

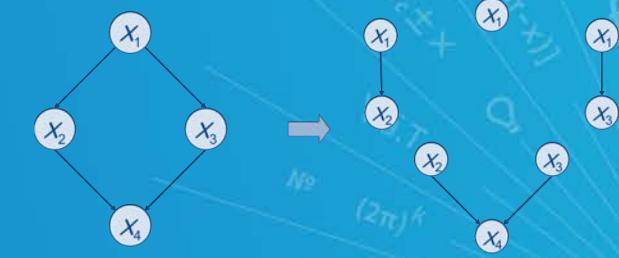
# Probabilistic Graphical Models

## Parameter Est. in fully observed BNs

Eric Xing

Lecture 5, January 30, 2019

Reading: see class homepage

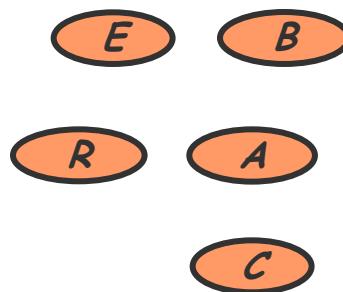




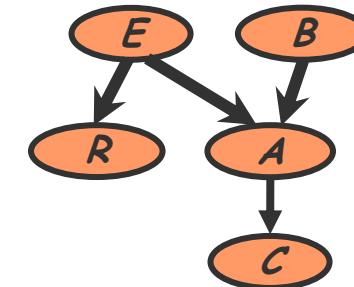
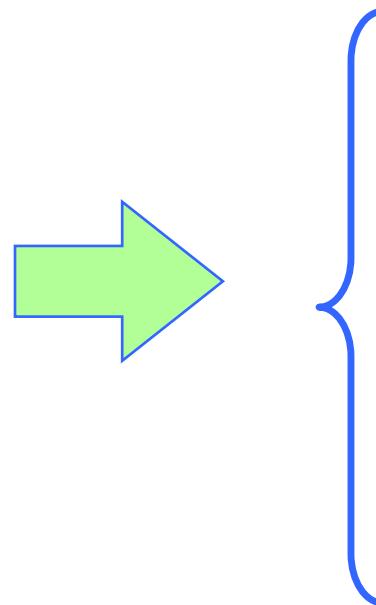
# Learning Graphical Models

The goal:

- Given set of independent samples (*assignments* of random variables), find the *best* (the most likely?) Bayesian Network (both DAG and CPDs)



$(B, E, A, C, R) = (T, F, F, T, F)$   
 $(B, E, A, C, R) = (T, F, T, T, F)$   
.....  
 $(B, E, A, C, R) = (F, T, T, T, F)$



**Structural learning**

$E$	$B$	$P(A   E, B)$	
$e$	$b$	0.9	0.1
$e$	$\bar{b}$	0.2	0.8
$\bar{e}$	$b$	0.9	0.1
$\bar{e}$	$\bar{b}$	0.01	0.99

**Parameter learning**

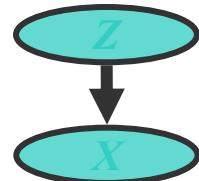




# Learning Graphical Models

- ❑ Scenarios:
  - ❑ completely observed GMs
    - ❑ directed
    - ❑ undirected
  - ❑ partially or unobserved GMs
    - ❑ directed
    - ❑ undirected (an open research topic)
- ❑ Estimation principles:
  - ❑ Maximal likelihood estimation (MLE)
  - ❑ Bayesian estimation
  - ❑ Maximal conditional likelihood
  - ❑ Maximal "Margin"
  - ❑ Maximum entropy
- ❑ We use **learning** as a name for the process of **estimating the parameters**, and in some cases, the topology of the network, from data.





# ML Parameter Est. for completely observed GMs of given structure

- The data:

$$\{(z_1, x_1), (z_2, x_2), (z_3, x_3), \dots (z_N, x_N)\}$$





# Parameter Learning

- Assume  $G$  is known and fixed,
  - from expert design
  - from an intermediate outcome of iterative structure learning
- Goal: estimate  $\theta$  from a dataset of  $N$  independent, identically distributed (*iid*) training cases  $D = \{x_1, \dots, x_N\}$ .
- In general, each training case  $\mathbf{x}_n = (x_{n,1}, \dots, x_{n,M})$  is a vector of  $M$  values, one per node,
  - the model can be completely observable, i.e., every element in  $x_n$  is known (no missing values, no hidden variables),
  - or, partially observable, i.e.,  $\exists i$ , s.t.  $x_{n,i}$  is not observed.
- In this lecture we consider learning parameters for a BN with given structure and is completely observable

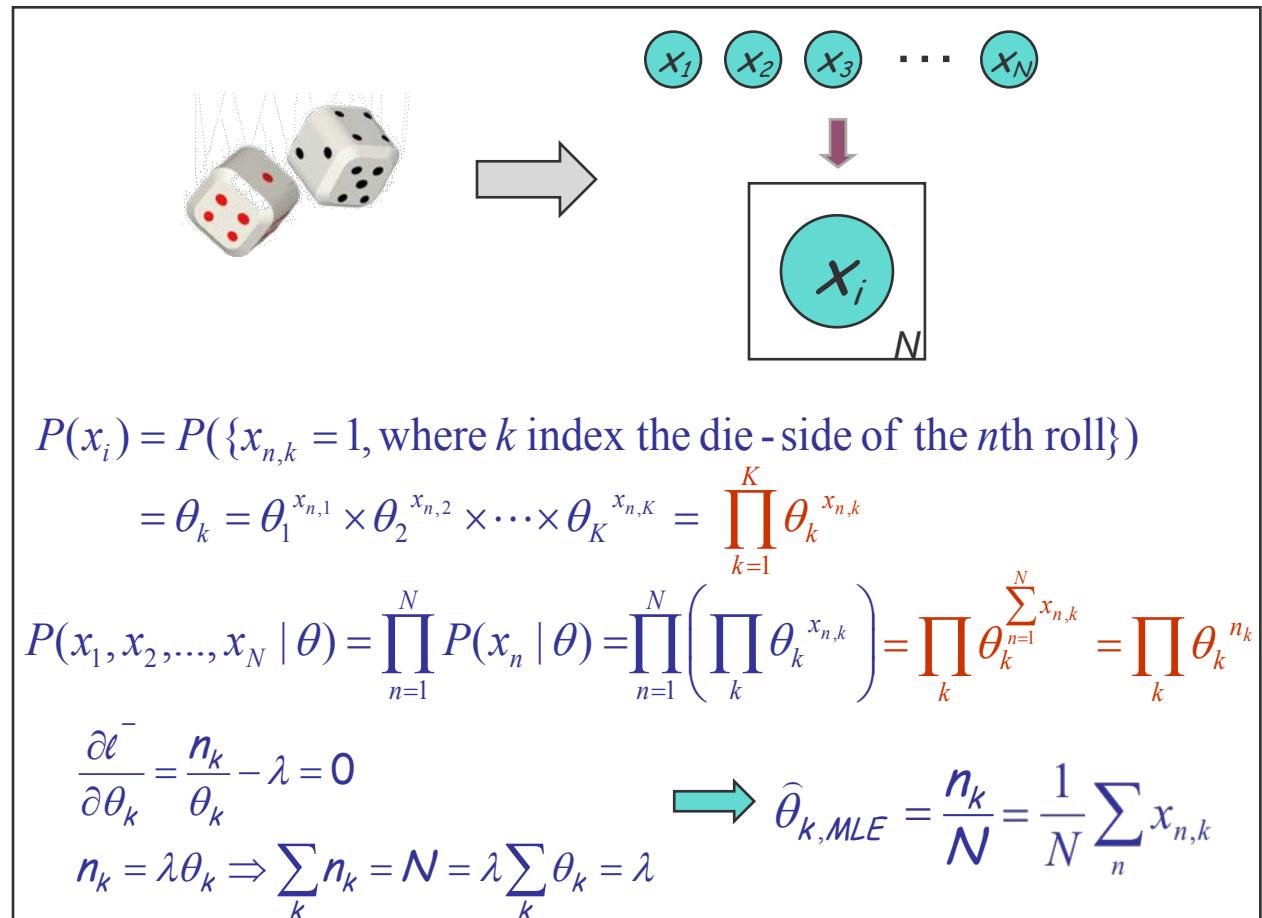
$$\ell(\theta; D) = \log p(D | \theta) = \log \prod_n \left( \prod_i p(x_{n,i} | \mathbf{x}_{n,\pi_i}, \theta_i) \right) = \sum_i \left( \sum_n \log p(x_{n,i} | \mathbf{x}_{n,\pi_i}, \theta_i) \right)$$





# Review of density estimation

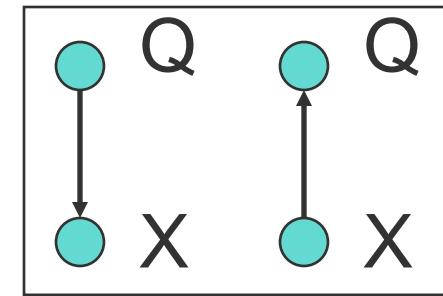
- Can be viewed as single-node GMs
- Instances of  
Exponential Family Dist.
- Building blocks of general GM
- MLE and Bayesian estimate
- See supplementary slides





# Estimation of conditional density

- ❑ Can be viewed as two-node graphical models
- ❑ Instances of GLIM (Generalized Linear Models)
- ❑ Building blocks of general GM
- ❑ MLE and Bayesian estimate
- ❑ See supplementary slides





# Exponential family, a basic building block

- For a numeric random variable  $X$

$$\begin{aligned} p(x | \eta) &= h(x) \exp\{\eta^T T(x) - A(\eta)\} \\ &= \frac{1}{Z(\eta)} h(x) \exp\{\eta^T T(x)\} \end{aligned}$$

is an **exponential family distribution** with natural (canonical) parameter  $\eta$

- Function  $T(x)$  is a *sufficient statistic*.
- Function  $A(\eta) = \log Z(\eta)$  is the log normalizer.
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma, ...





# Example: Multivariate Gaussian Distribution

- For a continuous vector random variable  $X \in \mathbb{R}^k$ :

$$\begin{aligned} p(x|\mu, \Sigma) &= \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\} \\ &= \frac{1}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} xx^T) + \mu^T \Sigma^{-1} x - \frac{1}{2} \mu^T \Sigma^{-1} \mu - \log|\Sigma|\right\} \end{aligned}$$

Moment parameter

- Exponential family representation

$$\begin{aligned} \eta &= [\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1})] = [\eta_1, \text{vec}(\eta_2)], \quad \eta_1 = \Sigma^{-1} \mu \text{ and } \eta_2 = -\frac{1}{2} \Sigma^{-1} \\ T(x) &= [x; \text{vec}(xx^T)] \\ A(\eta) &= \frac{1}{2} \mu^T \Sigma^{-1} \mu + \log|\Sigma| = -\frac{1}{2} \text{tr}(\eta_2 \eta_1 \eta_1^T) - \frac{1}{2} \log(-2\eta_2) \\ h(x) &= (2\pi)^{-k/2} \end{aligned}$$

Natural parameter

- Note: a  $k$ -dimensional Gaussian is a  $(d+d)$ -parameter distribution with a  $(d+d)$ -element vector of sufficient statistics (but because of symmetry and positivity, parameters are constrained and have lower degree of freedom)





# Example: Multinomial distribution

- For a binary vector random variable  $\boldsymbol{x} \sim \text{multi}(\boldsymbol{x} | \boldsymbol{\pi})$ ,

$$\begin{aligned} p(x|\boldsymbol{\pi}) &= \pi_1^{x_1} \pi_2^{x_2} \cdots \pi_K^{x_K} = \exp\left\{\sum_k x_k \ln \pi_k\right\} \\ &= \exp\left\{\sum_{k=1}^{K-1} x_k \ln \pi_k + \left(1 - \sum_{k=1}^{K-1} x_k\right) \ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right)\right\} \\ &= \exp\left\{\sum_{k=1}^{K-1} x_k \ln\left(\frac{\pi_k}{1 - \sum_{k=1}^{K-1} \pi_k}\right) + \ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right)\right\} \end{aligned}$$

- Exponential family representation

$$\begin{aligned} \boldsymbol{\eta} &= \left[ \ln\left(\frac{\pi_k}{\pi_K}\right); \mathbf{0} \right] \\ T(\boldsymbol{x}) &= [\boldsymbol{x}] \end{aligned}$$

$$A(\boldsymbol{\eta}) = -\ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right) = \ln\left(\sum_{k=1}^K e^{\eta_k}\right)$$

$$h(\boldsymbol{x}) = 1$$





# Why exponential family?

- Moment generating property

$$\begin{aligned}\frac{dA}{d\eta} &= \frac{d}{d\eta} \log Z(\eta) = \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= \frac{1}{Z(\eta)} \frac{d}{d\eta} \int h(x) \exp\{\eta^T T(x)\} dx \\ &= \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \\ &= E[T(x)]\end{aligned}$$

$$\begin{aligned}\frac{d^2 A}{d\eta^2} &= \int T^2(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx - \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= E[T^2(x)] - E^2[T(x)] \\ &= Var[T(x)]\end{aligned}$$





# Moment estimation

- We can easily compute moments of any exponential family distribution by taking the derivatives of the log normalizer  $A(\eta)$ .
- The  $q^{\text{th}}$  derivative gives the  $q^{\text{th}}$  centered moment.

$$\frac{dA(\eta)}{d\eta} = \text{mean}$$

$$\frac{d^2 A(\eta)}{d\eta^2} = \text{variance}$$

...

- When the sufficient statistic is a stacked vector, partial derivatives need to be considered.





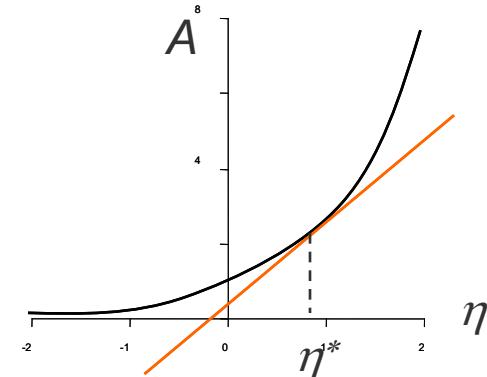
# Moment vs canonical parameters

- The moment parameter  $\mu$  can be derived from the natural (canonical) parameter

$$\frac{dA(\eta)}{d\eta} = E[T(x)] \stackrel{\text{def}}{=} \mu$$

- $A(\eta)$  is convex since

$$\frac{d^2 A(\eta)}{d\eta^2} = Var[T(x)] > 0$$



- Hence we can invert the relationship and infer the canonical parameter from the moment parameter (1-to-1):

$$\eta \stackrel{\text{def}}{=} \psi(\mu)$$

- A distribution in the exponential family can be parameterized not only by  $\eta$  – the canonical parameterization, but also by  $\mu$  – the moment parameterization.





# MLE for Exponential Family

- For *iid* data, the log-likelihood is

$$\begin{aligned}\ell(\eta; D) &= \log \prod_n h(x_n) \exp\left\{\eta^T T(x_n) - A(\eta)\right\} \\ &= \sum_n \log h(x_n) + \left( \eta^T \sum_n T(x_n) \right) - N A(\eta)\end{aligned}$$

- Take derivatives and set to zero:

$$\frac{\partial \ell}{\partial \eta} = \sum_n T(x_n) - N \frac{\partial A(\eta)}{\partial \eta} = 0$$

$$\Rightarrow \frac{\partial A(\eta)}{\partial \eta} = \frac{1}{N} \sum_n T(x_n)$$

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_n T(x_n)$$

- This amounts to **moment matching**.
- We can infer the canonical parameters using  $\hat{\eta}_{MLE} = \psi(\hat{\mu}_{MLE})$





# Sufficiency

- For  $p(x|\theta)$ ,  $T(x)$  is *sufficient* for  $\theta$  if there is no information in  $X$  regarding  $\theta$  beyond that in  $T(x)$ .
  - We can throw away  $X$  for the purpose of inference w.r.t.  $\theta$ .

