

Readings:

K&F: 3.4, 14.1, 14.2

Now it's personal!

# Parameter

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

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

1

## 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 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

10-708 – ©Carlos Guestrin 2006

3

## Obtaining a minimal I-map

- 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 –  $P(X_i | \mathbf{Pa}_{X_i})$

Flu, Allergy, SinusInfection, Headache

10-708 – ©Carlos Guestrin 2006

4

## 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 –  $P(X_i | \mathbf{Pa}_{X_i})$

Flu, Allergy, SinusInfection, Headache

10-708 – ©Carlos Guestrin 2006

5

## Perfect maps (P-maps)

- 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

10-708 – ©Carlos Guestrin 2006

6

## Inexistence of P-maps 1

- XOR (this is a hint for the homework)

10.708 – ©Carlos Guestrin 2006

7

## Inexistence of P-maps 2

- (Slightly un-PC) swinging couples example

10.708 – ©Carlos Guestrin 2006

8

## 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!

10-708 – ©Carlos Guestrin 2006

9

## 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?

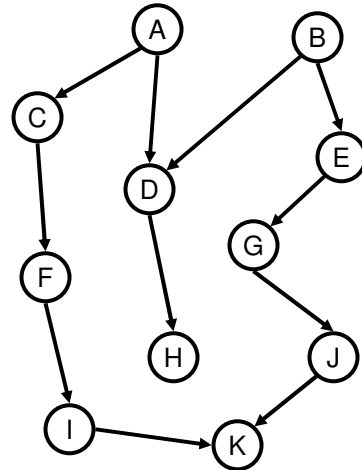
10-708 – ©Carlos Guestrin 2006

10

## 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 I-equivalent BN structures must have the same skeleton



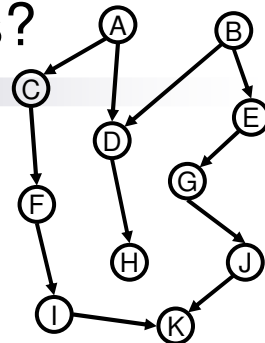
10-708 - ©Carlos Guestrin 2006

11

## 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



10-708 - ©Carlos Guestrin 2006

12

## 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

10-708 – ©Carlos Guestrin 2006

13

## Immoralities & I-Equivalence

- Key concept not V-structures, but “immoralities” (unmarried parents ☺)
  - $X \rightarrow Z \leftarrow Y$ , with no 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

10-708 – ©Carlos Guestrin 2006

14

## 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

10-708 – ©Carlos Guestrin 2006

15

## Identifying the skeleton 1

- When is there an edge between  $X$  and  $Y$ ?
- When is there no edge between  $X$  and  $Y$ ?

10-708 – ©Carlos Guestrin 2006

16



## 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  $\mathbf{U} \subseteq \mathbf{X} - \{X_i, X_j\}$ ,  $|\mathbf{U}| \leq 2d$ 
    - Is  $(X_i \perp X_j \mid \mathbf{U})$  ?
      - $E_{ij} \leftarrow \text{true}$
  - If  $E_{ij}$  is true
    - Add edge  $X - Y$  to skeleton

10-708 – ©Carlos Guestrin 2006

17

## Identifying immoralities

- Consider  $X - Z - Y$  in skeleton, when should it be an immorality?
- Must be  $X \rightarrow Z \leftarrow Y$  (immorality):
  - When  $X$  and  $Y$  are **never independent** given  $\mathbf{U}$ , if  $Z \in \mathbf{U}$
- 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  $\mathbf{U}$

10-708 – ©Carlos Guestrin 2006

18

## From immoralities and skeleton to BN structures

- Representing BN equivalence class as a **partially-directed acyclic graph** (PDAG)
- **Immoralities force direction on other BN edges**
- Full (polynomial-time) procedure described in reading

10-708 – ©Carlos Guestrin 2006

19

## 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

10-708 – ©Carlos Guestrin 2006

20

# Announcements

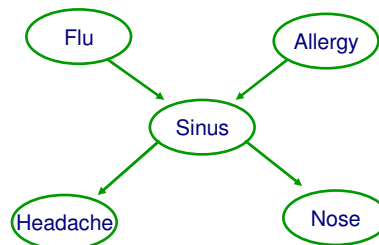
- I'll lead a special discussion session:
  - Today 2-3pm in NSH 1507
    - talk about homework, especially programming question

10-708 – ©Carlos Guestrin 2006

21

# Review

- Bayesian Networks
  - Compact representation for probability distributions
  - Exponential reduction in number of parameters
  - Exploits independencies
- Next – Learn BNs
  - parameters
  - structure



10-708 – ©Carlos Guestrin 2006

22

## 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  $D$  of  $\alpha_H$  Heads and  $\alpha_T$  Tails

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

10-708 – ©Carlos Guestrin 2006

23

## 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:

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

10-708 – ©Carlos Guestrin 2006

24

## Your first learning algorithm

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \ln P(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T}\end{aligned}$$

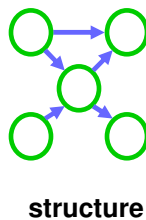
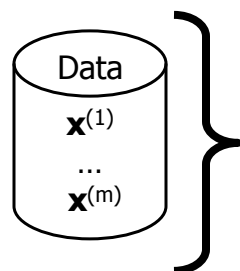
- Set derivative to zero:  $\frac{d}{d\theta} \ln P(\mathcal{D} | \theta) = 0$

