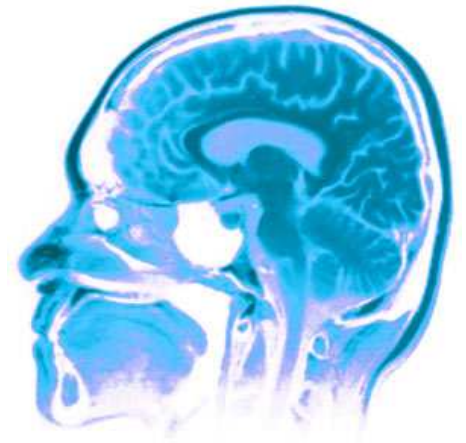




CPSC 540



Importance sampling &
Markov chain Monte Carlo (MCMC)



Nando de Freitas

March, 2013

University of British Columbia

Outline of the lecture

This is about Monte Carlo methods.

- We will revise importance sampling.
- Revise how Google works (Markov chains).
- Introduce Markov chain Monte Carlo (MCMC)

Bayesian logistic regression

The logistic regression model specifies the probability of a binary output $y_i \in \{0, 1\}$ given the input \mathbf{x}_i as follows:

$$\begin{aligned}
 p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) &= \prod_{i=1}^n \text{Ber}(y_i | \text{sigm}(\mathbf{x}_i \boldsymbol{\theta})) = \text{const} - \sum_{i=1}^n (y_i \log \pi_i + (1-y_i) \log(1-\pi_i)) \\
 &= \prod_{i=1}^n \left[\underbrace{\frac{1}{1 + e^{-\mathbf{x}_i \boldsymbol{\theta}}}}_{\pi_i} \right]^{y_i} \left[1 - \frac{1}{1 + e^{-\mathbf{x}_i \boldsymbol{\theta}}} \right]^{1-y_i} = \frac{1}{2\sigma^2} \|\boldsymbol{\theta}\|_2^2
 \end{aligned}$$

$\pi_i = -\log P(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}) = \text{const} - \log P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) - \log P(\boldsymbol{\theta})$

We also assume a Gaussian prior π_i

$$p(\boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (\boldsymbol{\theta} - \boldsymbol{\mu})' (\boldsymbol{\theta} - \boldsymbol{\mu})\right)$$

Posterior: $P(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}) = \frac{1}{Z} P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) P(\boldsymbol{\theta})$

$Z = \int P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) P(\boldsymbol{\theta}) d\boldsymbol{\theta}$ is unknown / hard

Importance sampling

$$z = \int P(y|\theta) P(\theta) d\theta$$

$$z = \int \frac{P(y|\theta) P(\theta) q(\theta)}{q(\theta)} d\theta$$

$$q(\theta) = \mathcal{N}(0, 1000)$$

$$z = \int w(\theta) q(\theta) d\theta$$

$$\theta^{(i)} \sim q(\theta), \quad i=1:N$$

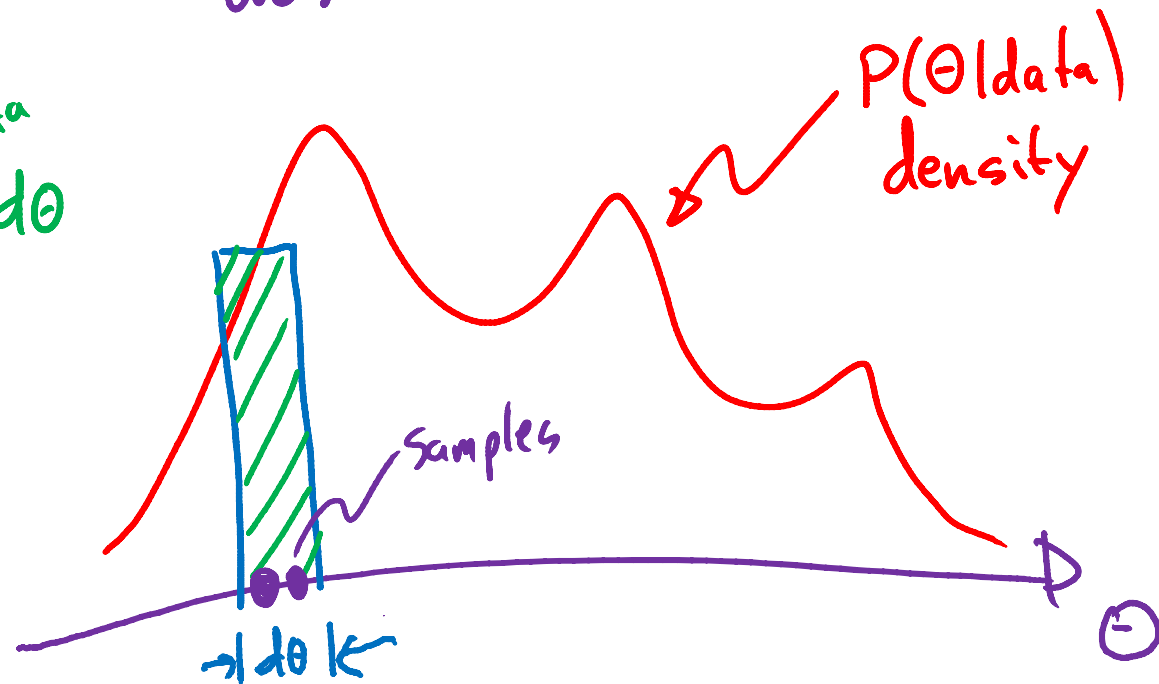
$$z \approx \frac{1}{N} \sum_{i=1}^N w(\theta^{(i)}) \quad \text{SLLN}$$