- Bayesian view

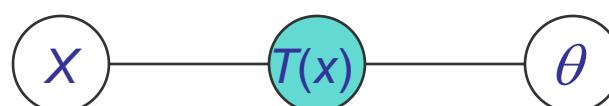


- Frequentist view



- The Neyman factorization theorem

- $T(x)$  is *sufficient* for  $\theta$  if



$$p(x, T(x), \theta) = \psi_1(T(x), \theta) \psi_2(x, T(x))$$

$$\Rightarrow p(x | \theta) = g(T(x), \theta) h(x, T(x))$$





# Examples

## □ Gaussian:

$$\begin{aligned}\eta &= \left[ \Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right] \\ T(x) &= \left[ x; \text{vec}(xx^T) \right] \\ A(\eta) &= \frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log |\Sigma| \\ h(x) &= (2\pi)^{-k/2}\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n T_1(x_n) = \frac{1}{N} \sum_n x_n$$

## □ Multinomial:

$$\begin{aligned}\eta &= \left[ \ln\left(\frac{\pi_k}{\pi_K}\right); \mathbf{0} \right] \\ T(x) &= [x] \\ A(\eta) &= -\ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right) = \ln\left(\sum_{k=1}^K e^{\eta_k}\right) \\ h(x) &= 1\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$

## □ Poisson:

$$\begin{aligned}\eta &= \log \lambda \\ T(x) &= x \\ A(\eta) &= \lambda = e^\eta \\ h(x) &= \frac{1}{x!}\end{aligned}$$

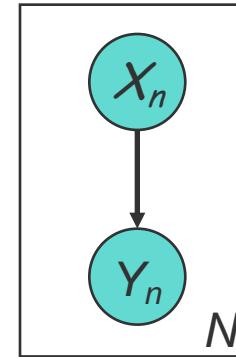
$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$





# Generalized Linear Models (GLIMs)

- The graphical model
  - Linear regression
  - Discriminative linear classification
  - Commonality:
    - model  $E_p(Y) = \mu = f(\theta^T X)$
    - What is  $p()$ ? the cond. dist. of  $Y$ .
    - What is  $f()$ ? the response function.
- GLIM
  - The observed input  $x$  is assumed to enter into the model via a linear combination of its elements
  - The conditional mean  $\mu$  is represented as a function  $f(\xi)$  of  $\xi$ , where  $f$  is known as the response function  $\xi = \theta^T x$
  - The observed output  $y$  is assumed to be characterized by an exponential family distribution with conditional mean  $\mu$ .





# Recall Linear Regression

- Let us assume that the target variable and the inputs are related by the equation:

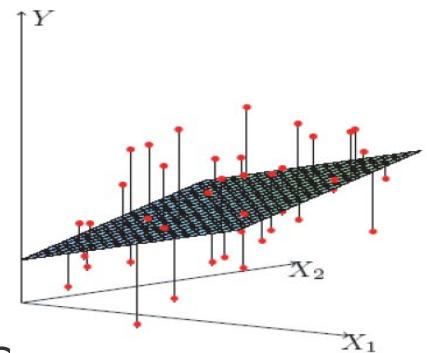
$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

where  $\varepsilon$  is an error term of unmodeled effects or random noise

- Now assume that  $\varepsilon$  follows a Gaussian  $N(0, \sigma)$ , then we have.

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

- We can use LMS algorithm, which is a gradient ascent/descent approach, to estimate the parameter





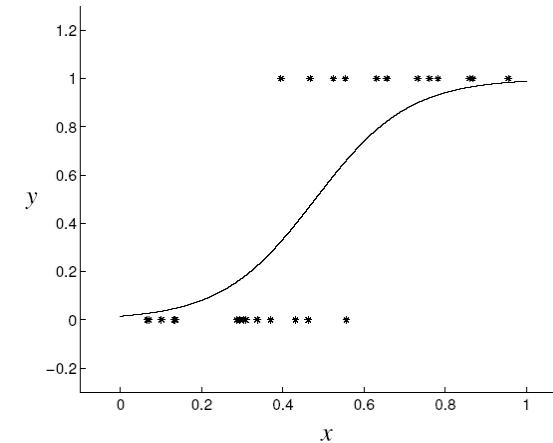
# Recall: Logistic Regression (sigmoid classifier, perceptron, etc.)

- The condition distribution: a Bernoulli

$$p(y|x) = \mu(x)^y (1 - \mu(x))^{1-y}$$

where  $\mu$  is a logistic function

$$\mu(x) = \frac{1}{1 + e^{-\theta^T x}}$$



- We can used the brute-force gradient method as in LR
- But we can also apply generic laws by observing the  $p(y|x)$  is an **exponential family function**, more specifically, a **generalized linear model!**

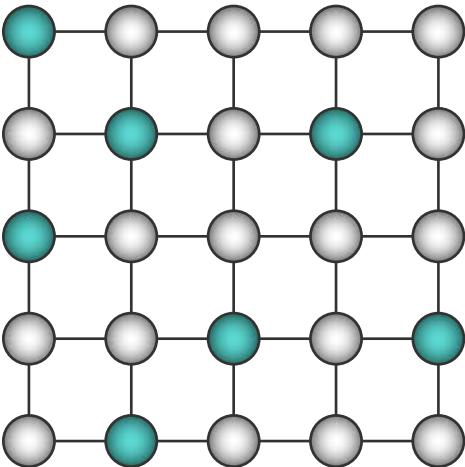




# More examples: parameterizing graphical models

- Markov random fields

$$p(\mathbf{x}) = \frac{1}{Z} \exp \left\{ - \sum_{c \in C} \phi_c(\mathbf{x}_c) \right\} = \frac{1}{Z} \exp \left\{ - H(\mathbf{x}) \right\}$$



$$p(X) = \frac{1}{Z} \exp \left\{ \sum_{i,j \in N_i} \theta_{ij} X_i X_j + \sum_i \theta_{i0} X_i \right\}$$

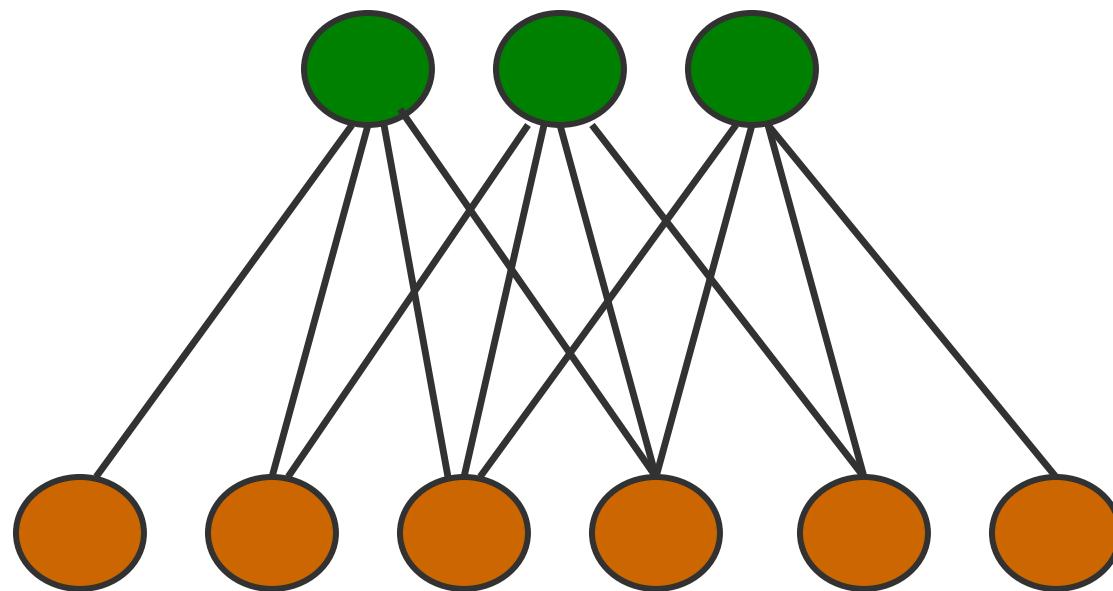




# Restricted Boltzmann Machines

hidden units

visible units

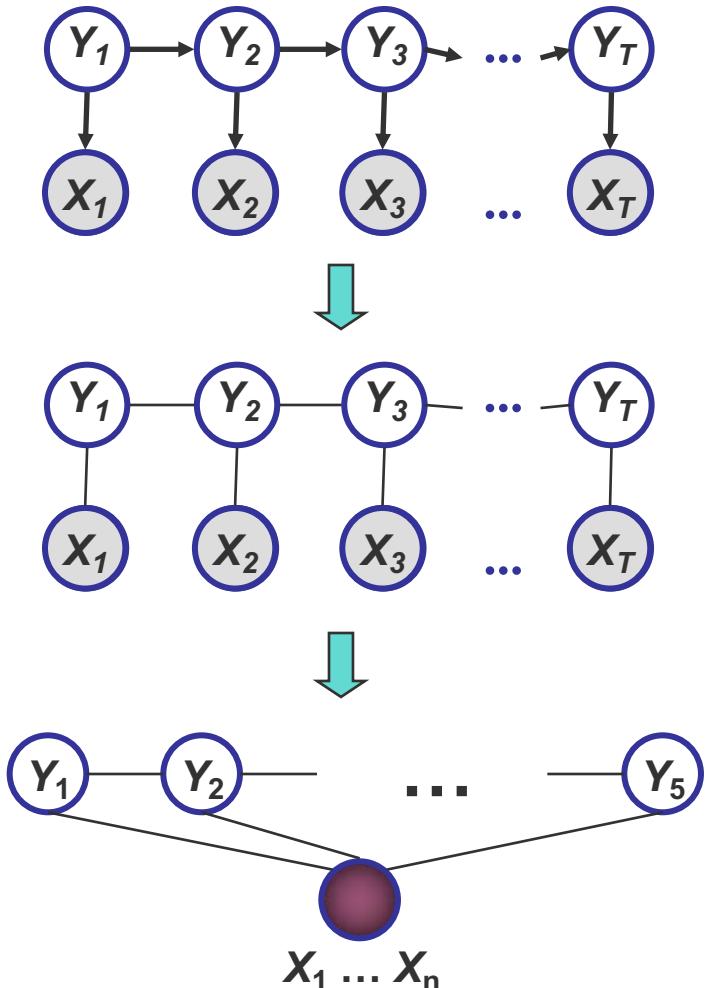


$$p(x, h | \theta) = \exp \left\{ \sum_i \theta_i \phi_i(x_i) + \sum_j \theta_j \phi_j(h_j) + \sum_{i,j} \theta_{i,j} \phi_{i,j}(x_i, h_j) - A(\theta) \right\}$$





# Conditional Random Fields



- Discriminative

$$p_\theta(y|x) = \frac{1}{Z(\theta, x)} \exp \left\{ \sum_c \theta_c f_c(x, y_c) \right\}$$

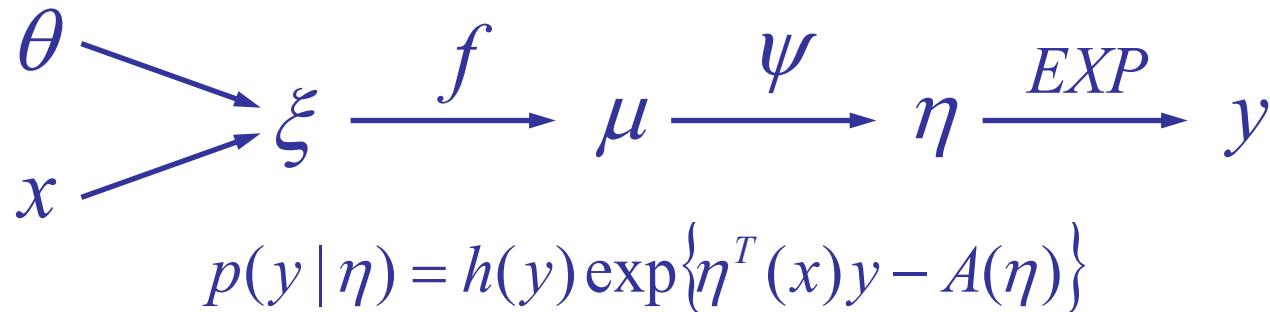
- $X_i$ 's are assumed as features that are inter-dependent

- When labeling  $X_i$  future observations are taken into account





# GLIM, cont.



$$\Rightarrow p(y|\eta, \phi) = h(y, \phi) \exp\left\{\frac{1}{\phi}(\eta^T(x)y - A(\eta))\right\}$$

- The choice of exp family is constrained by the nature of the data  $y$ 
  - Example:  $y$  is a continuous vector  $\rightarrow$  multivariate Gaussian
  - $y$  is a class label  $\rightarrow$  Bernoulli or multinomial
- The choice of the response function
  - Following some mild constraints, e.g.,  $[0,1]$ . Positivity ...
  - **Canonical response** function:
    - In this case  $\theta^T x$  directly corresponds to canonical parameter  $\eta$ .  $f = \psi^{-1}(\cdot)$





# Example canonical response functions

Model	Canonical response function
Gaussian	$\mu = \eta$
Bernoulli	$\mu = 1/(1 + e^{-\eta})$
multinomial	$\mu_i = \eta_i / \sum_j e^{\eta_j}$
Poisson	$\mu = e^\eta$
gamma	$\mu = -\eta^{-1}$





# MLE for GLIMs with natural response

- Log-likelihood

$$\ell = \sum_n \log h(y_n) + \sum_n (\theta^T x_n y_n - A(\eta_n))$$

- Derivative of Log-likelihood

$$\frac{d\ell}{d\theta} = \sum_n \left( x_n y_n - \frac{dA(\eta_n)}{d\eta_n} \frac{d\eta_n}{d\theta} \right)$$

$$= \sum_n (y_n - \mu_n) x_n$$

$$= X^T (y - \mu)$$

This is a fixed point function  
because  $\mu$  is a function of  $\theta$

- Online learning for canonical GLIMs

- Stochastic gradient ascent = least mean squares (LMS) algorithm:

$$\theta^{t+1} = \theta^t + \rho (y_n - \mu_n^t) x_n$$

where  $\mu_n^t = (\theta^t)^T x_n$  and  $\rho$  is a step size





# Batch learning for canonical GLIMs

- The Hessian matrix

$$\begin{aligned} H &= \frac{d^2\ell}{d\theta d\theta^T} = \frac{d}{d\theta^T} \sum_n (y_n - \mu_n) x_n = \sum_n x_n \frac{d\mu_n}{d\theta^T} \\ &= - \sum_n x_n \frac{d\mu_n}{d\eta_n} \frac{d\eta_n}{d\theta^T} \\ &= - \sum_n x_n \frac{d\mu_n}{d\eta_n} x_n^T \quad \text{since } \eta_n = \theta^T x_n \\ &= -X^T W X \end{aligned}$$

$$\begin{aligned} \mathbf{X} &= \begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_n & \cdots \end{bmatrix} \\ \vec{y} &= \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \end{aligned}$$

where  $X = [x_n^T]$  is the design matrix and

$$W = \text{diag}\left(\frac{d\mu_1}{d\eta_1}, \dots, \frac{d\mu_N}{d\eta_N}\right)$$

which can be computed by calculating the 2<sup>nd</sup> derivative of  $\mathcal{A}(\eta_n)$