10-708 – ©Carlos Guestrin 2006

25

## Learning Bayes nets

	Known structure	Unknown structure
Fully observable data		
Missing data		



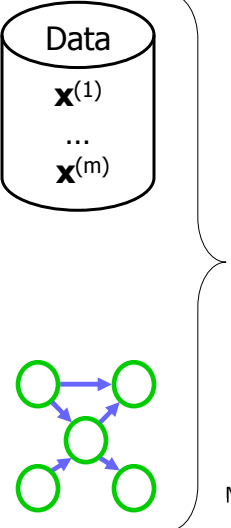
+

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

10-708 – ©Carlos Guestrin 2006

26

# Learning the CPTs

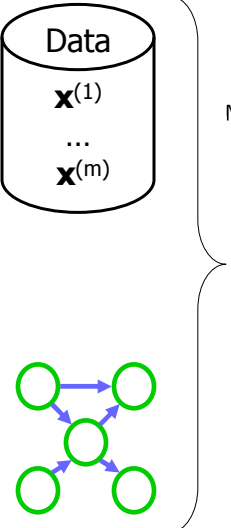


For each discrete variable  $X_i$

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

27

# Learning the CPTs



For each discrete variable  $X_i$

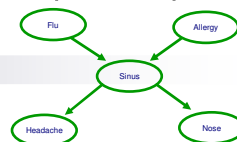
$$\text{MLE: } P(X_i = x_i \mid X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

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

28

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

- Given structure, log likelihood of data:  
 $\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$



10-708 – ©Carlos Guestrin 2006

29

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

- Data:  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$
- Restriction:  $\mathbf{x}^{(i)}[\mathbf{Pa}_{X_i}] \rightarrow$  assignment to  $\mathbf{Pa}_{X_i}$  in  $\mathbf{x}^{(i)}$
- Given structure, log likelihood of data:  
 $\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$

10-708 – ©Carlos Guestrin 2006

30

## Taking derivatives of MLE of BN parameters – General case

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

10-708 – ©Carlos Guestrin 2006

31

## General MLE for a CPT

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

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

10-708 – ©Carlos Guestrin 2006

32



# Parameter sharing

(basics now, more later in the semester)

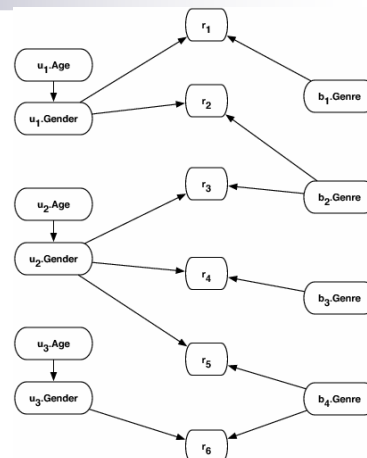
- Suppose we want to model customers' rating for books
- You know:
  - features of customers, e.g., age, gender, income,...
  - features of books, e.g., genre, awards, # of pages, has pictures,...
  - ratings: each user rates a few books
- A simple BN:

10-708 - ©Carlos Guestrin 2006

33

# Using recommender system

- Answer probabilistic question:

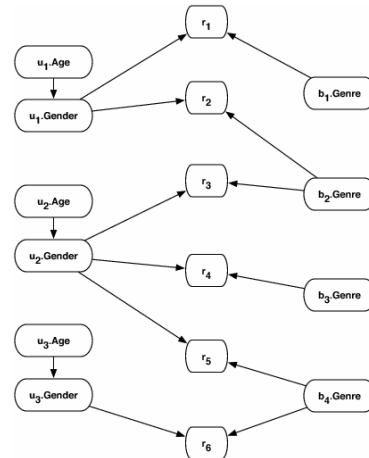


10-708 - ©Carlos Guestrin 2006

34

## Learning parameters of recommender system BN

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

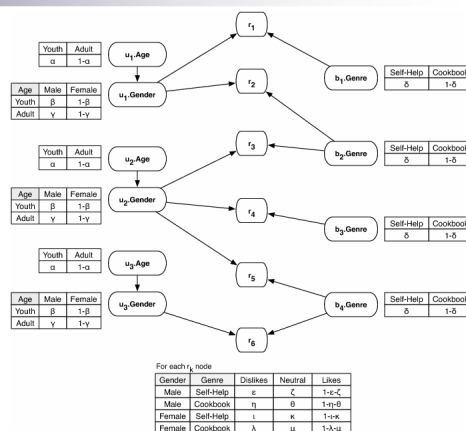


10.708 – ©Carlos Guestrin 2006

35

## 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?



10.708 – ©Carlos Guestrin 2006

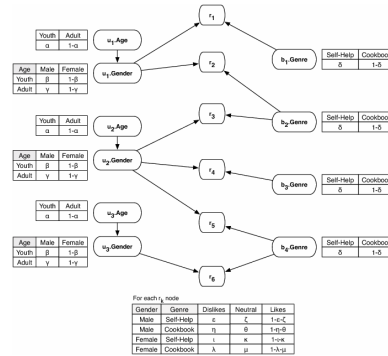
36

# MLE with simple parameter sharing

■ Estimating  $\alpha$ :

■ Estimating  $\beta$ :

■ Estimating  $\varepsilon$ :



10-708 – ©Carlos Guestrin 2006

37

## 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

10-708 – ©Carlos Guestrin 2006

38