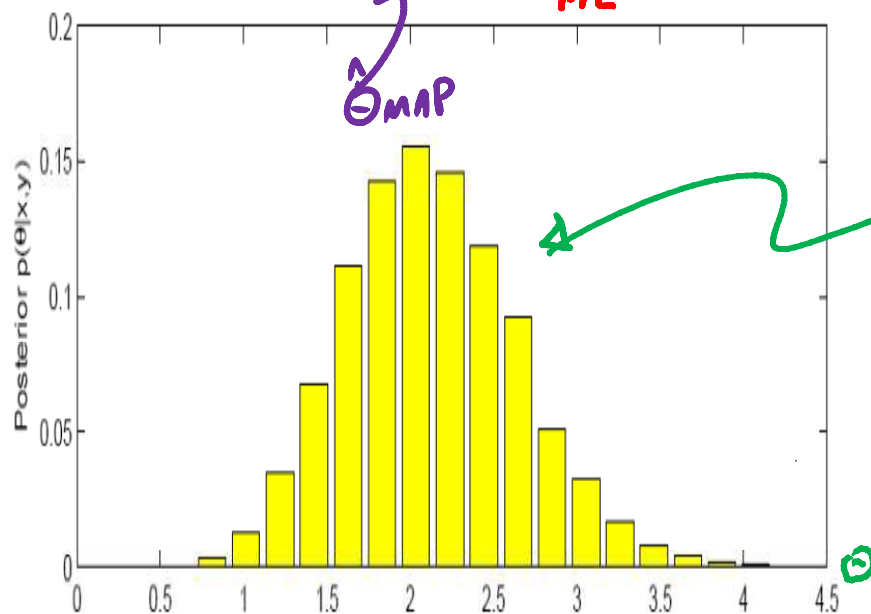
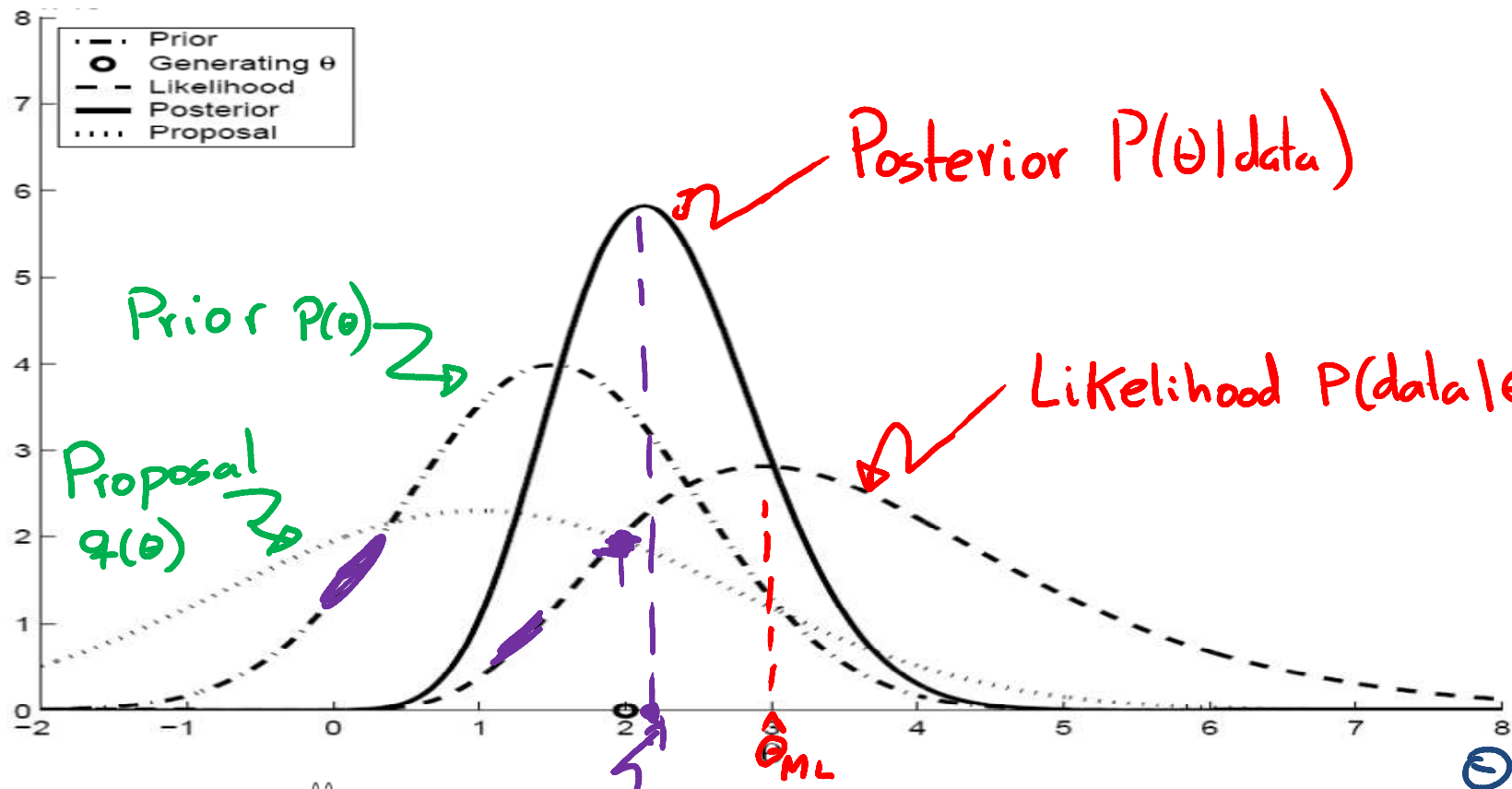
Importance sampling

$$P(\theta | \text{data}) = \frac{1}{N} \sum_{i=1}^N w(\theta^{(i)}) \delta_{\theta^{(i)}}(d\theta)$$

$\delta_{\theta^{(i)}}(d\theta)$ = Number of samples $\theta^{(i)}$ in the interval $d\theta$.

