

CS340 Machine learning

Final review

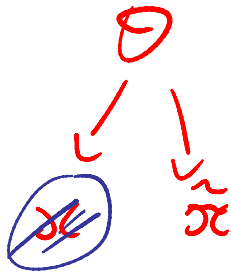
Covered in midterm review

- I – Basics:
 - Statistics (MLE, posteriors, Bayes factors, model selection, etc)
 - Info theory
 - Decision theory

Outline

- II- Models
 - Generative vs discriminative
 - Naïve Bayes
 - MVN
 - Markov chains
 - DGMs, including expert systems
 - UGMs, including Ising models
- III – Algorithms
 - Gibbs sampling

Unconditional density models



Eg $x \sim \text{bernoulli}$, $\theta \sim \text{beta}$

$x \sim \text{multinomial}$, $\theta \sim \text{dir}$

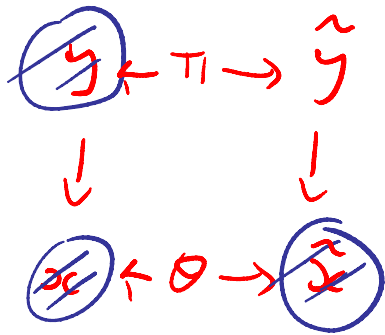
$x \sim \text{multinomial}$, $\theta \sim \text{mixture of dir}$

$x \sim \text{gaussian}$, $\theta = (\mu, \lambda) \sim \text{NormalGamma}$

$x \sim \text{MVN}$, $\theta = (\mu, \Lambda) \sim \text{NormalWishart}$

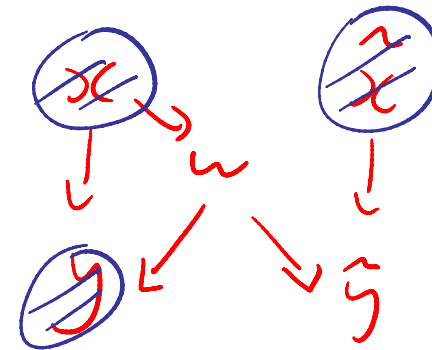
Generative vs discriminative models

Generative $y \rightarrow x$



$$p(\mathbf{x}, y | \boldsymbol{\pi}, \boldsymbol{\theta}) = p(y | \boldsymbol{\pi}) p(\mathbf{x} | y, \boldsymbol{\theta})$$

Discriminative $x \rightarrow y$



$$p(y | \mathbf{x}, \mathbf{w})$$

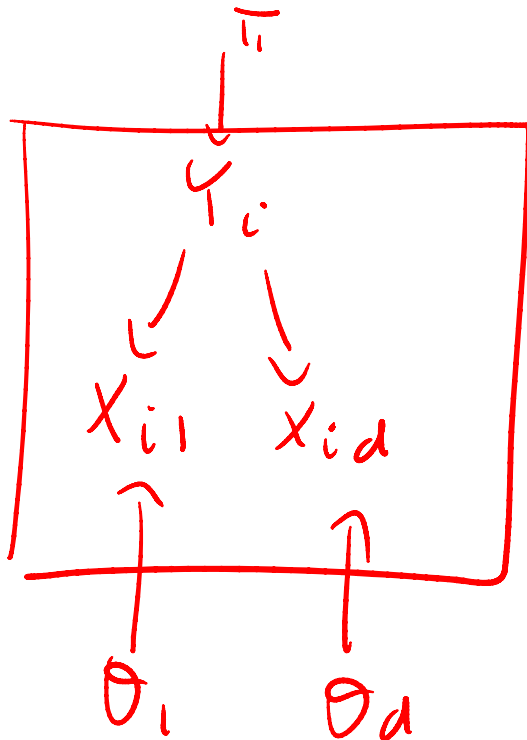
$P(x|y)$ = class conditional density

Eg fully factored (naïve Bayes)

Markov chain

full covariance Gaussian

Naïve Bayes



$P(x_j|Y=c)$ = bernoulli, gaussian, ...
Compute $p(\theta_{j|c} | D)$
Handle missing data
Log sum exp trick

