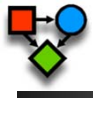


Readings: K&F 2.1, 2.2, 2.3, 3.1



# Introduction to Probabilistic Graphical Models

Lecture 1 – Mar 28, 2011  
CSE 515, Statistical Methods, Spring 2011

Instructor: Su-In Lee  
University of Washington, Seattle


## Logistics

- Teaching Staff
  - Instructor: Su-In Lee (suinlee@uw.edu, PAC 536)
    - Office hours: Fri 9-10am or by appointment (PAC 536)
  - TA: Andrew Guillory (guillory@cs.washington.edu)
    - Office hours: Wed 1:30-2:20 pm or by appointment (PAC 216)
- Course website
  - [cs.washington.edu/515](http://cs.washington.edu/515) ←
  - Discussion group: [course website](#) ↩
- Textbook
  - (required) Daphne Koller and Nir Friedman, Probabilistic Graphical Models: Principles and Techniques, MIT Press
  - Various research papers (copies available in class)

## Course requirement

- 4 homework assignments (60% of final grade)
  - Theory / implementation exercises
  - First one goes out next Monday!
  - 2 weeks to complete each
  - HW problems are long and hard
    - Please, please, please start early!
  - Late/collaboration policies are described on the website
- Final exam (35%)
  - Date will be announced later.
- Participation (5%)

## Probabilistic graphical models (PGMs)

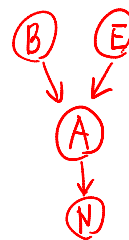
- One of the most exciting developments in machine learning (knowledge representation, AI, EE, Stats, ...) in the last two decades...
- Tool for representing complex systems and performing sophisticated reasoning tasks
- Why have a model? 
  - Compact and modular representation of complex systems
  - Ability to efficiently execute complex reasoning tasks
  - Make predictions
  - Generalize from particular problem

## Probabilistic graphical models (PGMs)

- Many classical probabilistic problems in statistics, information theory, pattern recognition, and statistical mechanics are special cases of the formalism
  - Graphical models provides a common framework
  - Advantage: specialized techniques developed in one field can be transferred between research communities
- PGMs are a marriage between graph theory and probability theory
  - Representation: graph
  - Reasoning: probability theory
  - Any simple example?

## A simple example

- We want to know/model whether our neighbor will inform us of the alarm being set off
- The alarm can set off (A) if
  - There is a burglary (B)
  - There is an earthquake (E)
- Whether our neighbor calls (N) depends on whether the alarm is set off (A)
- “Variables” in this system
  - Whether alarm being set off (A); burglary (B); earthquake (E); our neighbor calls (N)



# A simple

## Probabilistic Inference

**Task I** Say that the alarm is set off (A=True), then how likely is it to get a call from our neighbor (N=True)?

**Task II** Given that my neighbor calls (N=True), how likely it is that a burglary occurred (B=True)?

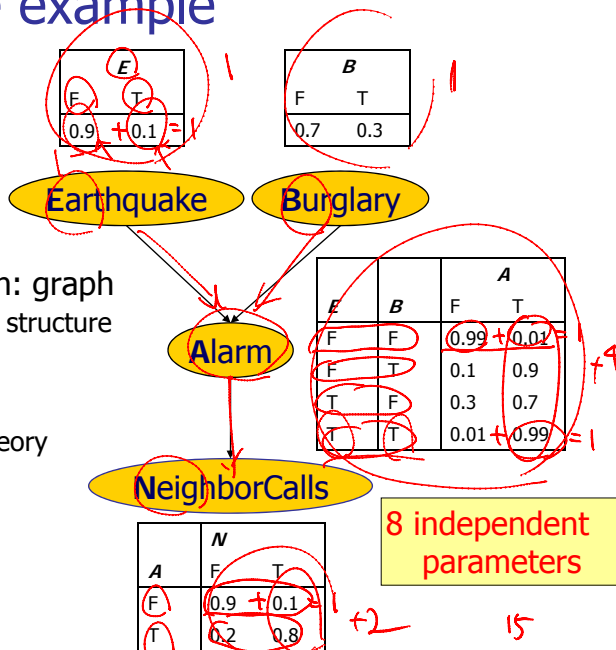
- Variable: Earthquake (E), Burglary (B), Alarm (A), NeighborCalls (N)