$$P(d\theta) = P(\theta) d\theta$$





Importance sampling

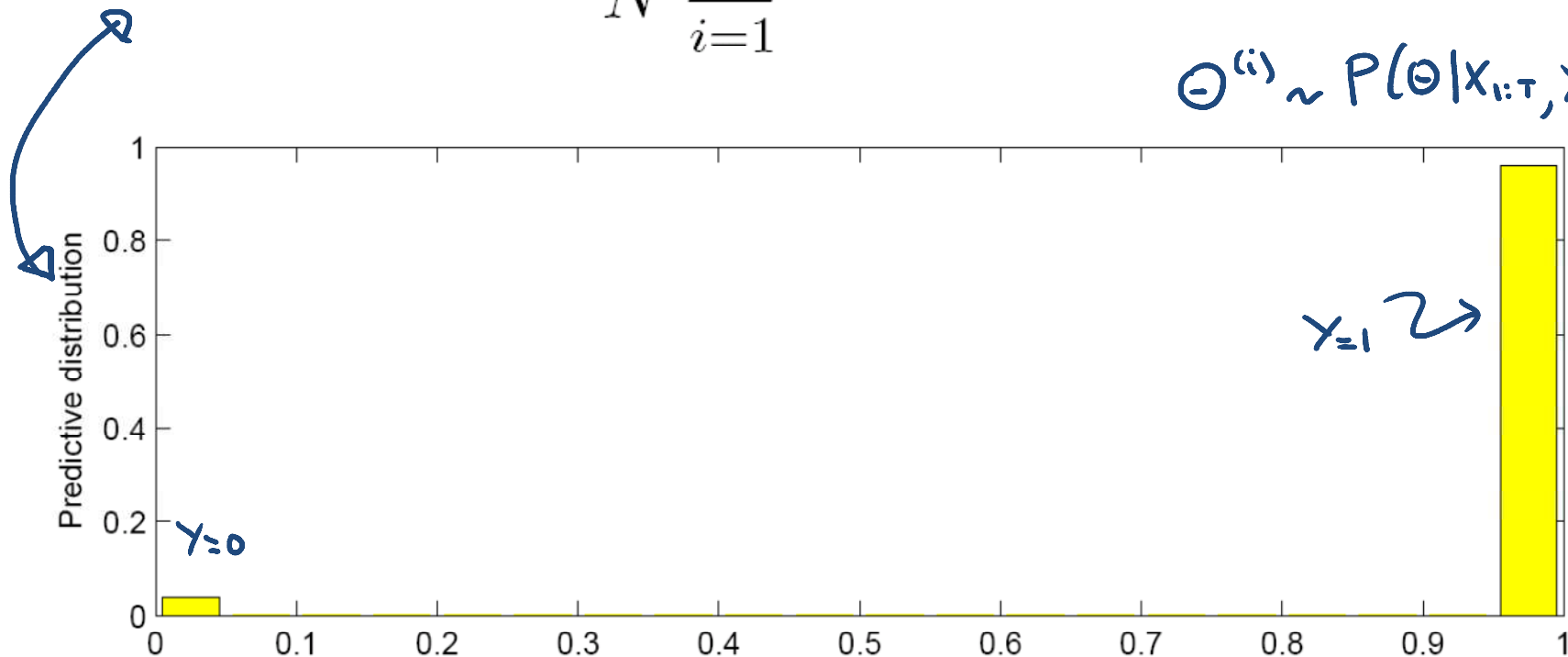
$$\begin{aligned} P(y_{t+1} | x_{t+1}, y_{1:t}, x_{1:t}) &= \int P(y_{t+1} | x_{t+1}, \theta) P(d\theta | x_{1:t}, y_{1:t}) \\ &= \int P(y_{t+1} | x_{t+1}, \theta) P(\theta | x_{1:t}, y_{1:t}) d\theta \\ &\stackrel{a.p.}{=} \int P(y_{t+1} | x_{t+1}, \theta) \frac{1}{N} \sum_{i=1}^N w(\theta^{(i)}) \delta_{\theta^{(i)}}(d\theta) \\ &\approx \frac{1}{N} \sum_{i=1}^N \int P(y_{t+1} | x_{t+1}, \theta) w(\theta^{(i)}) \delta_{\theta^{(i)}}(d\theta) \\ &\approx \frac{1}{N} \sum_{i=1}^N \underbrace{P(y_{t+1} | x_{t+1}, \theta^{(i)})}_{\text{likelihood}} w(\theta^{(i)}) \end{aligned}$$

Example: Logistic Regression

$$p(y_{T+1}|x_{1:T+1}) = \int_{\Theta} p(y_{T+1}|x_{T+1}, \theta) p(\theta|x_{1:T}, y_{1:T}) d\theta$$

$$p(y_{T+1}|x_{1:T+1}) = \frac{1}{N} \sum_{i=1}^N p(y_{T+1}|x_{T+1}, \theta^{(i)})$$

$$\Theta^{(i)} \sim P(\Theta|x_{1:T}, y_{1:T})$$



Un-normalized importance sampling

$D = \text{data}$

$$P(\theta | D) = \frac{1}{Z} P(D | \theta) P(\theta) = \frac{P(D | \theta) P(\theta)}{\int P(D | \theta) P(\theta) d\theta}$$

$$\begin{aligned} P(y_{t+1} | x_{t+1}, D) \\ = P(y_{t+1} | x_{1:t+1}, y_{1:t}) &= \frac{1}{Z} \int P(y_{t+1} | x_{t+1}, \theta) P(D | \theta) P(\theta) d\theta \\ &= \frac{\int P(y_{t+1} | x_{t+1}, \theta) P(D | \theta) P(\theta) \frac{q(\theta)}{q(\theta)} d\theta}{\int P(D | \theta) P(\theta) \frac{q(\theta)}{q(\theta)} d\theta} \end{aligned}$$



