

Readings:

K&F: 3.4, 14.1, 14.2

# BN Semantics 3 –

Now it's personal!

## Parameter Learning 1

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

September 22<sup>nd</sup>, 2006

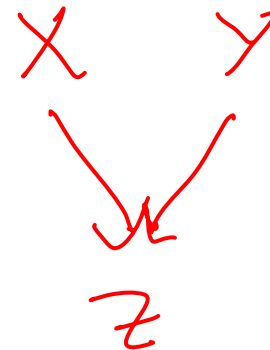
# Building BNs from independence properties

- From d-separation we learned:
  - Start from local Markov assumptions, obtain all independence assumptions encoded by graph
  - For most  $P$ 's that factorize over  $G$ ,  $I(G) = I(P)$
  - All of this discussion was for a given  $G$  that is an I-map for  $P$
- Now, give me a  $P$ , how can I get a  $G$ ?
  - i.e., give me the independence ~~assumptions~~ <sup>assertions</sup> entailed by  $P$
  - Many  $G$  are “equivalent”, how do I represent this?
  - Most of this discussion is not about practical algorithms, but useful concepts that will be used by practical algorithms
    - Practical algs next week

# Minimal I-maps

- One option:
  - $G$  is an I-map for  $P$
  - $G$  is as simple as possible
- $G$  is a minimal I-map for  $P$  if deleting any edges from  $G$  makes it no longer an I-map

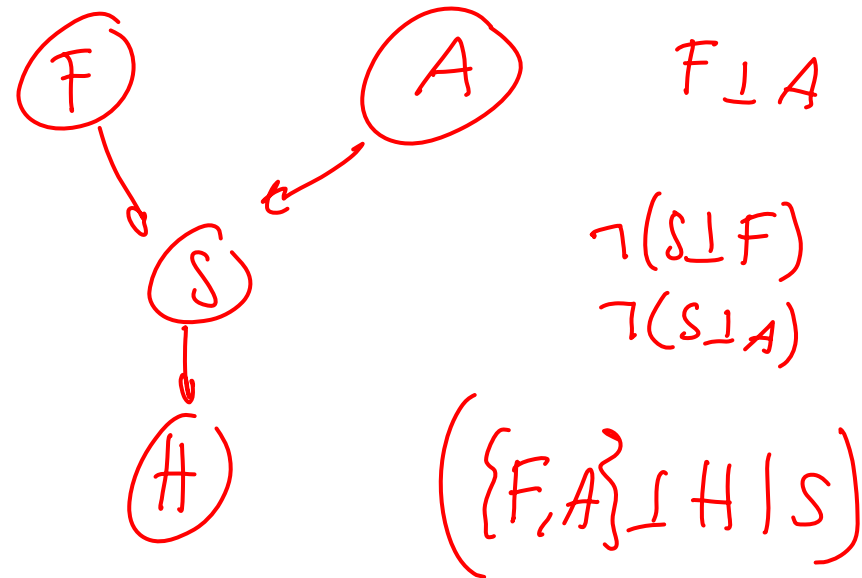
true  $P$   $X \perp Y$   
and nothing else  
vars  $X, Y, Z$



# Obtaining a minimal I-map

- Given a set of variables and conditional independence ~~assumptions~~ *assertions that P entails*
- Choose an ordering on variables, e.g.,  $X_1, \dots, X_n$
- For  $i = 1$  to  $n$ 
  - Add  $X_i$  to the network
  - Define parents of  $X_i$ ,  $\mathbf{Pa}_{X_i}$ , in graph as the minimal subset of  $\{X_1, \dots, X_{i-1}\}$  such that local Markov assumption holds –  $X_i$  independent of rest of  $\{X_1, \dots, X_{i-1}\}$ , given parents  $\mathbf{Pa}_{X_i}$
  - Define/learn CPT –  $P(X_i | \mathbf{Pa}_{X_i})$

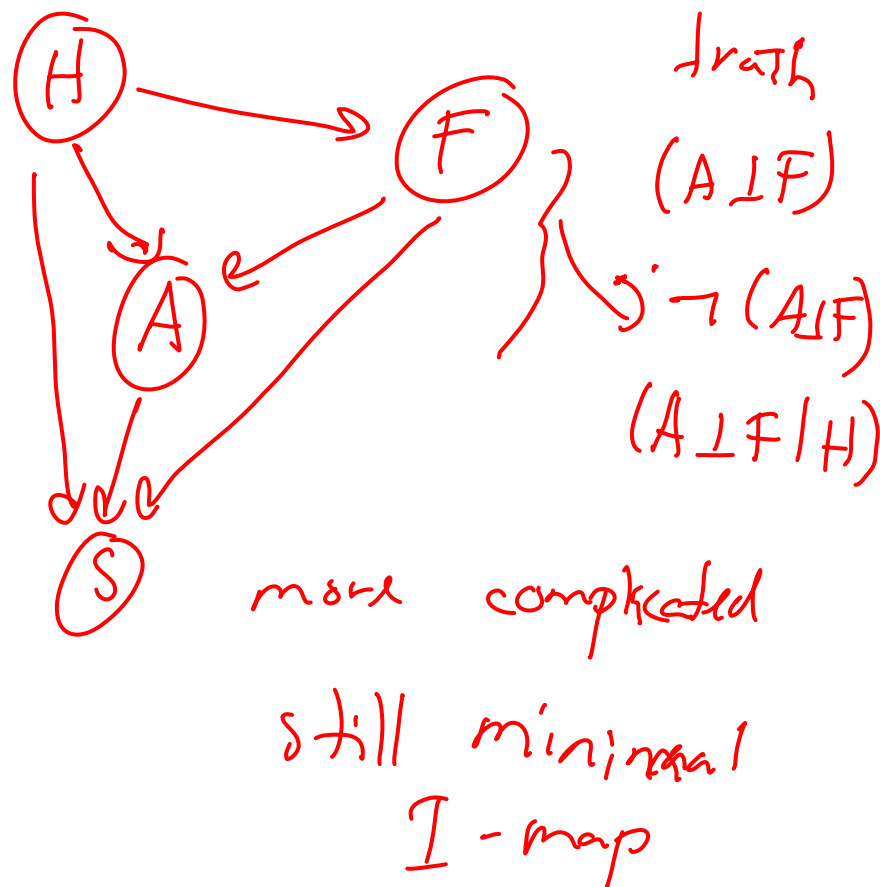
1 2 3 4  
Flu, Allergy, SinusInfection, Headache



