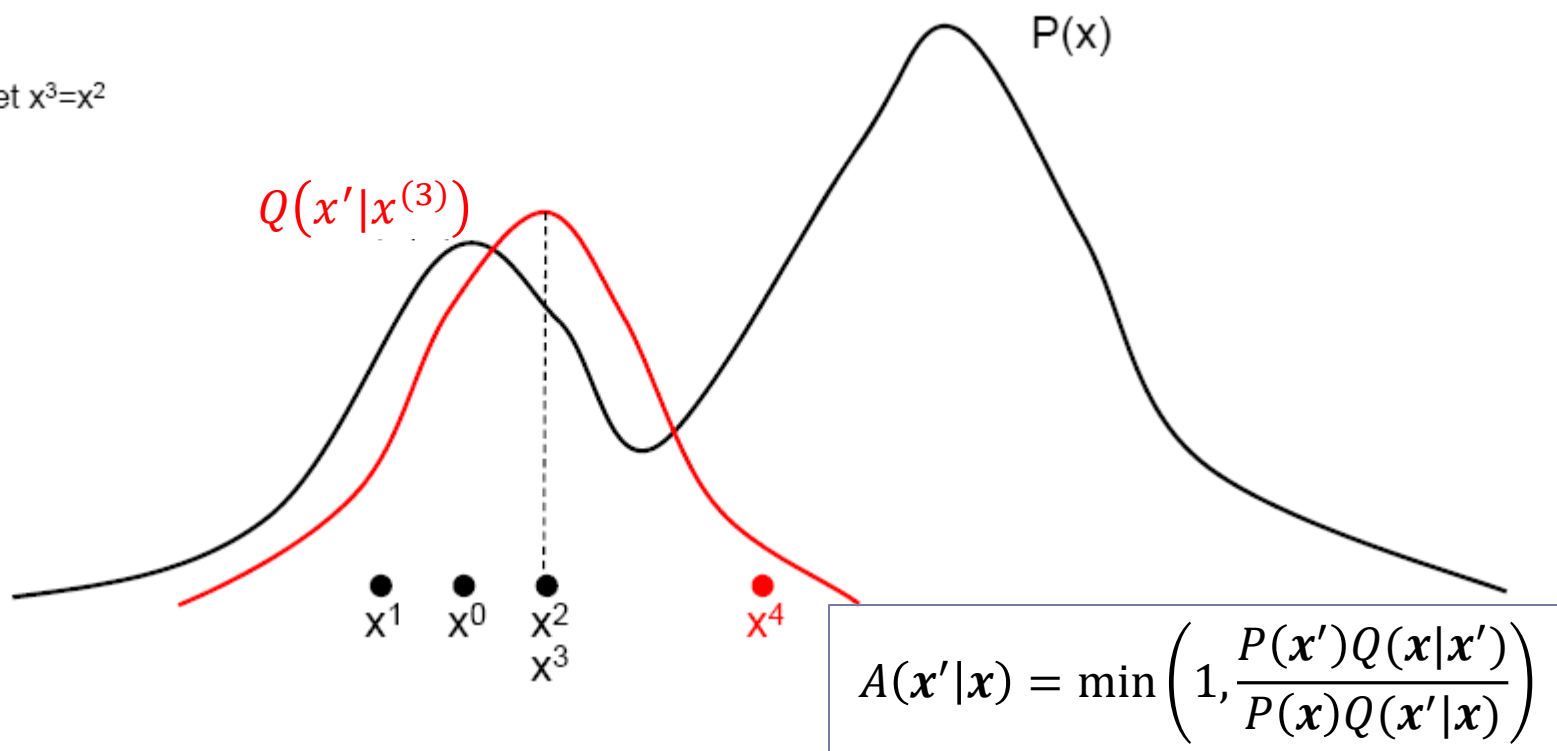
Initialize $x^{(0)}$

Draw, accept x^1

Draw, accept x^2

Draw but reject; set $x^3=x^2$

Draw, accept x^4



Metropolis-Hastings algorithm: Example

- ▶ Let $Q(x'|x)$ be a Gaussian centered on x
- ▶ We're trying to sample from a bimodal distribution $P(x)$