$Q(\mu, \sigma^2 I)$

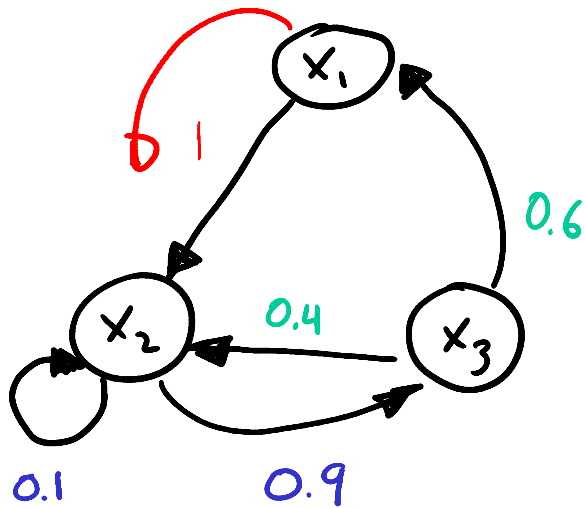
$$\tilde{w}^i = \frac{w^i}{\sum_j w^j}$$

$$\begin{aligned} &= \frac{\int P(y_{t+1} | x_{t+1}, \theta) w(\theta) q(\theta) d\theta}{\int w(\theta) q(\theta) d\theta} = \frac{\frac{1}{N} \sum_{i=1}^N w(\theta^{(i)}) P(y_{t+1} | x_{t+1}, \theta^{(i)})}{\frac{1}{N} \sum_{j=1}^N w(\theta^{(j)})} \\ &= \sum_{i=1}^N \tilde{w}(\theta^{(i)}) P(y_{t+1} | x_{t+1}, \theta^{(i)}) \end{aligned}$$

Markov Chain Monte Carlo

For simplicity, Let's consider only 3 states:

$$x_t \in \mathcal{X} = \{x_1, x_2, x_3\}$$



$$T = P(x_t | x_{t-1}) = \begin{matrix} & \begin{matrix} x_1 & x_2 & x_3 \end{matrix} \\ \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0.1 & 0.9 \\ 0.6 & 0.4 & 0 \end{bmatrix} \end{matrix}$$

Think of this as a webgraph. Our goal is to crawl it to find the "relevance" of each node.

Markov Chain Monte Carlo

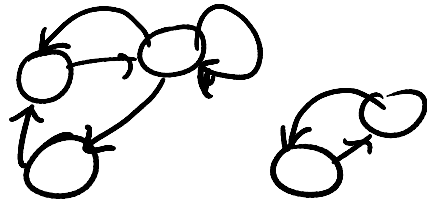
T is a stochastic matrix. As long as the graph (state space) is **aperiodic** and **irreducible**, we have that for any initial vector of Probabilities ν :

$$\underline{\nu} T^{(t)} \rightarrow \pi \quad \text{as } t \rightarrow \infty$$

Where π is the **invariant** or **Stationary** distribution of the chain. It is unique.

Markov Chain Monte Carlo

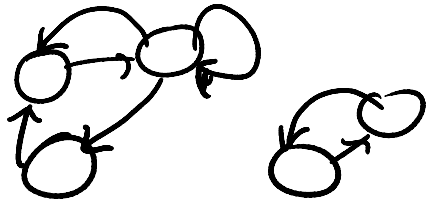
Need for irreducibility:



One cluster might
never be visited!

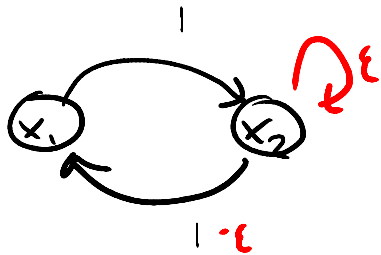
Markov Chain Monte Carlo

Need for irreducibility:



One cluster might never be visited!

Need for aperiodicity:



$$T = \begin{bmatrix} 0 & 1-\epsilon \\ 1-\epsilon & 0 \end{bmatrix}$$

$$\text{Let } \pi = \left[\frac{1}{3} \quad \frac{2}{3} \right]$$

$$\pi T = \left[\frac{2}{3} \quad \frac{1}{3} \right]$$

$$\pi T^2 = \left[\frac{1}{3} \quad \frac{2}{3} \right]$$

⋮

Oscillation!

Markov Chain Monte Carlo

In the limit:

$$\pi' T = \pi'$$

π is the left eigenvector of T with corresponding eigenvalue 1.

Markov Chain Monte Carlo

In the limit:

$$\underline{\pi' T = 1 \pi'}$$

π is the left eigenvector of T with corresponding eigenvalue 1. Componentwise, we have:

$$\sum_{i=1}^3 \pi_i T_{ij} = \pi_j$$

$$\frac{1}{\pi_1 \pi_2 \pi_3}$$

$$\sum \pi_i = 1$$

$$\pi' = [\pi_1 \pi_2 \pi_3]$$

Markov Chain Monte Carlo

In the limit:

$$\underline{\pi}' T = \underline{\pi}'$$

π is the left eigenvector of T with corresponding eigenvalue 1. Componentwise, we have:

$$\sum_{i=1}^3 \pi_i T_{ij} = \pi_j$$

As the state space grows:

$$\int \pi(x) \underbrace{P(y|x)}_{\text{Markov chain Kernel}} dx = \pi(y)$$

Markov Chain Monte Carlo

Detailed Balance:

$$\text{If } \int_{x_t} \pi(x_t) P(x_{t+1} | x_t) = \int_{x_t} \pi(x_{t+1}) P(x_t | x_{t+1})$$

Integrating over x_t yields $= \pi(x_{t+1}) \int_{x_t} P(x_t | x_{t+1})$

$$\int_{x_t} \pi(x_t) P(x_{t+1} | x_t) = 1 \pi(x_{t+1})$$

Which is the ergodic behaviour we want.
Now we have a sufficient condition for designing $P(x_{t+1} | x_t)$ so as to get samples from π \square

Metropolis-Hastings for logistic regression

$$P(\theta | x_{1:t}, y_{1:t}) = \frac{1}{Z} \prod_{i=1}^t \left[\pi_i^{y_i} (1-\pi_i)^{1-y_i} \right] e^{-\frac{1}{2S^2} \theta^T \theta}$$

want $\theta^{(i)} \sim P(\theta | y_{1:t}, x_{1:t})$

$$\rightarrow \pi(x)$$

$$x^i \sim \pi(x)$$

$$\tilde{\pi}(x) = \frac{1}{Z} \pi(x)$$

$$Z = \int \pi(x) dx$$

MCMC: Metropolis-Hastings

► Initialise $x^{(0)}$.