# Minimal I-map not unique (or minimal)

- Given a set of variables and conditional independence assumptions
- Choose an ordering on variables, e.g.,  $X_1, \dots, X_n$
- For  $i = 1$  to  $n$ 
  - Add  $X_i$  to the network
  - Define parents of  $X_i$ ,  $\mathbf{Pa}_{X_i}$ , in graph as the minimal subset of  $\{X_1, \dots, X_{i-1}\}$  such that local Markov assumption holds –  $X_i$  independent of rest of  $\{X_1, \dots, X_{i-1}\}$ , given parents  $\mathbf{Pa}_{X_i}$
  - Define/learn CPT –  $P(X_i | \mathbf{Pa}_{X_i})$

<sup>2</sup>  
~~Flu~~ <sup>3</sup> ~~Allergy~~ ~~SinusInfection~~ ~~Headache~~ <sup>1</sup>  
Flu, Allergy, SinusInfection, Headache



# Perfect maps (P-maps)

- <sup>minimal</sup> I-maps are not unique and often not simple enough
- Define “simplest”  $G$  that is I-map for  $P$ 
  - A BN structure  $G$  is a perfect map for a distribution  $P$  if  $I(P) = I(G)$
- Our goal:
  - Find a perfect map!
  - Must address equivalent BNs

# Inexistence of P-maps 1

- XOR (this is a hint for the homework)