E	B	A	N	Prob.
F	F	F	F	0.01
F	F	F	T	0.04
F	F	T	F	0.05
F	F	T	T	0.01
F	T	F	F	0.02
F	T	F	T	0.07
F	T	T	F	0.2
F	T	T	T	0.1
T	F	F	F	0.01
T	F	F	T	0.07
T	F	T	F	0.13
T	F	T	T	0.04
T	T	F	F	0.06
T	T	F	T	0.05
T	T	T	F	0.1
T	T	T	T	0.05

Handwritten notes on the table:

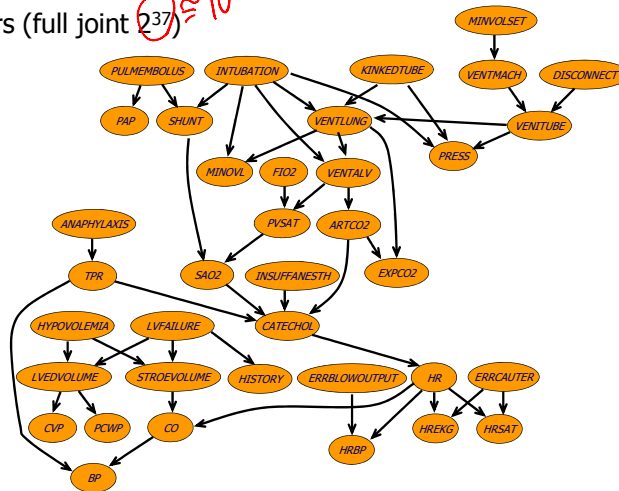
- Red circles around E, B, A, N columns.
- Red arrow pointing to the 0.01 probability cell with "1%".
- Red arrow pointing to the 0.05 probability cell with "5%".
- Equation:  $2^4 = 16$
- Equation:  $2^4 - 1 = 15$
- Text: "2^4 - 1 independent parameters" (circled in red).
- Equation:  $\sum P_i = 1 \Rightarrow 1 - \sum_{i \neq j} P_i$

# A simple example



## Example Bayesian network

- The "Alarm" network for monitoring intensive care patients
  - 37 variables
  - 509 parameters (full joint  $\approx 10^{37}$ )

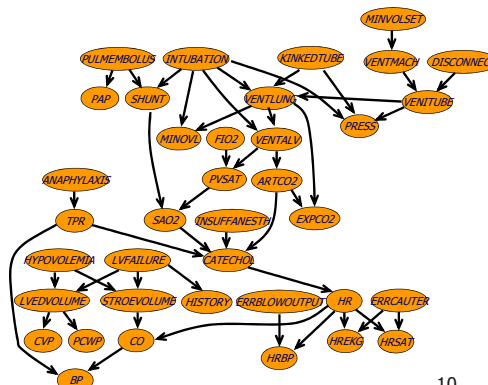


CSE 515 – Statistical Methods – Spring 2011

9

## Representation: graphs

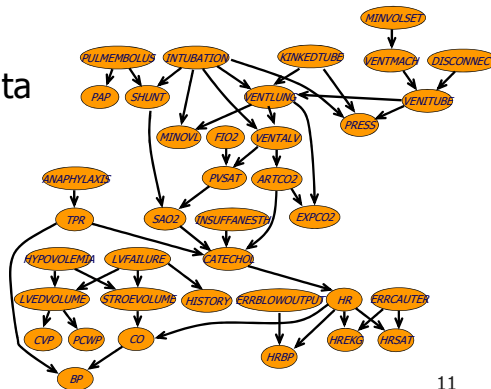
- Intuitive data structure for modeling highly-interacting sets of variables
  - Compact representation
  - Explicit model of modularity
- Data structure that allows for design of efficient general-purpose algorithms



10

## Reasoning: probability theories

- Well understood framework for **modeling uncertainty**
  - Partial knowledge of the state of the world
  - Noisy observations
  - Phenomenon not covered by our model
  - Inherent stochasticity
- Clear semantics
- Can be learned from data



11

## Probabilistic reasoning

- This course covers:
  - Probabilistic graphical model (PGM) **representation**
    - Bayesian networks (directed graph)
    - Markov networks (undirected graph)
  - Answering queries in PGMs ("**inference**")
    - What is the probability of X given some observations?
    - What is the most likely explanation for what is happening?
  - Learning PGMs from data ("**learning**")
    - What are the right/good parameters/structure of the model?
  - Application & special topics
    - Modeling **temporal processes** with PGMs
      - Hidden Markov Models (HMMs) as a special case
    - Modeling **decision-making processes**
      - Markov Decision Processes (MDPs) as a special case

