



10-708 Probabilistic Graphical Models

Markov Chain Monte Carlo (MCMC)

Readings:

MacKay Ch. 29

Jordan Ch. 21

Matt Gormley
Lecture 16
March 14, 2016

Housekeeping

- **Homework 2**
 - Due March 16, 12:00 noon (extended)
- **Midway Project Report**
 - Due March 23, 12:00 noon

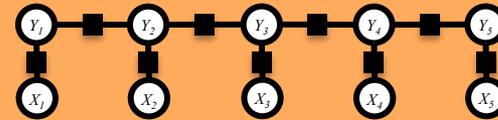
1. Data

$$\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$$

Sample 1:	n time	v flies	p like	d an	n arrow
Sample 2:	n time	n flies	v like	d an	n arrow
Sample 3:	n flies	v fly	p with	n their	n wing
Sample 4:	p with	n time	n you	v will	v see

2. Model

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$



3. Objective

$$\ell(\boldsymbol{\theta}; \mathcal{D}) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)} | \boldsymbol{\theta})$$

5. Inference

1. Marginal Inference

$$p(x_C) = \sum_{\mathbf{x}': \mathbf{x}'_C = x_C} p(\mathbf{x}' | \boldsymbol{\theta})$$

2. Partition Function

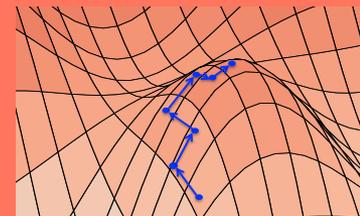
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$

3. MAP Inference

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta})$$

4. Learning

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D})$$



It's Pi Day 2016...

...so let's compute π .

Properties of Monte Carlo

$$\text{Estimator: } \int f(x)P(x) dx \approx \hat{f} \equiv \frac{1}{S} \sum_{s=1}^S f(x^{(s)}), \quad x^{(s)} \sim P(x)$$

Estimator is unbiased:

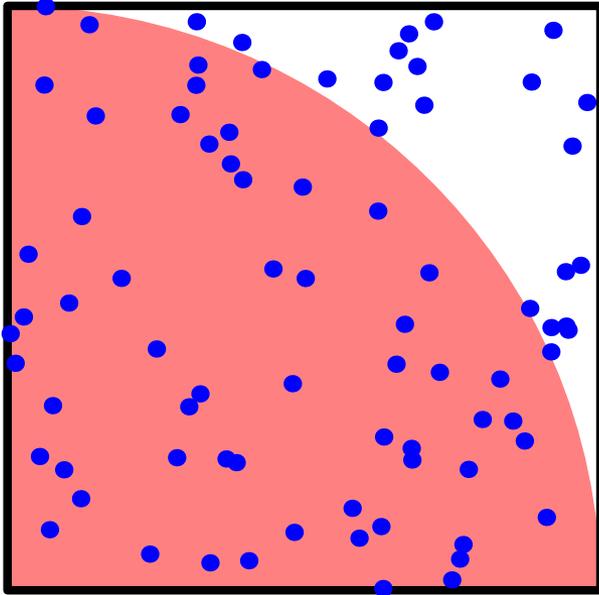
$$\mathbb{E}_{P(\{x^{(s)}\})} [\hat{f}] = \frac{1}{S} \sum_{s=1}^S \mathbb{E}_{P(x)} [f(x)] = \mathbb{E}_{P(x)} [f(x)]$$

Variance shrinks $\propto 1/S$:

$$\text{var}_{P(\{x^{(s)}\})} [\hat{f}] = \frac{1}{S^2} \sum_{s=1}^S \text{var}_{P(x)} [f(x)] = \text{var}_{P(x)} [f(x)] / S$$

“Error bars” shrink like \sqrt{S}

A dumb approximation of π



$$P(x, y) = \begin{cases} 1 & 0 < x < 1 \text{ and } 0 < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\pi = 4 \iint \mathbb{I}((x^2 + y^2) < 1) P(x, y) \, dx \, dy$$

```
octave:1> S=12; a=rand(S,2); 4*mean(sum(a.*a,2)<1)
```

```
ans = 3.3333
```

```
octave:2> S=1e7; a=rand(S,2); 4*mean(sum(a.*a,2)<1)
```

```
ans = 3.1418
```

Aside: don't always sample!

“Monte Carlo is an extremely bad method; it should be used only when all alternative methods are worse.”

— Alan Sokal, 1996

Example: numerical solutions to (nice) 1D integrals are fast

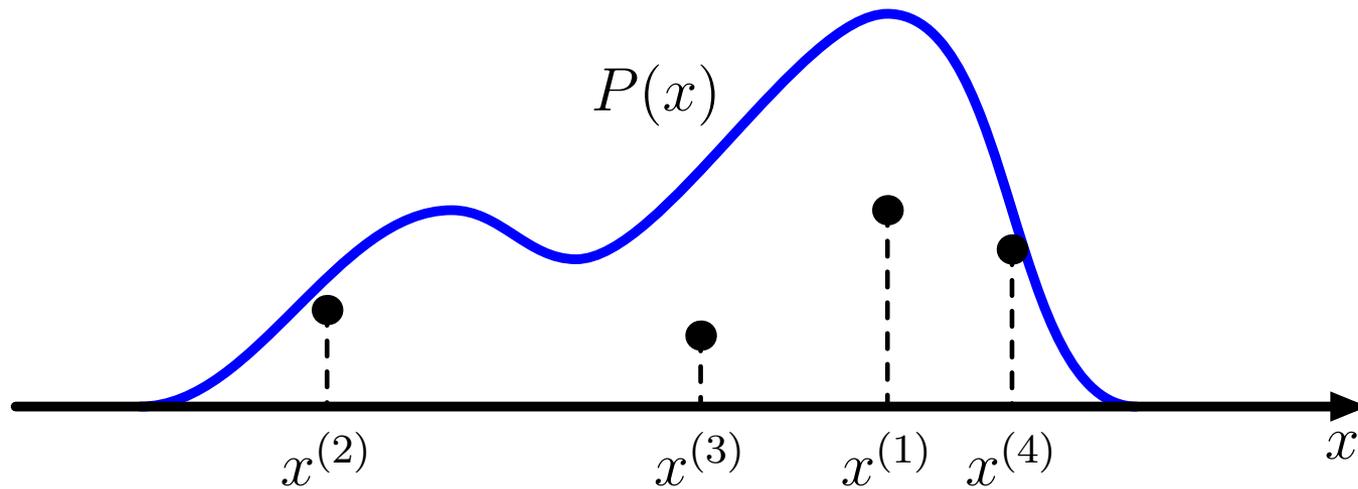
```
octave:1> 4 * quad1(@(x) sqrt(1-x.^2), 0, 1, tolerance)
```

Gives π to 6 dp's in 108 evaluations, machine precision in 2598.

(NB Matlab's `quad1` fails at zero tolerance)

Sampling from distributions

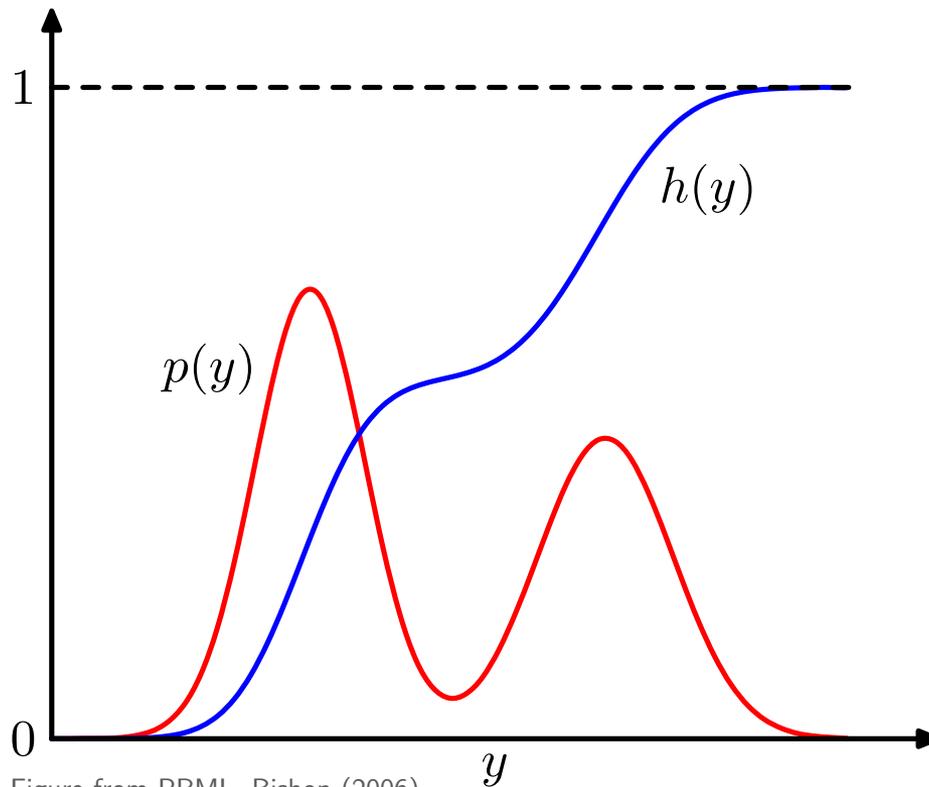
Draw points uniformly under the curve:



Probability mass to left of point \sim Uniform[0,1]