# Recall LMS

- Cost function in matrix form:

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2 \\ &= \frac{1}{2} (\mathbf{X}\theta - \bar{\mathbf{y}})^T (\mathbf{X}\theta - \bar{\mathbf{y}}) \end{aligned}$$

$$\begin{aligned} \mathbf{X} &= \begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_n & \cdots \end{bmatrix} \\ \bar{\mathbf{y}} &= \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \end{aligned}$$

- To minimize  $J(\theta)$ , take derivative and set to zero:

$$\begin{aligned} \nabla_{\theta} J &= \frac{1}{2} \nabla_{\theta} \text{tr}(\theta^T X^T X \theta - \theta^T X^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T X \theta + \bar{\mathbf{y}}^T \bar{\mathbf{y}}) \\ &= \frac{1}{2} (\nabla_{\theta} \text{tr} \theta^T X^T X \theta - 2 \nabla_{\theta} \text{tr} \bar{\mathbf{y}}^T X \theta + \nabla_{\theta} \text{tr} \bar{\mathbf{y}}^T \bar{\mathbf{y}}) \\ &= \frac{1}{2} (X^T X \theta + X^T X \theta - 2 X^T \bar{\mathbf{y}}) \\ &= X^T X \theta - X^T \bar{\mathbf{y}} = 0 \end{aligned}$$

$$\Rightarrow \boxed{X^T X \theta = X^T \bar{\mathbf{y}}} \\ \text{The normal equations}$$

$$\downarrow \\ \theta^* = (X^T X)^{-1} X^T \bar{\mathbf{y}}$$





# Iteratively Reweighted Least Squares (IRLS)

- Recall Newton-Raphson methods with cost function  $\mathcal{J}$

$$\theta^{t+1} = \theta^t - H^{-1} \nabla_{\theta} J$$

- We now have

$$\nabla_{\theta} J = X^T (y - \mu)$$

$$H = -X^T W X$$

$$\theta^* = (X^T X)^{-1} X^T \bar{y}$$

- Now:

$$\begin{aligned}\theta^{t+1} &= \theta^t + H^{-1} \nabla_{\theta} \ell \\ &= (X^T W^t X)^{-1} [X^T W^t X \theta^t + X^T (y - \mu^t)] \\ &= (X^T W^t X)^{-1} X^T W^t z^t\end{aligned}$$

where the adjusted response is  $z^t = X \theta^t + (W^t)^{-1} (y - \mu^t)$

- This can be understood as solving the following "Iteratively reweighted least squares" problem

$$\theta^{t+1} = \arg \min_{\theta} (z - X \theta)^T W (z - X \theta)$$





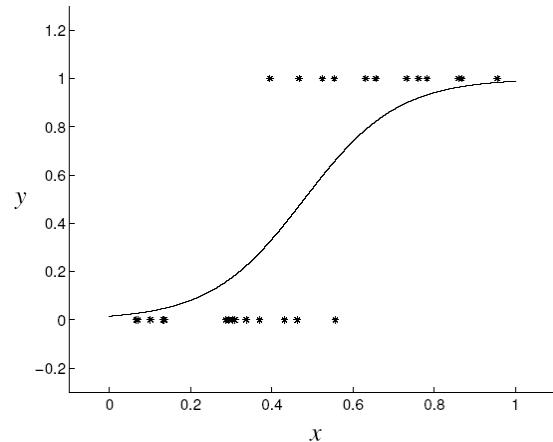
# Example 1: logistic regression (sigmoid classifier)

- The condition distribution: a Bernoulli

$$p(y|x) = \mu(x)^y (1 - \mu(x))^{1-y}$$

where  $\mu$  is a logistic function

$$\mu(x) = \frac{1}{1 + e^{-\eta(x)}}$$



- $p(y|x)$  is an exponential family function, with

- mean:  $E[y|x] = \mu = \frac{1}{1 + e^{-\eta(x)}}$

- and canonical response function  $\eta = \xi = \theta^T x$

- IRLS

$$\frac{d\mu}{d\eta} = \mu(1 - \mu)$$

$$W = \begin{pmatrix} \mu_1(1 - \mu_1) & & \\ & \ddots & \\ & & \mu_N(1 - \mu_N) \end{pmatrix}$$





# Logistic regression: practical issues

- It is very common to use *regularized* maximum likelihood.

$$p(y = \pm 1 | x, \theta) = \frac{1}{1 + e^{-y\theta^T x}} = \sigma(y\theta^T x)$$

$$p(\theta) \sim \text{Normal}(0, \lambda^{-1} I)$$

$$l(\theta) = \sum_n \log(\sigma(y_n \theta^T x_n)) - \frac{\lambda}{2} \theta^T \theta$$

- IRLS takes  $\mathcal{O}(Nd^3)$  per iteration, where  $N$  = number of training cases and  $d$  = dimension of input  $x$ .
- Quasi-Newton methods, that approximate the Hessian, work faster.
- Conjugate gradient takes  $\mathcal{O}(Nd)$  per iteration, and usually works best in practice.
- Stochastic gradient descent can also be used if  $N$  is large c.f. perceptron rule:

$$\nabla_{\theta} \ell = (1 - \sigma(y_n \theta^T x_n)) y_n x_n - \lambda \theta$$





## Example 2: linear regression

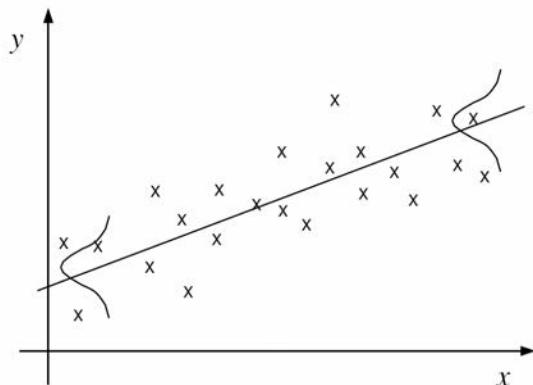
- The condition distribution: a Gaussian

$$p(y|x, \theta, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(y - \mu(x))^T \Sigma^{-1} (y - \mu(x))\right\}$$

**Rescale**  $\Rightarrow h(x) \exp\left\{-\frac{1}{2}\Sigma^{-1}(\eta^T(x)y - A(\eta))\right\}$

where  $\mu$  is a linear function

$$\mu(x) = \theta^T x = \eta(x)$$



- $p(y|x)$  is an exponential family function, with

- mean:

$$E[y|x] = \mu = \theta^T x$$

- and canonical response function

$$\eta_1 = \xi = \theta^T x$$

- IRLS

$$\begin{aligned} \frac{d\mu}{d\eta} &= 1 & \theta^{t+1} &= (X^T W^t X)^{-1} X^T W^t z^t \\ \frac{d\eta}{d\theta} &= 1 & &= (X^T X)^{-1} X^T (X\theta^t + (y - \mu^t)) \\ W &= I & &= \theta^t + (X^T X)^{-1} X^T (y - \mu^t) \end{aligned} \quad \stackrel{t \rightarrow \infty}{\Rightarrow} \quad \theta = (X^T X)^{-1} X^T Y$$

Steepest descent

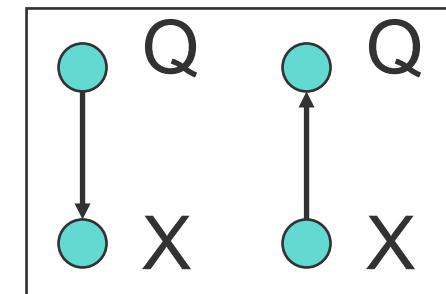
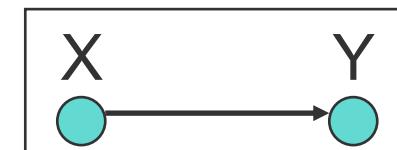
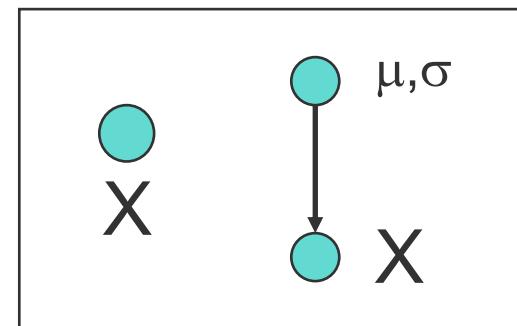
Normal equation





# Simple GMs are the building blocks of complex GMs

- ❑ Density estimation
  - ❑ Parametric and nonparametric methods
- ❑ Regression
  - ❑ Linear, conditional mixture, nonparametric
- ❑ Classification
  - ❑ Generative and discriminative approach
- ❑ Clustering

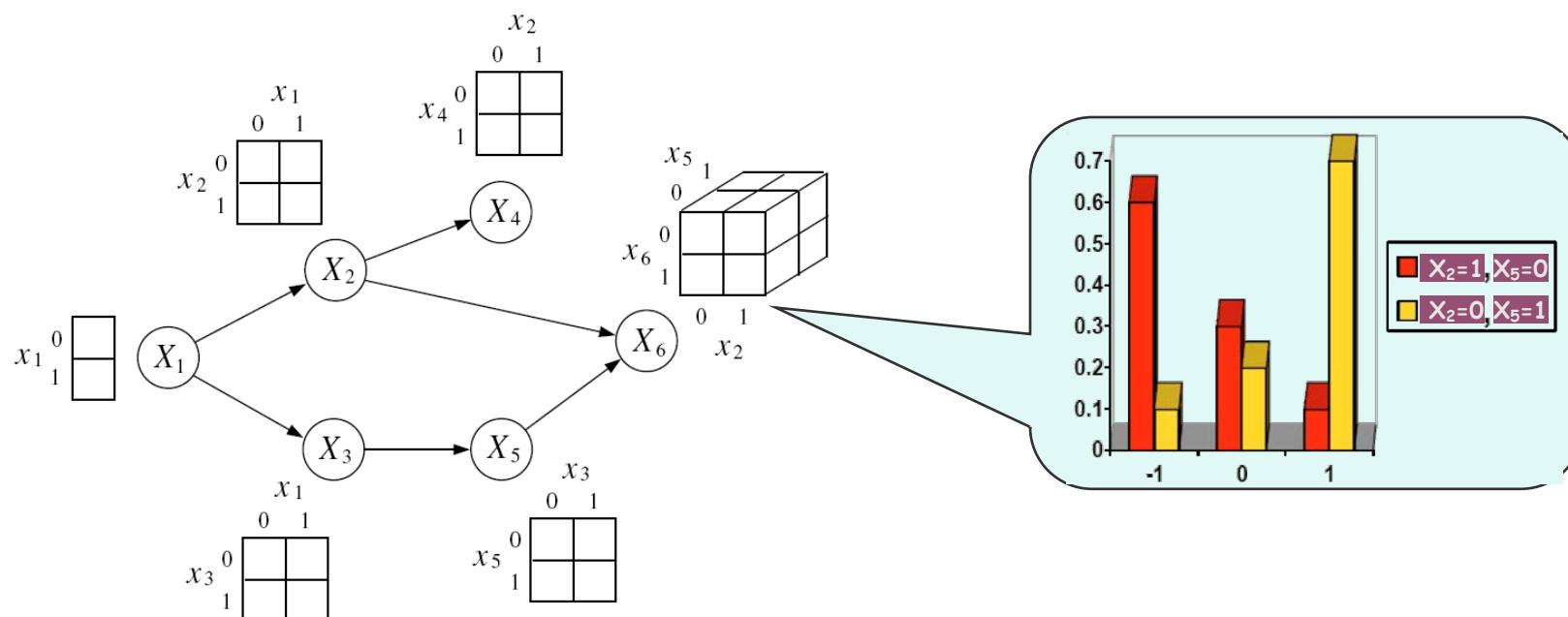




# MLE for general BNs

- If we assume the parameters for each CPD are globally independent, and all nodes are fully observed, then the log-likelihood function decomposes into a sum of local terms, one per node:

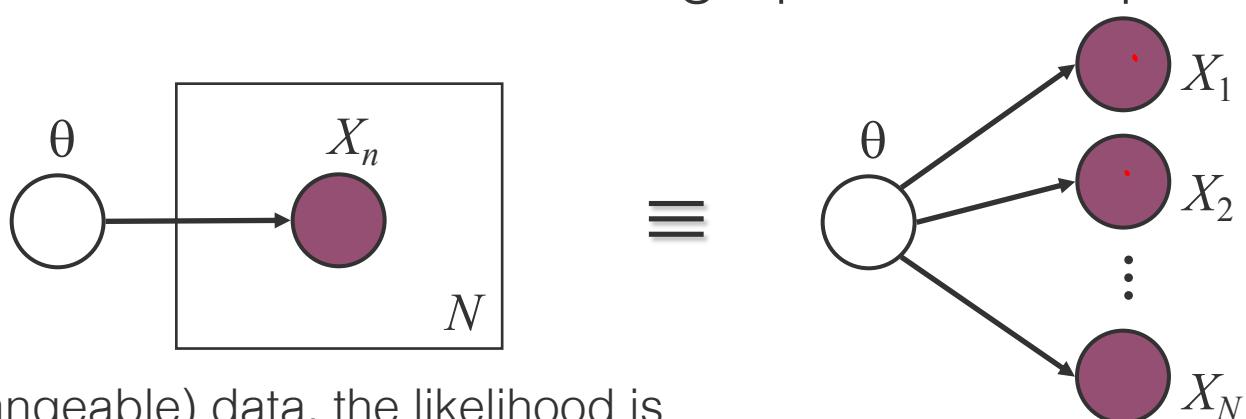
$$\ell(\theta; D) = \log p(D | \theta) = \log \prod_n \left( \prod_i p(x_{n,i} | \mathbf{x}_{n,\pi_i}, \theta_i) \right) = \sum_i \left( \sum_n \log p(x_{n,i} | \mathbf{x}_{n,\pi_i}, \theta_i) \right)$$





# Plates

- A plate is a “macro” that allows subgraphs to be replicated



- For iid (exchangeable) data, the likelihood is

$$p(D | \theta) = \prod_n p(x_n | \theta)$$

- We can represent this as a Bayes net with  $N$  nodes.
  - The rules of plates are simple: repeat every structure in a box a number of times given by the integer in the corner of the box (e.g.  $N$ ), updating the plate index variable (e.g.  $n$ ) as you go.
  - Duplicate every arrow going into the plate and every arrow leaving the plate by connecting the arrows to each copy of the structure.