## Course outline

Week	Topic	Reading
1	Introduction, Bayesian network representation	2.1-3, 3.1
	Bayesian network representation cont.	3.1-3
2	Local probability models	5
	Undirected graphical models	4
3	Exact inference	9.1-4
	Exact inference cont.	10.1-2
4	Approximate inference	12.1-3
	Approximate inference cont.	12.1-3
5	Parameter estimation	17
	Parameter estimation cont.	17
6	Partially observed data (EM algorithm)	19.1-3
	Structure learning BNs	18
7	Structure learning BNs cont.	18
	Partially observed data	19.4-5
8	Learning undirected graphical models	20.1-3
	Learning undirected graphical models cont.	20.1-3
9	Hidden Markov Models	TBD
	HMMs cont. and Kalman filter	TBD
10	Markov decision processes	TBD

## Application: recommendation systems

- Given user preferences, suggest recommendations
- **Example:** Amazon.com
  
- Input: movie preferences of many users
- Solution: model correlations between movie features
  - Users that like comedy, often like drama →
  - Users that like action, often do not like cartoons →
  - Users that like Robert Deniro films often like Al Pacino films
  - Given user preferences, can predict probability that new movies match preferences

# Diagnostic systems

- Diagnostic indexing for home health site at microsoft
- Enter symptoms → recommend multimedia content

**Describe the child**  
in the drop-down boxes at the right. Relevant information will appear below.

Age:  Sex:

Complaint:

---

Localized pain: Can the child localize, or point to, the site of the pain?

No, unable to localize

Below the navel to the child's left

Above the child's navel

Either of the child's sides

Below the navel to the child's right

Above the navel to the child's right

Above the navel to the child's left

Don't Know

**Results so far**

Disorder	Relevance
Viral gastroenteritis	<div style="width: 80%; height: 10px; background-color: red;"></div>
Psychosomatic pain	<div style="width: 70%; height: 10px; background-color: red;"></div>
Urinary tract infection	<div style="width: 20%; height: 10px; background-color: red;"></div>
Other	<div style="width: 10%; height: 10px; background-color: red;"></div>

# Many research areas in CS

- Full of tasks that require reasoning under uncertainty

**Speech recognition**

**HMM**

**computer vision**

**Tracking and robot localization**

**Kalman filter**

[Fox et al]

**Modeling sensor data**

[Barnard et al]

**evolutionary biology**

**Bayesian network**

[Friedman et al]

**Undirected graphical model**

[Guestrin et al]

**planning under uncertainty**

**Dynamic Bayesian network**

[Guestrin et al]

**medical diagnosis**



## Enjoy!

- Probabilistic graphical models are having significant impact in science, engineering and beyond
- This class should give you the basic foundation for applying PGMs and developing new methods
- The fun begins ...

## Today

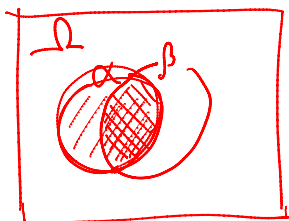
- **Basics of probability**
  - Conditional probabilities
  - Statistical independence
  - Random variable
- Simple Bayesian networks ←
  - Two nodes make a BN
  - Naïve Bayes
- Should be a review for everyone – Setting up notation for the class

## Sample spaces, events and probabilities

- Probability
  - A degree of confidence that an "event" of an uncertain nature will occur.
- Begin with a set  $\Omega$  -- the sample space
  - Space of possible outcomes
  - e.g. if we consider dice, we might have a set  $\Omega = \{1, 2, 3, 4, 5, 6\}$
  - $\alpha \in \Omega$  is a sample point / atomic event.
- A probability space is a sample space with an assignment  $P(\alpha)$  for every  $\alpha \in \Omega$  s.t.
  - $0 \leq P(\alpha) \leq 1$
  - $\sum_{\alpha} P(\alpha) = 1$
  - e.g.  $P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = 1/6$
- An event  $A$  is any subset of  $\Omega$ 
  - $P(A) = \sum_{\{\alpha \in A\}} P(\alpha)$
- E.g.,  $P(\text{die roll} < 4) = P(1) + P(2) + P(3) = 0.5$