Initialize $x^{(0)}$

Draw, accept x^1

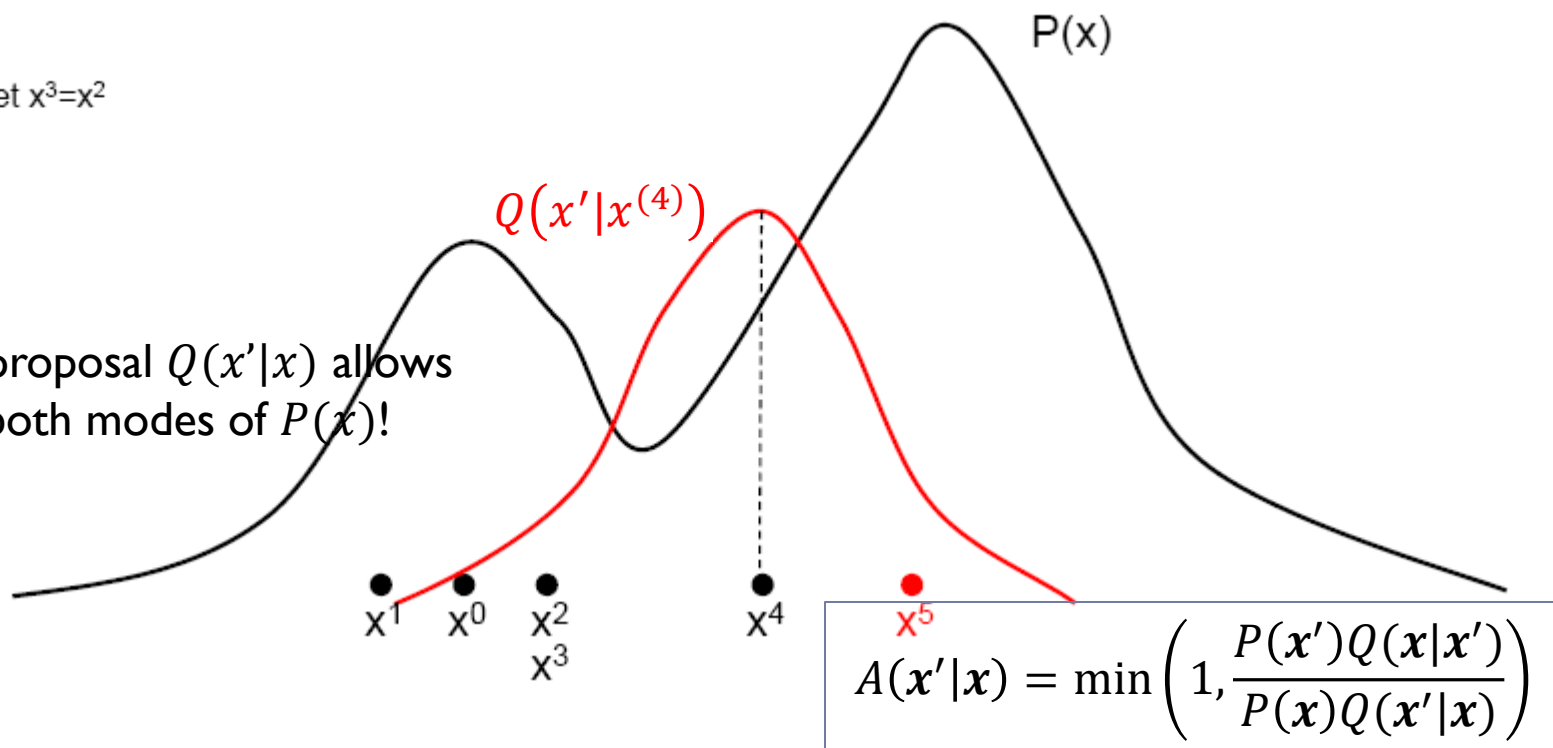
Draw, accept x^2

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 $P(x)$!



Metropolis-Hastings algorithm

- ▶ Initialize starting state $\mathbf{x}^{(0)}$, set $t = 0$
- ▶ Burn-in: while samples have “not converged”
 - $\mathbf{x} = \mathbf{x}^{(t)}$
 - sample $\mathbf{x}^* \sim Q(\mathbf{x}^* | \mathbf{x})$
 - $A(\mathbf{x}^* | \mathbf{x}) = \min \left(1, \frac{P(\mathbf{x}^*)Q(\mathbf{x} | \mathbf{x}^*)}{P(\mathbf{x})Q(\mathbf{x}^* | \mathbf{x})} \right)$
 - sample $u \sim \text{Uniform}(0,1)$
 - if $u < A(\mathbf{x}^* | \mathbf{x})$
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^*$
 - else
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)}$
 - $t = t + 1$
- ▶ For $n = 1 \dots N$
 - ▶ Draw sample from $Q(\mathbf{x} | \mathbf{x}^{(n-1)})$ with acceptance probability $A(\mathbf{x} | \mathbf{x}^{(n-1)})$