Sampling from distributions

How to convert samples from a Uniform[0,1] generator:



$$h(y) = \int_{-\infty}^y p(y') dy'$$

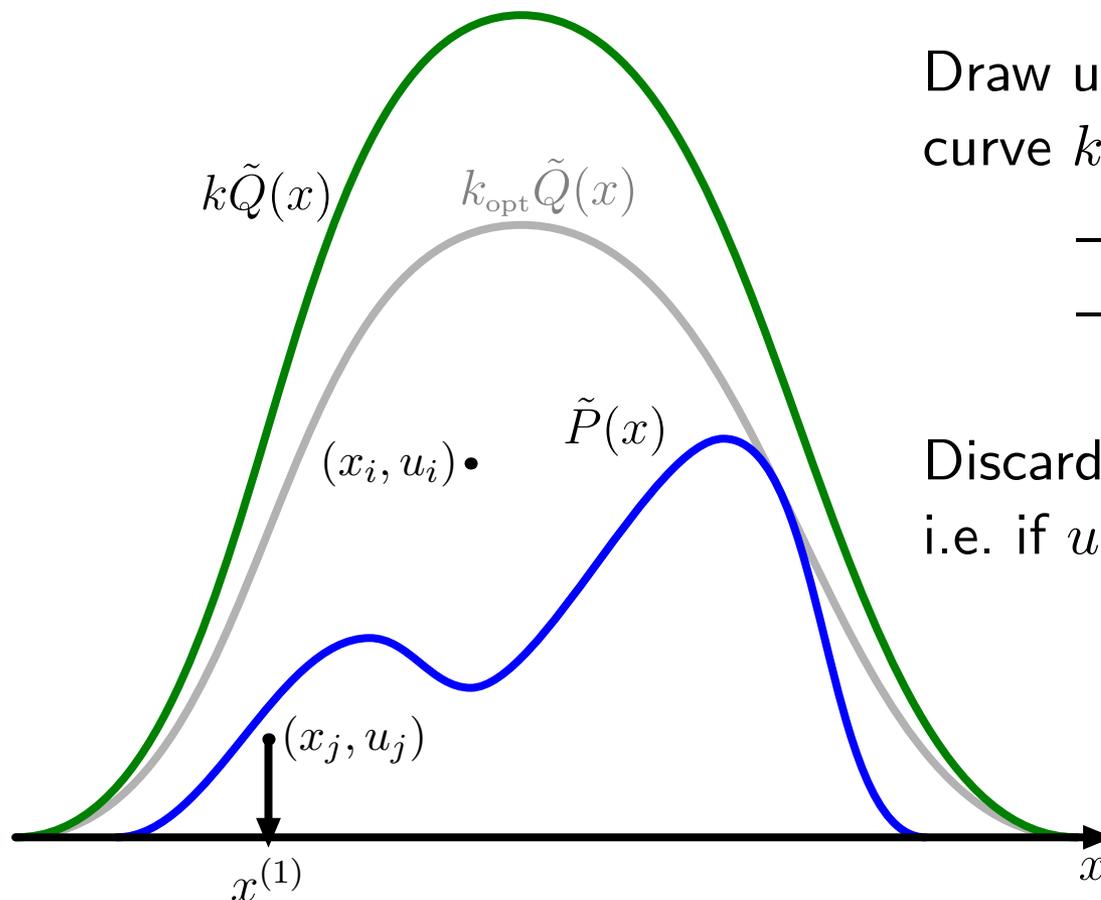
Draw mass to left of point:
 $u \sim \text{Uniform}[0,1]$

Sample, $y(u) = h^{-1}(u)$

Although we can't always compute and invert $h(y)$

Rejection sampling

Sampling underneath a $\tilde{P}(x) \propto P(x)$ curve is also valid



Draw underneath a simple curve $k\tilde{Q}(x) \geq \tilde{P}(x)$:

- Draw $x \sim Q(x)$
- height $u \sim \text{Uniform}[0, k\tilde{Q}(x)]$

Discard the point if above \tilde{P} ,
i.e. if $u > \tilde{P}(x)$

Importance sampling

Computing $\tilde{P}(x)$ and $\tilde{Q}(x)$, then *throwing x away* seems wasteful
Instead rewrite the integral as an **expectation under Q** :

$$\int f(x)P(x) dx = \int f(x)\frac{P(x)}{Q(x)}Q(x) dx, \quad (Q(x) > 0 \text{ if } P(x) > 0)$$
$$\approx \frac{1}{S} \sum_{s=1}^S f(x^{(s)})\frac{P(x^{(s)})}{Q(x^{(s)})}, \quad x^{(s)} \sim Q(x)$$

This is just simple Monte Carlo again, so it is unbiased.

Importance sampling applies when the integral is not an expectation.
Divide and multiply any integrand by a convenient distribution.

Importance sampling (2)

Previous slide assumed we could evaluate $P(x) = \tilde{P}(x)/\mathcal{Z}_P$

$$\int f(x)P(x) dx \approx \frac{\mathcal{Z}_Q}{\mathcal{Z}_P} \frac{1}{S} \sum_{s=1}^S f(x^{(s)}) \underbrace{\frac{\tilde{P}(x^{(s)})}{\tilde{Q}(x^{(s)})}}_{\tilde{r}^{(s)}}, \quad x^{(s)} \sim Q(x)$$

$$\approx \frac{1}{S} \sum_{s=1}^S f(x^{(s)}) \frac{\tilde{r}^{(s)}}{\frac{1}{S} \sum_{s'} \tilde{r}^{(s')}} \equiv \sum_{s=1}^S f(x^{(s)}) w^{(s)}$$

This estimator is **consistent** but **biased**

Exercise: Prove that $\mathcal{Z}_P/\mathcal{Z}_Q \approx \frac{1}{S} \sum_s \tilde{r}^{(s)}$

Summary so far

- Sums and integrals, often expectations, occur frequently in statistics
- **Monte Carlo** approximates expectations with a sample average
- **Rejection sampling** draws samples from complex distributions
- **Importance sampling** applies Monte Carlo to 'any' sum/integral

Pitfalls of Monte Carlo

Rejection & importance sampling scale badly with dimensionality

Example:

$$P(x) = \mathcal{N}(0, \mathbb{I}), \quad Q(x) = \mathcal{N}(0, \sigma^2 \mathbb{I})$$

Rejection sampling:

Requires $\sigma \geq 1$. Fraction of proposals accepted = σ^{-D}

Importance sampling:

Variance of importance weights = $\left(\frac{\sigma^2}{2-1/\sigma^2}\right)^{D/2} - 1$

Infinite / undefined variance if $\sigma \leq 1/\sqrt{2}$

Outline

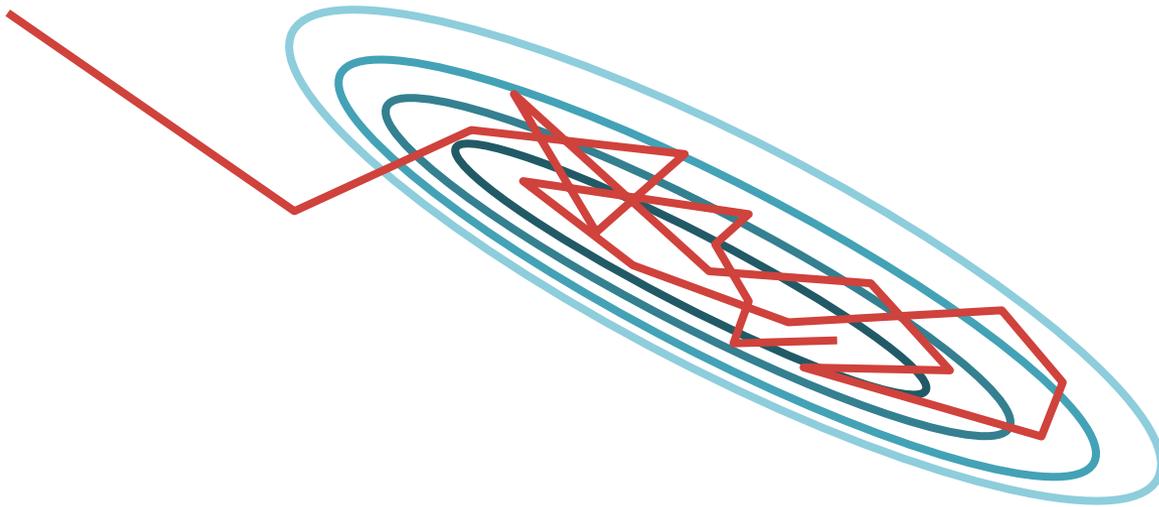
- **Review: Monte Carlo**
- **MCMC (Basic Methods)**
 - Metropolis algorithm
 - Metropolis-Hastings (M-H) algorithm
 - Gibbs Sampling
- **Markov Chains**
 - Transition probabilities
 - Invariant distribution
 - Equilibrium distribution
 - Markov chain as a WFSM
 - Constructing Markov chains
 - Why does M-H work?
- **MCMC (Auxiliary Variable Methods)**
 - Slice Sampling
 - Hamiltonian Monte Carlo

Metropolis, Metropolis-Hastings, Gibbs Sampling

MCMC (BASIC METHODS)

MCMC

- **Goal:** Draw approximate, correlated samples from a target distribution $p(x)$
- **MCMC:** Performs a biased random walk to explore the distribution