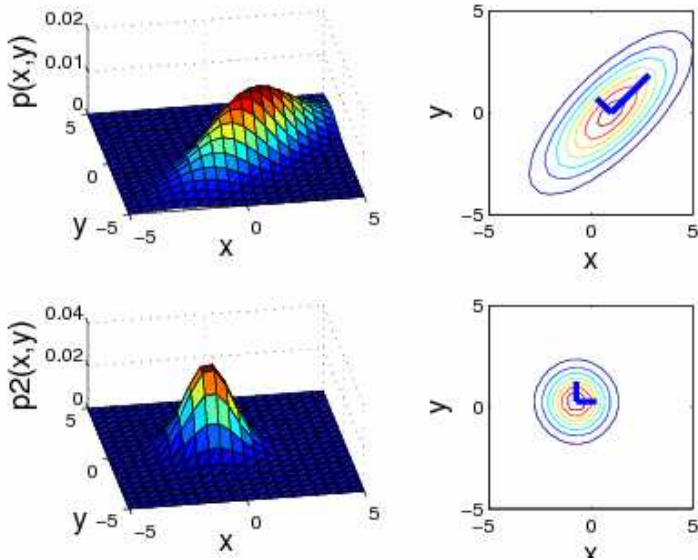
$$p(Y = c|x, \theta, \pi) \propto \exp \left[\log \pi_c + \sum_i I(x_i = 1) \log \theta_{ic} + I(x_i = 0) \log(1 - \theta_{ic}) \right]$$

$$x' = [1, I(x_1 = 1), I(x_1 = 0), \dots, I(x_d = 1), I(x_d = 0)]$$

$$\beta_c = [\log \pi_c, \log \theta_{1c}, \log(1 - \theta_{1c}), \dots, \log \theta_{dc}, \log(1 - \theta_{dc})]$$

$$p(Y = c|x, \beta) = \frac{\exp[\beta_c^T x']}{\sum_{c'} \exp[\beta_{c'}^T x']} \quad \text{Becomes sigmoid in 2-class case}$$

Multivariate normal



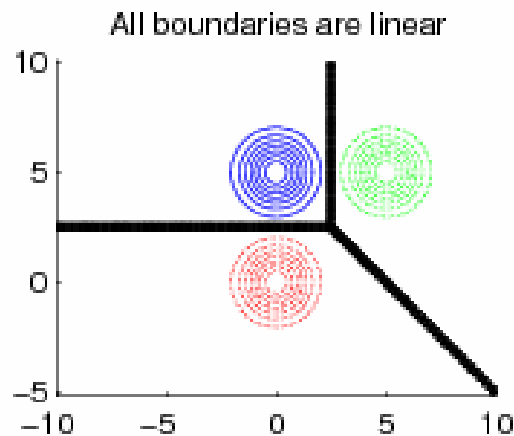
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \stackrel{\text{def}}{=} \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

$$\text{MLE: } \hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_i \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