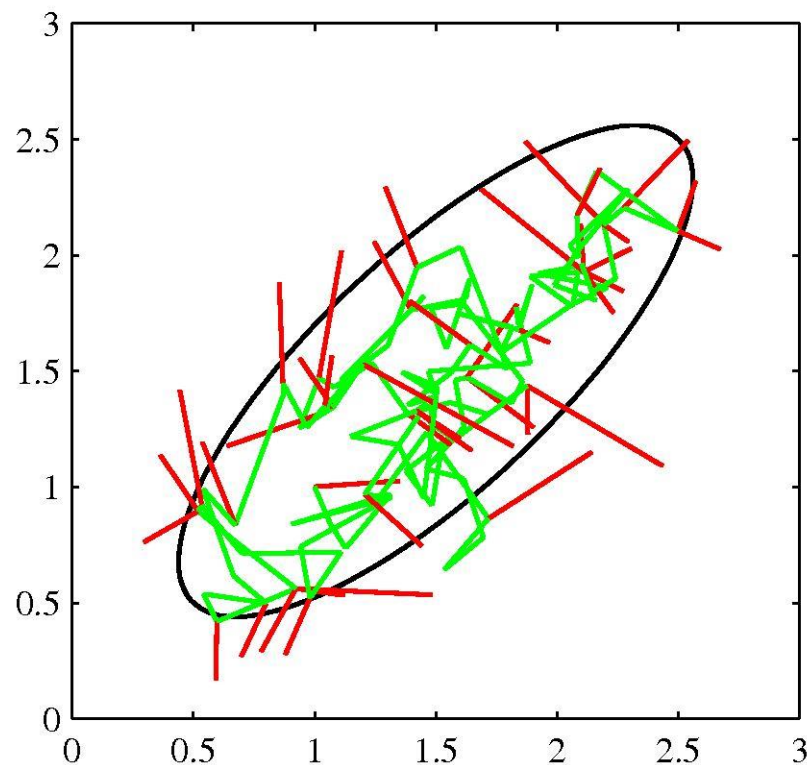
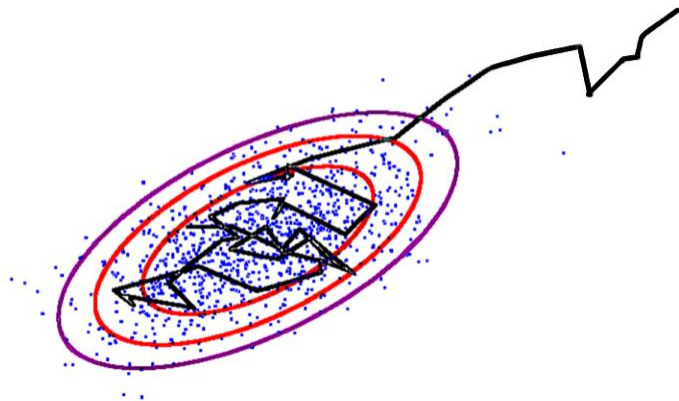
\mathbf{x}^* is accepted with
probability $A(\mathbf{x}^* | \mathbf{x})$

Metropolis algorithm: example

Let $Q(\mathbf{x}'|\mathbf{x})$ be a Gaussian centered on \mathbf{x} :

$$Q(\mathbf{x}'|\mathbf{x}) = N(\mathbf{x}'|\mathbf{x}, \sigma^2 \mathbf{I})$$

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



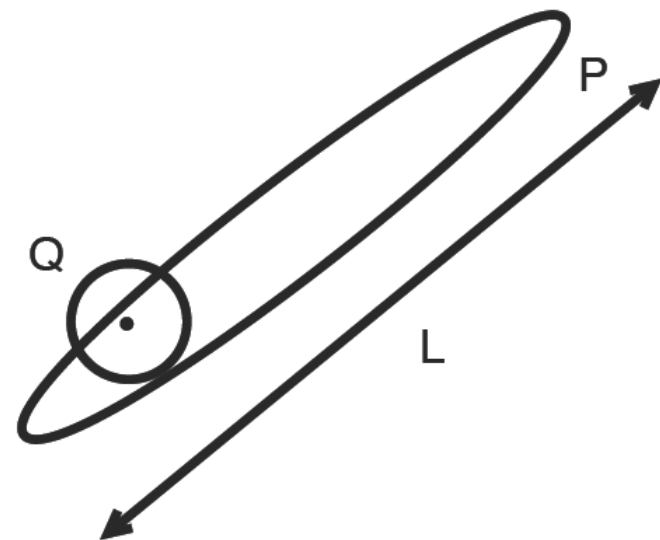
A biased random walk that explores the target distribution $P(\mathbf{x})$

[Bishop book]

Metropolis algorithm: example

- ▶ $Q(\mathbf{x}'|\mathbf{x}) = N(\mathbf{x}'|\mathbf{x}, \sigma^2 \mathbf{I})$
- ▶ large $\sigma^2 \Rightarrow$ many rejections
- ▶ small $\sigma^2 \Rightarrow$ slow exploration
 - ▶ $\left(\frac{L}{\sigma}\right)^2$ iterations are required to reach states

In general, finding a good proposal distribution is not always easy



Proposal distribution

- ▶ low-variance proposals:
 - ▶ high probability of acceptance
 - ▶ many iterations are required to explore $P(x)$
 - ▶ results in more correlated samples
- ▶ high-variance proposals
 - ▶ low probability of acceptance
 - ▶ have the potential to explore much of $P(x)$
 - ▶ Results in less correlated samples

Theoretical foundation of MH

- ▶ Why are the samples generated by MH method will eventually come from $P(\boldsymbol{x})$?
- ▶ Why does the MH algorithm have a “burn-in” period?