Simulations of MCMC

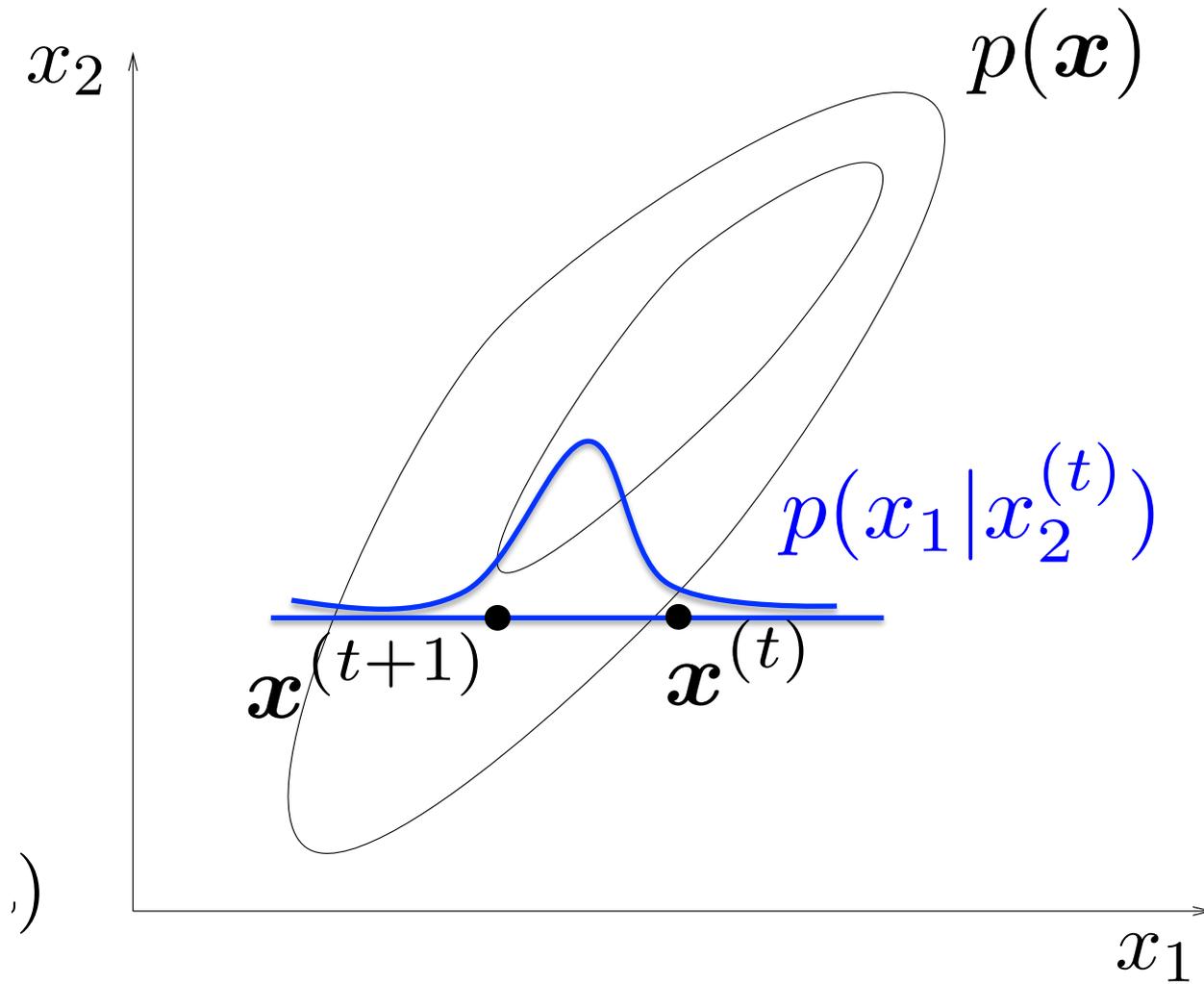
Visualization of Metropolis-Hastings, Gibbs Sampling, and Hamiltonian MCMC:

<http://twiecki.github.io/blog/2014/01/02/visualizing-mcmc/>

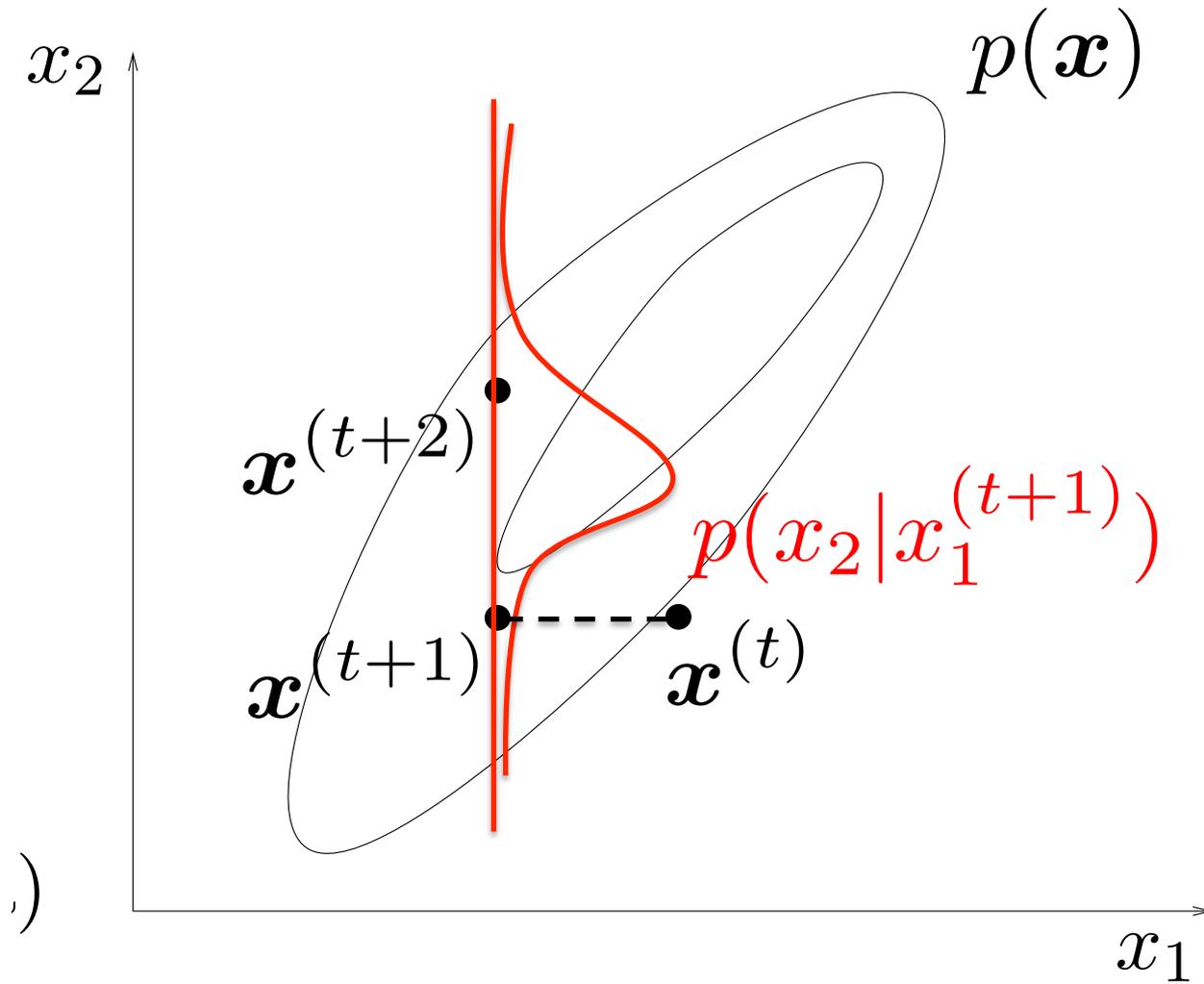
Whiteboard

- Metropolis Algorithm
- Metropolis-Hastings Algorithm
- Gibbs Sampling

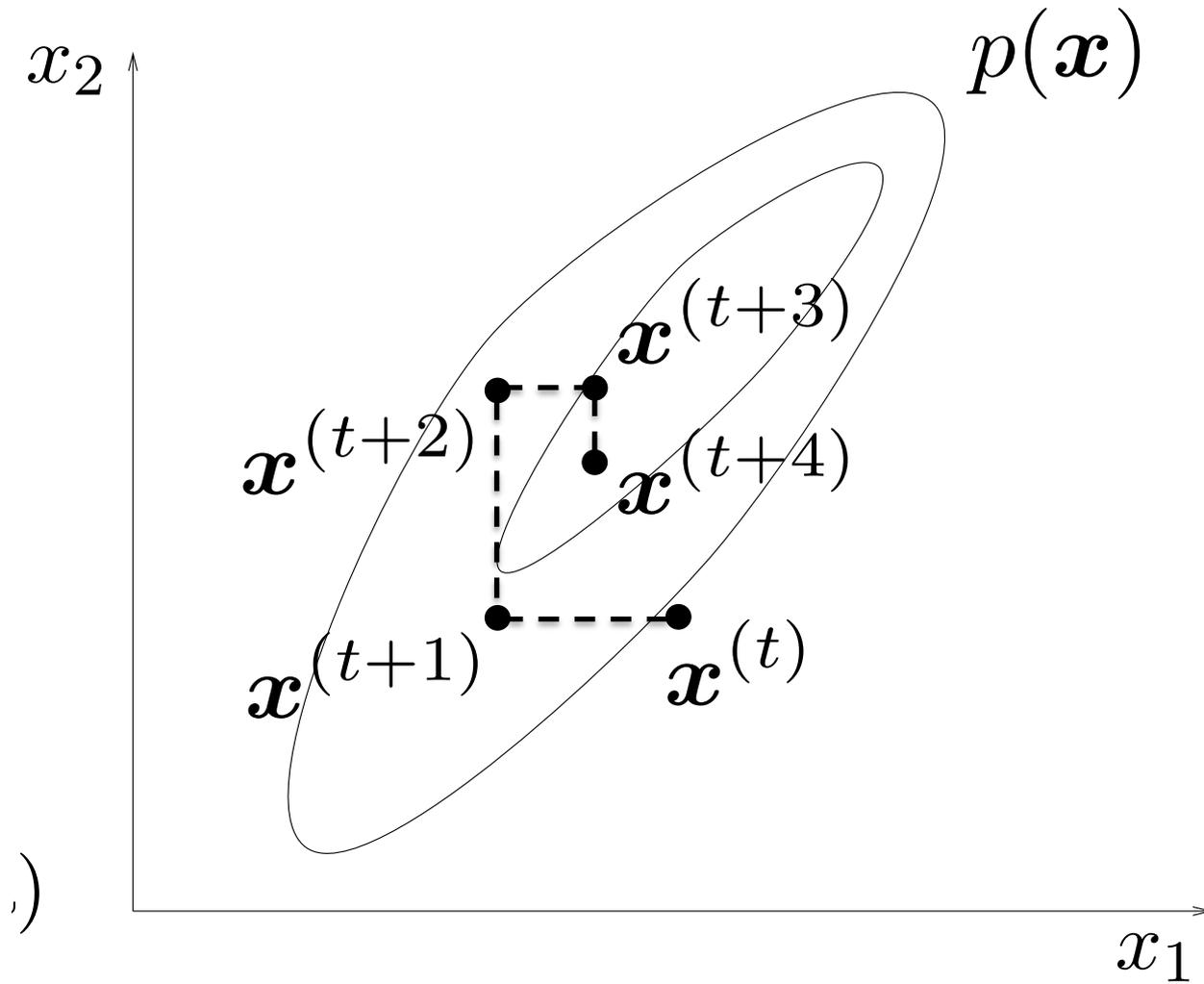
Gibbs Sampling



Gibbs Sampling



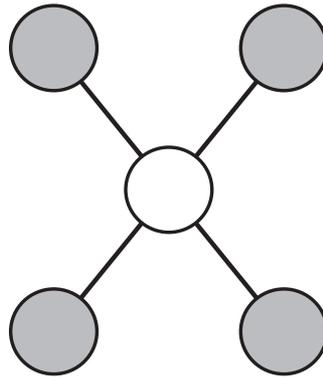
Gibbs Sampling



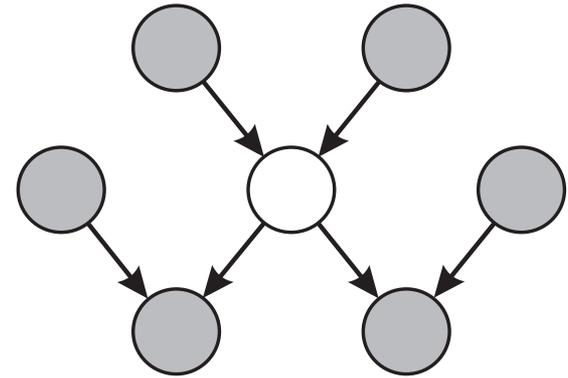
Gibbs Sampling

Full conditionals only need to condition on the Markov Blanket

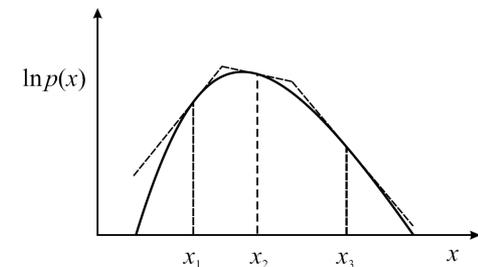
MRF



Bayes Net



- Must be “easy” to sample from conditionals
- Many conditionals are log-concave and are amenable to adaptive rejection sampling

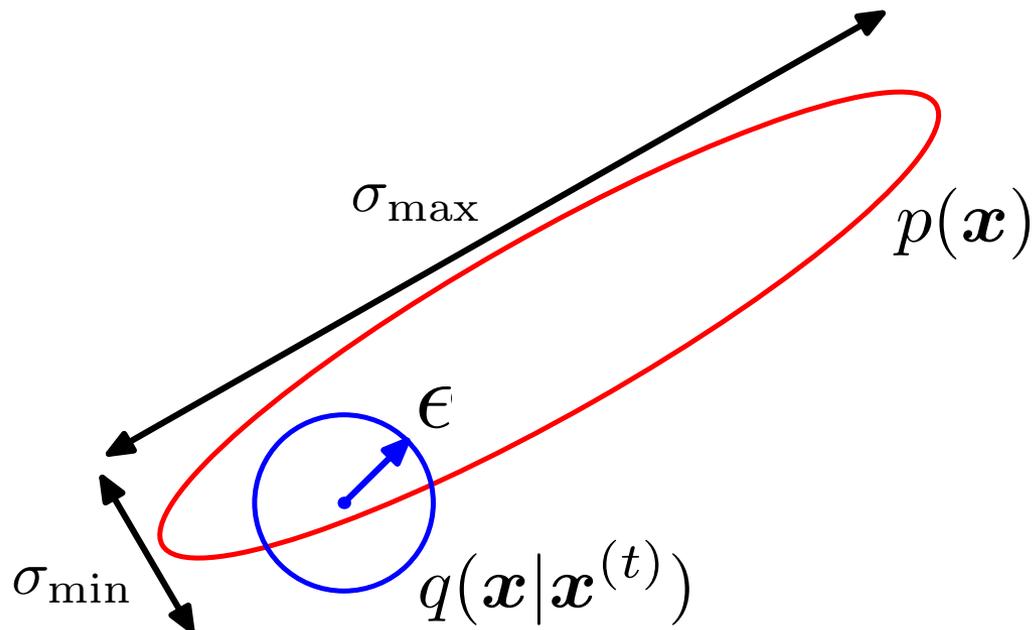


Whiteboard

- Gibbs Sampling as M-H

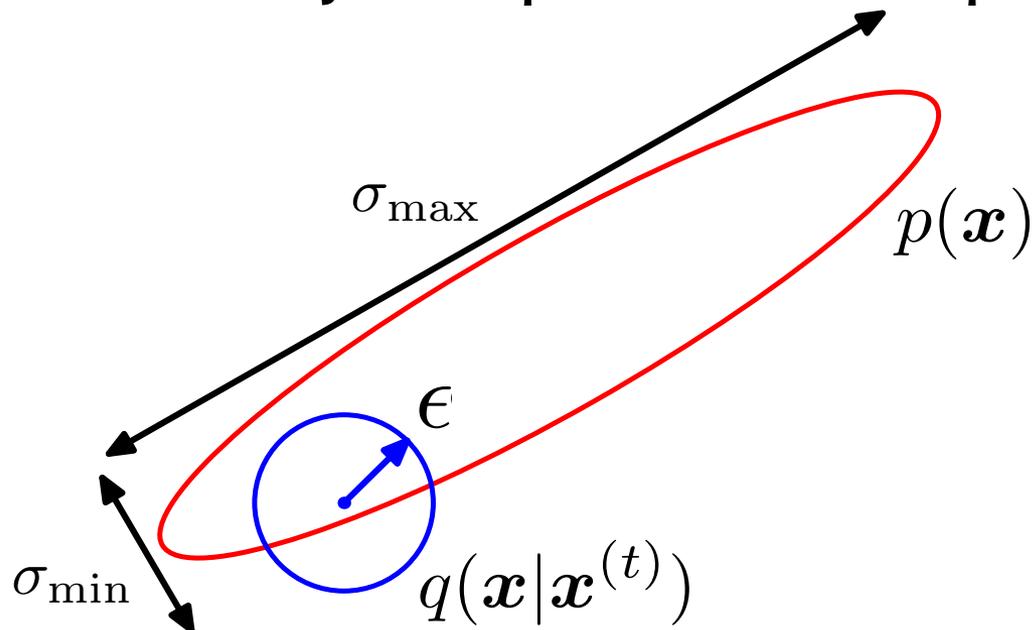
Random Walk Behavior of M-H

- For **Metropolis-Hastings**, a generic proposal distribution is: $q(\mathbf{x}|\mathbf{x}^{(t)}) = \mathcal{N}(\mathbf{0}, \epsilon^2)$
- If ϵ is large, many rejections
- If ϵ is small, slow mixing



Random Walk Behavior of M-H

- For **Rejection Sampling**, the accepted samples are **independent**
- But for **Metropolis-Hastings**, the samples are **correlated**
- **Question:** How long must we wait to get effectively independent samples?



A: independent states in the M-H random walk are separated by roughly $(\sigma_{\max}/\sigma_{\min})^2$ steps

Definitions and Theoretical Justification for MCMC

MARKOV CHAINS

Whiteboard

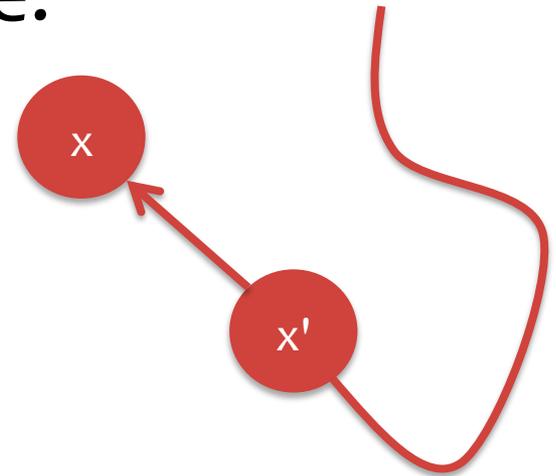
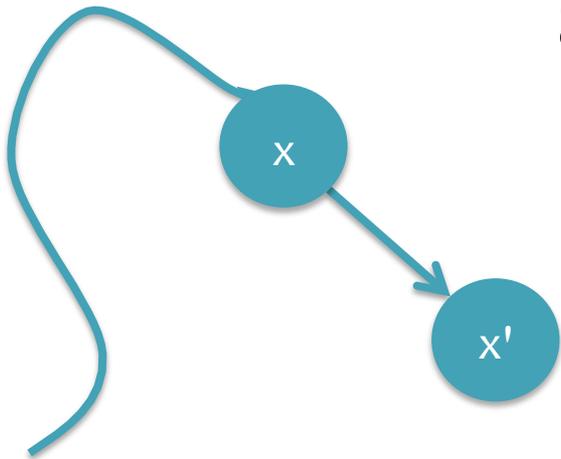
- Markov chains
- Transition probabilities
- Invariant distribution
- Equilibrium distribution
- Sufficient conditions for MCMC
- Markov chain as a WFSM

Detailed Balance

$$S(x' \leftarrow x)p(x) = S(x \leftarrow x')p(x')$$

Detailed balance means that, for each pair of states x and x' ,

arriving at x then x' and arriving at x' then x are equiprobable.



