

# Importance sampling (IS) & Markov chain Monte Carlo (MCMC)



Nando de Freitas  
*March, 2013*  
*University of British Columbia*

# Bayesian logistic regression

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

We also assume a Gaussian prior  $\pi(\theta)$

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

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

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

IS

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

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

$$q(\theta) = N(0, 1000)$$

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

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

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

SLLN

$$P(Y_{t+1} | X_{t+1}, Y_{1:t}, X_{1:t}) = \int P(Y_{t+1} | X_{t+1}, \theta) P(d\theta | X_{1:t}, Y_{1:t})$$

$$= \int P(Y_{t+1} | X_{t+1}, \theta) P(\theta | X_{1:t}, Y_{1:t}) d\theta$$

$$\approx \int P(Y_{t+1} | X_{t+1}, \theta) \frac{1}{N} \sum_{i=1}^N w(\theta^{(i)}) S_{\theta^{(i)}}(d\theta)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \int P(Y_{t+1} | X_{t+1}, \theta) w(\theta^{(i)}) S_{\theta^{(i)}}(d\theta)$$

$$\approx \frac{1}{N} \sum_{i=1}^N P(Y_{t+1} | X_{t+1}, \theta^{(i)}) w(\theta^{(i)})$$

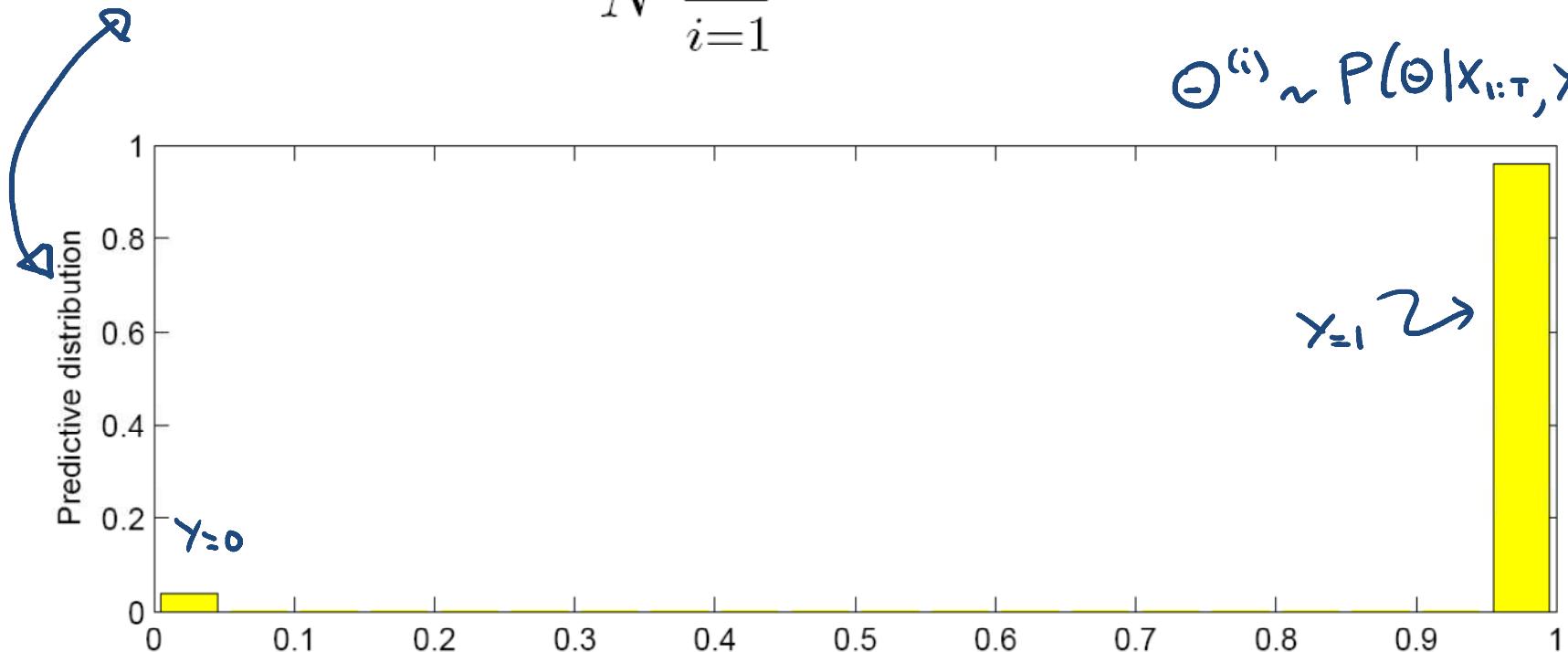
likelihood

# Example: Logistic Regression

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

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

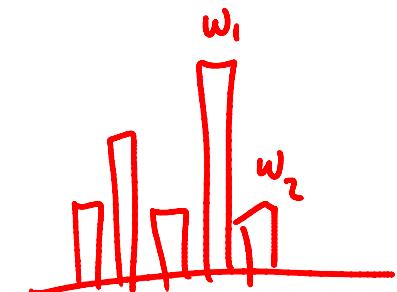
$\theta^{(i)} \sim P(\theta|x_{1:T}, y_{1:T})$



# Un-normalized IS

$D = \text{data}$

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



$$\begin{aligned}
 & P(Y_{t+1} | X_{t+1}, D) \\
 &= P(Y_{t+1} | X_{1:t+1}, \theta) = \frac{1}{Z} \int P(Y_{t+1} | X_{t+1}, \theta) P(D|\theta) P(\theta) d\theta \\
 &= \frac{\int P(Y_{t+1} | X_{t+1}, \theta) P(D|\theta) P(\theta) \frac{q(\theta)}{q(\theta)} d\theta}{\int P(D|\theta) P(\theta) \frac{q(\theta)}{q(\theta)} d\theta} \\
 &= \frac{\int P(Y_{t+1} | X_{t+1}, \theta) \omega(\theta) \frac{q(\theta)}{q(\theta)} d\theta}{\int \omega(\theta) q(\theta) d\theta} = \frac{1}{N} \sum_{i=1}^N \frac{\omega(\theta^{(i)}) P(Y_{t+1} | X_{t+1}, \theta^{(i)})}{\sum_{j=1}^N \omega(\theta^{(j)})}
 \end{aligned}$$

# MCMC

Detailed Balance:

If

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

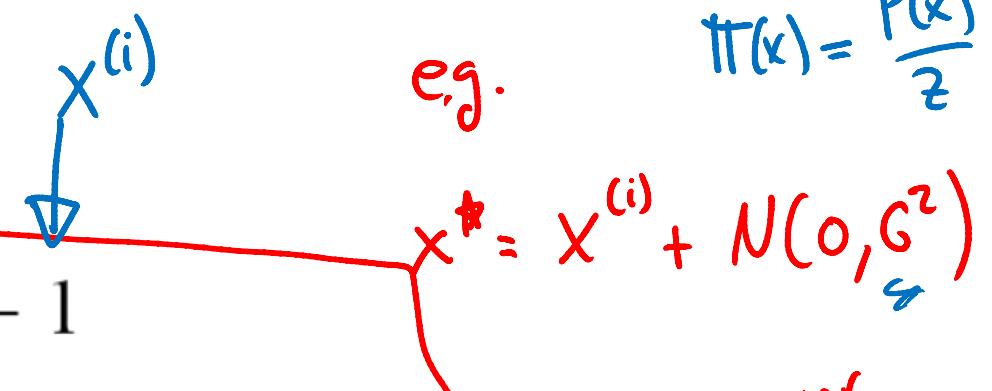
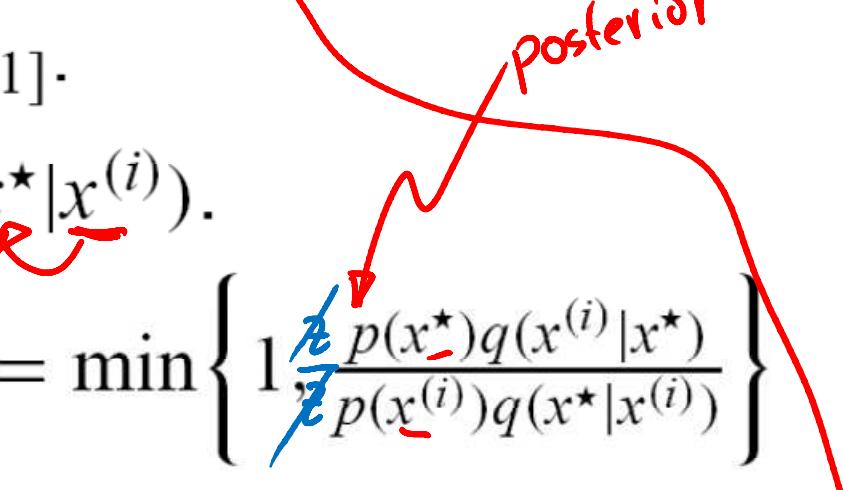
Integrating over  $x_t$  yields