$$Z = X \oplus Y$$

$Z$  true if exactly one of  $X$  or  $Y$  true

$$(X \perp Y)$$

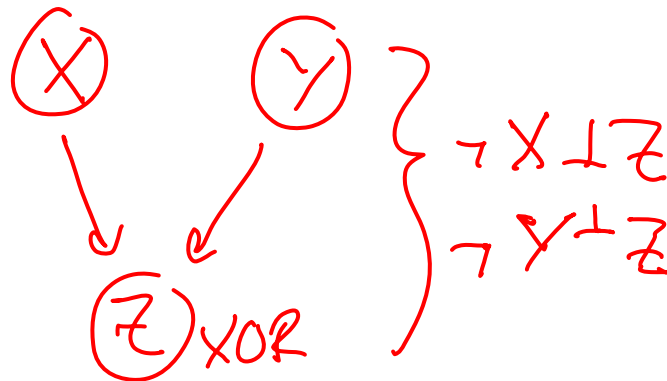
$$(Z \perp Y)$$

$$(Z \perp X)$$

$$\neg (Z \perp Y \mid X)$$

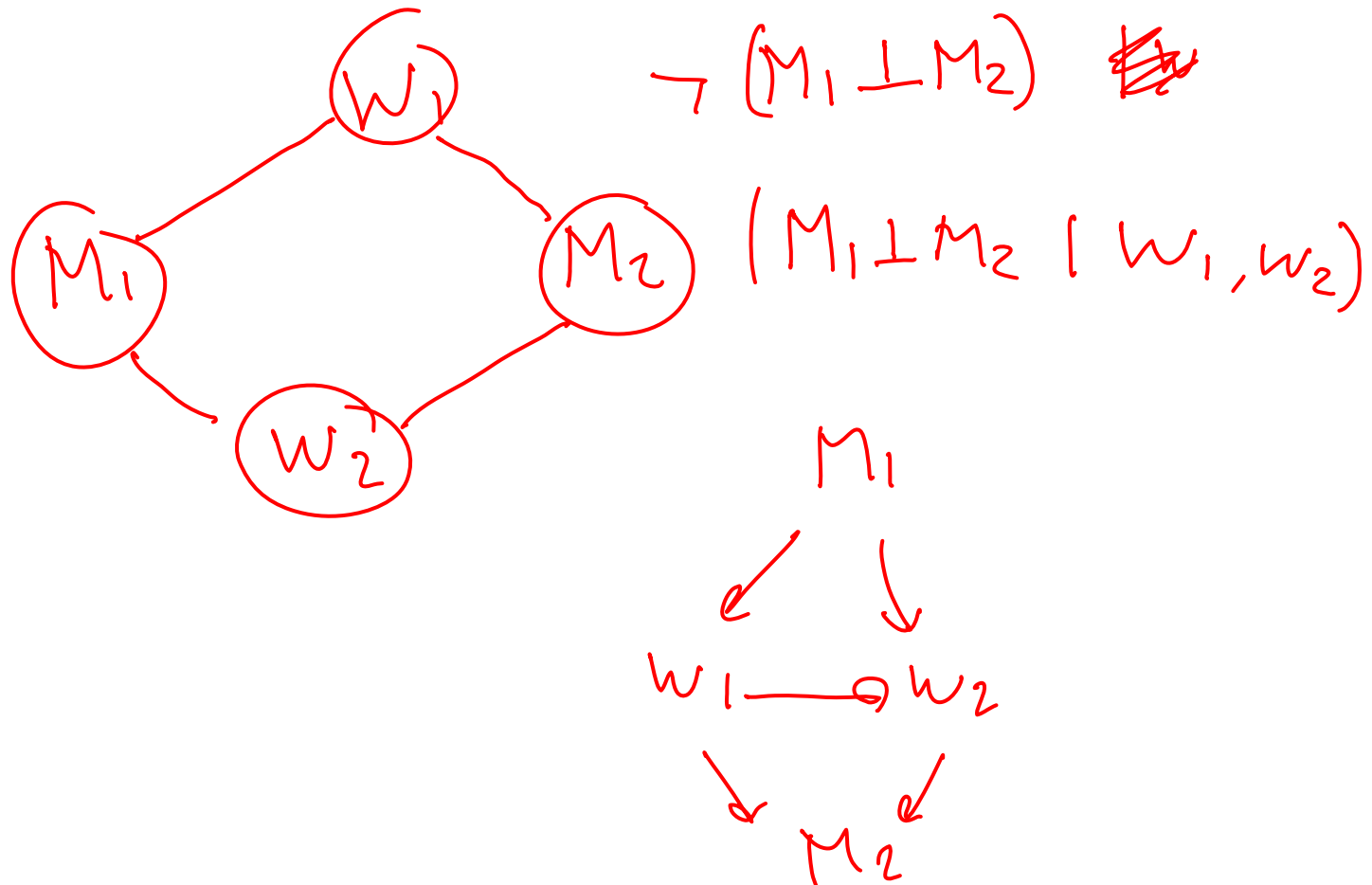
$$\neg (X \perp Y \mid Z)$$

$$\neg (Z \perp X \mid Y)$$



# Inexistence of P-maps 2

- (Slightly un-PC) swinging couples example



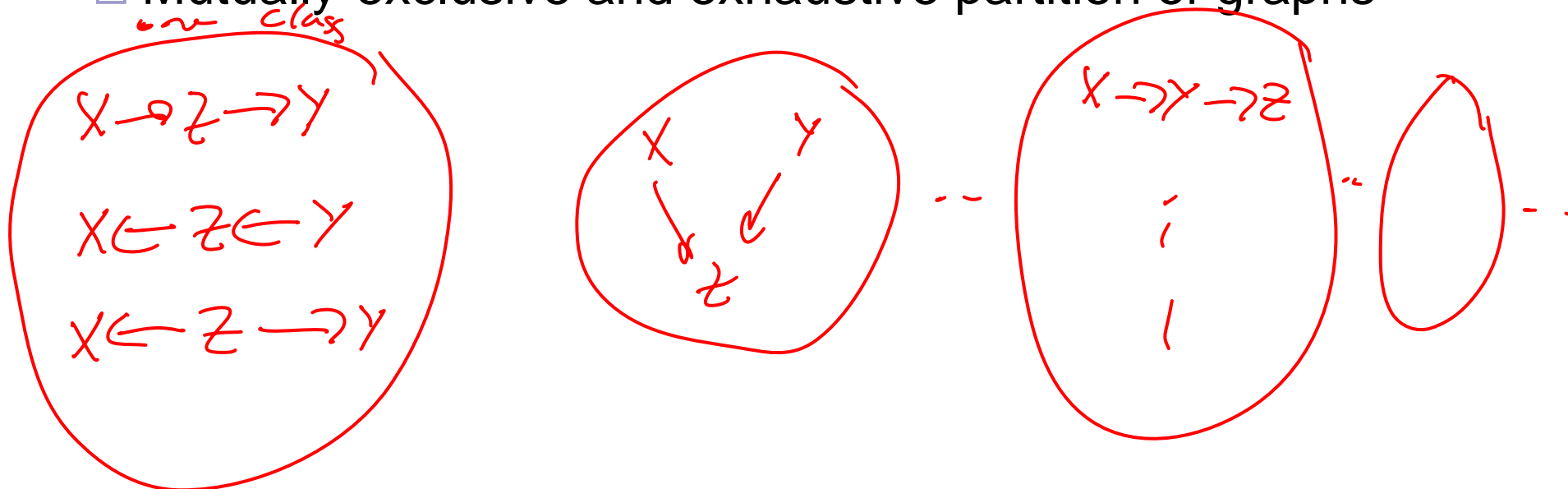


# Obtaining a P-map

- Given the independence assertions that are true for  $P$
- Assume that there exists a perfect map  $G^*$ 
  - Want to find  $G^*$
- Many structures may encode same independencies as  $G^*$ , when are we done?
  - Find all equivalent structures simultaneously!

# I-Equivalence

- Two graphs  $G_1$  and  $G_2$  are **I-equivalent** if  $I(G_1) = I(G_2)$
- Equivalence class** of BN structures
  - Mutually-exclusive and exhaustive partition of graphs



- How do we characterize these equivalence classes?

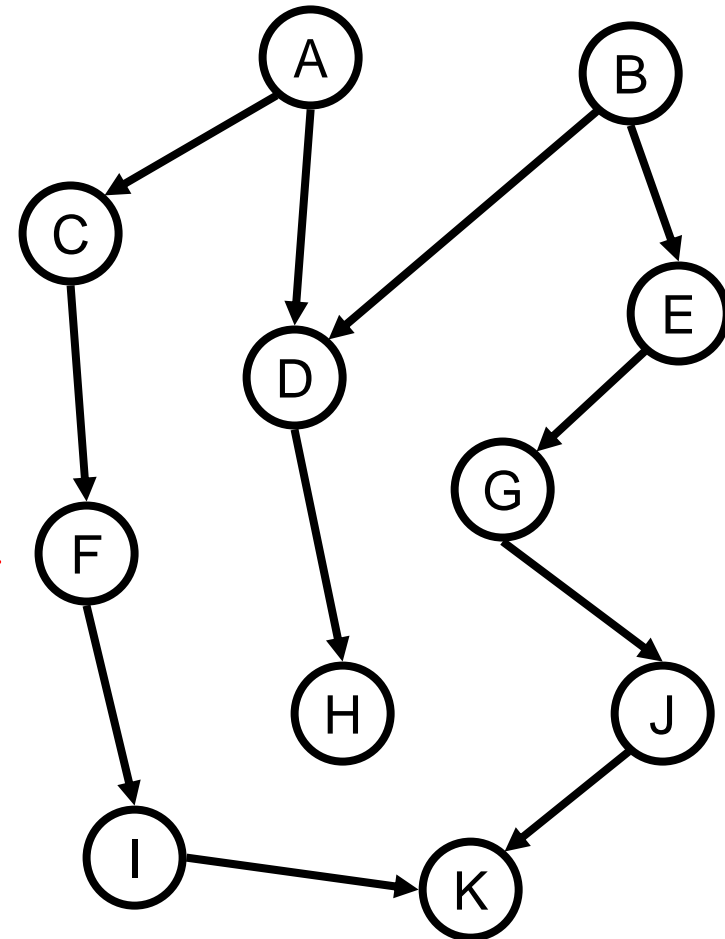
# Skeleton of a BN

- **Skeleton** of a BN structure  $G$  is an **undirected graph** over the same variables that has an edge  $X-Y$  for every  $X \rightarrow Y$  or  $Y \rightarrow X$  in  $G$