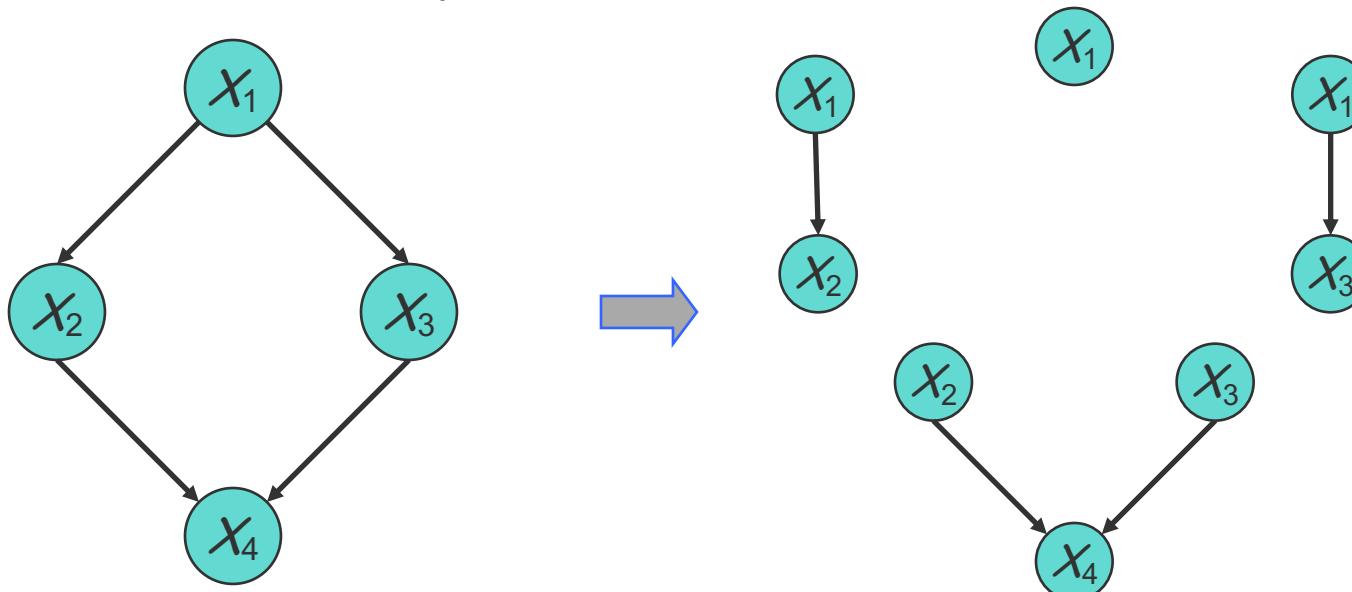


# Decomposable likelihood of a BN

- Consider the distribution defined by the directed acyclic GM:

$$p(x|\theta) = p(x_1|\theta_1)p(x_2|x_1,\theta_2)p(x_3|x_1,\theta_3)p(x_4|x_2,x_3,\theta_4)$$

- This is exactly like learning four separate small BNs, each of which consists of a node and its parents.





# MLE for BNs with tabular CPDs

- Assume each CPD is represented as a table (multinomial) where

$$\theta_{ijk} \stackrel{\text{def}}{=} p(X_i = j | X_{\pi_i} = k)$$

- Note that in case of multiple parents,  $\mathbf{X}_{\pi_i}$  will have a composite state, and the CPD will be a high-dimensional table
- The sufficient statistics are counts of family configurations

$$n_{ijk} \stackrel{\text{def}}{=} \sum_n x_{n,i}^j x_{n,\pi_i}^k$$

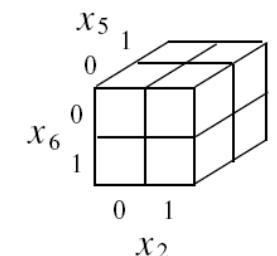
- The log-likelihood is

$$\ell(\theta; \mathcal{D}) = \log \prod_{i,j,k} \theta_{ijk}^{n_{ijk}} = \sum_{i,j,k} n_{ijk} \log \theta_{ijk}$$

- Using a Lagrange multiplier to enforce  $\sum_j \theta_{ijk} = 1$ , we get:

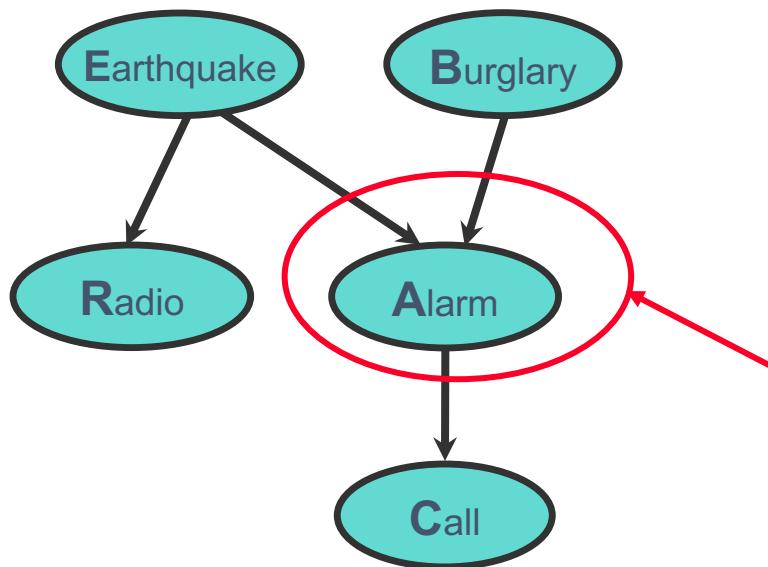
$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{j'} n_{ij'k}}$$

			$x_1$
	0	1	
$x_2$	0		





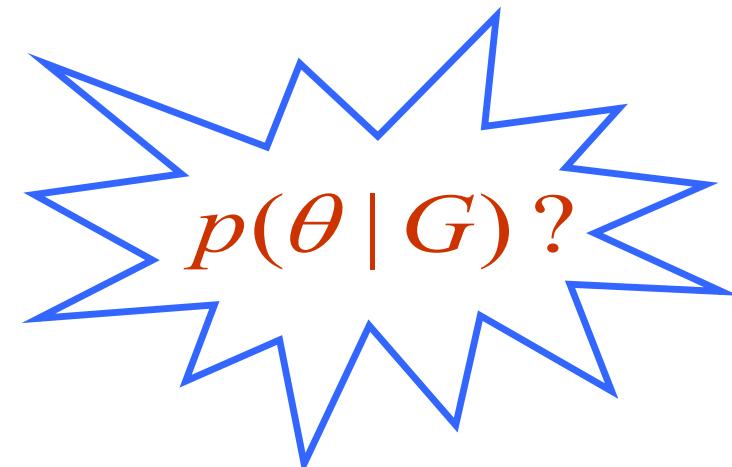
# How to define parameter prior?



Factorization:  $p(\mathbf{X} = \mathbf{x}) = \prod_{i=1}^M p(x_i | \mathbf{x}_{\pi_i})$

Local Distributions  
defined by, e.g., multinomial parameters:

$$p(x_i^k | \mathbf{x}_{\pi_i}^j) = \theta_{x_i^k | \mathbf{x}_{\pi_i}^j}$$



## Assumptions (Geiger & Heckerman 97,99):

- Complete Model Equivalence
- Global Parameter Independence
- Local Parameter Independence
- Likelihood and Prior Modularity





# Global & Local Parameter Independence

## ■ Global Parameter Independence

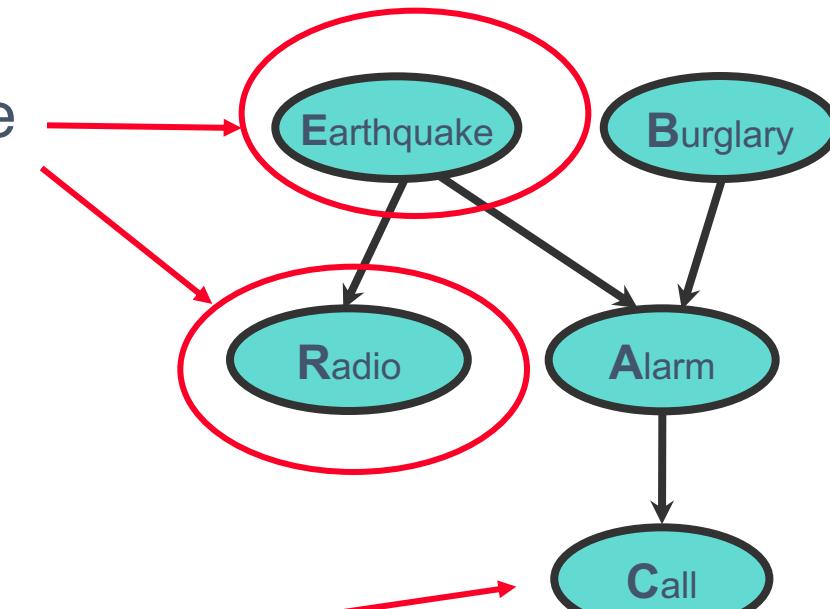
For every DAG model:

$$p(\theta_m | G) = \prod_{i=1}^M p(\theta_i | G)$$

## ■ Local Parameter Independence

For every node:

$$p(\theta_i | G) = \prod_{j=1}^{q_i} p(\theta_{x_i^k | \mathbf{x}_{\pi_i}^j} | G)$$



$$P(\theta_{Call|Alarm=YES})$$

independent of

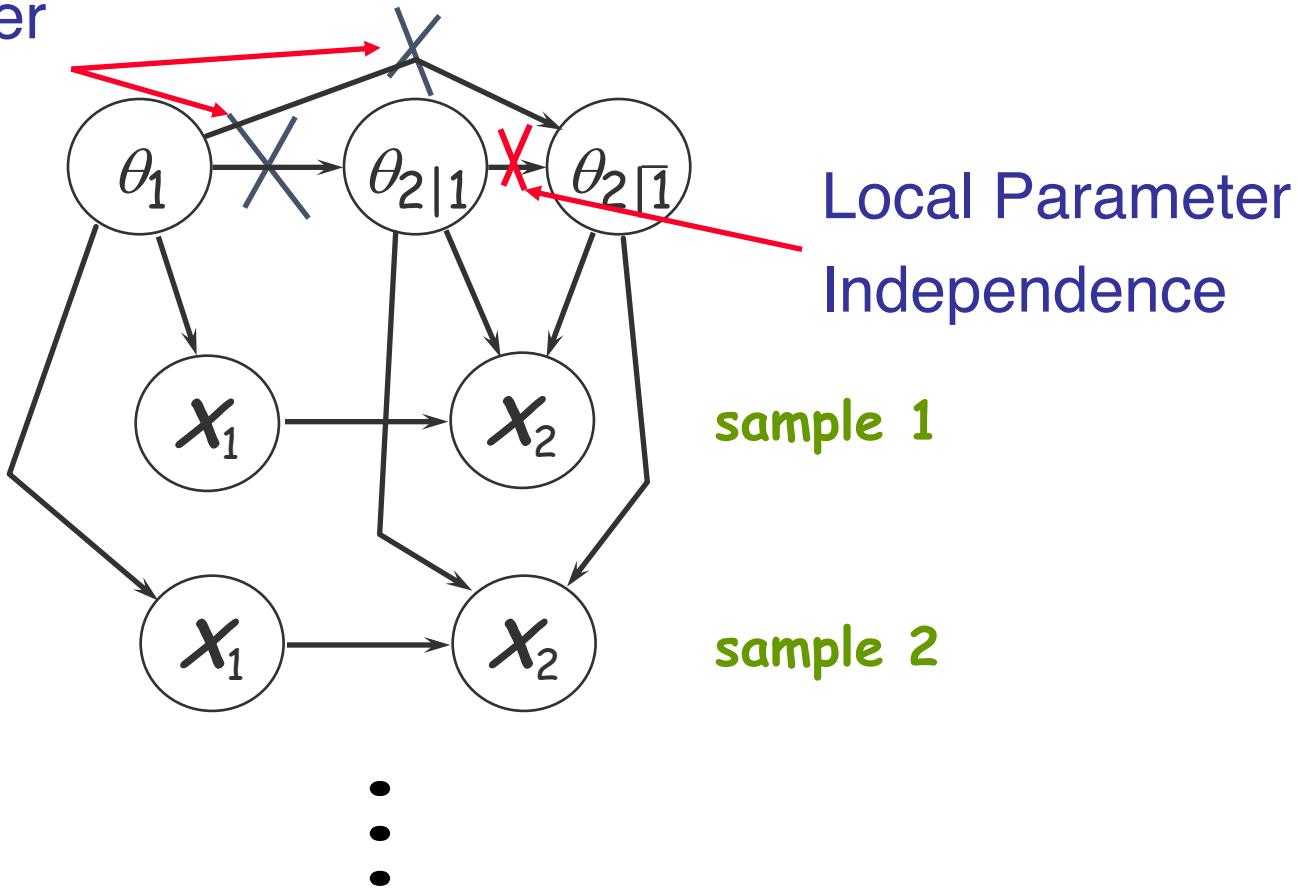
$$P(\theta_{Call|Alarm=NO})$$





# Parameter Independence, Graphical View

Global Parameter  
Independence



Local Parameter  
Independence

sample 1

sample 2

Provided all variables are observed in all cases, we can perform Bayesian update each parameter independently !!!





# Which PDFs Satisfy Our Assumptions? (Geiger & Heckerman 97,99)

## □ Discrete DAG Models:

Dirichlet prior:

$$x_i | \pi_{x_i}^j \sim \text{Multi}(\theta)$$

$$P(\theta) = \frac{\Gamma(\sum \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1} = C(\alpha) \prod_k \theta_k^{\alpha_k - 1}$$

## □ Gaussian DAG Models:

Normal prior:

$$x_i | \pi_{x_i}^j \sim \text{Normal}(\mu, \Sigma)$$

$$p(\mu | \nu, \Psi) = \frac{1}{(2\pi)^{n/2} |\Psi|^{1/2}} \exp\left\{-\frac{1}{2}(\mu - \nu)' \Psi^{-1} (\mu - \nu)\right\}$$

Normal-Wishart prior:

$$p(\mu | \nu, \alpha_\mu, \mathbf{W}) = \text{Normal}\left(\nu, (\alpha_\mu \mathbf{W})^{-1}\right),$$

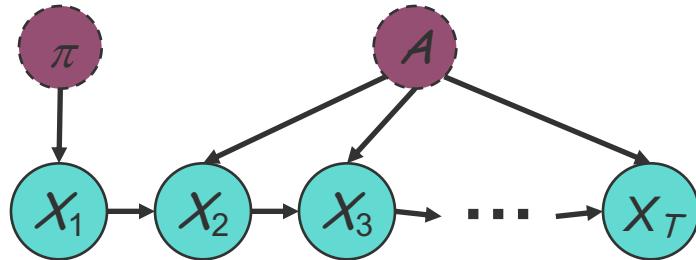
$$p(\mathbf{W} | \alpha_w, \mathbf{T}) = c(n, \alpha_w) |\mathbf{T}|^{\alpha_w/2} |\mathbf{W}|^{(\alpha_w - n - 1)/2} \exp\left\{\frac{1}{2} \text{tr}\{\mathbf{T}\mathbf{W}\}\right\},$$

$$\text{where } \mathbf{W} = \Sigma^{-1}.$$





# Parameter sharing



- Consider a time-invariant (stationary) 1<sup>st</sup>-order Markov model
  - Initial state probability vector:  $\pi_k \stackrel{\text{def}}{=} p(X_1^k = 1)$
  - State transition probability matrix:  $A_{ij} \stackrel{\text{def}}{=} p(X_t^j = 1 | X_{t-1}^i = 1)$
- The joint: 
$$p(X_{1:T} | \theta) = p(x_1 | \pi) \prod_{t=2}^T \prod_{i=1}^I p(X_t | X_{t-1})$$
- The log-likelihood: 
$$\ell(\theta; D) = \sum_n \log p(x_{n,1} | \pi) + \sum_n \sum_{t=2}^T \log p(x_{n,t} | x_{n,t-1}, A)$$
- Again, we optimize each parameter separately
  - $\pi$  is a multinomial frequency vector, and we've seen it before
  - What about  $A$ ?