Whiteboard

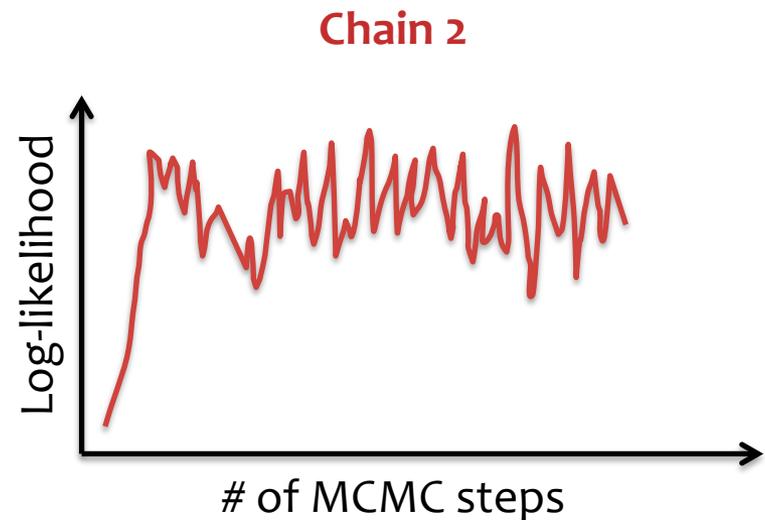
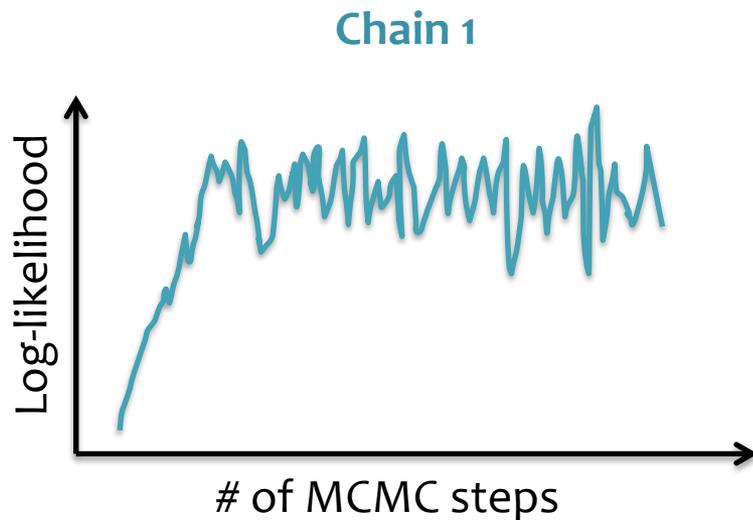
- Simple Markov chain example
- Constructing Markov chains
- Transition Probabilities for MCMC

Practical Issues

- **Question:** Is it better to move along one dimension or many?
- **Answer:** For **Metropolis-Hastings**, it is sometimes better to sample one dimension at a time
 - Q: Given a sequence of 1D proposals, compare rate of movement for **one-at-a-time** vs. **concatenation**.
- **Answer:** For **Gibbs Sampling**, sometimes better to sample a block of variables at a time
 - Q: When is it tractable to sample a block of variables?

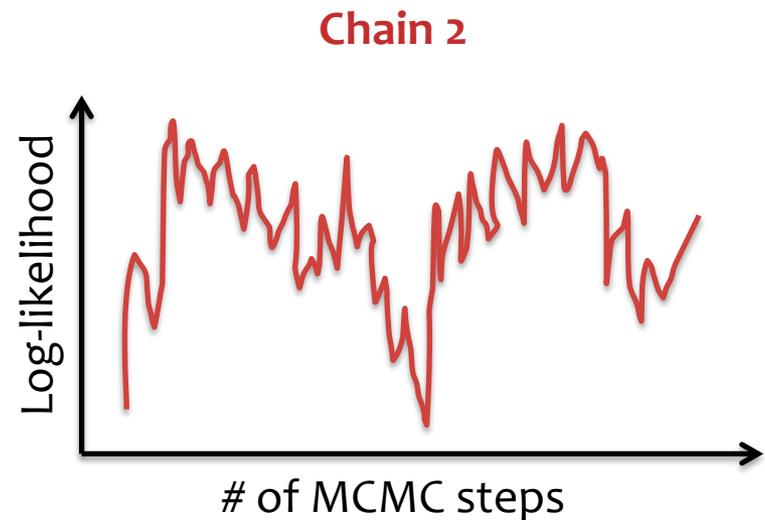
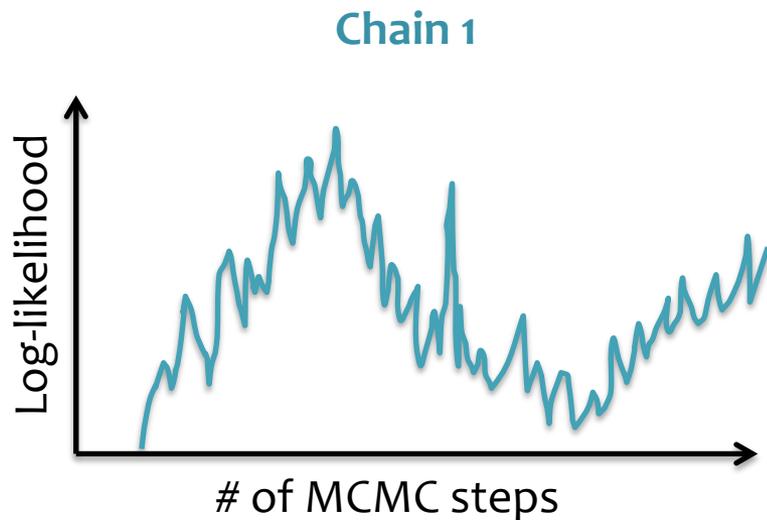
Practical Issues

- **Question:** How do we assess convergence of the Markov chain?
- **Answer:** It's not easy!
 - Compare statistics of multiple independent chains
 - Ex: Compare log-likelihoods



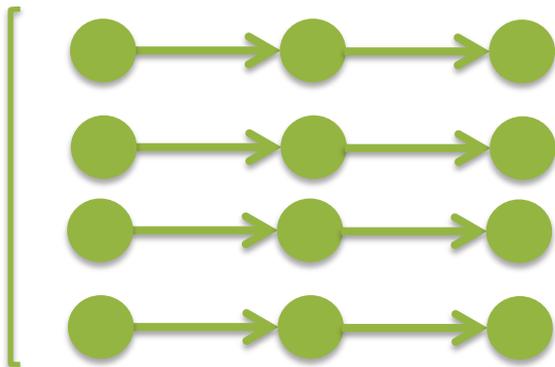
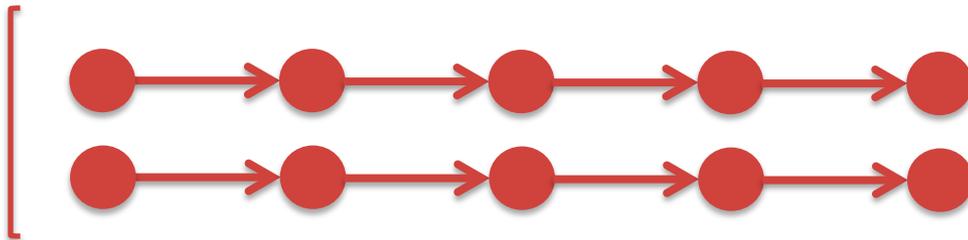
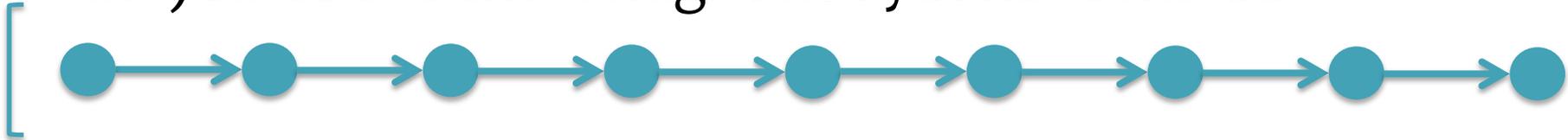
Practical Issues

- **Question:** How do we assess convergence of the Markov chain?
- **Answer:** It's not easy!
 - Compare statistics of multiple independent chains
 - Ex: Compare log-likelihoods



Practical Issues

- **Question:** Is one long Markov chain better than many short ones?
- **Note:** typical to discard initial samples (aka. “burn-in”) since the chain might not yet have mixed



- **Answer:** Often a balance is best:
 - Compared to one long chain: More independent samples
 - Compared to many small chains: Less samples discarded for burn-in
 - We can still parallelize
 - Allows us to assess mixing by comparing chains

Whiteboard

- Blocked Gibbs Sampling

Slice Sampling, Hamiltonian Monte Carlo

MCMC (AUXILIARY VARIABLE METHODS)

Auxiliary variables

The point of MCMC is to marginalize out variables, but one can introduce more variables:

$$\int f(x)P(x) dx = \int f(x)P(x, v) dx dv$$
$$\approx \frac{1}{S} \sum_{s=1}^S f(x^{(s)}), \quad x, v \sim P(x, v)$$

We might want to do this if

- $P(x|v)$ and $P(v|x)$ are simple
- $P(x, v)$ is otherwise easier to navigate

Slice Sampling

- Motivation:
 - Want **samples** from $p(x)$ and don't know the normalizer Z
 - Choosing a proposal at the correct **scale** is difficult
- Properties:
 - Similar to *Gibbs Sampling*: **one-dimensional** transitions in the state space
 - Similar to *Rejection Sampling*: (asymptotically) draws samples from the **region under the curve**

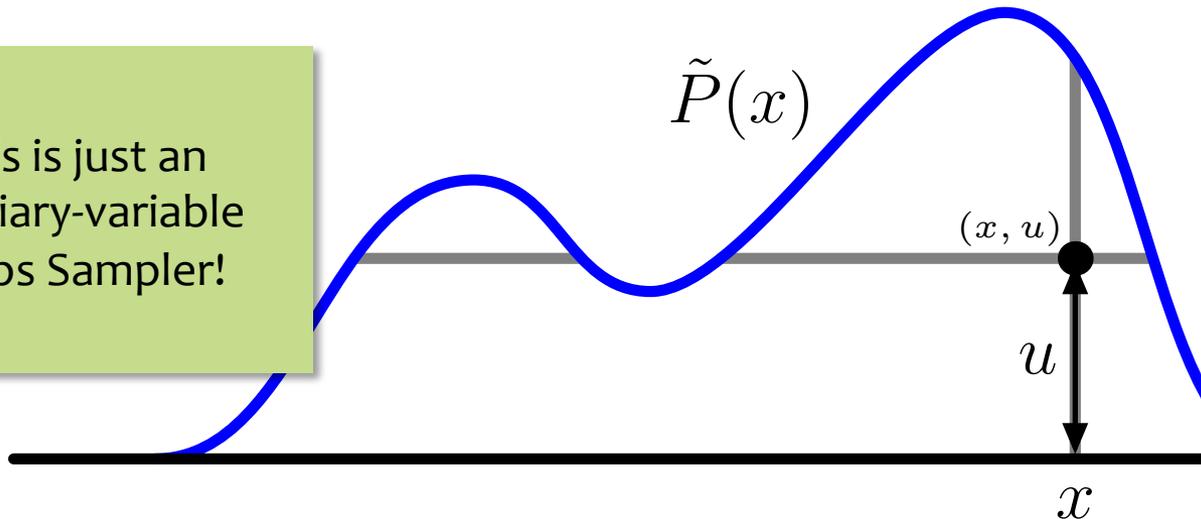


- An MCMC method with an **adaptive proposal**

Slice sampling idea

Sample point uniformly under curve $\tilde{P}(x) \propto P(x)$

This is just an
auxiliary-variable
Gibbs Sampler!