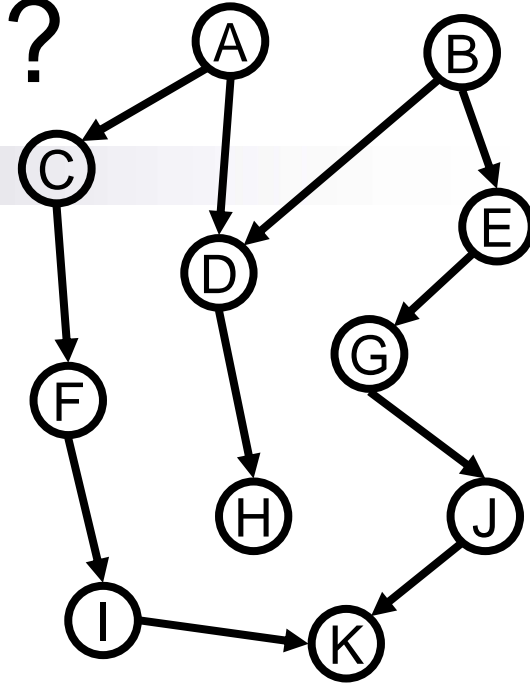
- (Little) **Lemma**: Two equivalent BN structures must have the same skeleton

counter example



# What about V-structures?

- V-structures are key property of BN structure

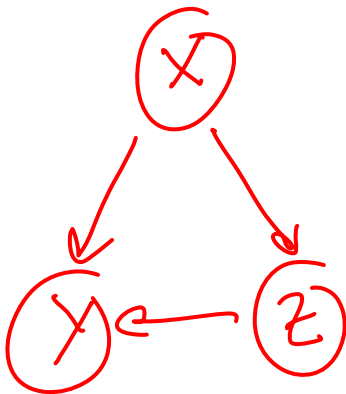


- **Theorem:** If  $G_1$  and  $G_2$  have the same skeleton and V-structures, then  $G_1$  and  $G_2$  are I-equivalent

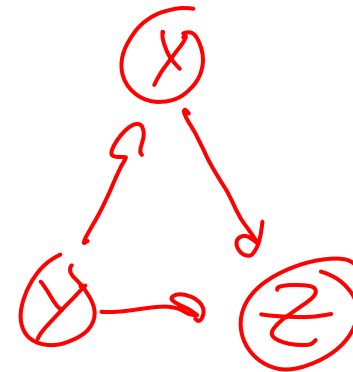
*not if and only if*

# Same V-structures not necessary

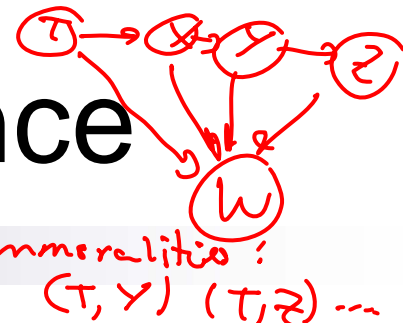
- **Theorem:** If  $G_1$  and  $G_2$  have the same skeleton and V-structures, then  $G_1$  and  $G_2$  are I-equivalent
- Though sufficient, same V-structures not necessary



diff. V-structures  
but  
I-equiv.



# Immoralities & I-Equivalence



- Key concept not V-structures, but “immoralities” (unmarried parents 😊)
  - $X \rightarrow Z \leftarrow Y$ , with no <sup>edge</sup> arrow between X and Y
  - Important pattern: X and Y independent given their parents, but not given Z
  - (If edge exists between X and Y, we have covered the V-structure)
- **Theorem**:  $G_1$  and  $G_2$  have the same skeleton and immoralities if and only if  $G_1$  and  $G_2$  are I-equivalent

# Obtaining a P-map

- Given the independence assertions that are true for  $P$ 
  - Obtain skeleton
  - Obtain immoralities
- From skeleton and immoralities, obtain every (and any) BN structure from the equivalence class

# Identifying the skeleton 1

- When is there an edge between X and Y?

Edge  $X \rightarrow Y$  if  $\nexists Z \subseteq \{x_1, \dots, x_n\} / (X, Y)$   
 $(X \perp Y | Z)$

- When is there no edge between X and Y?

$\exists Z \therefore (X \perp Y | Z)$



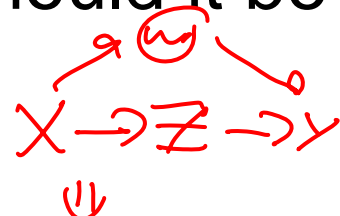
# Identifying the skeleton 2

- Assume d is max number of parents (d could be n)
- For each  $X_i$  and  $X_j$ 
  - $E_{ij} \leftarrow \text{true}$
  - For each  $\mathbf{U} \subseteq \mathbf{X} - \{X_i, X_j\}$ ,  $|\mathbf{U}| \leq \cancel{2}d$ 
    - Is  $(X_i \perp X_j \mid \mathbf{U})$  ?
      - $E_{ij} \leftarrow \text{~~true~~ false}$
  - If  $E_{ij}$  is true
    - Add edge  $X - Y$  to skeleton

# Identifying immoralities

- Consider  $X - Z - Y$  in skeleton, when should it be an immorality?