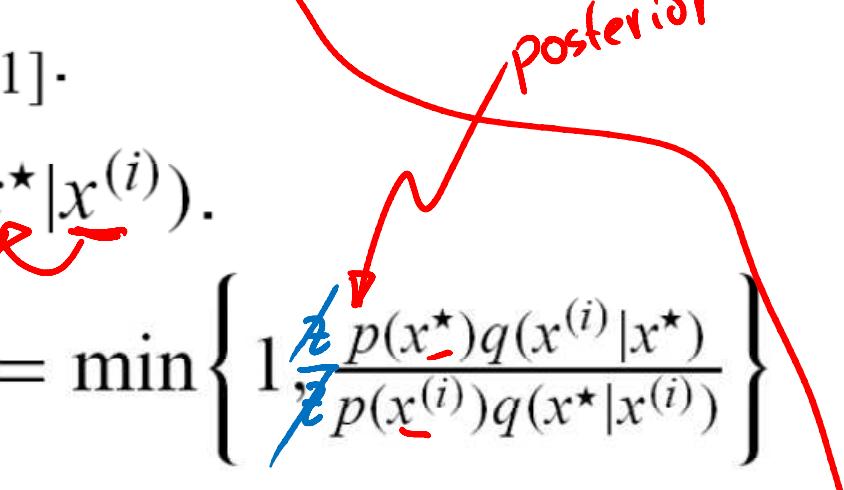
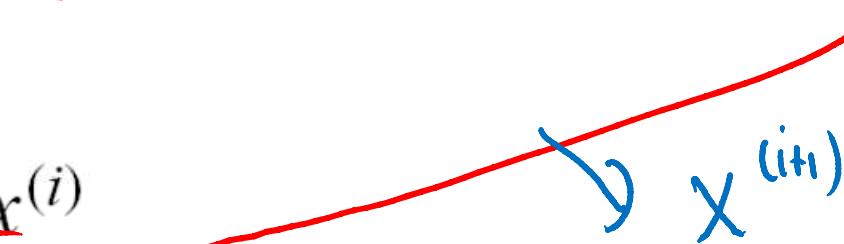
$$= \pi(x_{t+1}) \int_{x_t} P(x_t | x_{t+1})$$