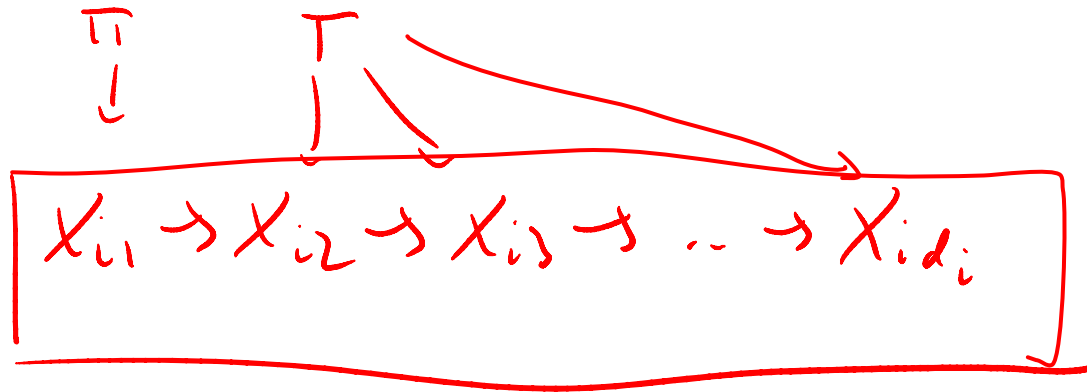
Gaussian classifiers

Tied Sigma, many classes

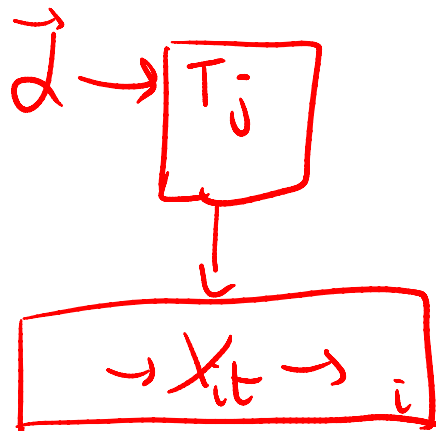
$$\begin{aligned} p(Y = c|\mathbf{x}) &= \frac{\pi_c \exp \left[-\frac{1}{2}(\mathbf{x} - \mu_c)^T \Sigma_c^{-1} (\mathbf{x} - \mu_c) \right]}{\sum_{c'} \pi_{c'} \exp \left[-\frac{1}{2}(\mathbf{x} - \mu_{c'})^T \Sigma_{c'}^{-1} (\mathbf{x} - \mu_{c'}) \right]} \\ &= \frac{\exp \left[\mu_c^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \log \pi_c \right]}{\sum_{c'} \exp \left[\mu_{c'}^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{c'}^T \Sigma^{-1} \mu_{c'} + \log \pi_{c'} \right]} \\ \theta_c &\stackrel{\text{def}}{=} \begin{pmatrix} -\mu_c^T \Sigma^{-1} \mu_c + \log \pi_c \\ \Sigma^{-1} \mu_c \end{pmatrix} = \begin{pmatrix} \gamma_c \\ \beta_c \end{pmatrix} \\ p(Y = c|\mathbf{x}) &= \frac{e^{\theta_c^T \mathbf{x}}}{\sum_{c'} e^{\theta_{c'}^T \mathbf{x}}} = \frac{e^{\beta_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\beta_{c'}^T \mathbf{x} + \gamma_{c'}}} \end{aligned}$$



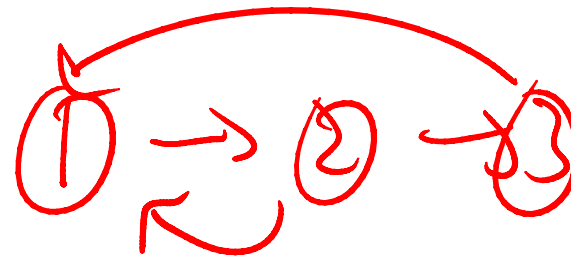
Markov chains



Language Models:
Empirical Bayes on rows of T
leads to backoff smoothing



Theory

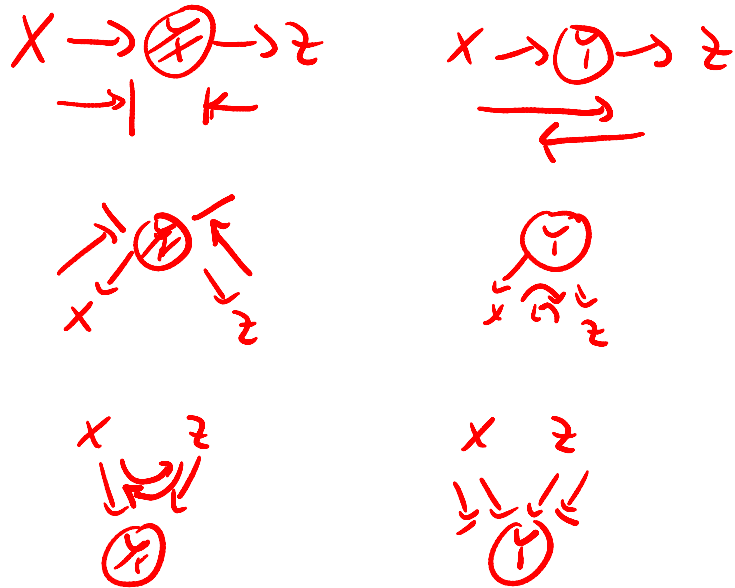


PageRank

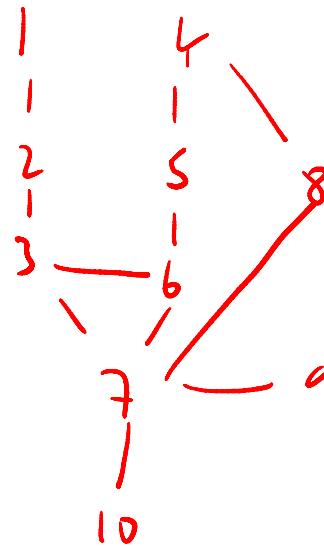
$$T_{ij} = \begin{cases} pG_{ij}/c_j + \delta & \text{if } c_j \neq 0 \\ 1/n & \text{if } c_j = 0 \end{cases}$$

