

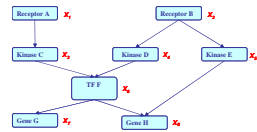
# Representation of directed GM

## Probabilistic Graphical Models (10-708)

Lecture 1, Sep 12, 2007

Eric Xing

Reading: MJ-Chap 2, KF-Chap. 3



1

- Recitation?
- Exam dates, poster dates, etc.
- Mailing list
- Questions?



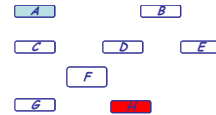
# Representing Multivariate Distribution



- Representation: what is the joint probability dist. on multiple variables?

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8,)$$

- How many state configurations in total? ---  $2^8$
- Are they all needed to be represented?
- Do we get any scientific/medical insight?



- Factored representation: the chain-rule

$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ &= P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2)P(X_4 | X_1, X_2, X_3)P(X_5 | X_1, X_2, X_3, X_4)P(X_6 | X_1, X_2, X_3, X_4, X_5) \\ &P(X_7 | X_1, X_2, X_3, X_4, X_5, X_6)P(X_8 | X_1, X_2, X_3, X_4, X_5, X_6, X_7) \end{aligned}$$

- This factorization is true for any distribution and any variable ordering
- Do we save any parameterization cost?

- If  $X_i$ 's are **independent**: ( $P(X_i | \cdot) = P(X_i)$ )

$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ &= P(X_1)P(X_2)P(X_3)P(X_4)P(X_5)P(X_6)P(X_7)P(X_8) = \prod_i P(X_i) \end{aligned}$$

- What do we gain?
- What do we lose?

Eric Xing

3

- Even in the simplest case where these variables are binary-valued, a joint distribution requires the specification of  $2^n$  numbers — the probabilities of the  $2^n$  different assignments of values  $x_1, \dots, x_n$

- Today's lecture is about ...

- how independence properties in the distribution can be used to represent such high-dimensional distributions much more compactly.
- how a combinatorial data structure — a directed acyclic graph — can provide us with a general-purpose modeling language for exploiting this type of structure in our representation.



Eric Xing

4

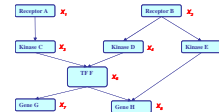
## Two types of GMs

- Directed edges give causality relationships (Bayesian Network or Directed Graphical Model):

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$= P(X_1) P(X_2) P(X_3/X_1) P(X_4/X_2) P(X_5/X_2)$$

$$P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)$$

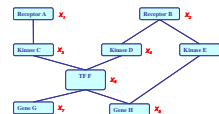


- Undirected edges simply give correlations between variables (Markov Random Field or Undirected Graphical model):

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$= \frac{1}{Z} \exp\{E(X_1)+E(X_2)+E(X_3, X_1)+E(X_4, X_2)+E(X_5, X_2)$$

$$+E(X_6, X_3, X_4)+E(X_7, X_6)+E(X_8, X_5, X_6)\}$$

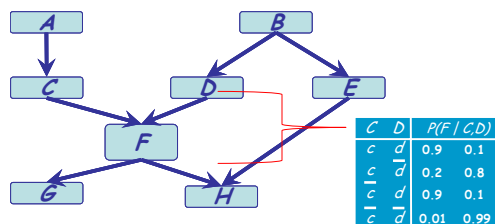


Eric Xing

5

## Specification of a directed GM

- There are two components to any GM:
  - the qualitative specification
  - the quantitative specification



Eric Xing

6

## Bayesian Network:

- A BN is a directed graph whose nodes represent the random variables and whose edges represent direct influence of one variable on another.
- It is a data structure that provides the skeleton for representing a **joint distribution** compactly in a **factorized** way;
- It offers a compact representation for a **set of conditional independence assumptions** about a distribution;
- We can view the graph as encoding a **generative sampling process** executed by nature, where the value for each variable is selected by nature using a distribution that depends only on its parents. In other words, each variable is a stochastic function of its parents.



$$P(x_1, \dots, x_n)$$

Eric Xing

7

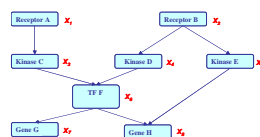
## Bayesian Network: Factorization Theorem

- Theorem:**

Given a DAG, The most general form of the probability distribution that is **consistent with** the graph factors according to “node given its parents”:

$$P(\mathbf{X}) = \prod_{i=1:d} P(X_i | \mathbf{X}_{\pi_i})$$

where  $\mathbf{X}_{\pi_i}$  is the set of parents of  $X_i$ ,  $d$  is the number of nodes (variables) in the graph.



$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ &= P(X_1) P(X_2) P(X_3 | X_1) P(X_4 | X_2) P(X_5 | X_2) \\ &\quad P(X_6 | X_3, X_4) P(X_7 | X_6) P(X_8 | X_5, X_6) \end{aligned}$$

Eric Xing

8

# Qualitative Specification



- Where does the qualitative specification come from?
  - Prior knowledge of causal relationships
  - Prior knowledge of modular relationships
  - Assessment from experts
  - Learning from data
  - We simply link a certain architecture (e.g. a layered graph)
  - ...