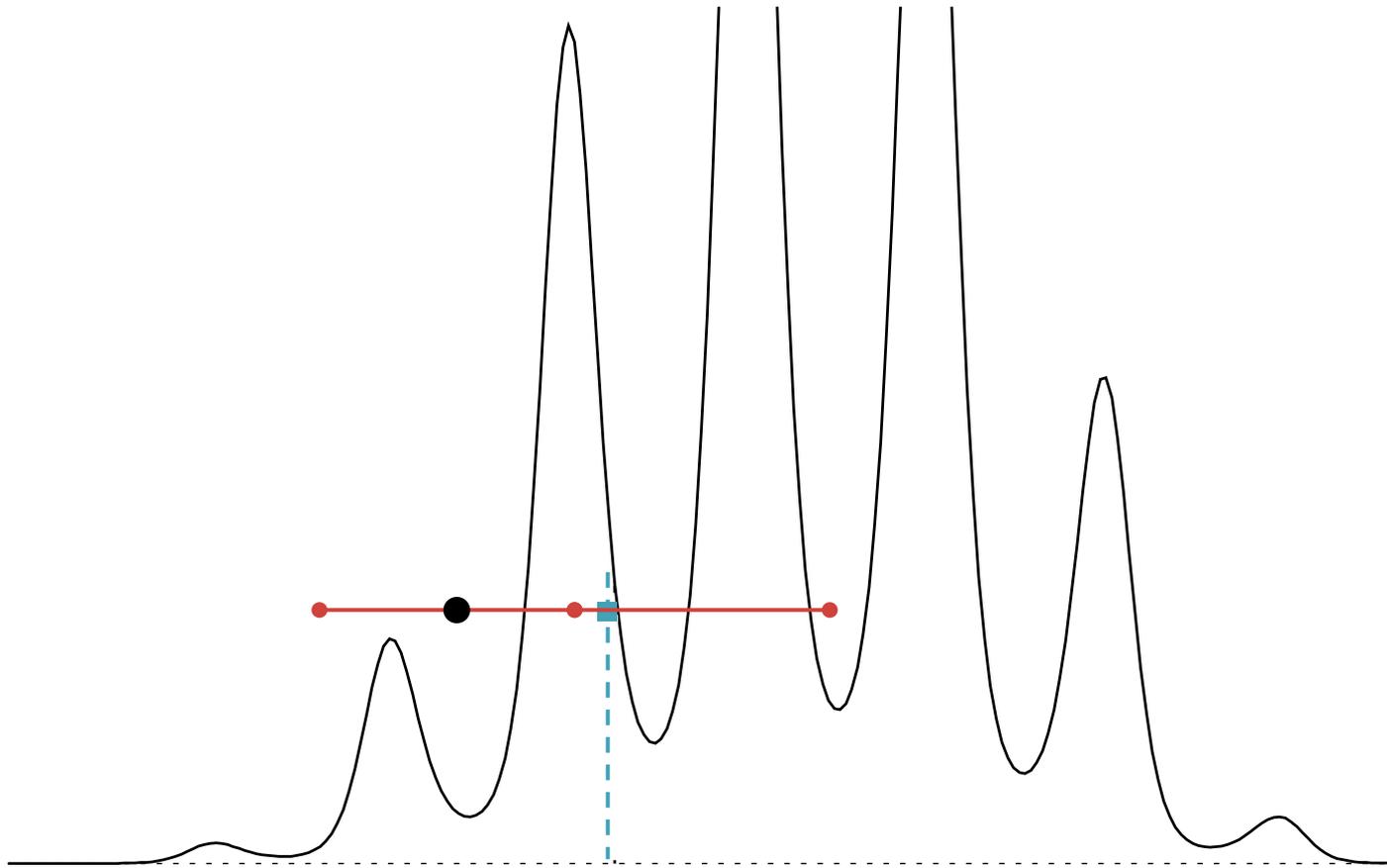


Problem: Sampling
from the conditional
 $p(x | u)$ might be
infeasible.

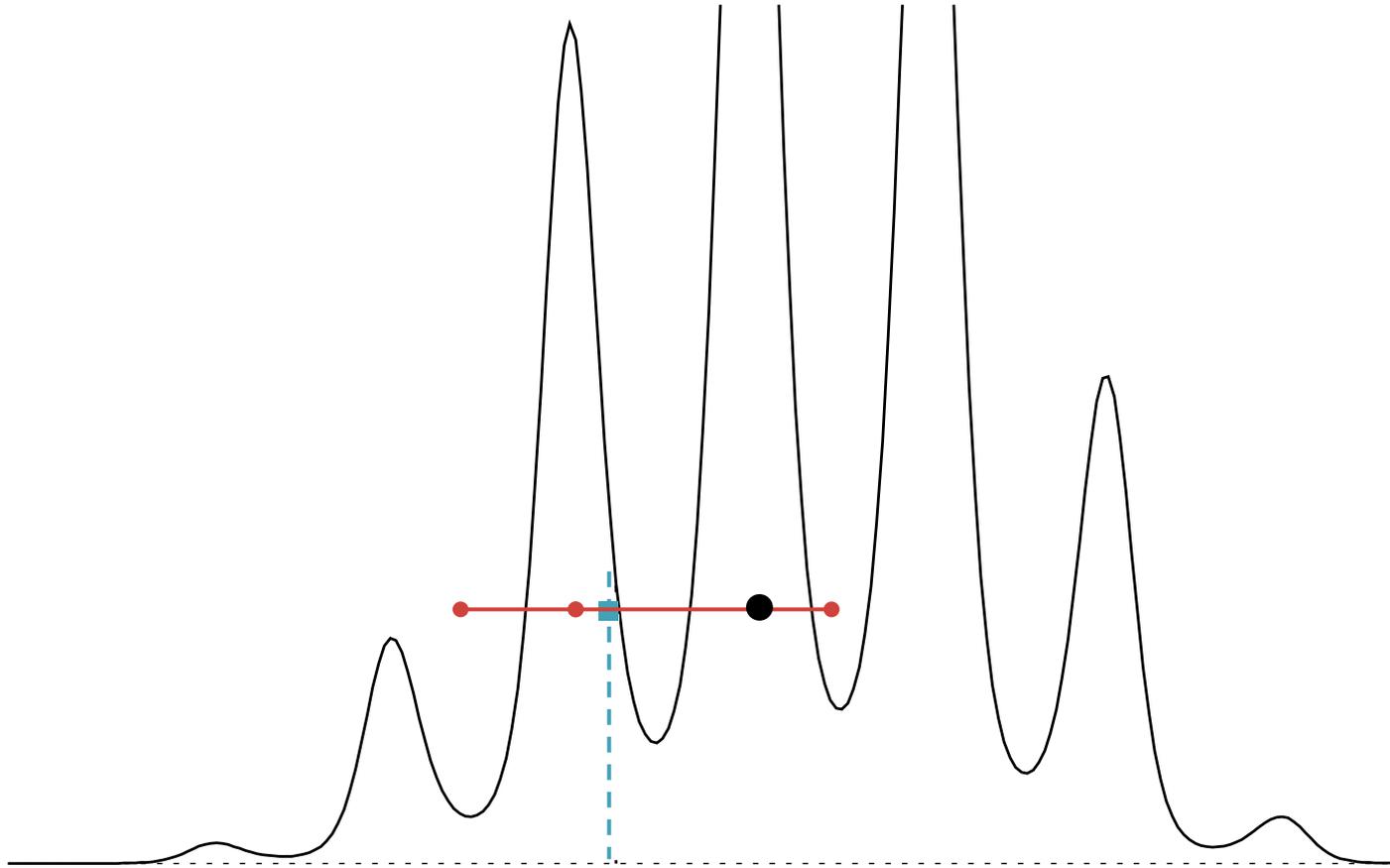
$$p(u|x) = \text{Uniform}[0, \tilde{P}(x)]$$

$$p(x|u) \propto \begin{cases} 1 & \tilde{P}(x) \geq u \\ 0 & \text{otherwise} \end{cases} = \text{"Uniform on the slice"}$$