Directed Graphical models

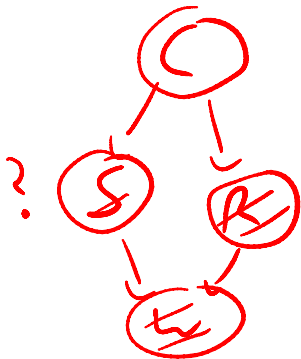
Bayes Ball



Moralization, ancestral graphs



State estimation



$$p(S = 1 | W = 1, R = 1) = \frac{p(S = 1, W = 1, R = 1)}{p(W = 1, R = 1)} = 0.19$$

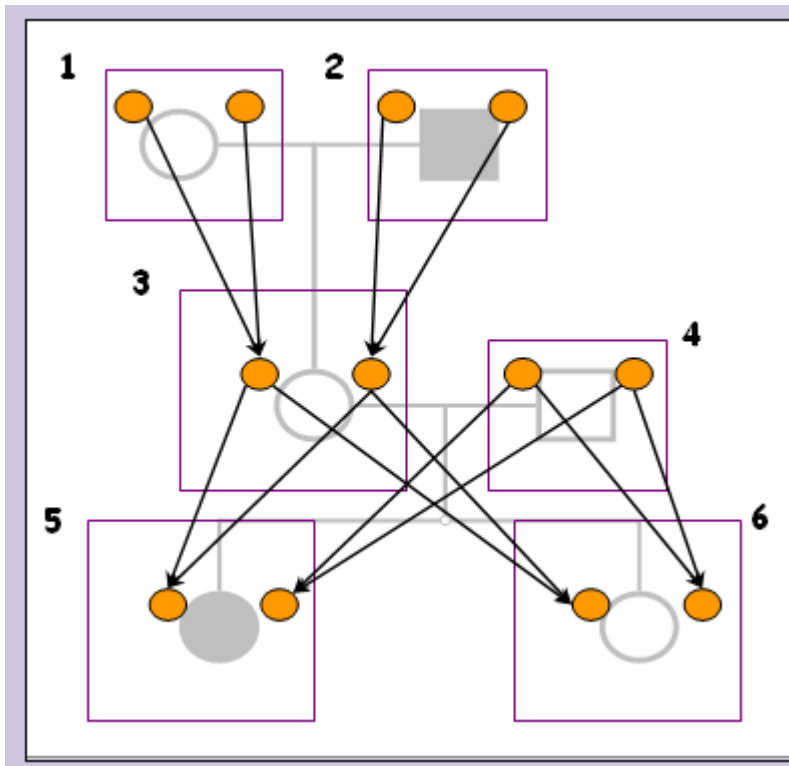
Parameter estimation



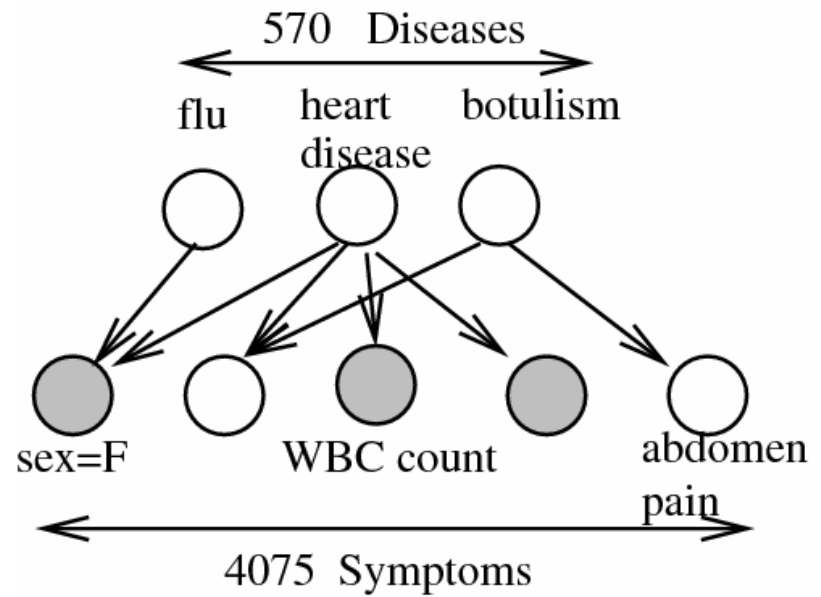
i	C	S	R	W	$p(\theta_C)$	$p(\theta_{R C=0})$	$p(\theta_{R C=1})$
1	0	0	0	0	1 1	1 1	1 1
2	0	0	1	1	2 1	2 1	1 1
3	1	1	1	1	3 1	2 2	1 1
					3 2	2 2	1 2

