

# Lecture 19. Gibbs Sampler.

Note Title

5/12/2006

The Gibbs sampler is easy to compute and requires no free parameters.

$$\underline{x} = (x_1, \dots, x_d)$$

Notation.  $\underline{x}_{/i} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d)$

Marginal Distribution  $\pi(x_i | \underline{x}_{/i})$

usually there is a Markov assumption.

so that  $\pi(x_i | \underline{x}_{/i}) = \pi(x_i | \underline{x}_{N(i)})$

where  $N(i)$  is the neighbourhood of  $i$ .

Gibbs sampler.

$$K_i(\underline{y} | \underline{x}) = \pi(y_i | \underline{x}_{/i}) \delta_{\underline{y}_{/i}, \underline{x}_{/i}}$$

Check detailed balance:

$$\begin{aligned} \pi(\underline{x}) K_i(\underline{y} | \underline{x}) &= \pi(\underline{x}) \pi(y_i | \underline{x}_{/i}) \delta_{\underline{y}_{/i}, \underline{x}_{/i}} \\ &= \pi(x_i | \underline{x}_{/i}) \pi(y_i | \underline{x}_{/i}) \pi(\underline{x}_{/i}) \delta_{\underline{y}_{/i}, \underline{x}_{/i}} \end{aligned}$$

symmetric in  $\underline{y}$  &  $\underline{x}$ ,  $= \pi(\underline{y}) K_i(\underline{y} | \underline{x})$

But this is not irreducible.

Page 2

To make it irreducible, must sample all  $x_i$ .

$$K(\underline{y}|\underline{x}) = \sum_{i=1}^n \alpha_i K_i(\underline{y}|\underline{x})$$

with  $\alpha_i > 0, \forall i$

$$\sum_{i=1}^n \alpha_i = 1.$$

This obeys detailed balance (linearity)

Algorithm: At each time step  $t$  at state  $\underline{x}^t$ .  
select  $i$  with probability  $\alpha_i$   
then select  $\underline{y}$  from  $K_i(\underline{y}|\underline{x}^t)$ .

Typically random scan.  $\alpha_i = 1/n, \forall i$ .

Can also do systematic scan:

$$\text{Let } \underline{x}^{(t)} = (x_1^{(t)}, \dots, x_d^{(t)}).$$

- Draw  $x_i^{t+1}$  from  $\pi(x_i | x_1^{t+1}, \dots, x_{i-1}^{t+1}, x_{i+1}^t, \dots, x_d^t)$

Scan through  $\bar{i}$ .

Gibbs Sampler converges geometrically (like M-H). The convergence rate depends on how well variables correlate with each other.

Example: Ising-Model

$$\pi(x_1, \dots, x_d) = \frac{1}{Z} e^{-\mu \sum_{i=1}^{d-1} x_i x_{i+1}} \quad x_i \in \{-1, 1\}$$

$$\pi(x_i | \underline{x}_{-i}) = \pi(x_i | x_{i-1}, x_{i+1})$$

Markov property.

Compute.

$$\pi(x_i | x_{i-1}, x_{i+1}) = \frac{\pi(x_{i-1}, x_i, x_{i+1})}{\pi(x_{i-1}, x_{i+1})}$$

$$\pi(\underline{x}) = \frac{1}{Z} e^{-\mu(x_{i-1}x_i + x_i x_{i+1})} \cdot f(x_{i+1}, \dots, x_d) \cdot g(x_1, \dots, x_{i-1})$$

$$\pi(x_{i-1}, x_i, x_{i+1}) = \frac{1}{Z} e^{-\mu(x_{i-1}x_i + x_i x_{i+1})} \sum_{x_{i+2}, \dots, x_d} f(x_{i+1}, \dots, x_d) \sum_{x_1, \dots, x_{i-2}} g(x_1, \dots, x_{i-1})$$

Hence  $\pi(x_i | x_{i-1}, x_{i+1}) = \frac{e^{-\mu(x_{i-1}x_i + x_i x_{i+1})}}{\sum_{x_i} e^{-\mu(x_{i-1}x_i + x_i x_{i+1})}}$

$$\pi(x_i | x_{i-1}, x_{i+1}) = \frac{e^{-\mu(x_{i-1}x_i + x_i x_{i+1})}}{e^{\mu(x_{i-1} + x_{i+1})} + e^{\mu(x_{i+1} + x_{i-1})}}$$

Moral: the conditional is usually easy to compute for Markov Random Fields (MRF)

Page 4:

Example.