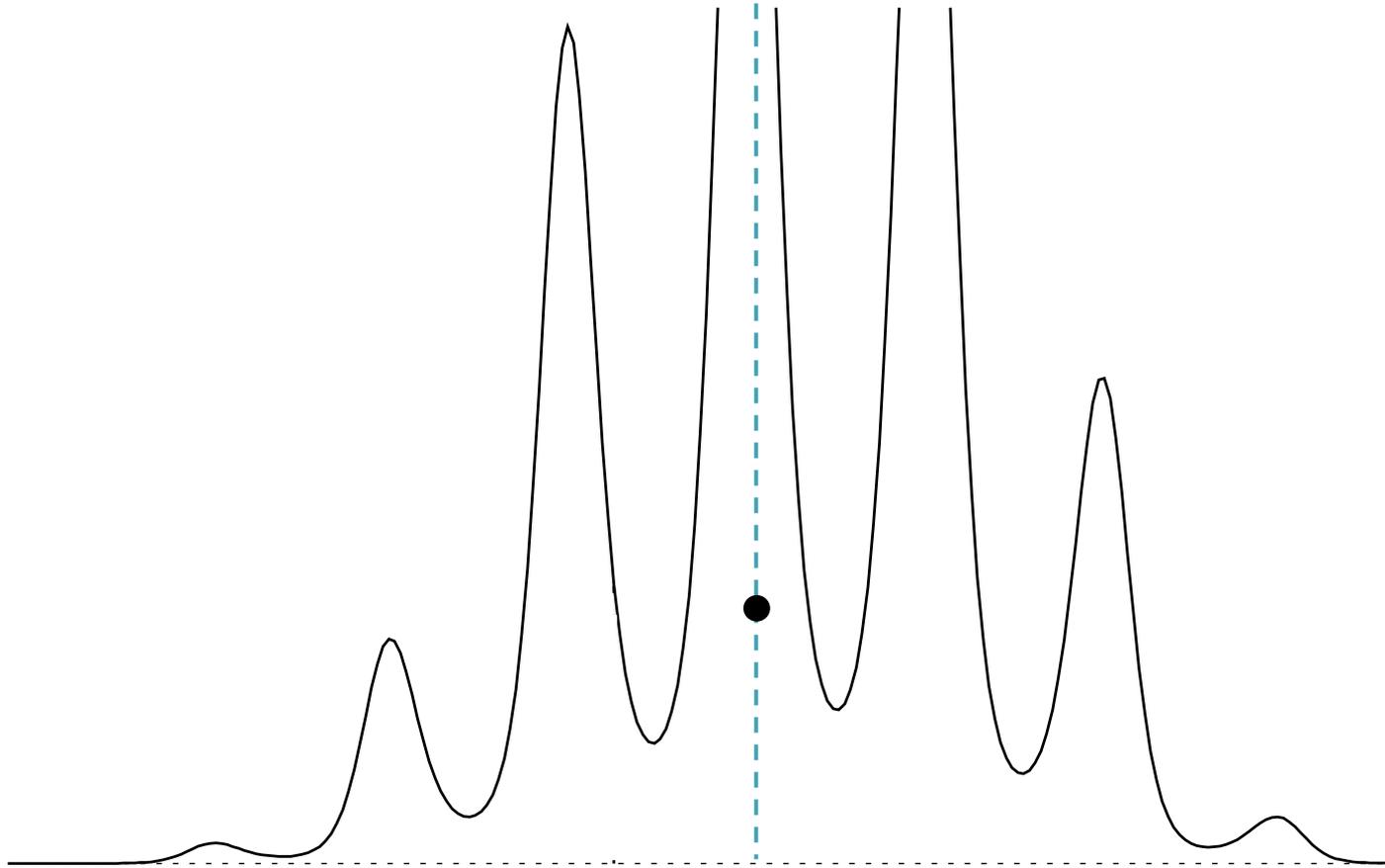
Slice Sampling



Slice Sampling



Slice Sampling



Slice Sampling

Goal: sample (x, u) given $(u^{(t)}, x^{(t)})$.

Part 1: Stepping Out

Sample interval (x_l, x_r) enclosing $x^{(t)}$.

Expand until endpoints are "outside" region under curve.

Part 2: Sample x (Shrinking)

Draw x from within the interval (x_l, x_r) , then accept or shrink.

Algorithm:

Slice Sampling

Goal: sample (x, u) given $(u^{(t)}, x^{(t)})$.

$u \sim \text{Uniform}(0, p(x^{(t)}))$

Part 1: Stepping Out

Sample interval (x_l, x_r) enclosing $x^{(t)}$.

$r \sim \text{Uniform}(u, w)$

$(x_l, x_r) = (x^{(t)} - r, x^{(t)} + w - r)$

Expand until endpoints are "outside" region under curve.

while($\tilde{p}(x_l) > u$) $\{x_l = x_l - w\}$

while($\tilde{p}(x_r) > u$) $\{x_r = x_r + w\}$

Part 2: Sample x (Shrinking)

Draw x from within the interval (x_l, x_r) , then accept or shrink.

Algorithm:

Slice Sampling

Goal: sample (x, u) given $(u^{(t)}, x^{(t)})$.

$u \sim \text{Uniform}(0, p(x^{(t)}))$

Part 1: Stepping Out

Sample interval (x_l, x_r) enclosing $x^{(t)}$.

$r \sim \text{Uniform}(u, w)$

$(x_l, x_r) = (x^{(t)} - r, x^{(t)} + w - r)$

Expand until endpoints are "outside" region under curve.

while($\tilde{p}(x_l) > u$) { $x_l = x_l - w$ }

while($\tilde{p}(x_r) > u$) { $x_r = x_r + w$ }

Part 2: Sample x (Shrinking)

while(true) {

Draw x from within the interval (x_l, x_r) , then accept or shrink.

$x \sim \text{Uniform}(x_l, x_r)$

if($\tilde{p}(x) > u$) { break }

else if($x > x^{(t)}$) { $x_r = x$ }

else { $x_l = x$ }

}

$x^{(t+1)} = x, u^{(t+1)} = u$

Algorithm:

Slice Sampling

Multivariate Distributions

- Resample each variable x_i **one-at-a-time** (just like Gibbs Sampling)
- Does not require sampling from
$$p(x_i | \{x_j\}_{j \neq i})$$
- Only need to evaluate a quantity **proportional** to the conditional

$$p(x_i | \{x_j\}_{j \neq i}) \propto \tilde{p}(x_i | \{x_j\}_{j \neq i})$$

Hamiltonian Monte Carlo

- Suppose we have a distribution of the form:

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

where $\mathbf{x} \in \mathcal{R}^N$

- We could use **random-walk M-H** to draw samples, but it seems a shame to **discard gradient information** $\nabla_{\mathbf{x}}E(\mathbf{x})$
- If we can evaluate it, the gradient tells us where to look for **high-probability regions!**

Background: Hamiltonian Dynamics

Applications:

- Following the motion of atoms in a fluid through time
- Integrating the motion of a solar system over time
- Considering the evolution of a galaxy (i.e. the motion of its stars)
- “molecular dynamics”
- “N-body simulations”

Properties:

- Total energy of the system $H(x,p)$ stays constant
- Dynamics are reversible

Important for
detailed balance

Background: Hamiltonian Dynamics

Let $\mathbf{x} \in \mathcal{R}^N$ be a position

$\mathbf{p} \in \mathcal{R}^N$ be a momentum

Potential energy: $E(\mathbf{x})$

Kinetic energy: $K(\mathbf{p}) = \mathbf{p}^T \mathbf{p} / 2$

Total energy: $H(\mathbf{x}, \mathbf{p}) = E(\mathbf{x}) + K(\mathbf{p})$



Hamiltonian function

Given a starting position $x^{(1)}$ and a starting momentum $p^{(1)}$ we can simulate the Hamiltonian dynamics of the system via:

1. Euler's method
2. Leapfrog method
3. etc.

Background: Hamiltonian Dynamics

Parameters to tune:

1. Step size, ϵ
2. Number of iterations, L

Leapfrog Algorithm:

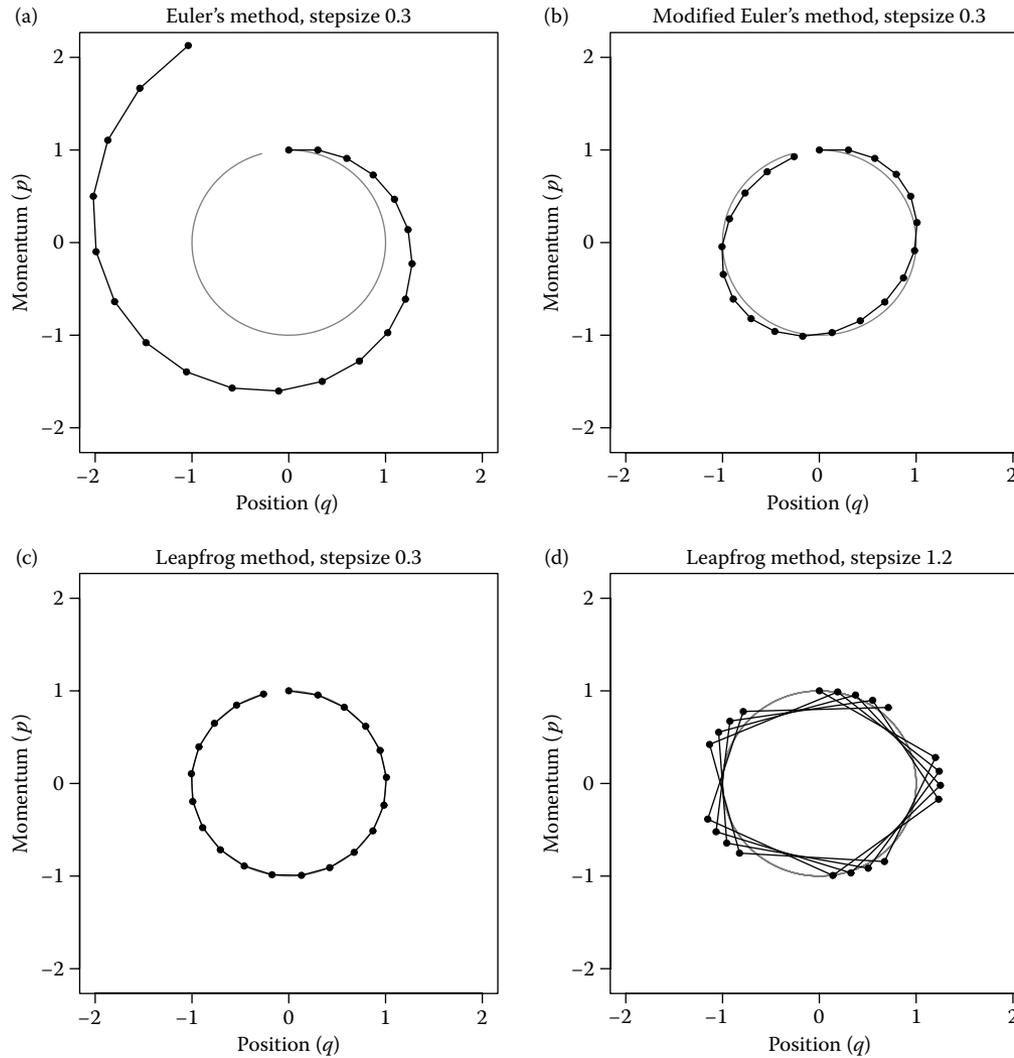
for τ in $1 \dots L$:

$$\mathbf{p} = \mathbf{p} - \frac{\epsilon}{2} \nabla_{\mathbf{x}} E(\mathbf{x})$$

$$\mathbf{x} = \mathbf{x} + \epsilon \mathbf{p}$$

$$\mathbf{p} = \mathbf{p} - \frac{\epsilon}{2} \nabla_{\mathbf{x}} E(\mathbf{x})$$