*other active path*



- Must be  $X \rightarrow Z \leftarrow Y$  (immorality):

*must check larger subsets*

- When  $X$  and  $Y$  are **never independent** given  $\mathbf{U}$ , if  $Z \in \mathbf{U}$

*$\nexists U \subseteq \{X_1, \dots, X_n\} - \{X, Y\} \text{ s.t. } Z \in U \text{ and } (X \perp Y | U)$*

- Must not be  $X \rightarrow Z \leftarrow Y$  (not immorality):

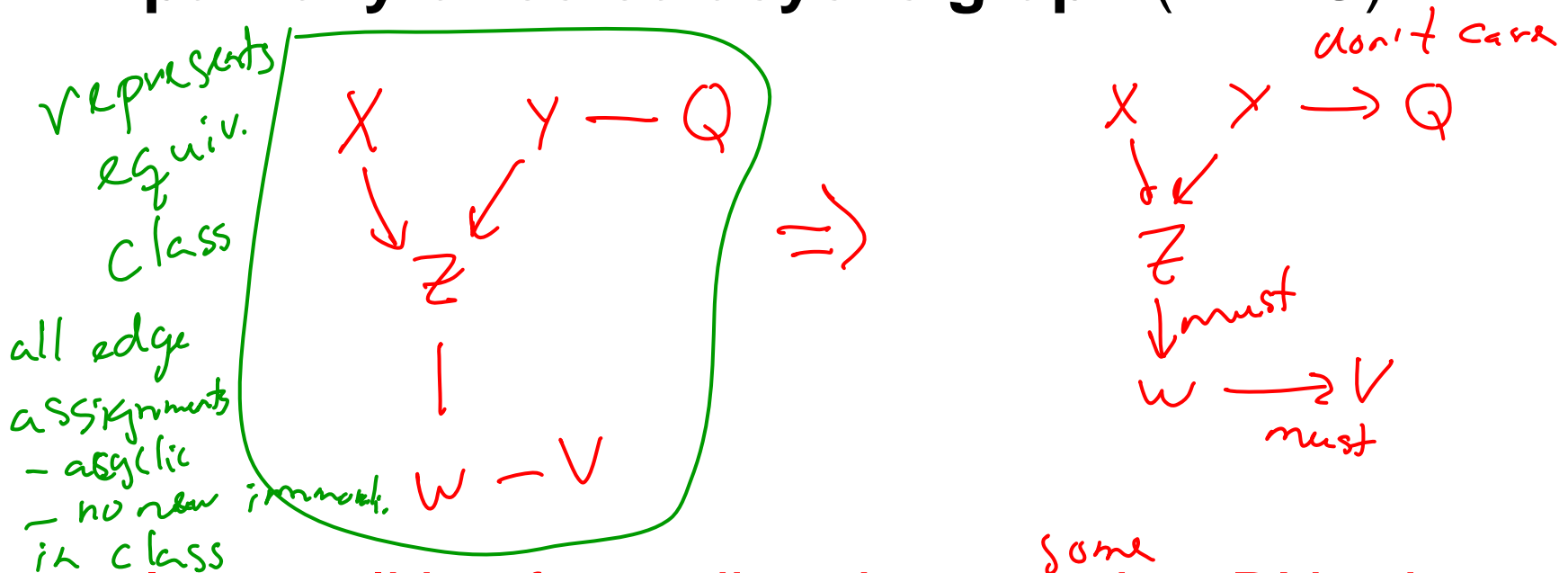
- When there exists  $\mathbf{U}$  with  $Z \in \mathbf{U}$ , such that  $X$  and  $Y$  are **independent** given  $\mathbf{U}$

*possible dirs:*

*$X \rightarrow Z \rightarrow Y$   
 $X \leftarrow Z \leftarrow Y$   
 $X \leftarrow Z \rightarrow Y$*

# From immoralities and skeleton to BN structures

- Representing BN equivalence class as a **partially-directed acyclic graph (PDAG)**



- Immoralities force direction on <sup>some</sup> other BN edges
- Full (polynomial-time) procedure described in reading

# What you need to know

- Minimal I-map

- every  $P$  has one, but usually many

- Perfect map

- better choice for BN structure
  - not every  $P$  has one
  - can find one (if it exists) by considering I-equivalence
  - Two structures are I-equivalent if they have same skeleton and immoralities

# Announcements

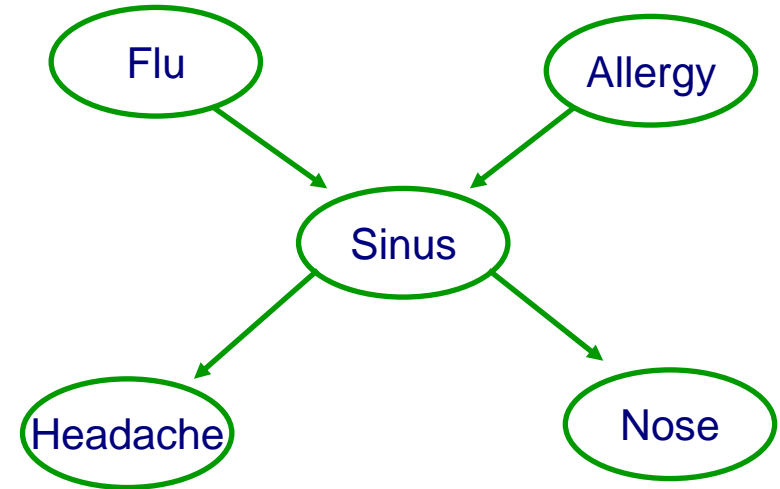


- I'll lead a special discussion session:
  - Today 2-3pm in NSH 1507
    - talk about homework, especially programming question

# Review

## ■ Bayesian Networks

- Compact representation for probability distributions
- Exponential reduction in number of parameters
- Exploits independencies



## ■ Next – Learn BNs

- parameters
- structure

# Thumbtack – Binomial Distribution

- $P(\text{Heads}) = \theta$ ,  $P(\text{Tails}) = 1 - \theta$



- Flips are i.i.d.:

- ☐ Independent events
- ☐ Identically distributed according to Binomial distribution

- Sequence  $\mathcal{D}$  of  $\alpha_H$  Heads and  $\alpha_T$  Tails

$$P(\underline{\mathcal{D}} \mid \underline{\theta}) = \underline{\theta}^{\alpha_H} (1 - \theta)^{\alpha_T}$$

# Maximum Likelihood Estimation

- **Data:** Observed set  $D$  of  $\alpha_H$  Heads and  $\alpha_T$  Tails
- **Hypothesis:** Binomial distribution
- Learning  $\theta$  is an optimization problem
  - What's the objective function?
- MLE: Choose  $\theta$  that maximizes the probability of observed data:

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} P(\mathcal{D} \mid \theta) \\ &= \arg \max_{\theta} \ln P(\mathcal{D} \mid \theta)\end{aligned}$$



# Your first learning algorithm

$$\ln a^5 = 5 \ln a$$

$$\hat{\theta} = \arg \max_{\theta} \ln P(\mathcal{D} | \theta)$$

$$\ln a \cdot b = \ln a + \ln b$$

$$= \arg \max_{\theta} \ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

$$\frac{\partial}{\partial \theta} \ln \theta = \frac{1}{\theta}$$

$$\frac{\partial}{\partial \theta} \ln(1 - \theta) = \frac{-1}{1 - \theta}$$

■ Set derivative to zero:

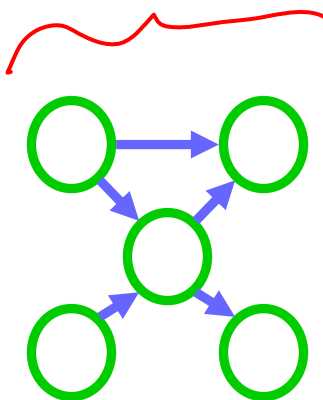
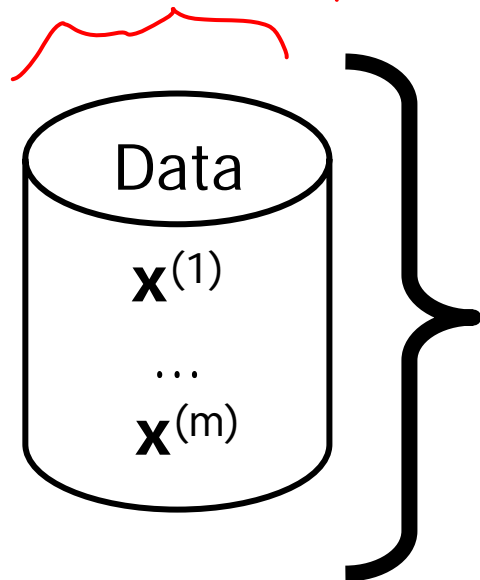
$$\frac{d}{d\theta} \ln P(\mathcal{D} | \theta) = 0$$

$$\begin{aligned} & \frac{\partial}{\partial \theta} \ln[\theta^{\alpha_H} (1 - \theta)^{\alpha_T}] \\ &= \frac{\partial}{\partial \theta} [\alpha_H \ln \theta + \alpha_T \ln(1 - \theta)] = \frac{\partial}{\partial \theta} \alpha_H \ln \theta + \frac{\partial}{\partial \theta} \alpha_T \ln(1 - \theta) \\ &= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1 - \theta} = 0 \quad \Rightarrow \quad \theta = \frac{\alpha_H}{\alpha_H + \alpha_T} \end{aligned}$$

# Learning Bayes nets

	Known structure	Unknown structure
Fully observable data	<i>easy</i>	<i>NP-hard but not solved</i>
Missing data	<i>hard (things like EM)</i>	<i>very hard but → in a few weeks</i>

*→  $\langle x_1=t, x_2=?, x_3=f \rangle$*



**structure**

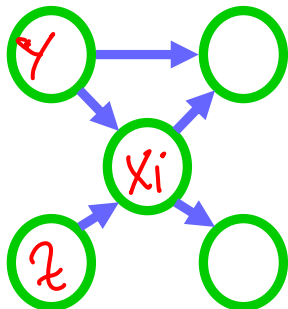
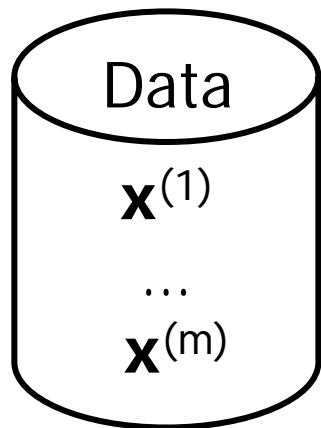
+



CPTs –  
 $P(X_i | \mathbf{Pa}_{X_i})$

**parameters**

# Learning the CPTs



For each discrete variable  $X_i$

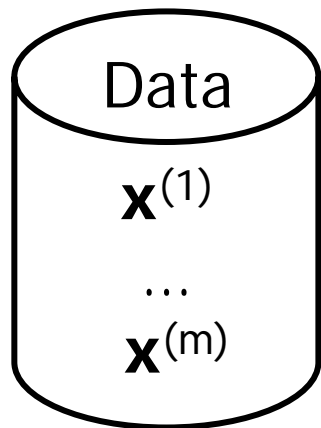
$$P(X_i | \text{Pa}_{X_i}) = P(X_i | Y, Z)$$

$\stackrel{\text{MLE}}{\approx}$

$$P(X_i = x_i | Y = y, Z = z) \stackrel{\text{MLE}}{\approx} \frac{\text{Count}(X_i = x_i, Y = y, Z = z)}{\text{Count}(Y = y, Z = z)}$$

$$\text{MLE: } P(X_i = x_i | X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

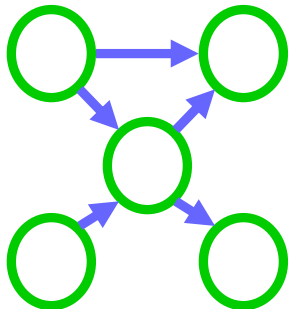
# Learning the CPTs



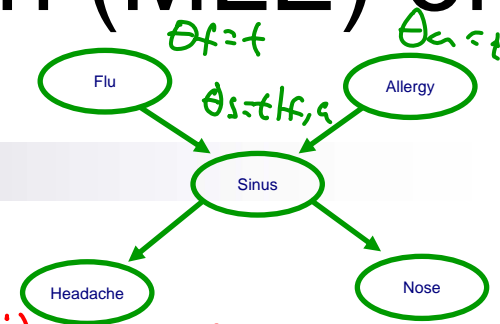
For each discrete variable  $X_i$

$$\text{MLE: } P(X_i = x_i \mid X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