# Learning a Markov chain transition matrix

- $A$  is a stochastic matrix:  $\sum_j A_{ij} = 1$
- Each row of  $A$  is multinomial distribution.
- So **MLE** of  $A_{ij}$  is the fraction of transitions from  $i$  to  $j$

$$A_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T x_{n,t-1}^i x_{n,t}^j}{\sum_n \sum_{t=2}^T x_{n,t-1}^i}$$

- Application:
  - if the states  $X_t$  represent words, this is called a *bigram language model*
- Sparse data problem:
  - If  $i \rightarrow j$  did not occur in data, we will have  $A_{ij} = 0$ , then any future sequence with word pair  $i \rightarrow j$  will have zero probability.
  - A standard hack: *backoff smoothing* or *deleted interpolation*

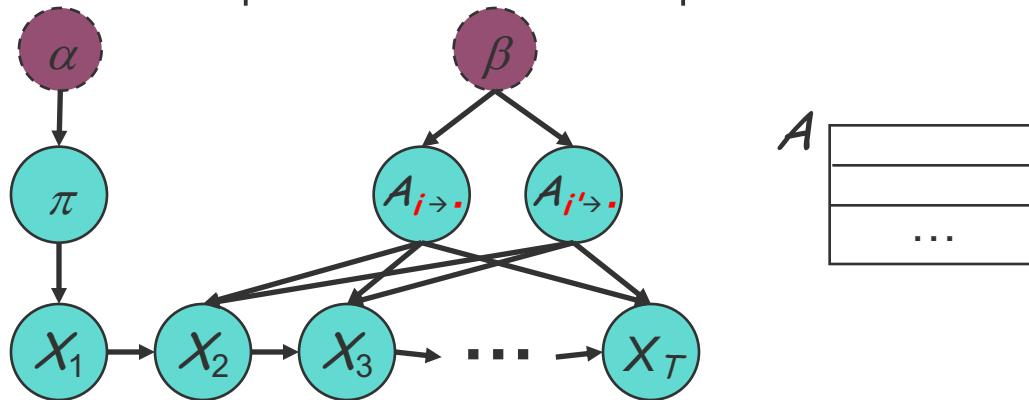
$$\tilde{A}_{i \rightarrow \bullet} = \lambda \eta_t + (1 - \lambda) A_{i \rightarrow \bullet}^{ML}$$





# Bayesian language model

- Global and local parameter independence



- The posterior of  $A_{i \rightarrow \cdot}$  and  $A_{i' \rightarrow \cdot}$  is factorized despite v-structure on  $X_t$ , because  $X_{t-1}$  acts like a **multiplexer**
- Assign a Dirichlet prior  $\beta_i$  to each row of the transition matrix:

$$A_{ij}^{Bayes} \stackrel{\text{def}}{=} p(j | i, D, \beta_i) = \frac{\#(i \rightarrow j) + \beta_{i,k}}{\#(i \rightarrow \bullet) + |\beta_i|} = \lambda_i \beta_{i,k}' + (1 - \lambda_i) A_{ij}^{ML}, \text{ where } \lambda_i = \frac{|\beta_i|}{|\beta_i| + \#(i \rightarrow \bullet)}$$

- We could consider more realistic priors, e.g., mixtures of Dirichlets to account for types of words (adjectives, verbs, etc.)





# Example: HMM: two scenarios

- Supervised learning: estimation when the “right answer” is known
  - Examples:  
**GIVEN**: a genomic region  $x = x_1 \dots x_{1,000,000}$  where we have good (experimental) annotations of the CpG islands  
**GIVEN**: the casino player allows us to observe him one evening, as he changes dice and produces 10,000 rolls
- Unsupervised learning: estimation when the “right answer” is unknown
  - Examples:  
**GIVEN**: the porcupine genome; we don’t know how frequent are the CpG islands there, neither do we know their composition  
**GIVEN**: 10,000 rolls of the casino player, but we don’t see when he changes dice
- **QUESTION**: Update the parameters  $\theta$  of the model to maximize  $P(x|\theta)$  --- Maximal likelihood (ML) estimation





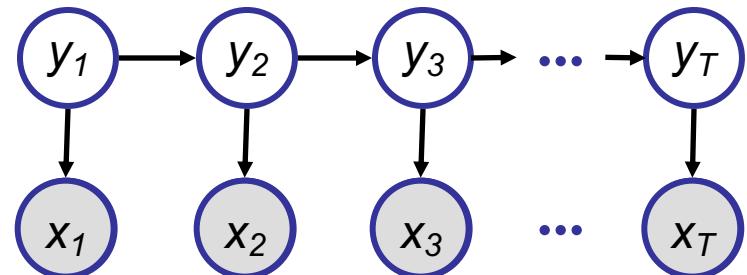
# Recall definition of HMM

- Transition probabilities between any two states

$$p(y_t^j = 1 | y_{t-1}^i = 1) = a_{i,j},$$

or

$$p(y_t | y_{t-1}^i = 1) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in I.$$



- Start probabilities

$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

- Emission probabilities associated with each state

$$p(x_t | y_t^i = 1) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in I.$$

or in general:

$$p(x_t | y_t^i = 1) \sim f(\cdot | \theta_i), \forall i \in I.$$





# Supervised ML estimation

- Given  $x = x_1 \dots x_N$  for which the true state path  $y = y_1 \dots y_N$  is known,
  - Define:

$A_{ij}$  = # times state transition  $i \rightarrow j$  occurs in  $y$

$B_{ik}$  = # times state  $i$  in  $y$  emits  $k$  in  $x$

- We can show that the **maximum likelihood** parameters  $\theta$  are:

$$a_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T y_{n,t-1}^i y_{n,t}^j}{\sum_n \sum_{t=2}^T y_{n,t-1}^i} = \frac{A_{ij}}{\sum_{j'} A_{ij'}}$$

$$b_{ik}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^T y_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T y_{n,t}^i} = \frac{B_{ik}}{\sum_{k'} B_{ik'}}$$

- What if  $x$  is continuous? We can treat  $\{(x_{n,t}, y_{n,t}) : t = 1:T, n = 1:N\}$  as  $N T$  observations of, e.g., a Gaussian, and apply learning rules for Gaussian ...





# Supervised ML estimation, ctd.

- Intuition:
  - When we know the underlying states, the best estimate of  $\theta$  is the average frequency of transitions & emissions that occur in the training data
- Drawback:
  - Given little data, there may be overfitting:
    - $P(x|\theta)$  is maximized, but  $\theta$  is unreasonable:**0 probabilities – VERY BAD**
- Example:
  - Given 10 casino rolls, we observe
    - $x = 2, 1, 5, 6, 1, 2, 3, 6, 2, 3$
    - $y = F, F, F, F, F, F, F, F, F, F$
  - Then:  $a_{FF} = 1; a_{FL} = 0$   
 $b_{F1} = b_{F3} = .2;$   
 $b_{F2} = .3; b_{F4} = 0; b_{F5} = b_{F6} = .1$





# Pseudocounts

- Solution for small training sets:
  - Add pseudocounts
$$A_{ij} = \text{\# times state transition } i \rightarrow j \text{ occurs in } \mathbf{y} + R_{ij}$$
$$B_{ik} = \text{\# times state } i \text{ in } \mathbf{y} \text{ emits } k \text{ in } \mathbf{x} + S_{ik}$$
  - $R_{ij}, S_{ij}$  are pseudocounts representing our prior belief
  - Total pseudocounts:  $R_i = \sum_j R_{ij}$ ,  $S_i = \sum_k S_{ik}$ ,
    - --- "strength" of prior belief,
    - --- total number of imaginary instances in the prior
- Larger total pseudocounts  $\Rightarrow$  strong prior belief
- Small total pseudocounts: just to avoid 0 probabilities --- smoothing
- This is equivalent to Bayesian est. under a uniform prior with "parameter strength" equals to the pseudocounts





# Summary: Learning GM

- ❑ For fully observed BN, the log-likelihood function decomposes into a sum of local terms, one per node; thus learning is also factored
  - ❑ Structural learning
    - ❑ Chow liu
    - ❑ Neighborhood selection
  - ❑ Learning single-node GM – density estimation: exponential family dist.
    - ❑ Typical discrete distribution
    - ❑ Typical continuous distribution
    - ❑ Conjugate priors
  - ❑ Learning two-node BN: GLIM
    - ❑ Conditional Density Est.
    - ❑ Classification
  - ❑ Learning BN with more nodes
    - ❑ Local operations



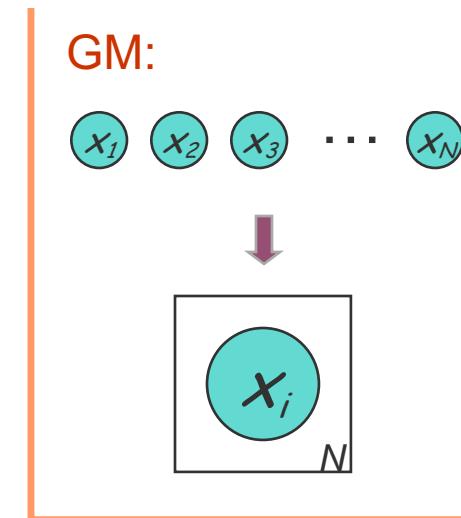
# Supplementary





# Review of density estimation

- Can be viewed as single-node graphical models
- Instances of exponential family dist.
- Building blocks of general GM
- MLE and Bayesian estimate





# Discrete Distributions

- Bernoulli distribution:  $\text{Ber}(p)$

$$P(x) = \begin{cases} 1-p & \text{for } x=0 \\ p & \text{for } x=1 \end{cases} \Rightarrow P(x) = p^x (1-p)^{1-x}$$



- Multinomial distribution:  $\text{Mult}(1, \theta)$

- Multinomial (indicator) variable:

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix}, \quad \text{where} \quad \begin{aligned} X_j &= [0,1], \quad \text{and} \quad \sum_{j \in [1, \dots, 6]} X_j = 1 \\ X_j &= 1 \text{ w.p. } \theta_j, \quad \sum_{j \in [1, \dots, 6]} \theta_j = 1 \end{aligned}$$



$$\begin{aligned} p(x(j)) &= P(\{X_j = 1, \text{ where } j \text{ index the dice - face}\}) \\ &= \theta_j = \theta_A^{x_A} \times \theta_C^{x_C} \times \theta_G^{x_G} \times \theta_T^{x_T} = \prod_k \theta_k^{x_k} = \theta^x \end{aligned}$$





# Discrete Distributions

- Multinomial distribution:  $\text{Mult}(n, \theta)$
- Count variable:

$$n = \begin{bmatrix} n_1 \\ \vdots \\ n_K \end{bmatrix}, \quad \text{where } \sum_j n_j = N$$

$$p(n) = \frac{N!}{n_1! n_2! \cdots n_K!} \theta_1^{n_1} \theta_2^{n_2} \cdots \theta_K^{n_K} = \frac{N!}{n_1! n_2! \cdots n_K!} \theta^n$$

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.





# Example: multinomial model

- Data:

- We observed  $N$  *iid* die rolls ( $K$ -sided):  $D=\{5, 1, K, \dots, 3\}$

- Representation:

Unit basis vectors:

$$x_n = \begin{pmatrix} x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,K} \end{pmatrix}, \text{ where } x_{n,k} = \{0,1\}, \text{ and } \sum_{k=1}^K x_{n,k} = 1$$

- Model:

$$X_{n,k} = 1 \text{ w.p. } \theta_k, \text{ and } \sum_{k \in \{1, \dots, K\}} \theta_k = 1$$

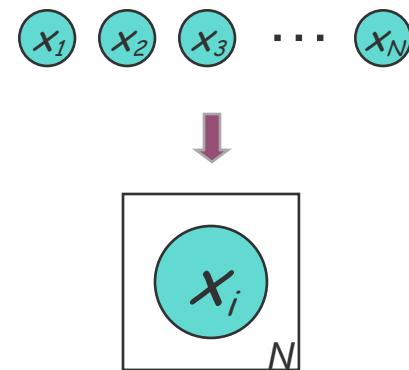
- How to write the likelihood of a single observation  $x_n$ ?

$$\begin{aligned} P(x_n) &= P(\{x_{n,k} = 1, \text{ where } k \text{ index the die - side of the } n \text{th roll}\}) \\ &= \theta_k = \theta_1^{x_{n,1}} \times \theta_2^{x_{n,2}} \times \dots \times \theta_K^{x_{n,K}} = \prod_{k=1}^K \theta_k^{x_{n,k}} \end{aligned}$$

- The likelihood of dataset  $D=\{x_1, \dots, x_N\}$ :

$$P(x_1, x_2, \dots, x_N | \theta) = \prod_{n=1}^N P(x_n | \theta) = \prod_{n=1}^N \left( \prod_k \theta_k^{x_{n,k}} \right) = \prod_k \theta_k^{\sum_{n=1}^N x_{n,k}} = \prod_k \theta_k^{n_k}$$

GM:





# MLE: constrained optimization with Lagrange multipliers

- Objective function:

$$\ell(\theta; \mathcal{D}) = \log P(\mathcal{D} | \theta) = \log \prod_k \theta_k^{n_k} = \sum_k n_k \log \theta_k$$

- We need to maximize this subject to the constrain  $\sum_{k=1}^K \theta_k = 1$

- Constrained cost function with a Lagrange multiplier

$$\ell^- = \sum_k n_k \log \theta_k + \lambda \left( 1 - \sum_{k=1}^K \theta_k \right)$$

- Take derivatives wrt  $\theta_k$

$$\frac{\partial \ell}{\partial \theta_k} = \frac{n_k}{\theta_k} - \lambda = 0$$

$$n_k = \lambda \theta_k \Rightarrow \sum_k n_k = N = \lambda \sum_k \theta_k = \lambda$$

$$\rightarrow \hat{\theta}_{k, MLE} = \frac{n_k}{N} \quad \text{or} \quad \hat{\theta}_{k, MLE} = \frac{1}{N} \sum_n x_{n,k}$$

Frequency as  
sample mean

- Sufficient statistics

- The counts,  $\bar{n} = (n_1, \dots, n_K)$ ,  $n_k = \sum_n x_{n,k}$ , are **sufficient statistics** of data  $\mathcal{D}$

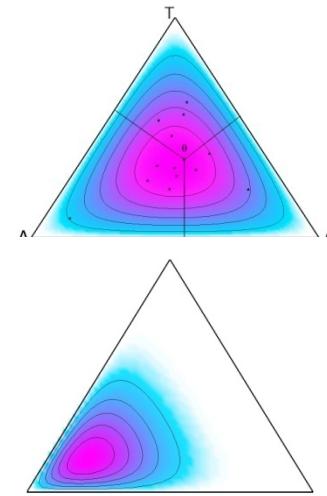




# Bayesian estimation:

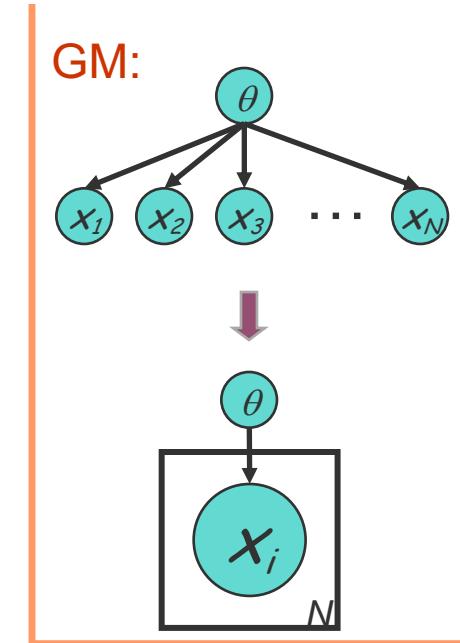
- Dirichlet distribution:

$$P(\theta) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k-1} = C(\alpha) \prod_k \theta_k^{\alpha_k-1}$$