$x^{(i)}$

e.g.

$$\pi(x) = \frac{P(x)}{Z}$$

► For $i = 0$ to $N - 1$

► Sample $u \sim U_{[0,1]}$.

► Sample $x^* \sim q(x^* | x^{(i)})$.

► If $u < A(x^{(i)}, x^*) = \min \left\{ 1, \frac{p(x^*)q(x^{(i)} | x^*)}{p(x^{(i)})q(x^* | x^{(i)})} \right\}$

0.8

$$x^{(i+1)} = x^*$$

else

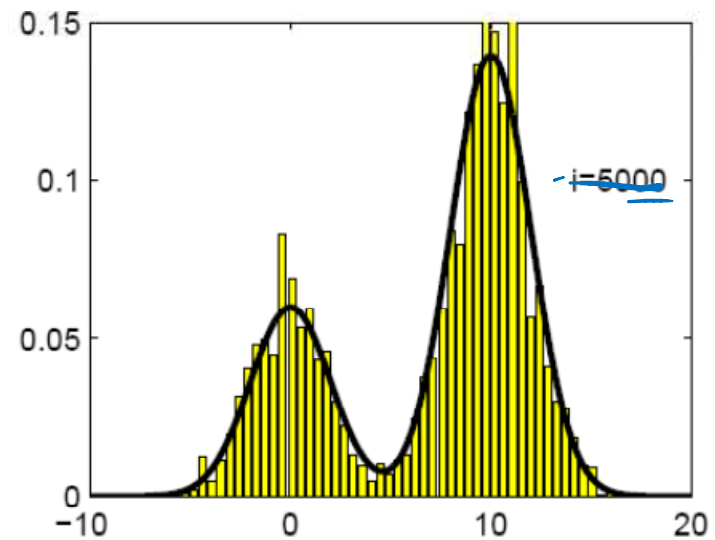
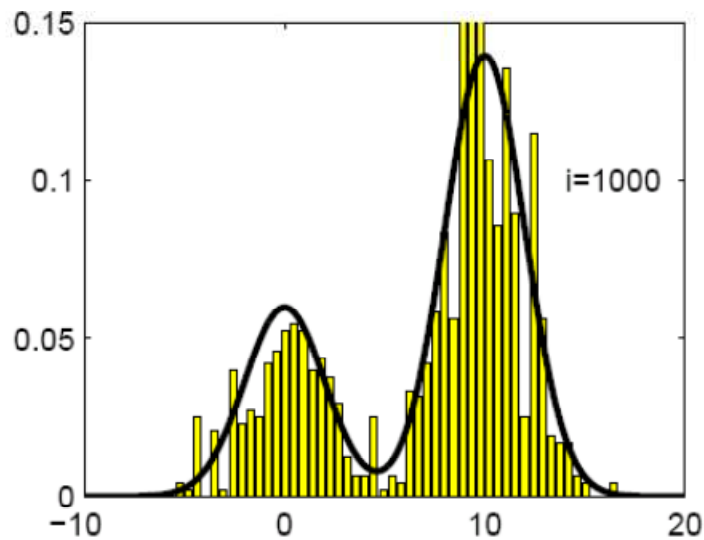
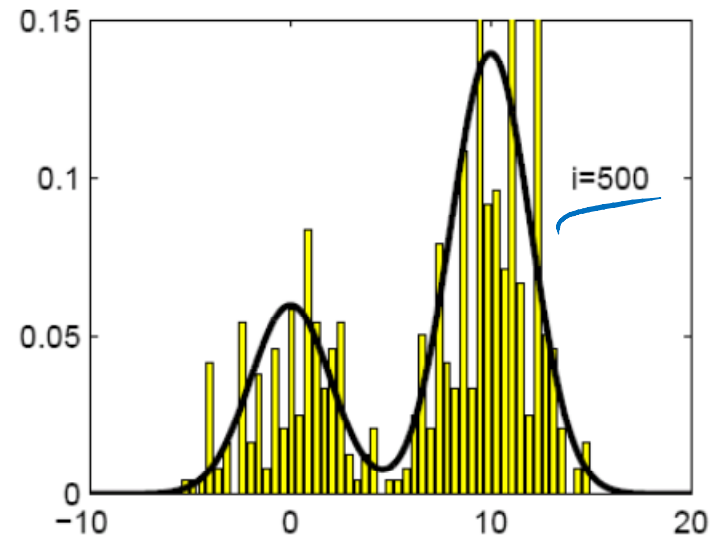
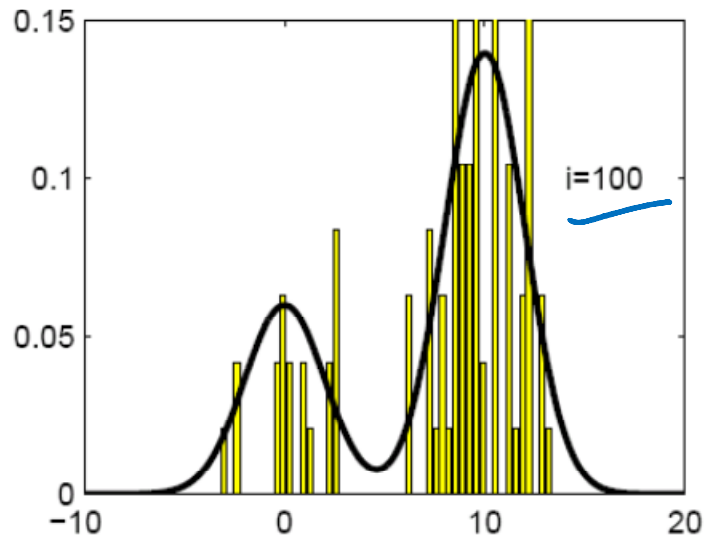
$$x^{(i+1)} = x^{(i)}$$