$$\underline{x} = (x_1, x_2)$$

$$\pi(x) = N \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right\}$$

$$\pi(x_1 | x_2) = N(\rho x_2, (1-\rho^2))$$

$$\pi(x_2 | x_1) = N(\rho x_1, (1-\rho^2))$$

Systematic scan.

$$\begin{pmatrix} x_1^{(t)} \\ x_2^{(t)} \end{pmatrix} \sim N \left\{ \begin{pmatrix} \rho^{2t-1} x_2^{(0)} \\ \rho^{2t} x_2^{(0)} \end{pmatrix}, \begin{pmatrix} 1-\rho^{4t-2} & \rho-\rho^{4t-1} \\ \rho-\rho^{4t-1} & 1-\rho^{4t} \end{pmatrix} \right\}$$

As  $t \rightarrow \infty$ , the joint distribution of  $(x_1^{(t)}, x_2^{(t)})$  converges to the target distribution.

Also, the rate of convergence is exponential.

Rate of convergence is equal to the maximal correlation between  $x_i^{(t)}$  and  $x_i^{(t+1)}$  which is  $\rho^2$ .

The Gibbs Sampler can be thought of as a special case of Metropolis-Hastings.

Gibbs is M-H with a proposal which is automatically accepted.

Recall M-H

$$K(\underline{y}|\underline{x}) = T(\underline{y}|\underline{x}) \min \left\{ 1, \frac{\pi(\underline{y}) T(\underline{x}|\underline{y})}{\pi(\underline{x}) T(\underline{y}|\underline{x})} \right\}$$

Suppose  $T(\underline{y}|\underline{x}) = K(\underline{y}|\underline{x})$

Gibbs Sampler.

Now 
$$\frac{\pi(\underline{y}) T(\underline{x}|\underline{y})}{\pi(\underline{x}) T(\underline{y}|\underline{x})} = 1$$

because  $K(\underline{y}|\underline{x})$  obeys detailed balance.

### Metropolized Gibbs Sampler.

$$\underline{x} = (x_1, \dots, x_d)$$

Each  $x_i$  takes  $m_i$  possible values.

- select  $i$  at random.

- draw  $y_i (\neq x_i)$  with probability  $\frac{\pi(y_i | \underline{x}_{-i})}{1 - \pi(x_i | \underline{x}_{-i})}$

then replace  $x_i$  by  $y_i$  with probability  $\min \left\{ 1, \frac{1 - \pi(x_i | \underline{x}_{-i})}{1 - \pi(y_i | \underline{x}_{-i})} \right\}$   
 (statistically more efficient than Gibbs)

Page 6.

Data Augmentation (DA)

A stochastic alternative to the EM algorithm.

$\underline{y}_{obs} \sim$  observed data,  $\underline{y}_{mis} \sim$  missing data.

$$p(\underline{\theta} | \underline{y}_{obs}, \underline{y}_{mis}) \quad \& \quad p(\underline{y}_{mis} | \underline{y}_{obs})$$

Given  $p(\underline{\theta}, \underline{y}_{mis} | \underline{y}_{obs})$

Sample  $\underline{\theta}$  &  $\underline{y}_{mis}$  in turn.

Initialize:  $\underline{\theta}^0$  &  $\underline{y}_{mis}^0$

Sample  $\underline{\theta}^t$  from  $p(\underline{\theta} | \underline{y}_{mis}^{t-1}, \underline{y}_{obs})$

$\underline{y}_{mis}^t$  from  $p(\underline{y}_{mis} | \underline{\theta}^t, \underline{y}_{obs})$

This is a form of Gibbs sampling and so is guaranteed to converge to samples from  $p(\underline{\theta}, \underline{y}_{mis} | \underline{y}_{obs})$

First Example: Hierarchical model.

$$y_i | \theta_i \sim f_i(y_i | \theta_i)$$

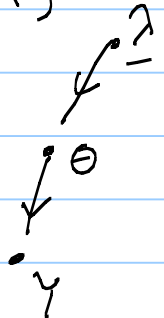
$$\theta_i \sim G(\theta | \lambda)$$

Prior  $\lambda \sim f_0(\mu, \sigma^2)$

Want to estimate  $\theta_i, \lambda$  and quantify their uncertainties:  $\lambda = (\mu, \sigma^2)$ .

$$P(Y, \theta, \lambda) = P(Y | \theta) P(\theta | \lambda) P(\lambda)$$

Need.  $P(\theta | Y, \lambda)$   
 $P(\lambda | Y, \theta)$



$$P(\theta | Y, \lambda) = \frac{P(Y | \theta) P(\theta | \lambda)}{\sum_{\theta} P(Y | \theta) P(\theta | \lambda)}$$

$$P(\lambda | Y, \theta) = \frac{P(\theta | \lambda) P(\lambda)}{\sum_{\lambda} P(\theta | \lambda) P(\lambda)}$$

It may be impossible to compute the denominator  
 $\sum_{\theta} P(Y | \theta) P(\theta | \lambda)$  &  $\sum_{\lambda} P(\theta | \lambda) P(\lambda)$   
 if not, doing weighted sampling.