- Posterior distribution of  $\theta$ :

$$P(\theta | x_1, \dots, x_N) = \frac{p(x_1, \dots, x_N | \theta)p(\theta)}{p(x_1, \dots, x_N)} \propto \prod_k \theta_k^{n_k} \prod_k \theta_k^{\alpha_k-1} = \prod_k \theta_k^{\alpha_k + n_k - 1}$$



- Notice the isomorphism of the posterior to the prior,- such a prior is called a **conjugate prior**

- Posterior mean estimation:

$$\theta_k = \int \theta_k p(\theta | D) d\theta = C \int \theta_k \prod_k \theta_k^{\alpha_k + n_k - 1} d\theta = \frac{n_k + \alpha_k}{N + |\alpha|}$$

Dirichlet parameters  
can be understood  
as pseudo-counts





# More on Dirichlet Prior:

- Where is the normalize constant  $C(\alpha)$  come from?

$$\frac{1}{C(\alpha)} = \int \cdots \int \theta_1^{\alpha_1-1} \cdots \theta_K^{\alpha_K-1} d\theta_1 \cdots d\theta_K = \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$$

- Integration by parts
- $\Gamma(\alpha)$  is the gamma function:  $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$
- For integers,  $\Gamma(n+1) = n!$

- Marginal likelihood:

$$p(\{x_1, \dots, x_N\} | \vec{\alpha}) = p(\bar{n} | \vec{\alpha}) = \int p(\bar{n} | \vec{\theta}) p(\vec{\theta} | \vec{\alpha}) d\vec{\theta} = \frac{C(\vec{\alpha})}{C(\bar{n} + \vec{\alpha})}$$

- Posterior in closed-form:

$$P(\vec{\theta} | \{x_1, \dots, x_N\}, \vec{\alpha}) = \frac{p(\bar{n} | \theta) p(\theta | \vec{\alpha})}{p(\bar{n} | \vec{\alpha})} = C(\bar{n} + \vec{\alpha}) \prod_k \theta_k^{\alpha_k + n_k - 1} = \text{Dir}(\bar{n} + \vec{\alpha})$$

- Posterior predictive rate:

$$p(x_{N+1} = i | \{x_1, \dots, x_N\}, \vec{\alpha}) = \int C(\bar{n} + \vec{\alpha}) \prod_k \theta_k^{\alpha_k + n_k - 1} \times \theta_i^{\alpha_i + n_i} d\vec{\theta} = \frac{C(\bar{n} + \vec{\alpha})}{C(\bar{n} + \vec{\alpha} + \vec{x}_N)} = \frac{n_i + \alpha_i}{|\bar{n}| + |\vec{\alpha}|}$$





# Sequential Bayesian updating

- ❑ Start with Dirichlet prior  $P(\bar{\theta} | \bar{\alpha}) = \text{Dir}(\bar{\theta} : \bar{\alpha})$
- ❑ Observe  $N'$  samples with sufficient statistics  $\bar{n}'$ . Posterior becomes:

$$P(\bar{\theta} | \bar{\alpha}, \bar{n}') = \text{Dir}(\bar{\theta} : \bar{\alpha} + \bar{n}')$$

- ❑ Observe another  $N''$  samples with sufficient statistics  $\bar{n}''$ . Posterior becomes:

$$P(\bar{\theta} | \bar{\alpha}, \bar{n}', \bar{n}'') = \text{Dir}(\bar{\theta} : \bar{\alpha} + \bar{n}' + \bar{n}'')$$

- ❑ So sequentially absorbing data in any order is equivalent to batch update.

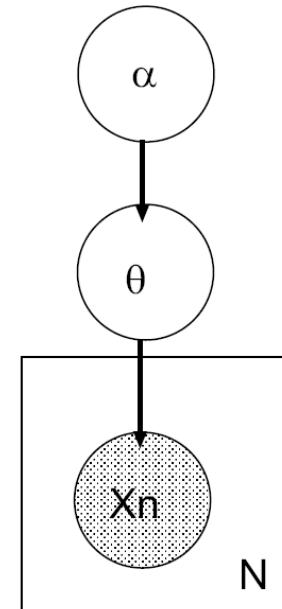




# Hierarchical Bayesian Models

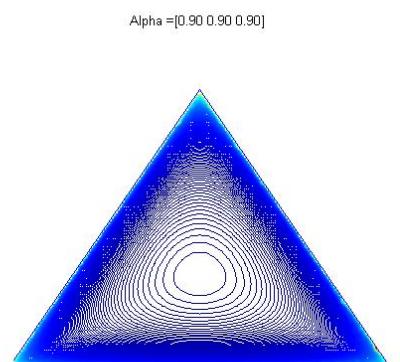
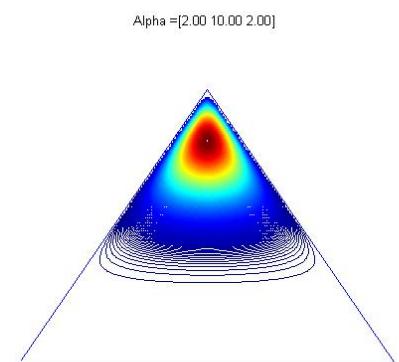
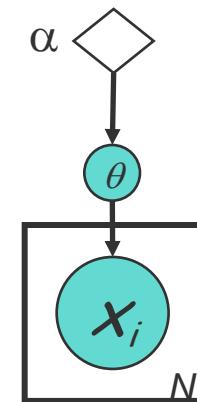
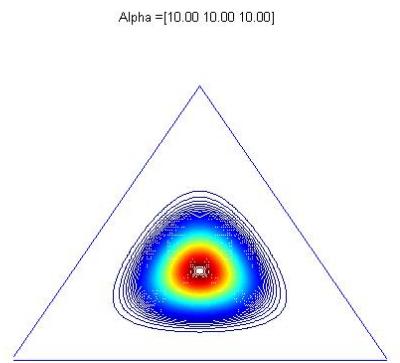
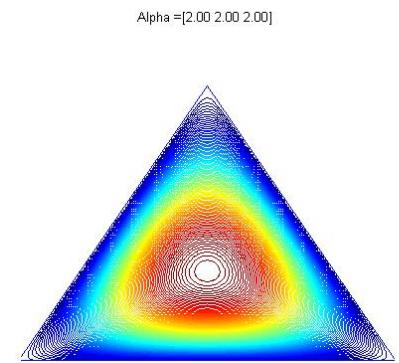
- $\theta$  are the parameters for the likelihood  $p(x|\theta)$
- $\alpha$  are the parameters for the prior  $p(\theta|\alpha)$ .
- We can have hyper-hyper-parameters, etc.
- We stop when the choice of hyper-parameters makes no difference to the marginal likelihood; typically make hyper-parameters constants.
- Where do we get the prior?
  - Intelligent guesses
  - Empirical Bayes (Type-II maximum likelihood)  
→ computing point estimates of  $\alpha$ :

$$\hat{\alpha}_{MLE} = \arg \max_{\bar{\alpha}} = p(\bar{n} | \bar{\alpha})$$





# Limitation of Dirichlet Prior:





# The Logistic Normal Prior

$$\theta \sim LN_K(\mu, \Sigma)$$

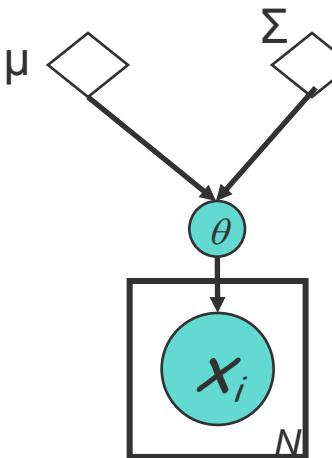
$$\gamma \sim N_{K-1}(\mu, \Sigma)$$

$$\theta_i = \exp\left\{\gamma_i - \log\left(1 + \sum_{i=1}^{K-1} e^{\gamma_i}\right)\right\}$$

$$C(\gamma) = \log\left(1 + \sum_{i=1}^{K-1} e^{\gamma_i}\right)$$

Problem

- Log Partition Function
- Normalization Constant

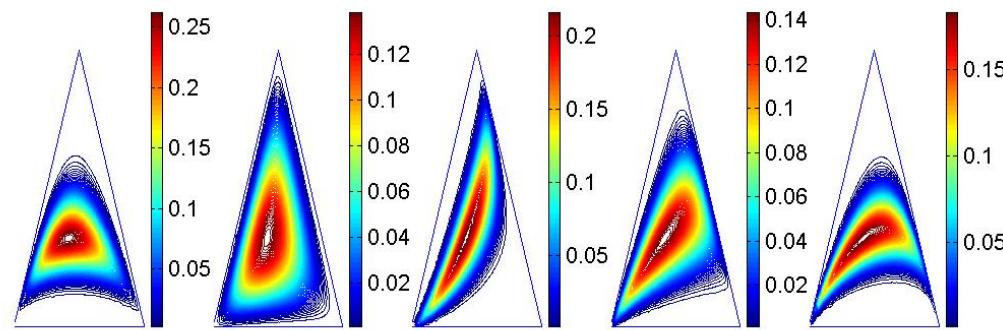


- **Pro: co-variance structure**
- **Con: non-conjugate (we will discuss how to solve this later)**

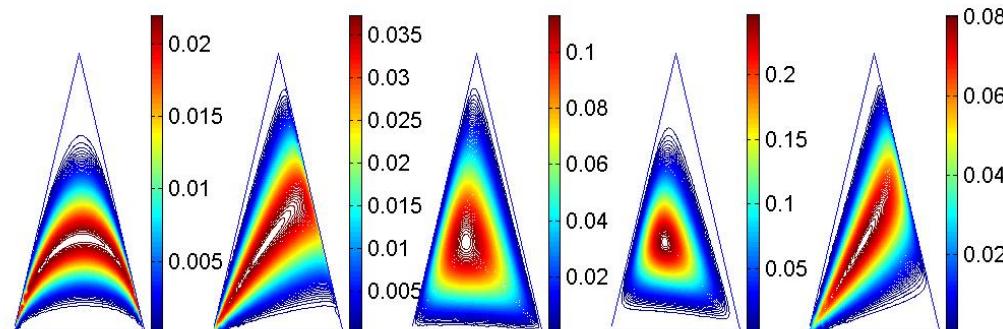




# Logistic Normal Densities



Logistic  
Normal

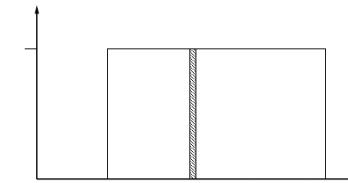




# Continuous Distributions

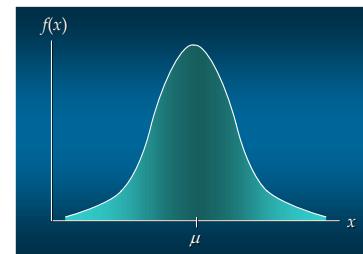
- Uniform Probability Density Function

$$p(x) = \begin{cases} 1/(b-a) & \text{for } a \leq x \leq b \\ 0 & \text{elsewhere} \end{cases}$$



- Normal (Gaussian) Probability Density Function

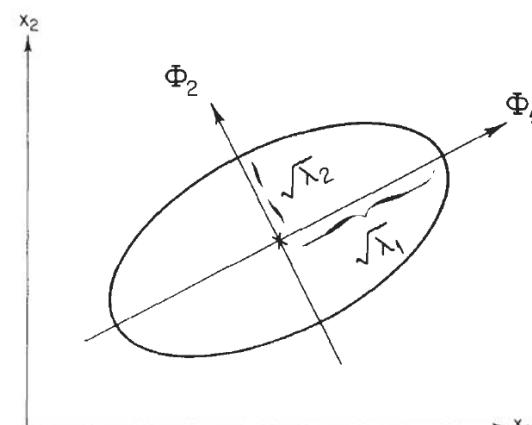
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$



- The distribution is symmetric, and is often illustrated as a bell-shaped curve.
- Two parameters,  $\mu$  (mean) and  $\sigma$  (standard deviation), determine the location and shape of the distribution.
- The highest point on the normal curve is at the mean, which is also the median and mode.
- The mean can be any numerical value: negative, zero, or positive.

- Multivariate Gaussian

$$p(X; \vec{\mu}, \Sigma) = \frac{1}{(\sqrt{2\pi})^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (X - \vec{\mu})^T \Sigma^{-1} (X - \vec{\mu})\right\}$$





# MLE for a multivariate-Gaussian

- It can be shown that the MLE for  $\mu$  and  $\Sigma$  is

$$\boldsymbol{\mu}_{MLE} = \frac{1}{N} \sum_n (\mathbf{x}_n)$$

$$\boldsymbol{\Sigma}_{MLE} = \frac{1}{N} \sum_n (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T = \frac{1}{N} S$$

where the scatter matrix is

$$S = \sum_n (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T = \left( \sum_n \mathbf{x}_n \mathbf{x}_n^T \right) - N \boldsymbol{\mu}_{ML} \boldsymbol{\mu}_{ML}^T$$

$$\mathbf{x}_n = \begin{pmatrix} x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,K} \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} \cdots \mathbf{x}_1^T \cdots \\ \cdots \mathbf{x}_2^T \cdots \\ \vdots \\ \cdots \mathbf{x}_N^T \cdots \end{pmatrix}$$

- The sufficient statistics are  $\sum_n \mathbf{x}_n$  and  $\sum_n \mathbf{x}_n \mathbf{x}_n^T$ .
- Note that  $\mathbf{X}^T \mathbf{X} = \sum_n \mathbf{x}_n \mathbf{x}_n^T$  may not be full rank (eg. if  $N < D$ ), in which case  $\boldsymbol{\Sigma}_{ML}$  is not invertible





# Bayesian parameter estimation for a Gaussian

- There are various reasons to pursue a Bayesian approach
  - We would like to update our estimates sequentially over time.
  - We may have prior knowledge about the expected magnitude of the parameters.
  - The MLE for  $\Sigma$  may not be full rank if we don't have enough data.
- We will restrict our attention to conjugate priors.
- We will consider various cases, in order of increasing complexity:
  - Known  $\sigma$ , unknown  $\mu$
  - Known  $\mu$ , unknown  $\sigma$
  - Unknown  $\mu$  and  $\sigma$





# Bayesian estimation: unknown $\mu$ , known $\sigma$

- Normal Prior:

$$P(\mu) = (2\pi\tau^2)^{-1/2} \exp\left\{-(\mu - \mu_0)^2 / 2\tau^2\right\}$$

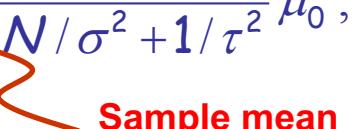
- Joint probability:

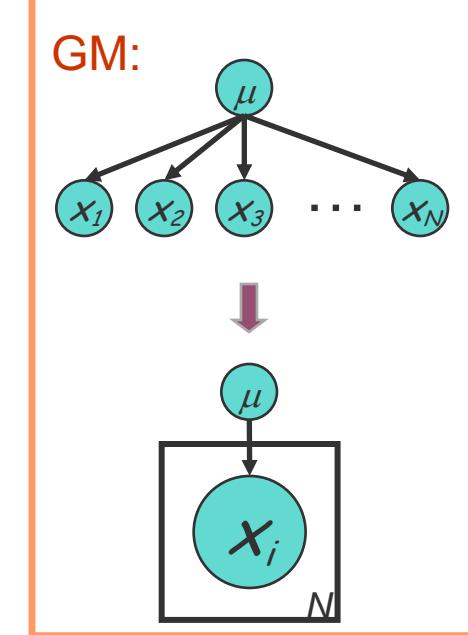
$$\begin{aligned} P(x, \mu) &= (2\pi\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right\} \\ &\quad \times (2\pi\tau^2)^{-1/2} \exp\left\{-(\mu - \mu_0)^2 / 2\tau^2\right\} \end{aligned}$$