$$x^* = x^{(i)} + N(0, \sigma^2)$$

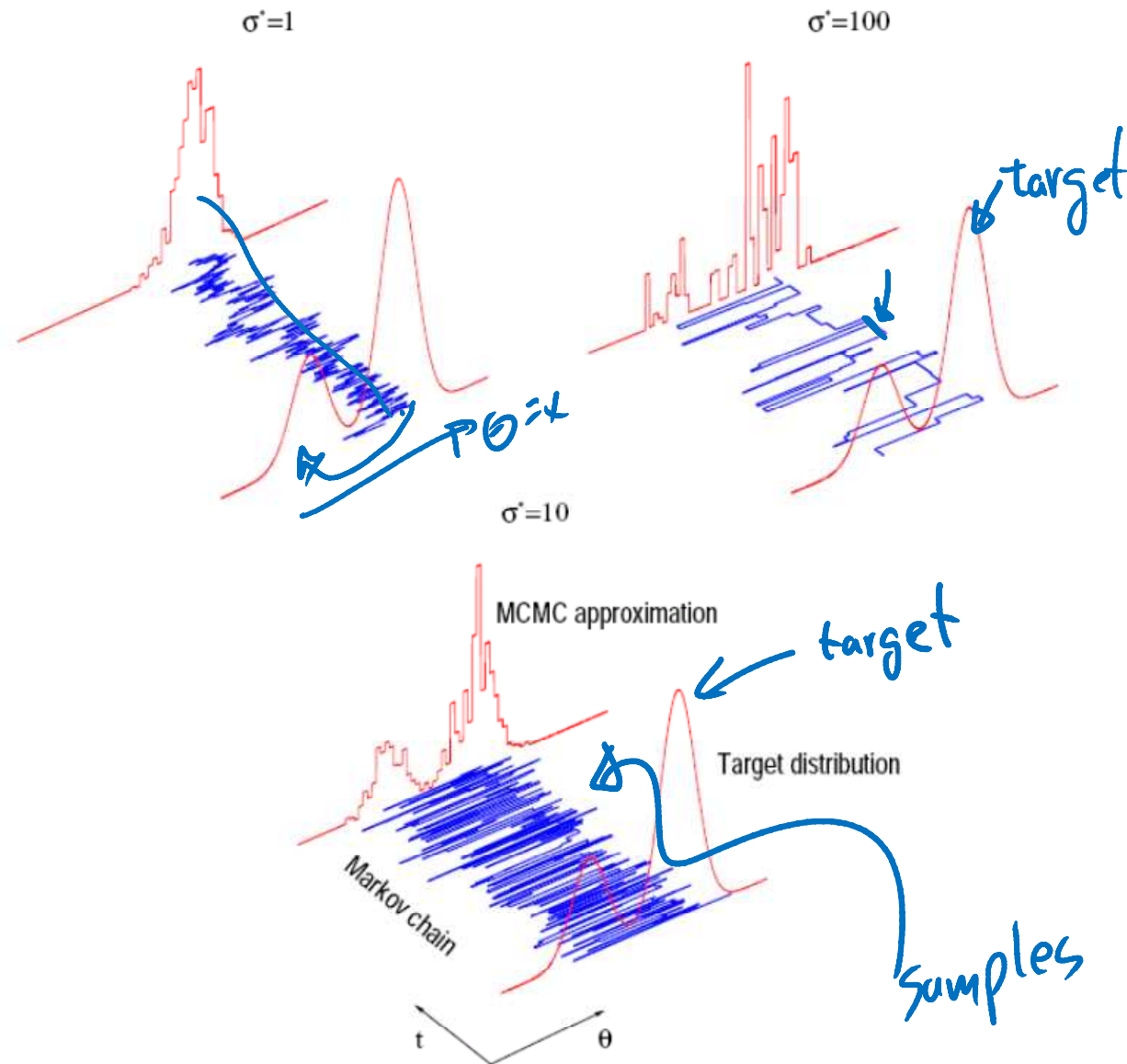
posterior

$x^{(i+1)}$

MCMC: Metropolis-Hastings



MCMC: Choosing the Right Proposal



MCMC: Theory

Kernel:

$$T = K(x, B) = \begin{cases} \gamma(B|x) A(x, B) & x \notin B \\ 1 - \int_{x' \in \{X \setminus B\}} \gamma(x'|x) A(x, x') & x \in B \end{cases}$$

\swarrow Prob of going from x to interval B .
 \nwarrow all space \setminus minus B

$$\therefore K(x, B) = \gamma(B|x) A(x, B) + \prod_{x \in B} \left\{ 1 - \gamma(B|x) A(x, B) - \int_{x' \in \{X \setminus B\}} \gamma(x'|x) A(x, x') \right\}$$

$$K(x, B) = \gamma(B|x) A(x, B) + \prod_{x \in B} \left\{ 1 - \int_{x' \in X} \gamma(x'|x) A(x, x') \right\}$$