**WHY??????????**



# Maximum likelihood estimation (MLE) of BN parameters – example



■ Given structure, log likelihood of data:

$$\begin{aligned}
 \log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) &\stackrel{\text{iid}}{=} \log \prod_i P(F=f^{(i)}, A=a^{(i)}, S=s^{(i)}, H=h^{(i)}, N=n^{(i)}) \\
 &\stackrel{\text{structure}}{=} \log \prod_i P(F=f^{(i)}) P(A=a^{(i)}) \cdot P(S=s^{(i)} \mid a^{(i)}, f^{(i)}) \cdot P(H=h^{(i)} \mid s^{(i)}) \cdot P(N=n^{(i)} \mid s^{(i)}) \\
 &= \sum_i [\log P(f^{(i)}) + \log P(a^{(i)}) + \log P(s^{(i)} \mid a^{(i)}, f^{(i)}) + \log P(h^{(i)} \mid s^{(i)}) + \log P(n^{(i)} \mid s^{(i)})] \\
 &\stackrel{\text{argmax}}{=} \left[ \sum_i \log P(f^{(i)}) \right] + \sum_i \log P(a^{(i)}) + \sum_i \log P(s^{(i)} \mid a^{(i)}, f^{(i)}) + \dots \\
 &= \left[ \argmax_{\theta_f=t} \sum_i \log P(f^{(i)} \mid \theta_f=t) \right] + \dots \argmax_{\theta_a=t} \dots + \argmax_{\theta_{s=t,f,a}} \dots
 \end{aligned}$$

# Maximum likelihood estimation (MLE) of BN parameters – General case

- Data:  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$
- Restriction:  $\mathbf{x}^{(j)}[\mathbf{Pa}_{x_i}] \rightarrow$  assignment to  $\mathbf{Pa}_{x_i}$  in  $\mathbf{x}^{(j)}$
- Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \log \prod_i \prod_j P(x_i = x_i^{(j)} \mid \mathbf{Pa}_{x_i} = \mathbf{x}^{(j)}[\mathbf{Pa}_{x_i}])$$

# Taking derivatives of MLE of BN parameters – General case

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log P \left( X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)} [\mathbf{Pa}_{X_i}] \right)$$

# General MLE for a CPT

- Take a CPT:  $P(X|U)$
- Log likelihood term for this CPT
- Parameter  $\theta_{X=x|U=u}$  :

$$\text{MLE: } P(X = x \mid U = u) = \theta_{X=x|U=u} = \frac{\text{Count}(X = x, U = u)}{\text{Count}(U = u)}$$



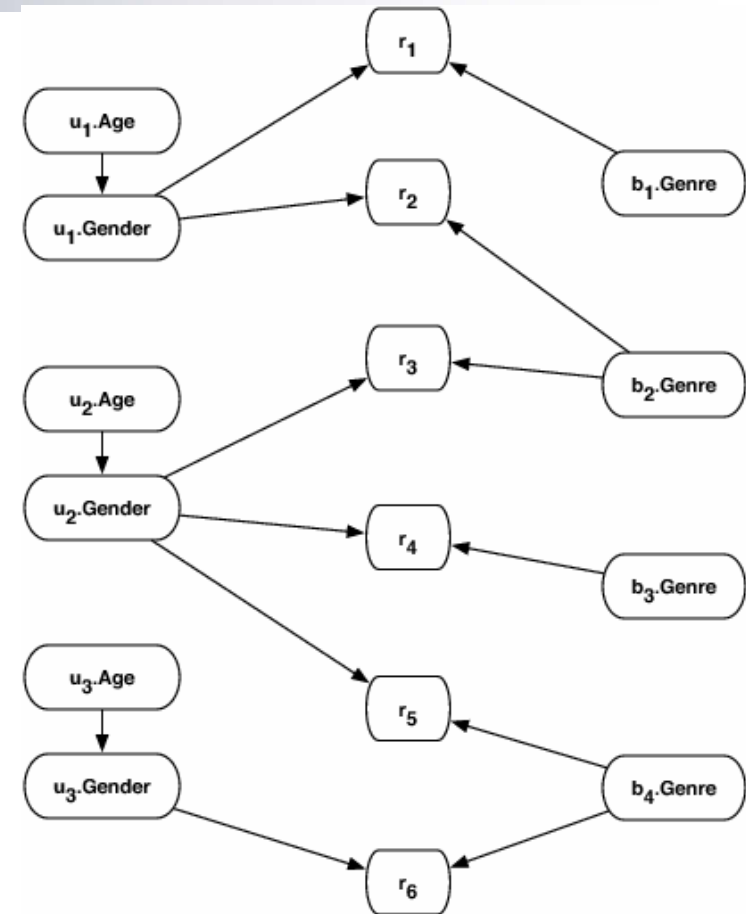
# Parameter sharing

(basics now, more later in the semester)

- Suppose we want to model customers' rating for books
- You know:
  - features of customers, e.g., age, gender, income,...
  - features of books, e.g., genre, awards, # of pages, has pictures,...
  - ratings: each user rates a few books
- A simple BN:

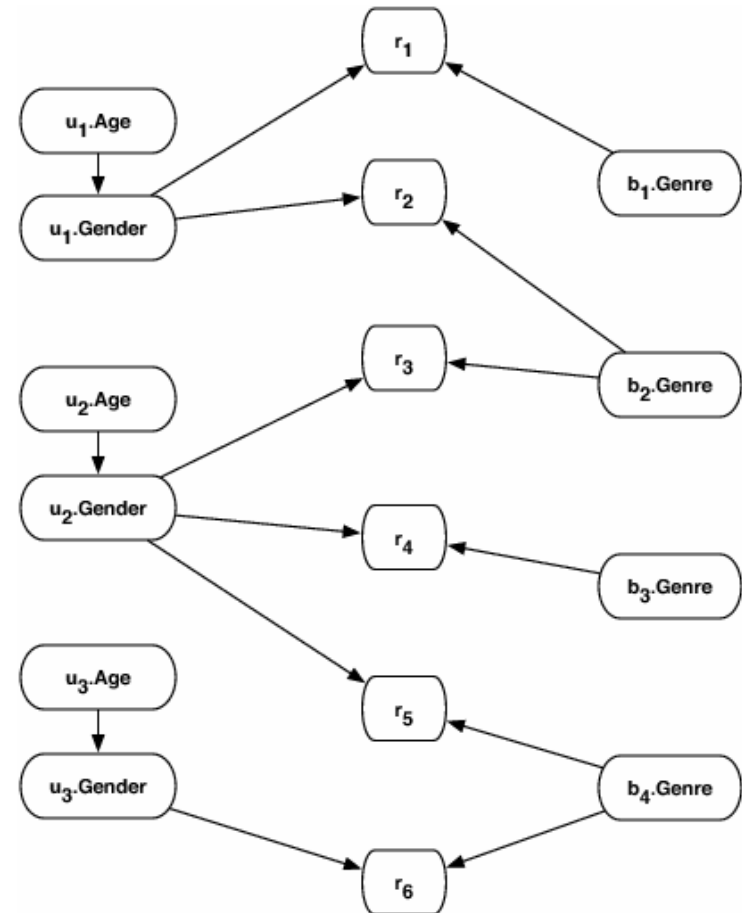
# Using recommender system

- Answer probabilistic question:



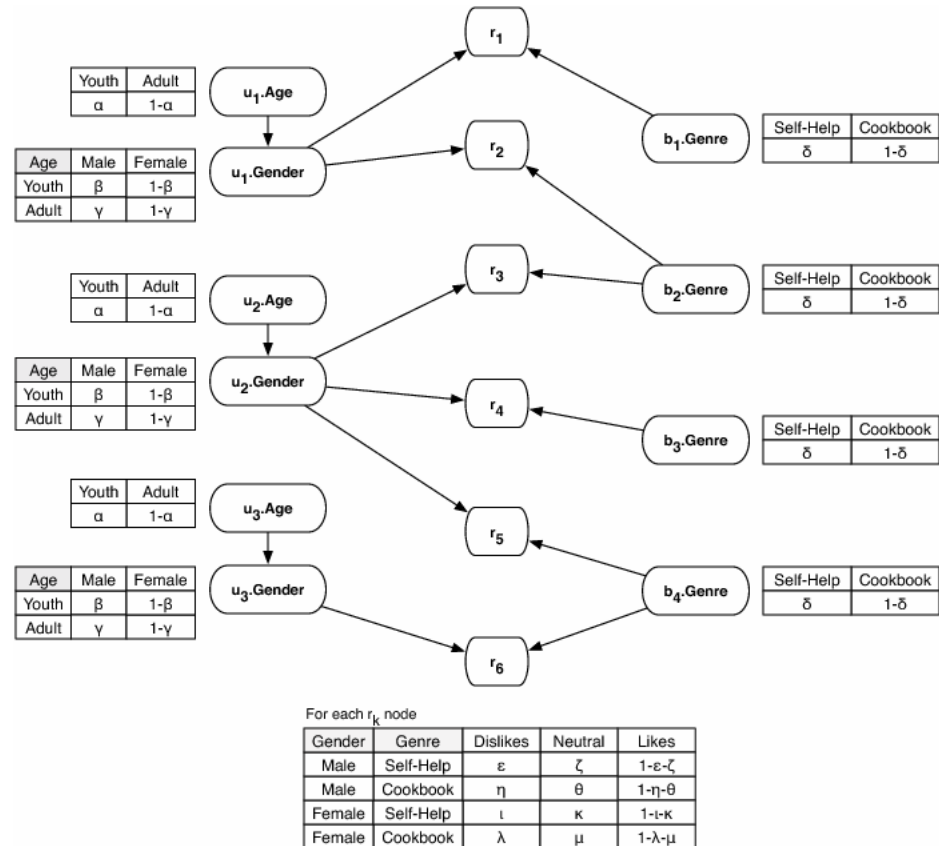
# Learning parameters of recommender system BN

- How many parameters do I have to learn?
- How many samples do I have?



# Parameter sharing for recommender system BN

- Use same parameters in many CPTs
- How many parameters do I have to learn?
- How many samples do I have?

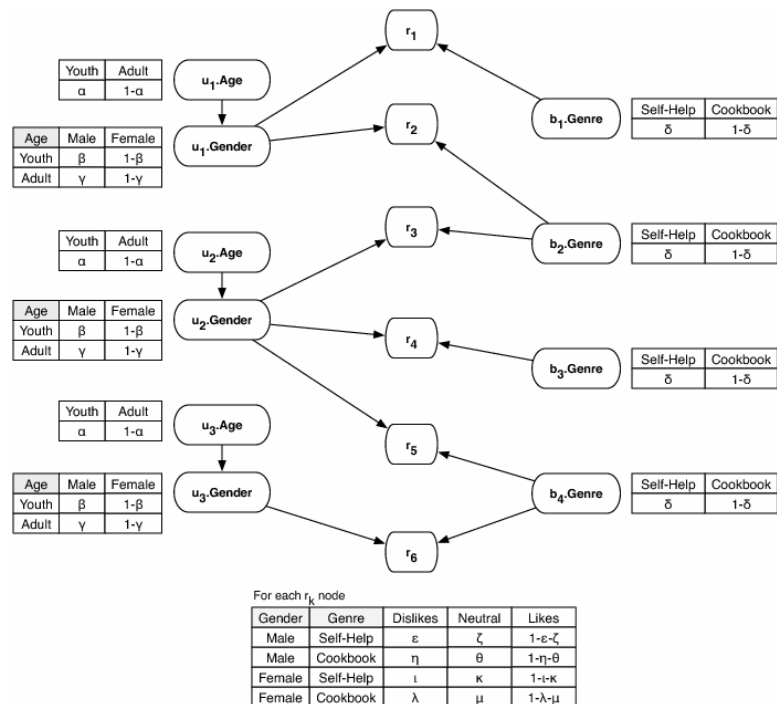


# MLE with simple parameter sharing

■ Estimating  $\alpha$ :

■ Estimating  $\beta$ :

■ Estimating  $\varepsilon$ :



# What you need to know about learning BNs thus far

- Maximum likelihood estimation
  - decomposition of score
  - computing CPTs
- Simple parameter sharing
  - why share parameters?
  - computing MLE for shared parameters