## Conditional probabilities

- Consider two events  $\alpha$  and  $\beta$ ,
  - e.g.  $\alpha$  = getting admitted to the UW CSE,  $\beta$  = getting a job offer from Microsoft.
- After learning that  $\alpha$  is true, how do we feel about  $\beta$ ?
  - $P(\beta|\alpha)$   ~~$P(\beta)$~~  ?



$$P(\beta|\alpha) = \frac{P(\alpha \cap \beta)}{P(\alpha)}$$

## Two of the most important rules of the quarter: 1. The chain rule

- From the definition of the conditional distribution, we immediately see that

- $$P(\alpha \cap \beta) = P(\alpha)P(\beta|\alpha)$$

$$P(\beta|\alpha) = \frac{P(\alpha \cap \beta)}{P(\alpha)}$$

- More generally:  $\alpha_1, \dots, \alpha_k \quad (k > 2)$

- $$P(\alpha_1 \cap \dots \cap \alpha_k) = P(\alpha_1)P(\alpha_2|\alpha_1) \dots P(\alpha_k|\alpha_1 \cap \dots \cap \alpha_{k-1})$$



## Two of the most important rules of the quarter: 2. Bayes rule

- Another immediate consequence of the definition of conditional probability is:

$$P(\alpha|\beta) = \frac{P(\alpha \cap \beta)}{P(\beta)} = \frac{P(\alpha)P(\beta|\alpha)}{P(\beta)}$$

"invert"

- A more general version of Bayes' rule, where all the probabilities are conditioned on some "background" event  $\gamma$ :

$$P(\alpha|\beta \cap \gamma) = \frac{P(\beta|\alpha \cap \gamma)P(\alpha|\gamma)}{P(\beta|\gamma)}$$

## Most important concept of the quarter: a) Independence

- $\alpha$  and  $\beta$  are **independent**, if  $P(\beta|\alpha) = P(\beta)$ 
  - Denoted  $P \models (\alpha \perp \beta)$
- **Proposition:**  $\alpha$  and  $\beta$  are **independent** if and only if  $P(\alpha \cap \beta) = P(\alpha)P(\beta)$

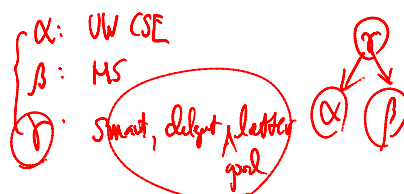
$$P(\alpha \cap \beta) = P(\alpha) \cdot P(\beta|\alpha) = P(\alpha)P(\beta)$$

## Most important concept of the quarter: b) Conditional independence

- Independence is rarely true, but conditionally...
- $\alpha$  and  $\beta$  **conditionally independent** given  $\gamma$  if

$$P(\beta|\alpha \cap \gamma) = P(\beta|\gamma)$$

- $P \models (\alpha \perp \beta | \gamma)$



**Proposition:**  $P \models (\alpha \perp \beta | \gamma)$  if and only if  $P(\alpha \cap \beta | \gamma) = P(\alpha | \gamma)P(\beta | \gamma)$

## Random variables

- Probability distributions are defined for events
  - Events are complicated – so, let's think about attributes
    - Age, Grade, HairColor
  - A random variable (such as *Grade*), is defined by a function that associates each outcome in  $\Omega$  (each person) with a value.
    - $\text{Grade} = A$  – shorthand for event  $\{w \in \Omega: f_{\text{Grade}}(w) = A\}$
    - $\text{Grade} = B$  – shorthand for event  $\{w \in \Omega: f_{\text{Grade}}(w) = B\}$
    - ⋮
  - Properties of a random variable  $X$ :
    - $\text{Val}(X)$  = a set of possible values of random variable  $X$
    - For discrete (categorical):  $\sum_{i=1, \dots, |\text{Val}(X)|} P(X=x_i) = 1$
    - $P(X) \geq 0$
- Handwritten notes:* Job Offer, Grade A,  $X: \text{Grade}$ ,  $\text{Val}(X) = \{A, B, C, \dots, F\}$ ,  $P(X=A)$