Why does Metropolis-Hastings work?

- ▶ MH is a general construction algorithm that allows us to build a reversible Markov chain with a particular stationary distribution $P(\mathbf{x})$
- ▶ If we draw a sample \mathbf{x}' according to $Q(\mathbf{x}'|\mathbf{x})$, and then accept/reject according to $A(\mathbf{x}'|\mathbf{x})$, we have a transition kernel:

$$T(\mathbf{x}'|\mathbf{x}) = A(\mathbf{x}'|\mathbf{x})Q(\mathbf{x}'|\mathbf{x})$$

$$T(\mathbf{x}'|\mathbf{x}) = A(\mathbf{x}'|\mathbf{x})Q(\mathbf{x}'|\mathbf{x}) \text{ if } \mathbf{x}' \neq \mathbf{x}$$

$$T(\mathbf{x}|\mathbf{x}) = Q(\mathbf{x}|\mathbf{x}) + \sum_{\mathbf{x}' \neq \mathbf{x}} Q(\mathbf{x}'|\mathbf{x})(1 - A(\mathbf{x}'|\mathbf{x}))$$

MH satisfies detailed balance

- ▶ Theorem: MH satisfies detailed balance
- ▶ Proof:

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

If $A(\mathbf{x}'|\mathbf{x}) \leq 1$ then $\frac{P(\mathbf{x})Q(\mathbf{x}'|\mathbf{x})}{P(\mathbf{x}')Q(\mathbf{x}|\mathbf{x}')} \geq 1$ then $A(\mathbf{x}|\mathbf{x}') = 1$

Suppose that $A(\mathbf{x}'|\mathbf{x}) < 1$

$$A(\mathbf{x}'|\mathbf{x}) = \frac{P(\mathbf{x}')Q(\mathbf{x}|\mathbf{x}')}{P(\mathbf{x})Q(\mathbf{x}'|\mathbf{x})}$$

$$\Rightarrow P(\mathbf{x})Q(\mathbf{x}'|\mathbf{x})A(\mathbf{x}'|\mathbf{x}) = P(\mathbf{x}')Q(\mathbf{x}|\mathbf{x}')$$

$$\xRightarrow{A(\mathbf{x}|\mathbf{x}')=1} P(\mathbf{x})Q(\mathbf{x}'|\mathbf{x})A(\mathbf{x}'|\mathbf{x}) = P(\mathbf{x}')Q(\mathbf{x}|\mathbf{x}')A(\mathbf{x}|\mathbf{x}')$$

$$\Rightarrow P(\mathbf{x})T(\mathbf{x}'|\mathbf{x}) = P(\mathbf{x}')T(\mathbf{x}|\mathbf{x}')$$

$T(\mathbf{x}'|\mathbf{x}) = A(\mathbf{x}'|\mathbf{x})Q(\mathbf{x}'|\mathbf{x})$

MH properties

- ▶ MH algorithm eventually converges to a stationary distribution $P(\mathbf{x})$ that is the true distribution
- ▶ However, we have no guarantees as to when this will occur
 - ▶ the burn-in period is a way to ignore the un-converged part of the Markov chain
 - ▶ but deciding when to halt burn-in is an art that needs experimentation.
- ▶ Q must be chosen to fulfill the technical requirements

Gibbs sampling algorithm

Suppose the graphical model contains variables x_1, \dots, x_M

Initialize starting values for x_1, \dots, x_M

Do until convergence:

Pick an ordering of the M variables (can be fixed or random)

For each variable x_i in order:

Sample x from $P(X_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_M)$

Update $x_i \leftarrow x$



the current values of all other variables

When we update x_i , we immediately use its new value for sampling other variables x_j

Gibbs sampling algorithm

Suppose the graphical model contains variables x_1, \dots, x_M

If the current sample is $\mathbf{x} = [x_1, \dots, x_M]$, the next sample $\mathbf{x}' = [x'_1, \dots, x'_M]$ is drawn as:

Sample x'_1 from $P(X_1 | x_2, \dots, x_M)$

Sample x'_2 from $P(X_2 | x'_1, x_3, \dots, x_M)$

...

Sample x'_i from $P(X_i | x'_1, \dots, x'_{i-1}, x_{i+1}, \dots, x_M)$

...

Sample x'_M from $P(X_M | x'_1, \dots, x'_{M-1})$