- Posterior:

$$P(\mu | x) = (2\pi\tilde{\sigma}^2)^{-1/2} \exp\left\{-(\mu - \tilde{\mu})^2 / 2\tilde{\sigma}^2\right\}$$

where  $\tilde{\mu} = \frac{N/\sigma^2}{N/\sigma^2 + 1/\tau^2} \bar{x} + \frac{1/\tau^2}{N/\sigma^2 + 1/\tau^2} \mu_0$ , and  $\tilde{\sigma}^2 = \left( \frac{N}{\sigma^2} + \frac{1}{\tau^2} \right)^{-1}$

 Sample mean

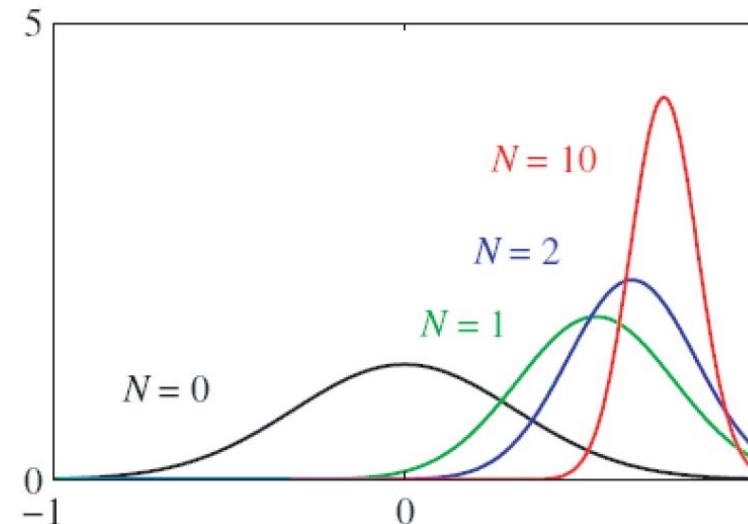




# Bayesian estimation: unknown $\mu$ , known $\sigma$

$$\mu_N = \frac{N/\sigma^2}{N/\sigma^2 + 1/\sigma_0^2} \bar{x} + \frac{1/\sigma_0^2}{N/\sigma^2 + 1/\sigma_0^2} \mu_0, \quad \tilde{\sigma}^2 = \left( \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right)^{-1}$$

- The posterior mean is a convex combination of the prior and the MLE, with weights proportional to the relative noise levels.
- The precision of the posterior  $1/\sigma_N^2$  is the precision of the prior  $1/\sigma_0^2$  plus one contribution of data precision  $1/\sigma^2$  for each observed data point.
- Sequentially updating the mean
  - $\mu^* = 0.8$  (unknown),  $(\sigma^2)^* = 0.1$  (known)
  - Effect of single data point  $\mu_1 = \mu_0 + (x - \mu_0) \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2} = x - (x - \mu_0) \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2}$
  - Uninformative (vague/ flat) prior,  $\sigma_0^2 \rightarrow \infty$





# Other scenarios

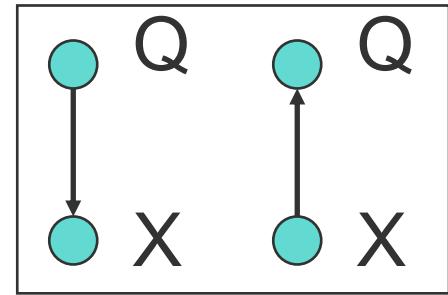
- ❑ Known  $\mu$ , unknown  $\lambda = 1/\sigma_2^2$ 
  - ❑ The conjugate prior for  $\lambda$  is a **Gamma** with shape  $a_0$  and rate (inverse scale)  $b_0$ 
$$p(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$
  - ❑ The conjugate prior for  $\sigma^2$  is  $IG(\sigma^2|a, b) = \frac{1}{\Gamma(a)} b^a (\sigma^2)^{-a-1} \exp(-b/(\sigma^2))$
- ❑ Unknown  $\mu$  and unknown  $\sigma_2^2$ 
  - ❑ The conjugate prior is **Normal-Inverse-Gamma**
$$\begin{aligned} P(\mu, \sigma^2) &= P(\mu|\sigma^2)P(\sigma^2) \\ &= \mathcal{N}(\mu|m, \sigma^2V) \cdot IG(\sigma^2|a, b) \end{aligned}$$
  - ❑ Semi conjugate prior
  - ❑ Multivariate case:
    - ❑ The conjugate prior is **Normal-Inverse-Wishart**
$$\begin{aligned} P(\mu, \Sigma) &= P(\mu|\Sigma)P(\Sigma) \\ &= \mathcal{N}(\mu|\mu_0, \frac{1}{\kappa_0}\Sigma) \cdot \mathcal{IW}(\Sigma|\Lambda_0^{-1}, \nu_0) \end{aligned}$$





# Two node fully observed BNs

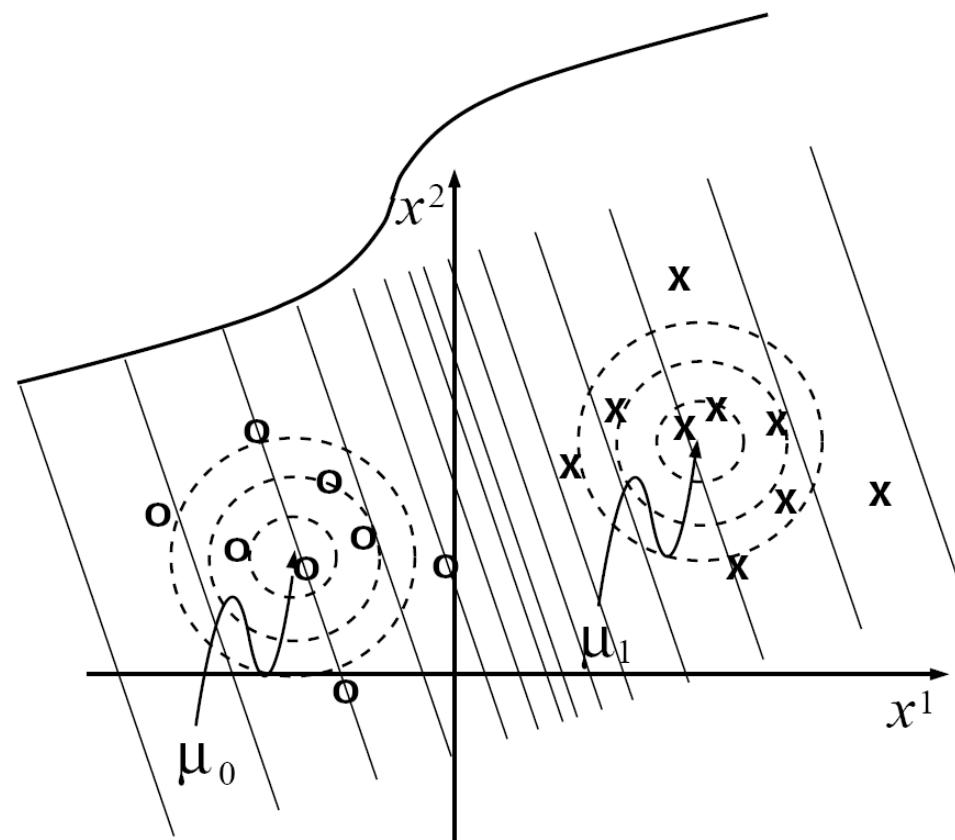
- Conditional mixtures
- Linear/Logistic Regression
- 
- Classification
  - Generative and discriminative approaches





# Classification:

- Goal: Wish to learn  $f: X \rightarrow Y$
- Generative:
  - Modeling the joint distribution of all data
- Discriminative:
  - Modeling only points at the boundary





# Conditional Gaussian

- The data:

$$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_N, y_N)\}$$

- Both nodes are observed:
  - $\mathbf{y}$  is a class indicator vector

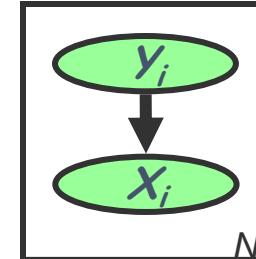
$$p(y_n) = \text{multi}(y_n : \pi) = \prod_k \pi_k^{y_{n,k}}$$

- $X$  is a conditional Gaussian variable with a class-specific mean

$$p(x_n | y_{n,k} = 1, \mu, \sigma) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2} (x_n - \mu_k)^2\right\}$$

$$p(x | y, \mu, \sigma) = \prod_n \left( \prod_k N(x_n : \mu_k, \sigma) \right)^{y_{n,k}}$$

GM:





# MLE of conditional Gaussian

- Data log-likelihood

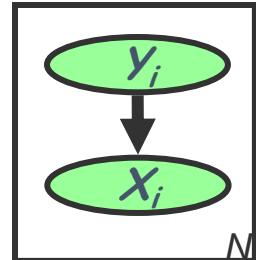
$$\ell(\boldsymbol{\theta}; D) = \log \prod_n p(x_n, y_n) = \log \prod_n p(y_n | \pi) p(x_n | y_n, \mu, \sigma)$$

- MLE

$$\hat{\pi}_{k,MLE} = \arg \max_{\pi} \ell(\boldsymbol{\theta}; D), \quad \hat{\pi}_{k,MLE} = \frac{\sum_n y_{n,k}}{N} = \frac{n_k}{N}$$

$$\hat{\mu}_{k,MLE} = \arg \max \ell(\boldsymbol{\theta}; D), \quad \hat{\mu}_{k,MLE} = \frac{\sum_n y_{n,k} x_n}{\sum_n y_{n,k}} = \frac{\sum_n y_{n,k} x_n}{n_k}$$

GM:



)

the fraction of  
samples of class  $m$

the average of  
samples of class  $m$





# Byesian estimation of conditional Gaussian

- Prior:

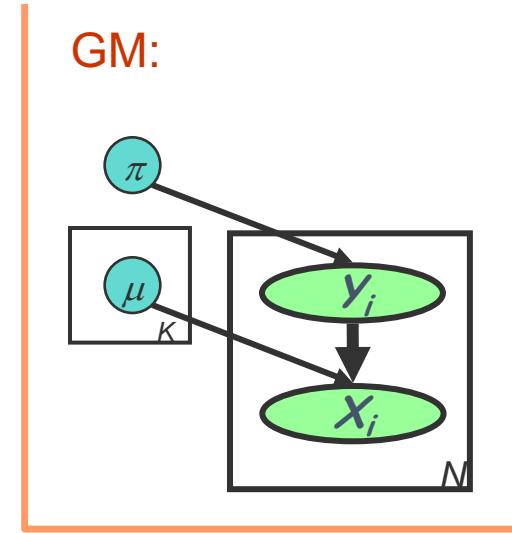
$$P(\bar{\pi} | \bar{\alpha}) = \text{Dir}(\bar{\pi} : \bar{\alpha})$$

$$P(\mu_k | \nu) = \text{Normal}(\mu_k : \nu, \tau)$$

- Posterior mean (Bayesian est.)

$$\pi_{k, Bayes} = \frac{N}{N+|\alpha|} \hat{\pi}_{k, ML} + \frac{|\alpha|}{N+|\alpha|} \frac{\alpha_k}{|\alpha|} = \frac{n_k + \alpha_k}{N+|\alpha|}$$

$$\mu_{k, Bayes} = \frac{n_k / \sigma^2}{n_k / \sigma^2 + 1 / \tau^2} \hat{\mu}_{k, ML} + \frac{1 / \tau^2}{n_k / \sigma^2 + 1 / \tau^2} \nu, \quad \text{and} \quad \sigma_{Bayes}^2 = \left( \frac{N}{\sigma^2} + \frac{1}{\tau^2} \right)^{-1}$$



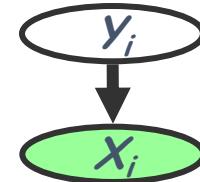


# Classification

- ❑ Gaussian Discriminative Analysis:

- ❑ The joint probability of a datum and its label is:

$$\begin{aligned} p(x_n, y_{n,k} = 1 | \mu, \sigma) &= p(y_{n,k} = 1) \times p(x_n | y_{n,k} = 1, \mu, \sigma) \\ &= \pi_k \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x_n - \mu_k)^2\right\} \end{aligned}$$



- ❑ Given a datum  $x_n$ , we predict its label using the conditional probability of the label given the datum:

$$p(y_{n,k} = 1 | x_n, \mu, \sigma) = \frac{\pi_k \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x_n - \mu_k)^2\right\}}{\sum_{k'} \pi_{k'} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x_n - \mu_{k'})^2\right\}}$$

- ❑ This is basic inference

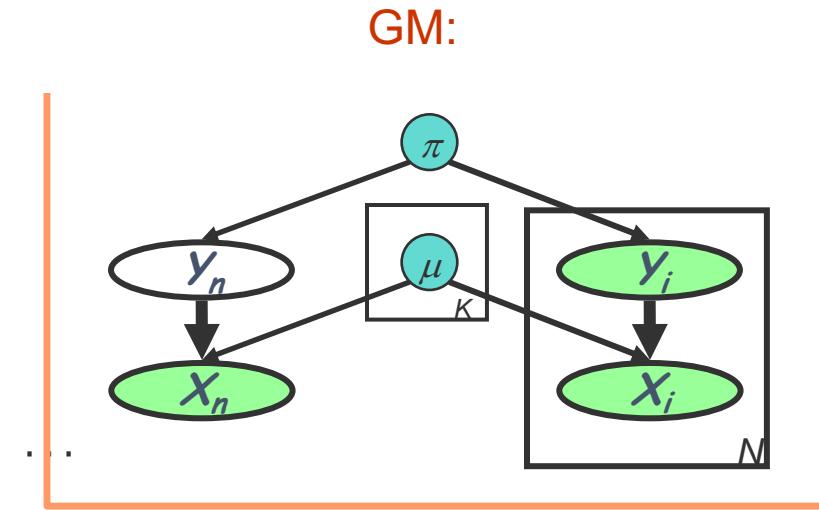
- ❑ introduce evidence, and then normalize





# Transductive classification

- ❑ Given  $X_n$ , what is its corresponding  $Y_n$  when we know the answer for a set of training data?
- ❑ Frequentist prediction:
  - ❑ we fit  $\pi$ ,  $\mu$  and  $\sigma$  from data first, and then ...



$$p(y_{n,k} = 1 | x_n, \mu, \sigma, \pi) = \frac{p(y_{n,k} = 1, x_n | \mu, \sigma, \pi)}{p(x_n | \mu, \sigma, \pi)} = \frac{\pi_k N(x_n, | \mu_k, \sigma)}{\sum_{k'} \pi_{k'} N(x_n, | \mu_{k'}, \sigma)}$$

- ❑ Bayesian:
  - ❑ we compute the posterior dist. of the parameters first ...