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

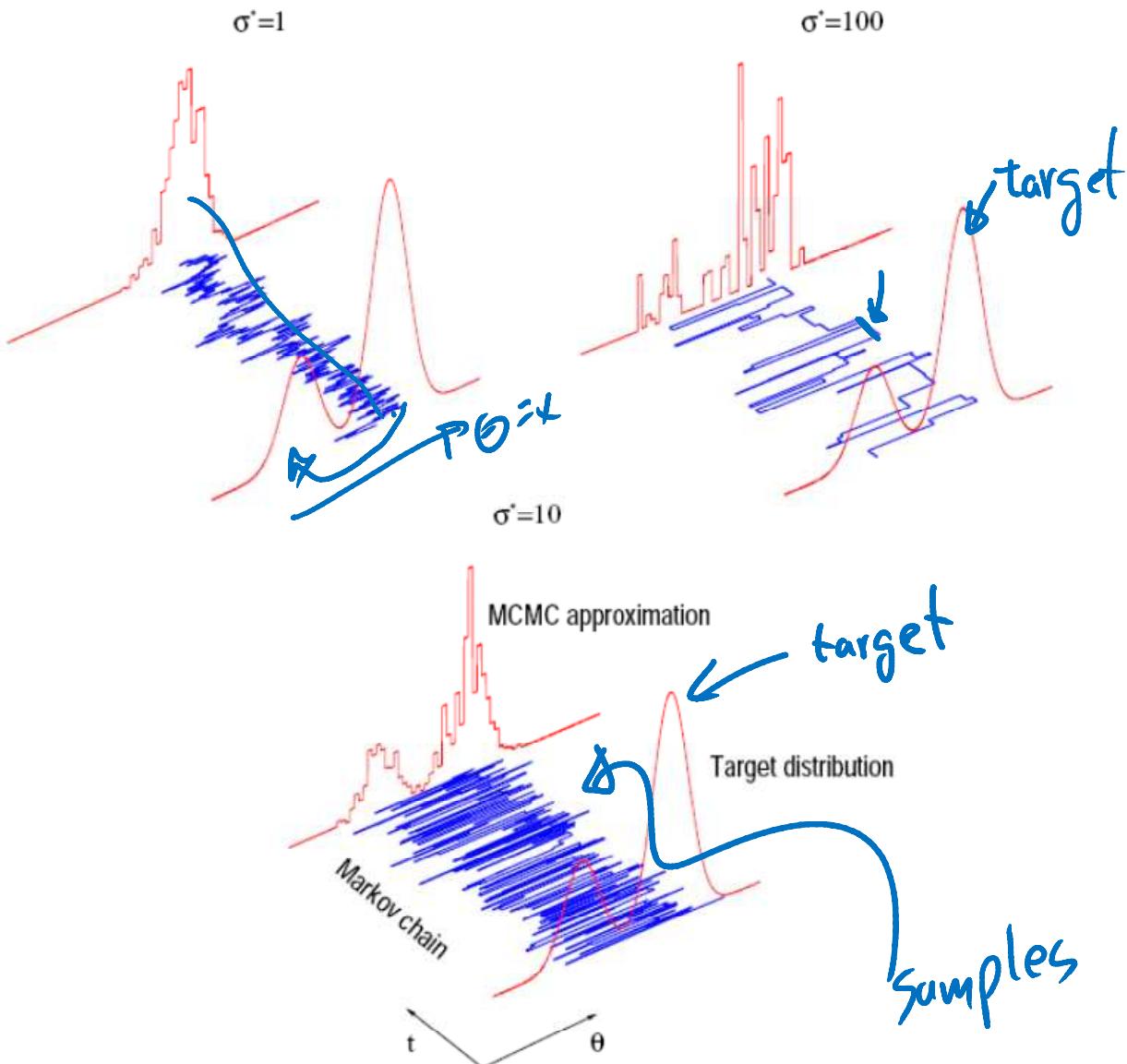
which is the ergodic behaviour we want.  
Now we have a sufficient condition for designing  
 $P(x_{t+1} | x_t)$  so as to get samples from  $\pi$

# MCMC: Metropolis-Hastings

- ▶ Initialise  $x^{(0)}$ .  

- ▶ For  $i = 0$  to  $N - 1$ 
  - ▶ Sample  $u \sim U_{[0,1]}$ .
  - ▶ Sample  $x^* \sim q(x^* | x^{(i)})$ .  


A red circle highlights the term  $q(x^* | x^{(i)})$  in the acceptance ratio formula.
  - ▶ If  $u < A(x^{(i)}, x^*) = \min\left\{1, \frac{p(x^*)q(x^{(i)}|x^*)}{p(x^{(i)})q(x^*|x^{(i)})}\right\}$   
 $x^{(i+1)} = x^*$   

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


# MCMC: Choosing the Right Proposal



# MCMC: Theory

Kernel:

$$T = K(x, B) =$$

Prob of going from  
x to interval B.