Background: Hamiltonian Dynamics



Hamiltonian Monte Carlo

Preliminaries

Goal: $p(\mathbf{x}) = \exp\{-E(\mathbf{x})\}/Z$ where $\mathbf{x} \in \mathcal{R}^N$

Define: $K(\mathbf{p}) = \mathbf{p}^T \mathbf{p}/2$

$$H(\mathbf{x}, \mathbf{p}) = E(\mathbf{x}) + K(\mathbf{p})$$

$$\begin{aligned} p(\mathbf{x}, \mathbf{p}) &= \exp\{-H(\mathbf{x}, \mathbf{p})\}/Z_H \\ &= \exp\{-E(\mathbf{x})\} \exp\{-K(\mathbf{p})\}/Z_H \end{aligned}$$

Note:

Since $p(\mathbf{x}, \mathbf{p})$ is separable...

$$\Rightarrow \sum_{\mathbf{p}} p(\mathbf{x}, \mathbf{p}) = \exp\{-E(\mathbf{x})\}/Z$$

Target dist.

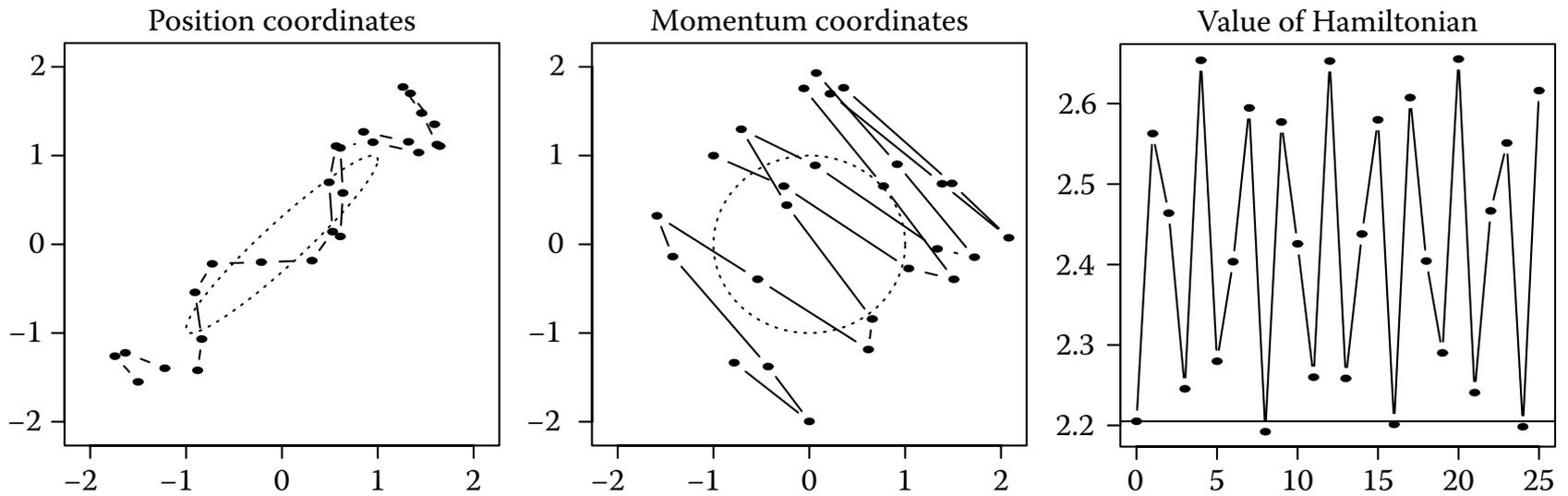
$$\Rightarrow \sum_{\mathbf{x}} p(\mathbf{x}, \mathbf{p}) = \exp\{-K(\mathbf{p})\}/Z_K$$

Gaussian

Whiteboard

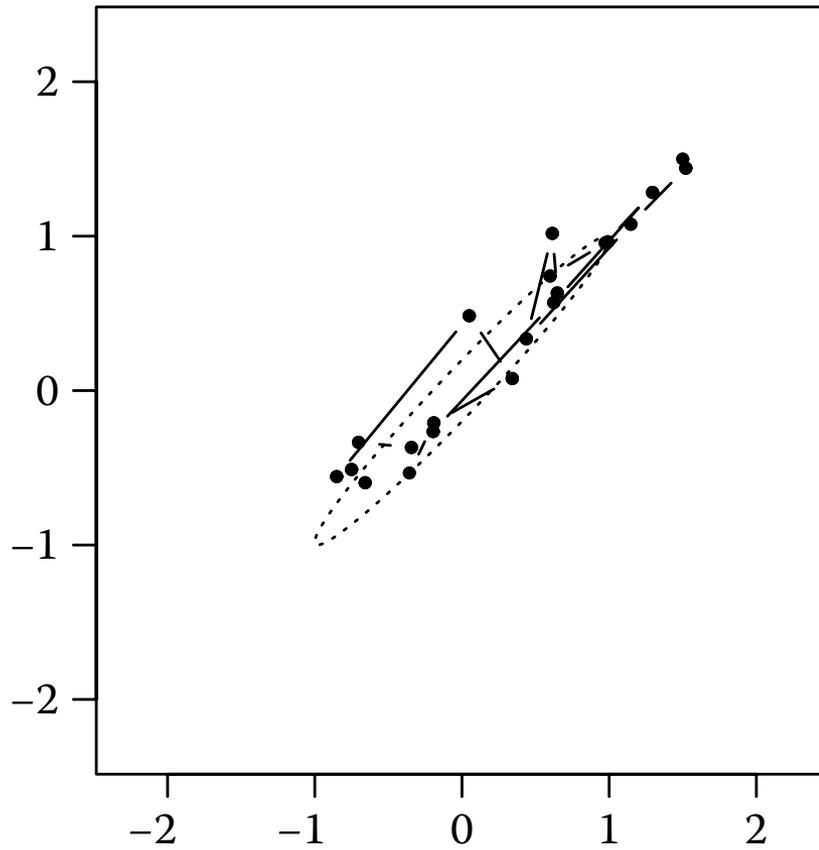
- Hamiltonian Monte Carlo algorithm
(aka. Hybrid Monte Carlo)

Hamiltonian Monte Carlo

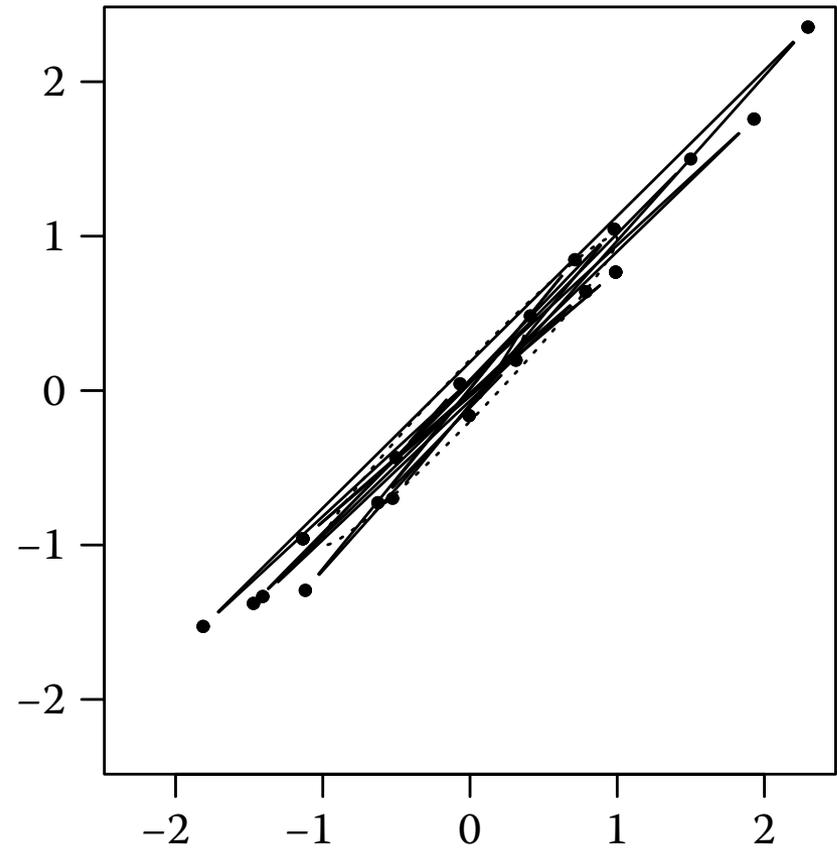


M-H vs. HMC

Random-walk Metropolis



Hamiltonian Monte Carlo



Simulations of MCMC

Visualization of Metropolis-Hastings, Gibbs Sampling, and Hamiltonian MCMC:

<http://twiecki.github.io/blog/2014/01/02/visualizing-mcmc/>

MCMC Summary

- **Pros**
 - Very general purpose
 - Often easy to implement
 - Good theoretical guarantees as $t \rightarrow \infty$
- **Cons**
 - Lots of tunable parameters / design choices
 - Can be quite slow to converge
 - Difficult to tell whether it's working