Page 8

## Second Example

Data is 1-dimensional and is generated by one of two Gaussian distributions

$$N(\mu_1, \sigma) \quad N(\mu_2, \sigma)$$

$\mu_1, \mu_2$  unknown random variables  
 $\sigma$  known.

$x^i$  data,  $v^i$  missing data

$$v^i = 1 \text{ means } x^i \sim N(\mu_1, \sigma)$$

$$= 0 \text{ means } x^i \sim N(\mu_2, \sigma)$$

This can be summarized by

$$P(x^i | v^i, \mu_1, \mu_2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} (x^i - v^i\mu_1 - (1-v^i)\mu_2)^2}$$

Prior for  $v^i$  is

$$P(v^i) = e^{v^i \log \alpha + (1-v^i) \log(1-\alpha)} \quad \alpha = \text{known}$$

Prior for  $\mu_1$  &  $\mu_2$  is

$$P(\mu_1, \mu_2) = \frac{1}{2\pi \sigma_m^2} e^{-\frac{(\mu_1 - \alpha_1)^2}{2\sigma_m^2}} e^{-\frac{(\mu_2 - \alpha_2)^2}{2\sigma_m^2}}$$

$\alpha_1, \alpha_2, \sigma_m^2$  known

For a set of data  $\{x^i : i=1 \text{ to } M\}$

Full distribution  $m$

$$P(\mu_1, \mu_2) \left\{ \prod_{i=1}^M P(x^i | v^i, \mu_1, \mu_2) P(v^i) \right\}$$



To do Data Augmentation, we need to compute

$$P(v^i | \underline{\mu}, x^i) \text{ for } i=1 \text{ to } M.$$

$$\& P(\underline{\mu} | \{v^i : i=1 \text{ to } M\}, \{x^i : i=1 \text{ to } M\})$$

$$\underline{\mu} = (\mu_1, \mu_2)$$

$$P(v^i | \underline{\mu}, x^i) = \frac{e^{-\frac{1}{2\sigma^2} (x^i - v^i \mu_1 - (1-v^i) \mu_2)^2} \times e^{v^i \log d + (1-v^i) \log(1-d)}}{Z[\underline{\mu}, x^i]} \leftarrow \text{normalization constant.}$$

This can be simplified (homework)

$$P(\underline{\mu} | \{v^i\}, \{x^i\}) = \frac{1}{(2\pi\sigma^2)^{M/2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^M (x^i - v^i \mu_1 - (1-v^i) \mu_2)^2} \times \frac{1}{2\pi\sigma_m^2} e^{-\frac{(\mu_1 - d_1)^2}{2\sigma_m^2}} e^{-\frac{(\mu_2 - d_2)^2}{2\sigma_m^2}}$$

normalization constant.  $\rightarrow Z[\{v^i\}, \{x^i\}]$

After some algebra.

$$P(\mu_1 | \{v^i\}, \{x^i\}) \sim N(\tilde{\mu}_1, \tilde{\sigma}_1^2)$$

$$P(\mu_2 | \{v^i\}, \{x^i\}) \sim N(\tilde{\mu}_2, \tilde{\sigma}_2^2)$$

with  $\tilde{\mu}_1 = \frac{\sum_{i=1}^M v^i x^i}{\sum_{i=1}^M v^i} + \sigma^2 d_1$

$$\frac{\sigma_m^2 \sum_{i=1}^M v^i + \sigma^2}{\sigma_m^2 \sum_{i=1}^M v^i + \sigma^2}$$

Page 10.

$$\hat{\sigma}_1^2 = \frac{\sigma^2 \sigma_m^2}{\sigma^2 + \sigma_m^2 \sum_{i=1}^M V_i}$$

$$\hat{\mu}_2 = \frac{\sigma_m^2 \sum_{i=1}^M (1-V_i) X_i + \sigma^2 \alpha_2}{\sigma_m^2 \sum_{i=1}^M (1-V_i) + \sigma^2}$$

$$\hat{\sigma}_2^2 = \frac{\sigma^2 \sigma_m^2}{\sigma^2 + \sigma_m^2 \sum_{i=1}^M (1-V_i)}$$