$$= \begin{cases} q(B|x) A(x, B) & x \notin B \\ 1 - \int_{x' \in \mathcal{X} \setminus B} q(x'|x) A(x, x') & x \in B \end{cases}$$

all space  $\mathcal{X}$  minus B

$$\therefore k(x, B) = q(B|x) A(x, B) + \mathbb{I}_{x \in B} \left\{ 1 - \int_{x' \in \mathcal{X} \setminus B} q(x'|x) A(x, x') \right\}$$

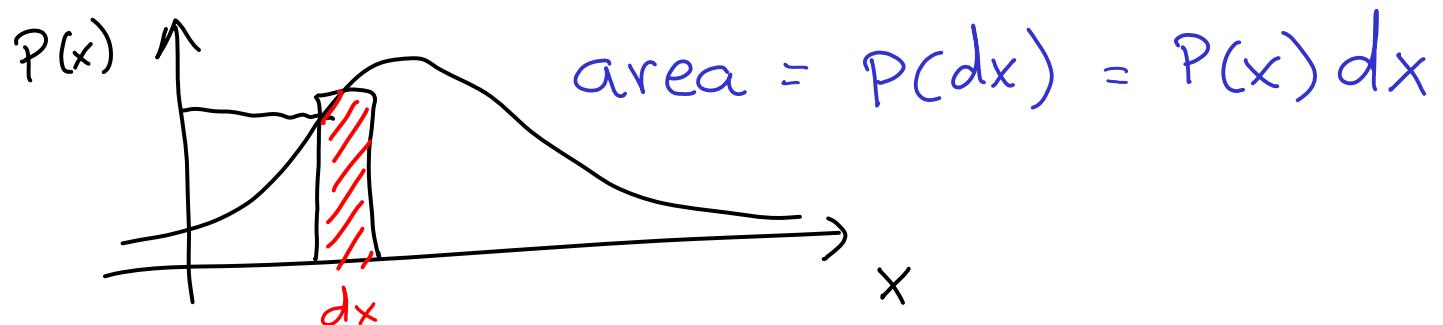
$$k(x, B) = q(B|x) A(x, B) + \mathbb{I}_{x \in B} \left\{ 1 - \int_{x' \in \mathcal{X}} q(x'|x) A(x, x') \right\}$$

Detailed balance :