## Basic concepts for random variables

- Atomic event: assignment  $x_1, \dots, x_n$  to  $X_1, \dots, X_n$
  - Conditional probability:  $P(Y|X) = P(X, Y) / P(X)$ 
    - For all values  $x \in \text{Val}(X), y \in \text{Val}(Y)$
  - Bayes rule:  $P(X|Y) = \frac{P(X)P(Y|X)}{P(Y)}$ 
    - $\frac{P(Y|X)}{P(X)}$
  - Chain rule:
    - $P(X_1, \dots, X_n) = P(X_1)P(X_2|X_1) \dots P(X_n|X_1, \dots, X_{n-1})$
    - $P(X, Y) = P(X)P(Y|X)$
    - $P(\text{Grade} = A, \text{Intell} = H) = P(\text{Grade} = A | \text{Intell} = H) P(\text{Intell} = H)$
    - $P(\text{Grade} = A | \text{Intell} = L)$
- Handwritten notes:* A Grade,  $P(X, Y)$ ,  $P(Y)$ ,  $P(X \cap Y)$ ,  $P(\text{Grade} | \text{Intell})$ ,  $P(\text{Grade} = A | \text{Intell} = H)$ ,  $P(\text{Grade} = A | \text{Intell} = L)$ ,  $P(X)P(Y|X)$ ,  $P(X \cap Y)$ ,  $0.5$

## Joint distribution, marginalization

- Two random variables – Grade & Intelligence

$$G: \{a, b\}$$

$$I: \{vh, h\}$$

$$P(G, I)$$

$G \backslash I$	$vh$	$h$
$a$	0.7	0.1
$b$	0.15	0.05

$\rightarrow P(G=a, I=h)$  (points to 0.1)  
 $\downarrow P(G=a, I=vh)$  (points to 0.7)

- Marginalization** – Compute marginal over single variable

$$P(G, I)$$

$$\downarrow$$

$$P(G) \text{ or } P(I)$$

$$P(G=b) = P(G=b, I=vh) + P(G=b, I=h)$$

$$P(G=a) = 1 - P(G=b)$$



## Marginalization – the general case

- Compute marginal distribution  $P(X_i)$  from joint distribution  $P(X_1, \dots, X_i, \dots, X_n)$ :

$$P(X_1, X_2, \dots, X_i) = \sum_{x_{i+1}, \dots, x_n} P(X_1, X_2, \dots, X_i, x_{i+1}, \dots, x_n)$$

$$P(X_i) = \sum_{x_1, \dots, x_{i-1}} P(x_1, \dots, x_{i-1}, X_i)$$

$\sum \sum \sum \sum$   
 $x_1 \ x_2 \ \dots \ x_n$   
 $2^{n-1}$

## Today

- Basics of probability
  - Conditional probabilities
  - Statistical independence
  - Random variable
- **Two nodes make a BN**
- Naïve Bayes
  
- Should be a review for everyone – Setting up notation for the class

## Representing joint distributions

- Random variables:  $X_1, \dots, X_n$
- $P$  is a joint distribution over  $X_1, \dots, X_n$



If  $X_1, \dots, X_n$  binary, need  $2^n$  parameters to describe  $P$

$p(x_1, \dots, x_n)$   
↑↑↑↑↑

$x: \in \{T, F\}$

Can we represent  $P$  more compactly?

- Key: Exploit independence properties

## Independent random variables

- If  $X_1, \dots, X_n$  are independent then:
  - $P(X_1, \dots, X_n) = P(X_1) \dots P(X_n)$
  - $O(n)$  parameters
  - All  $2^n$  probabilities are implicitly defined
  - Cannot represent many types of distributions
- X and Y are conditionally independent given Z if
  - $P(X=x|Y=y, Z=z) = P(X=x|Z=z)$  for all values  $x, y, z$
  - Equivalently, if we know Z, then knowing Y does not change predictions of X
  - Notation:  $(X \perp Y | Z)$

$X_1, X_2$   
 $P(X_2) = P(X_2|X_1)$   
 $P(X_1, X_2) = P(X_1)P(X_2)$   
 $= P(X_1)P(X_2)$   
 $X \perp Y | Z$

## Conditional parameterization

- S = SAT score,  $\text{Val}(S) = \{s^0, s^1\}$
- I = Intelligence,  $\text{Val}(I) = \{i^0, i^1\}$