Expert systems

Pedigree trees

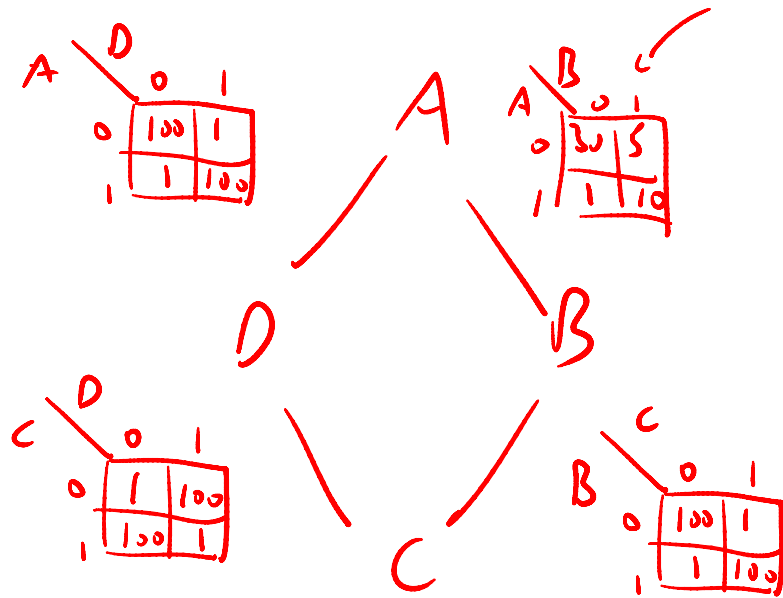


QMR



Noisy-or

Undirected graphical models



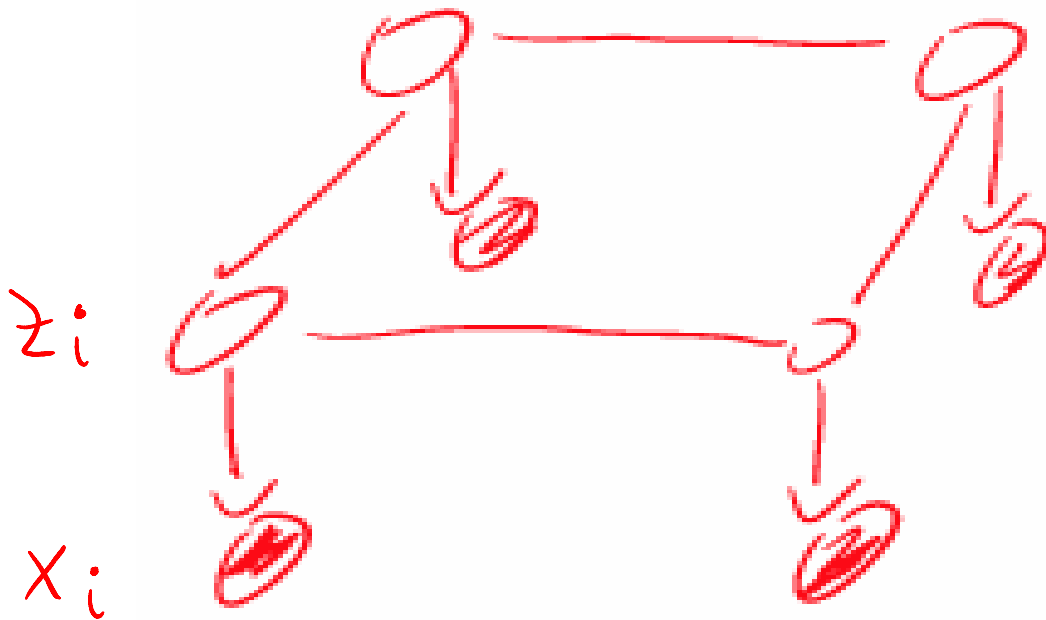
$$p(\mathbf{x}) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{x}_c)$$