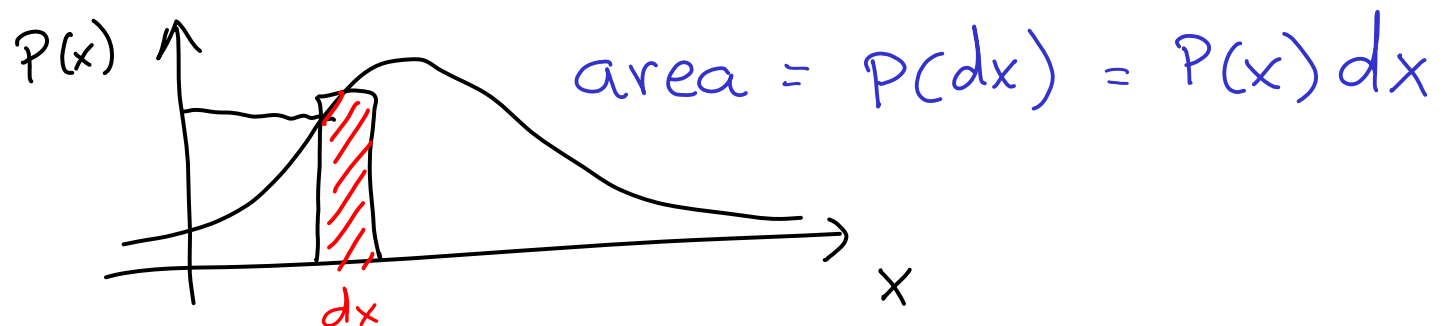
MCMC: Theory

Detailed balance:

$$\underline{\pi(A)} \underline{K(A,B)} = \pi(B) K(B,A)$$

$$\int_{x \in A} \pi(dx) K(x,B) = \int_{\gamma \in B} \pi(d\gamma) K(\gamma,A)$$

Note: $\int f(x) p(x) dx \equiv \int f(x) p(dx)$



Variations of Metropolis-Hastings

$$\min \left\{ 1, \frac{P(x^*)}{P(x^{(i)})} \frac{q(x^{(i)}|x^*)}{q(x^*|x^{(i)})} \right\}$$

$$x^* = x^{(i)} + N(0, \sigma^2)$$

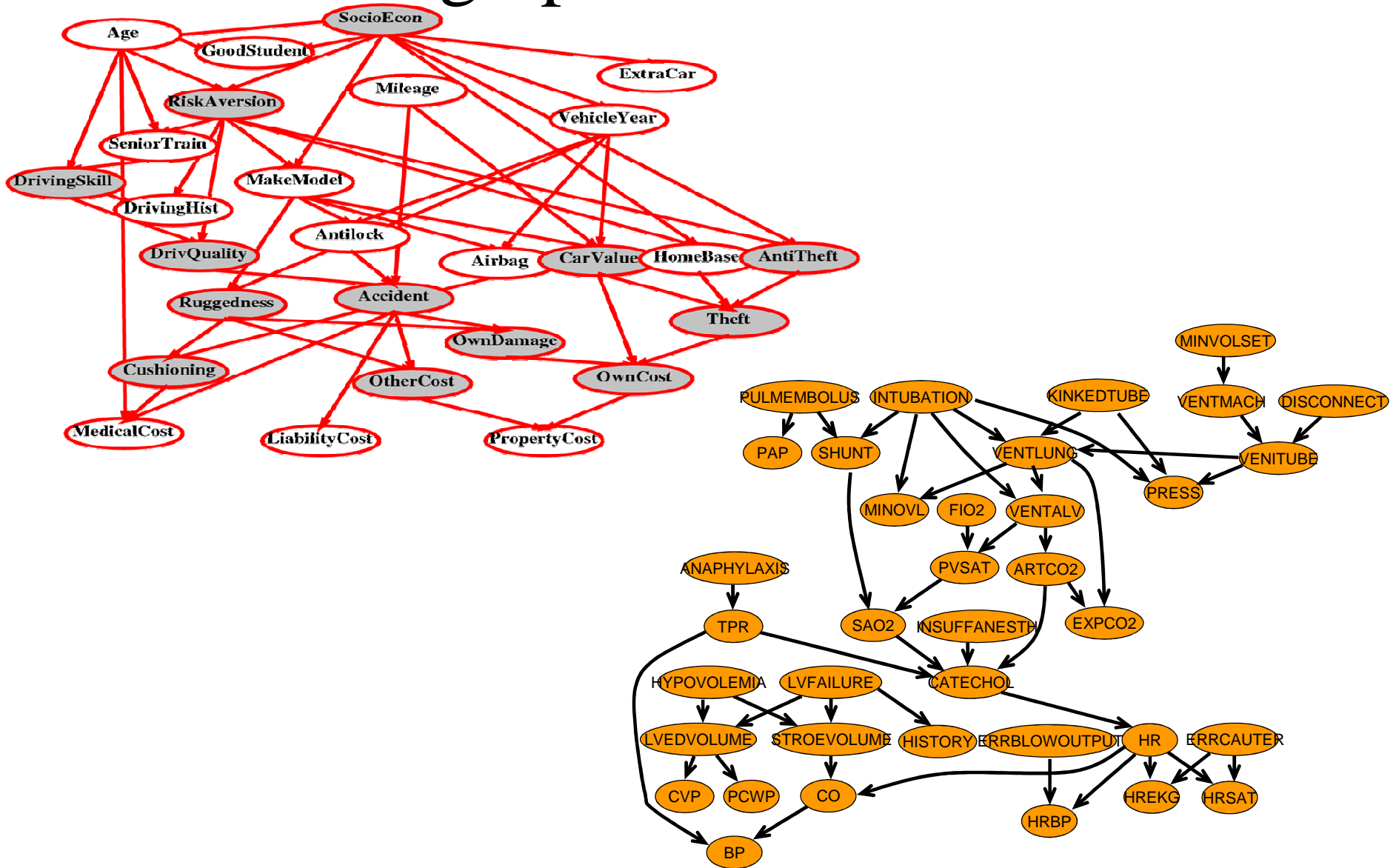
$$q(x^{(i)}|x^*) \propto e^{-\|x^* - x^{(i)}\|^2 / \sigma^2}$$