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

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

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



# Variations of Metropolis-Hastings

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

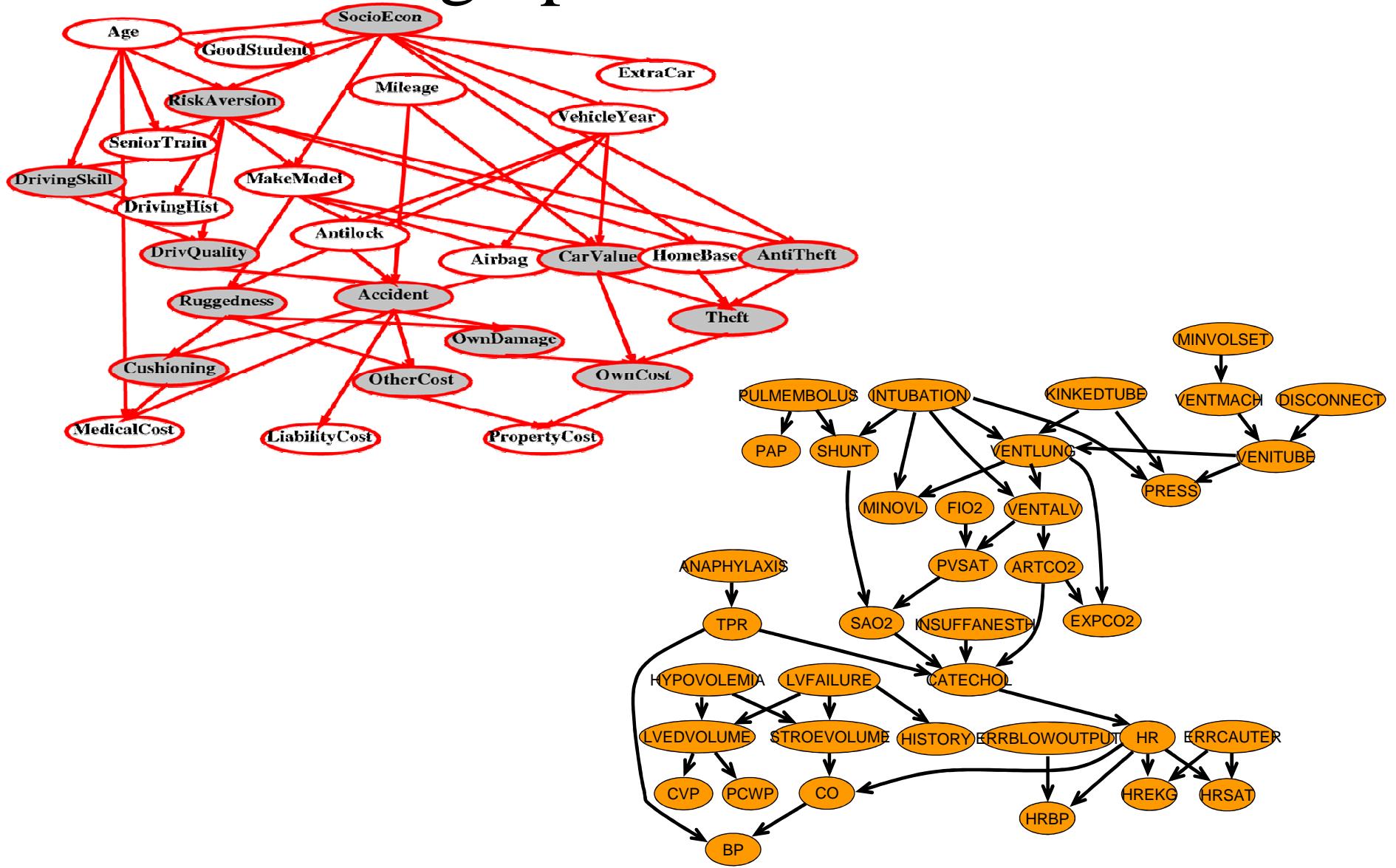
$$x^* = x^{(i)} + N(0, \sigma^2)$$

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

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

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

# Extending MH to directed probabilistic graphical models

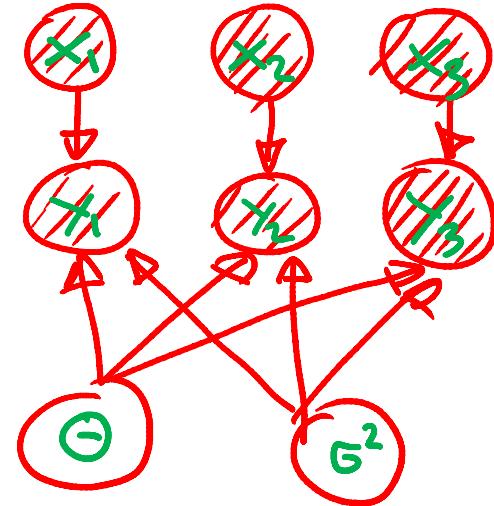


# Bayesian graphical models and Gibbs

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

$$P(\boldsymbol{\theta}) = N(0, \sigma^2 I) \quad \equiv$$

$$P(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}) = N(\mathbf{x}\boldsymbol{\theta}, \boldsymbol{\sigma}^2)$$



GIBBS:

FOR  $i=1$  to  $N_{\text{samples}}$

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

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

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:

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

- 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  $\rho$  and  $L$ .
- Next we “leapfrog”.

- Sample  $v \sim U_{[0,1]}$  and  $u^\star \sim N(0, I_{n_x})$ .
- Let  $x_0 = x^{(i)}$  and  $u_0 = u^\star + \rho \Delta(x_0)/2$ .
- 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  $v < A = \min \left\{ 1, \frac{p(x_L)}{p(x^{(i)})} \exp \left( -\frac{1}{2} (u_L^\top u_L - u^{\star \top} u^\star) \right) \right\}$
- $(x^{(i+1)}, u^{(i+1)}) = (x_L, u_L)$
- else  $(x^{(i+1)}, u^{(i+1)}) = (x^{(i)}, u^\star)$