$P(I, S)$

I	S	$P(I, S)$
$i^0$	$s^0$	0.665
$i^0$	$s^1$	0.035
$i^1$	$s^0$	0.06
$i^1$	$s^1$	0.24

Joint parameterization

3 parameters

$P(I, S) = P(I)P(S|I)$

P(I)		P(S I)	
I		S	
$i^0$	$i^1$	$s^0$	$s^1$
0.7	0.3	0.95	0.05
		$i^0$	$i^1$
		0.2	0.8

Conditional parameterization

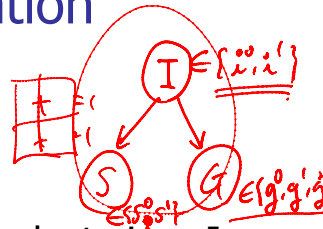
4 parameters

Alternative parameterization:  $P(S)$  and  $P(I|S)$



## Conditional parameterization

- S = SAT score, Val(S) = {s<sup>0</sup>, s<sup>1</sup>}
- I = Intelligence, Val(I) = {i<sup>0</sup>, i<sup>1</sup>}
- G = Grade, Val(G) = {g<sup>0</sup>, g<sup>1</sup>, g<sup>2</sup>}
- Assume that G and S are independent given I



- Joint parameterization

- 2·2·3=12-1=11 independent parameters

- Conditional parameterization has

- P(I, S, G) = P(I)P(S|I)P(G|I, S) = P(I)P(S|I)P(G|I)
  - P(I) – 1 independent parameter
  - P(S|I) – 2·1 independent parameters
  - P(G|I) – 2·2 independent parameters
  - 7 independent parameters



Handwritten derivations for joint and conditional probabilities:

$$P(I, S, G) = P(I)P(S|I)P(G|I, S)$$

$$= P(I)P(S|I)P(G|I)$$

$$= P(I)P(S|I)P(G|I)$$

Annotations: "PCS, I, G", "2x2", "2x2", "+1", "+2", "+1"

## Naïve Bayes model

- Class variable C, Val(C) = {c<sub>1</sub>, ..., c<sub>k</sub>}
- Evidence variables X<sub>1</sub>, ..., X<sub>n</sub>
- Naïve Bayes assumption: evidence variables are conditionally independent given C

$$P(C, X_1, \dots, X_n) = P(C) \prod_{i=1}^n P(X_i | C)$$

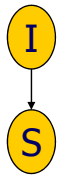
- Applications in medical diagnosis, text classification
- Used as a classifier:

$$\frac{P(C = c_1 | x_1, \dots, x_n)}{P(C = c_2 | x_1, \dots, x_n)} = \frac{P(C = c_1)}{P(C = c_2)} \prod_{i=1}^n \frac{P(x_i | C = c_1)}{P(x_i | C = c_2)}$$

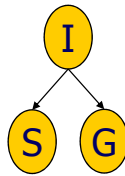
- Problem: Double counting correlated evidence

## Bayesian network (informal)

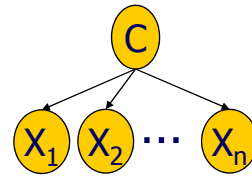
- Directed acyclic graph  $G$ 
  - Nodes represent random variables
  - Edges represent direct influences between random variables
- Local probability models



Example 1



Example 2



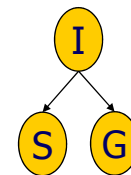
Naïve Bayes

CSE 515 – Statistical Methods – Spring 2011

35

## Bayesian network (informal)

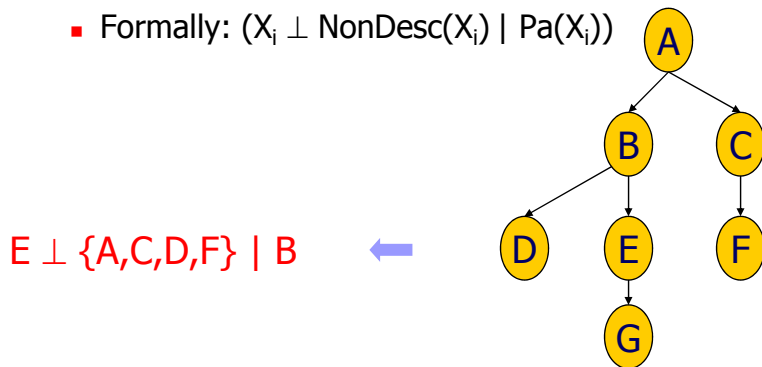
- Represent a joint distribution
  - Specifies the probability for  $P(\mathbf{X}=\mathbf{x})$
  - Specifies the probability for  $P(\mathbf{X}=\mathbf{x}|\mathbf{E}=\mathbf{e})$
- Allows for reasoning patterns
  - **Prediction** (e.g., intelligent  $\rightarrow$  high scores)
  - **Explanation** (e.g., low score  $\rightarrow$  not intelligent)
  - **Explaining away** (different causes for same effect interact)