$$q(x^*|x^{(i)}) \propto e^{-\|x^{(i)} - x^*\|^2 / \sigma^2}$$

$$\min \left\{ 1, \frac{P(x^*)}{P(x^{(i)})} \right\}$$

If anneal
with T ,
concentrate on
peaks of $P(x)$

Extending MH to directed probabilistic graphical models

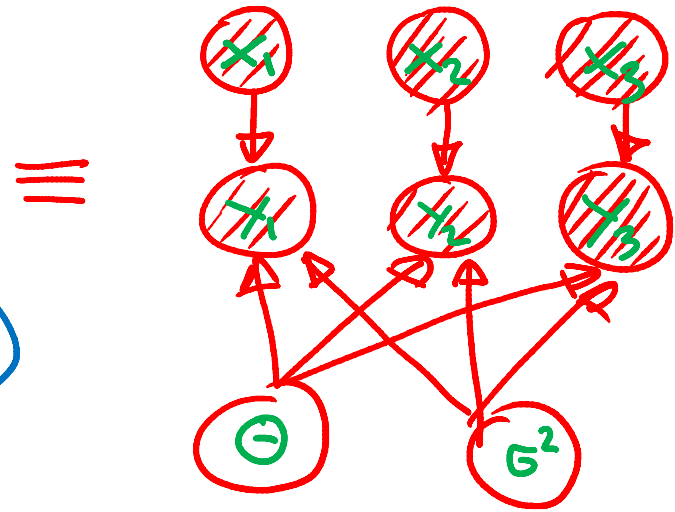


Bayesian graphical models and Gibbs

$$P(\sigma^2) = \text{IG}(a, b)$$

$$P(\theta) = \mathcal{N}(0, \sigma^2 \mathbf{I})$$

$$P(y_i | x_i, \theta) = \mathcal{N}(x_i \theta, \sigma^2)$$



GIBBS :

FOR $i=1$ to N_{samples}

$$\theta^{(i+1)} \sim P(\theta | \sigma^{2(i)}, x, y)$$

$$\sigma^{2(i+1)} \sim P(\sigma^2 | \theta^{(i)}, x, y)$$

END

Gibbs Sampling

Choose the following proposal:

$$q(x^\star | x^{(i)}) = \begin{cases} p(x_j^\star | x_{-j}^{(i)}) & \text{If } x_{-j}^\star = x_{-j}^{(i)} \\ 0 & \text{Otherwise.} \end{cases}$$

where $x_{-j} = \{x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n\}$.

Then the acceptance is:

$$A(x^{(i)}, x^\star) = \min \left\{ 1, \frac{p(x^\star) q(x^{(i)} | x^\star)}{p(x^{(i)}) q(x^\star | x^{(i)})} \right\} = 1.$$

Gibbs Sampling

- ▶ Initialise $x_{1:n}^{(0)}$.
- ▶ For $i = 0$ to $N - 1$
 - ▶ Sample $x_1^{(i+1)} \sim p(x_1 | x_2^{(i)}, x_3^{(i)}, \dots, x_n^{(i)})$.
 - ▶ Sample $x_2^{(i+1)} \sim p(x_2 | x_1^{(i+1)}, x_3^{(i)}, \dots, x_n^{(i)})$.
 - ▶ \vdots
 - ▶ Sample $x_j^{(i+1)} \sim p(x_j | x_1^{(i+1)}, \dots, x_{j-1}^{(i+1)}, x_{j+1}^{(i)}, \dots, x_n^{(i)})$.
 - ▶ \vdots
 - ▶ Sample $x_n^{(i+1)} \sim p(x_n | x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_{n-1}^{(i+1)})$.

Gibbs Sampling For Graphical models

A large-dimensional joint distribution is factored into a directed graph that encodes the conditional independencies in the model. In particular, if $x_{pa(j)}$ denotes the parent nodes of node x_j , we have

$$p(x) = \prod_j p(x_j | x_{pa(j)}).$$

It follows that the full conditionals simplify as follows

$$p(x_j | x_{-j}) = p(x_j | x_{pa(j)}) \prod_{k \in ch(j)} p(x_k | x_{pa(k)})$$

where $ch(j)$ denotes the children nodes of x_j .



Auxiliary Variable Samplers

- ▶ It is often easier to sample from an augmented distribution $p(x, u)$, where u is an auxiliary variable, than from $p(x)$.
- ▶ It is possible to obtain marginal samples $x^{(i)}$ by sampling $(x^{(i)}, u^{(i)})$ according to $p(x, u)$ and, then, ignoring the samples $u^{(i)}$.
- ▶ This very useful idea was proposed in the physics literature (Swendsen and Wang, 1987).

Hybrid (Hamiltonian) Monte Carlo

➤ The idea is to exploit gradient information.

➤ Define the extended target distribution:

$$p(x, u) = p(x)N(u; 0, I_{n_x}).$$

➤ Introduce the gradient vector: $\Delta(x) = \partial \log p(x) / \partial x$

➤ Introduce the parameters ρ and L .

➤ Next we “leapfrog”.

Hybrid Monte Carlo

- ▶ Sample $v \sim U_{[0,1]}$ and $u^* \sim N(0, I_{n_x})$.
- ▶ Let $x_0 = x^{(i)}$ and $u_0 = u^* + \rho \Delta(x_0)/2$.
- ▶ For $l = 1, \dots, L$, take steps

$$x_l = x_{l-1} + \rho u_{l-1}$$

$$u_l = u_{l-1} + \rho_l \Delta(x_l)$$

where $\rho_l = \rho$ for $l < L$ and $\rho_L = \rho/2$.

- ▶ If $v < A = \min \left\{ 1, \frac{p(x_L)}{p(x^{(i)})} \exp \left(-\frac{1}{2} (u_L^\top u_L - u^{*\top} u^*) \right) \right\}$

$$(x^{(i+1)}, u^{(i+1)}) = (x_L, u_L)$$

else $(x^{(i+1)}, u^{(i+1)}) = (x^{(i)}, u^*)$

Next lecture

In the next lecture, we look at constrained optimization and sparse methods in learning.