State estimation

Assignment				Unnormalized	Normalized
a^0	b^0	c^0	d^0	300000	0.04
a^0	b^0	c^0	d^1	300000	0.04
a^0	b^0	c^1	d^0	300000	0.04
a^0	b^0	c^1	d^1	30	$4.1 \cdot 10^{-6}$
a^0	b^1	c^0	d^0	500	$6.9 \cdot 10^{-5}$
a^0	b^1	c^0	d^1	500	$6.9 \cdot 10^{-5}$
a^0	b^1	c^1	d^0	5000000	0.69
a^0	b^1	c^1	d^1	500	$6.9 \cdot 10^{-5}$
a^1	b^0	c^0	d^0	100	$1.4 \cdot 10^{-5}$
a^1	b^0	c^0	d^1	1000000	0.14
a^1	b^0	c^1	d^0	100	$1.4 \cdot 10^{-5}$
a^1	b^0	c^1	d^1	100	$1.4 \cdot 10^{-5}$
a^1	b^1	c^0	d^0	10	$1.4 \cdot 10^{-6}$
a^1	b^1	c^0	d^1	100000	0.014
a^1	b^1	c^1	d^0	100000	0.014
a^1	b^1	c^1	d^1	100000	0.014

Ising models

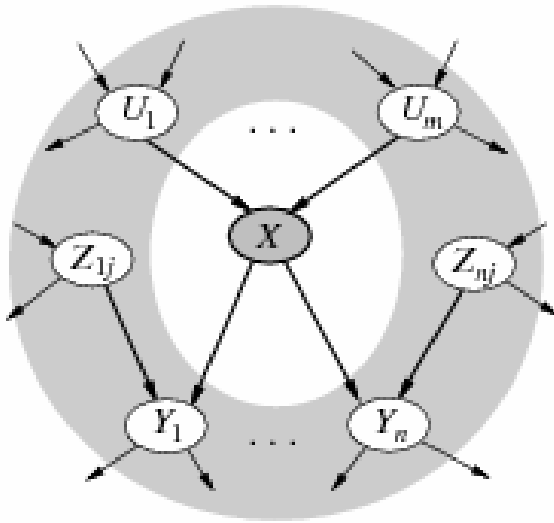
$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \left[\frac{1}{Z} \prod_{\langle ij \rangle} \psi_{ij}(z_i, z_j) \right] \left[\prod_i p(x_i|z_i) \right]$$



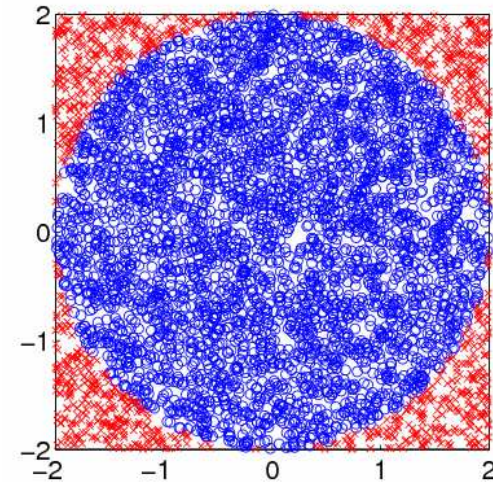
Gibbs sampling

1. $x_1^{s+1} \sim p(x_1|x_2^s, \dots, x_D^s)$
2. $x_2^{s+1} \sim p(x_2|x_1^{s+1}, x_3^s, \dots, x_D^s)$
3. $x_i^{s+1} \sim p(x_i|x_{1:i-1}^{s+1}, x_{i+1:D}^s)$
4. $x_D^{s+1} \sim p(x_D|x_1^{s+1}, \dots, x_{D-1}^{s+1})$

Markov blanket



Monte Carlo integration



Full conditional

$$p(X_i|X_{-i}) \propto p(X_i|Pa(X_i)) \prod_{Y_j \in ch(X_i)} p(Y_j|Pa(Y_j))$$