# Linear Regression

- The data:

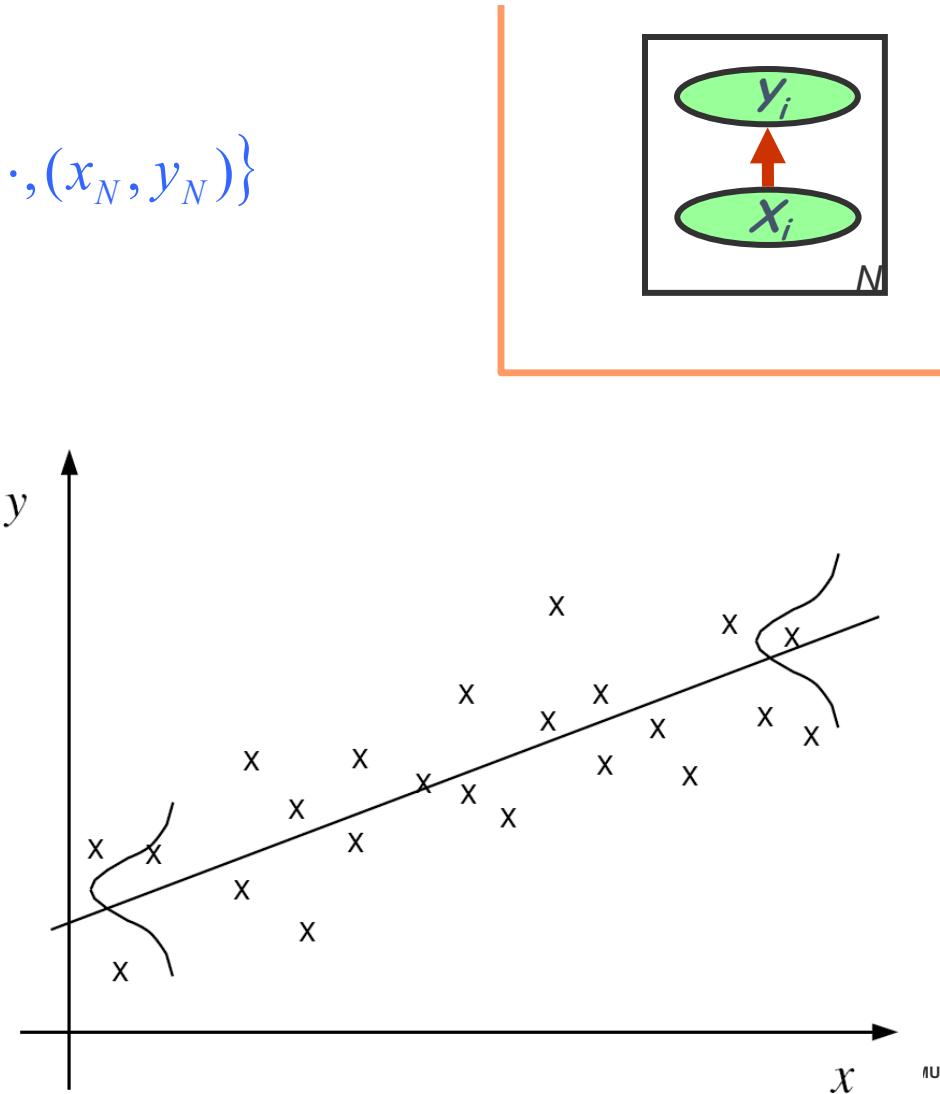
$$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_N, y_N)\}$$

- Both nodes are observed:

- $X$  is an input vector
- $Y$  is a response vector

(we first consider  $y$  as a generic continuous response vector, then we consider the special case of classification where  $y$  is a discrete indicator)

- A regression scheme can be used to model  $p(y|x)$  directly, rather than  $p(x,y)$





# A discriminative probabilistic model

- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

where  $\varepsilon$  is an error term of unmodeled effects or random noise

- Now assume that  $\varepsilon$  follows a Gaussian  $N(0, \sigma)$ , then we have:

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

- By independence assumption:

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$





# Linear regression

- ❑ Hence the log-likelihood is:

$$l(\theta) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

- ❑ Do you recognize the last term?

Yes it is:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

- ❑ It is same as the MSE!





# A recap:

- LMS update rule

$$\theta^{t+1} = \theta^t + \alpha(y_n - \mathbf{x}_n^T \theta^t) \mathbf{x}_n$$

- Pros: on-line, low per-step cost
- Cons: coordinate, maybe slow-converging

- Steepest descent

$$\theta^{t+1} = \theta^t + \alpha \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta^t) \mathbf{x}_i$$

- Pros: fast-converging, easy to implement
- Cons: a batch,

- Normal equations

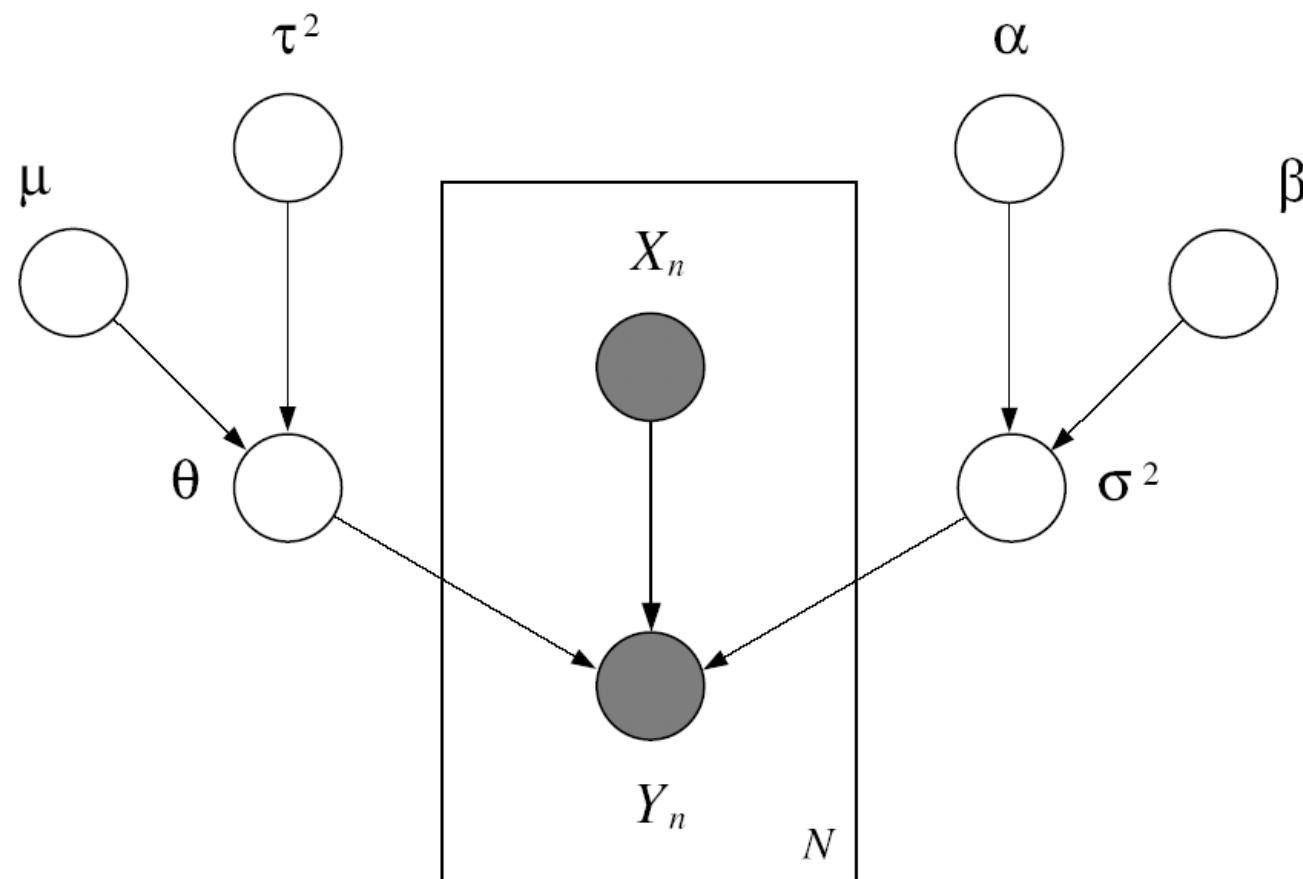
$$\theta^* = (X^T X)^{-1} X^T \bar{y}$$

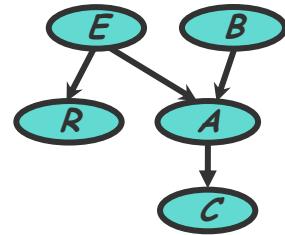
- Pros: a single-shot algorithm! Easiest to implement.
- Cons: need to compute pseudo-inverse ( $X^T X$ ) $^{-1}$ , expensive, numerical issues (e.g., matrix is singular ..)





# Bayesian linear regression





# ML Structural Learning for completely observed GMs



$(x_1^{(1)}, \dots, x_n^{(1)})$   
 $(x_1^{(2)}, \dots, x_n^{(2)})$

...

$(x_1^{(M)}, \dots, x_n^{(M)})$





# Two “Optimal” approaches

- ❑ “Optimal” here means the employed algorithms guarantee to return a structure that maximizes the objectives (e.g., LogLik)
  - ❑ Many heuristics used to be popular, but they provide no guarantee on attaining optimality, interpretability, or even do not have an explicit objective
  - ❑ E.g.: structured EM, Module network, greedy structural search, deep learning via auto-encoders, etc.
- ❑ We will learn two classes of algorithms for guaranteed structure learning, which are likely to be the only known methods enjoying such guarantee, but they only apply to certain families of graphs:
  - ❑ Trees: The Chow-Liu algorithm (this lecture)
  - ❑ Pairwise MRFs: covariance selection, neighborhood-selection (later)





# Structural Search

- How many graphs over  $n$  nodes?  $O(2^{n^2})$
- How many trees over  $n$  nodes?  $O(n!)$
- But it turns out that we can find exact solution of an optimal tree (under MLE)!
  - Trick: MLE score decomposable to edge-related elements
  - Trick: in a tree each node has only one parent!
  - Chow-liu algorithm





# Information Theoretic Interpretation of ML

$$\begin{aligned}\ell(\theta_G, G; D) &= \log p(D | \theta_G, G) \\ &= \log \prod_n \left( \prod_i p(x_{n,i} | \mathbf{x}_{n, \pi_i(G)}, \theta_{i|\pi_i(G)}) \right) \\ &= \sum_i \left( \sum_n \log p(x_{n,i} | \mathbf{x}_{n, \pi_i(G)}, \theta_{i|\pi_i(G)}) \right) \\ &= M \sum_i \left( \sum_{x_i, \mathbf{x}_{\pi_i(G)}} \frac{\text{count}(x_i, \mathbf{x}_{\pi_i(G)})}{M} \log p(x_i | \mathbf{x}_{\pi_i(G)}, \theta_{i|\pi_i(G)}) \right) \\ &= M \sum_i \left( \sum_{x_i, \mathbf{x}_{\pi_i(G)}} \hat{p}(x_i, \mathbf{x}_{\pi_i(G)}) \log p(x_i | \mathbf{x}_{\pi_i(G)}, \theta_{i|\pi_i(G)}) \right)\end{aligned}$$

**From sum over data points to sum over count of variable states**





# Information Theoretic Interpretation of ML (con'd)

$$\begin{aligned}\ell(\theta_G, G; D) &= \log \hat{p}(D | \theta_G, G) \\ &= M \sum_i \left( \sum_{x_i, \mathbf{x}_{\pi_i(G)}} \hat{p}(x_i, \mathbf{x}_{\pi_i(G)}) \log \hat{p}(x_i | \mathbf{x}_{\pi_i(G)}, \theta_{i|\pi_i(G)}) \right) \\ &= M \sum_i \left( \sum_{x_i, \mathbf{x}_{\pi_i(G)}} \hat{p}(x_i, \mathbf{x}_{\pi_i(G)}) \log \frac{\hat{p}(x_i, \mathbf{x}_{\pi_i(G)}, \theta_{i|\pi_i(G)})}{\hat{p}(\mathbf{x}_{\pi_i(G)})} \frac{\hat{p}(x_i)}{\hat{p}(x_i)} \right) \\ &= M \sum_i \left( \sum_{x_i, \mathbf{x}_{\pi_i(G)}} \hat{p}(x_i, \mathbf{x}_{\pi_i(G)}) \log \frac{\hat{p}(x_i, \mathbf{x}_{\pi_i(G)}, \theta_{i|\pi_i(G)})}{\hat{p}(\mathbf{x}_{\pi_i(G)}) \hat{p}(x_i)} \right) - M \sum_i \left( \sum_{x_i} \hat{p}(x_i) \log \hat{p}(x_i) \right) \\ &= M \sum_i \hat{I}(x_i, \mathbf{x}_{\pi_i(G)}) - M \sum_i \hat{H}(x_i)\end{aligned}$$

Decomposable score and a function of the graph structure





# Chow-Liu tree learning algorithm

- Objection function:

$$\begin{aligned}\ell(\theta_G, G; D) &= \log \hat{p}(D | \theta_G, G) \\ &= M \sum_i \hat{I}(x_i, \mathbf{x}_{\pi_i(G)}) - M \sum_i \hat{H}(x_i)\end{aligned} \Rightarrow C(G) = M \sum_i \hat{I}(x_i, \mathbf{x}_{\pi_i(G)})$$

- Chow-Liu:

- For each pair of variable  $x_i$  and  $x_j$ 
  - Compute empirical distribution:
  - Compute mutual information:
- Define a graph with node  $x_1, \dots, x_n$ 
  - Edge  $(i,j)$  gets weight

$$\hat{I}(X_i, X_j)$$

$$\begin{aligned}\hat{p}(X_i, X_j) &= \frac{\text{count}(x_i, x_j)}{M} \\ \hat{I}(X_i, X_j) &= \sum_{x_i, x_j} \hat{p}(x_i, x_j) \log \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}\end{aligned}$$





# Chow-Liu algorithm (con'd)

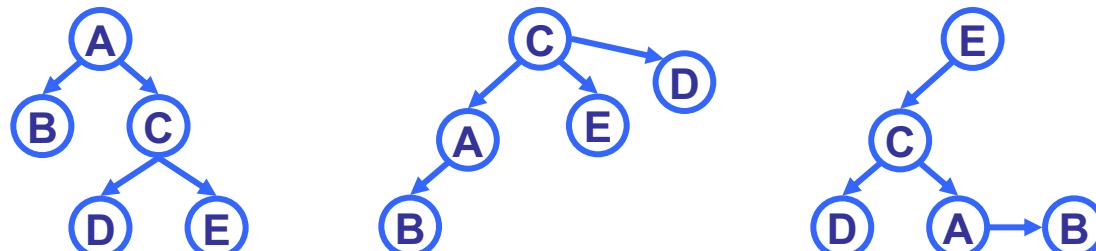
- Objection function:

$$\begin{aligned}\ell(\theta_G, G; D) &= \log \hat{p}(D | \theta_G, G) \\ &= M \sum_i \hat{I}(x_i, \mathbf{x}_{\pi_i(G)}) - M \sum_i \hat{H}(x_i)\end{aligned} \Rightarrow C(G) = M \sum_i \hat{I}(x_i, \mathbf{x}_{\pi_i(G)})$$

- Chow-Liu:

Optimal tree BN

- Compute maximum weight spanning tree
- Direction in BN: pick any node as root, do breadth-first-search to define directions
- I-equivalence:



$$C(G) = I(A, B) + I(A, C) + I(C, D) + I(C, E)$$





# Structure Learning for general graphs

- Theorem:
  - The problem of learning a BN structure with at most  $d$  parents is NP-hard for any (fixed)  $d \geq 2$
- Most structure learning approaches use heuristics
  - Exploit score decomposition
  - Two heuristics that exploit decomposition in different ways
    - Greedy search through space of node-orders
    - Local search of graph structures