Example 2

## Bayesian network structure

- Directed acyclic graph  $G$ 
  - Nodes  $X_1, \dots, X_n$  represent random variables
- $G$  encodes local Markov assumptions
  - $X_i$  is independent of its non-descendants given its parents
  - Formally:  $(X_i \perp \text{NonDesc}(X_i) \mid \text{Pa}(X_i))$



CSE 515 – Statistical Methods – Spring 2011

37

## Independency mappings (I-maps)

- Let  $P$  be a distribution over  $\mathbf{X}$
- Let  $I(P)$  be the independencies  $(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$  in  $P$
- A Bayesian network structure is an I-map (independency mapping) of  $P$  if  $I(G) \subseteq I(P)$

I

S

$I$	$S$	$P(I, S)$
$i^0$	$s^0$	0.25
$i^0$	$s^1$	0.25
$i^1$	$s^0$	0.25
$i^1$	$s^1$	0.25

$I(G) = \{I \perp S\}$      $I(P) = \{I \perp S\}$

$I$	$S$	$P(I, S)$
$i^0$	$s^0$	0.4
$i^0$	$s^1$	0.3
$i^1$	$s^0$	0.2
$i^1$	$s^1$	0.1

I

S

$I(P) = \emptyset$      $I(G) = \emptyset$

CSE 515 – Statistical Methods – Spring 2011

38

## Factorization Theorem

- If  $G$  is an I-Map of  $P$ , then  $P$  factorizes over  $G$ .

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa(X_i))$$

**Proof:**

- wlog. (without loss of generality)  
 $X_1, \dots, X_n$  is an ordering consistent with  $G$
- By chain rule:  $P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})$
- From assumption:  $Pa(X_i) \subseteq \{X_1, \dots, X_{i-1}\}$   
 $\{X_1, \dots, X_{i-1}\} - Pa(X_i) \subseteq NonDesc(X_i)$
- Since  $G$  is an I-Map  $\rightarrow (X_i; NonDesc(X_i) | Pa(X_i)) \in I(P)$



$$P(X_i | X_1, \dots, X_{i-1}) = P(X_i | Pa(X_i))$$

CSE 515 – Statistical Methods – Spring 2011

39

## Factorization implies I-Map

- $P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa(X_i)) \rightarrow G$  is an I-Map of  $P$

**Proof:**

- Need to show  $(X_i; NonDesc(X_i) | Pa(X_i)) \in I(P)$  or that  
 $P(X_i | NonDesc(X_i)) = P(X_i | Pa(X_i))$
- wlog.  $X_1, \dots, X_n$  is an ordering consistent with  $G$

$$\begin{aligned} P(X_i | NonDesc(X_i)) &= \frac{P(X_i, NonDesc(X_i))}{P(NonDesc(X_i))} \\ &= \frac{\prod_{k=1}^i P(X_k | Pa(X_k))}{\prod_{k=1}^{i-1} P(X_k | Pa(X_k))} \\ &= P(X_i | Pa(X_i)) \end{aligned}$$

CSE 515 – Statistical Methods – Spring 2011

40

## Bayesian network definition

- A Bayesian network is a pair  $(G,P)$ 
  - $P$  factorizes over  $G$
  - $P$  is specified as set of CPDs associated with  $G$ 's nodes (and its parents)
- Parameters
  - Joint distribution:  $2^n$
  - Bayesian network (bounded in-degree  $k$ ):  $n2^k$

## Today and next class

- Next class
  - Details on semantics of BNs, relate them to independence assumptions encoded by the graph.
- Today's To-Do List
  - Visit the course website.
  - Reading K&F 2.1-3, 3.1.

## Acknowledgement

- These lecture notes were generated based on the slides from Profs Eran Segal and Carlos Guestrin.