Eric Xing

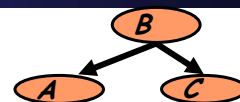
9

# Local Structures & Independencies



- Common parent
  - Fixing B decouples A and C

"given the level of gene B, the levels of A and C are independent"



- Cascade
  - Knowing B decouples A and C

"given the level of gene B, the level gene A provides no extra prediction value for the level of gene C"



- V-structure
  - Knowing C couples A and B
  - because A can "explain away" B w.r.t. C

"If A correlates to C, then chance for B to also correlate to B will decrease"



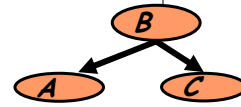
- The language is compact, the concepts are rich!

Eric Xing

10

## A simple justification

$$A \perp C | B.$$



$$\begin{aligned}
 P(A, C | B) &\neq P(A|B)P(C|B) \\
 \uparrow \\
 \frac{P(A, B, C)}{P(B)} &= \frac{P(B) \times P(A|B) \times P(C|B)}{P(B)} \\
 &= P(A|B)P(C|B)
 \end{aligned}$$

Eric Xing

11

## I-maps

- **Defn (3.2.2):** Let  $P$  be a distribution over  $\mathbf{X}$ . We define  $I(P)$  to be the set of independence assertions of the form  $(X \perp Y | Z)$  that hold in  $P$  (however how we set the parameter-values).

$$P. \rightarrow I_P$$

- **Defn (3.2.3):** Let  $K$  be any graph object associated with a set of independencies  $I(K)$ . We say that  $K$  is an **I-map** for a set of independencies  $I$ ,  $I(K) \subseteq I$ .

$$K \rightarrow I(K)$$

$$I = I_P \quad \text{with } I(K) \not\subseteq I_P$$

- We now say that  $G$  is an I-map for  $P$  if  $G$  is an I-map for  $I(P)$ , where we use  $I(G)$  as the set of independencies associated.

Eric Xing

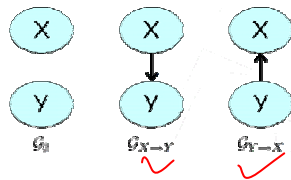
12

## Facts about I-map

- For  $G$  to be an I-map of  $P$ , it is necessary that  $G$  does not mislead us regarding independencies in  $P$ :

any independence that  $G$  asserts must also hold in  $P$ . Conversely,  $P$  may have additional dependencies that are not reflected in  $G$

- Example:



$$P_1$$

$X$	$Y$	$P(X, Y)$
$x^0$	$y^0$	0.08
$x^0$	$y^1$	0.32
$x^1$	$y^0$	0.12
$x^1$	$y^1$	0.48

$$P_2$$

$X$	$Y$	$P(X, Y)$
$x^0$	$y^0$	0.4
$x^0$	$y^1$	0.3
$x^1$	$y^0$	0.2
$x^1$	$y^1$	0.1

$$P(X) = \begin{matrix} 0.4 & 0 \\ 0.6 & 1 \end{matrix}$$

$$P(Y) = \begin{matrix} 0.2 & 0 \\ 0.8 & 1 \end{matrix}$$

$$P(X \neq Y = 1) = 0.48$$

$$= P(X=1)P(Y=1)$$

Eric Xing

13

## What is in $I(G)$ --- local Markov assumptions of BN

A Bayesian network structure  $G$  is a directed acyclic graph whose nodes represent random variables  $X_1, \dots, X_n$ .

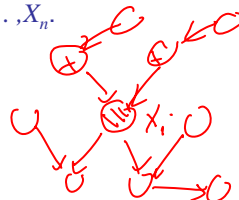
### local Markov assumptions

- Defn (3.2.1):

Let  $Pa_{X_i}$  denote the parents of  $X_i$  in  $G$ , and  $NonDescendants_{X_i}$  denote the variables in the graph that are not descendants of  $X_i$ . Then  $G$  encodes the following set of **local conditional independence assumptions**  $I_G(G)$ :

$$I_G(G): \{X_i \perp NonDescendants_{X_i} \mid Pa_{X_i} : \forall i\},$$

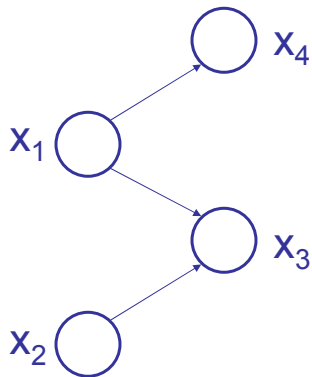
In other words, each node  $X_i$  is independent of its nondescendants given its parents.



Eric Xing

14

## Example



$$I_d(G) = \left\{ \begin{array}{l} x_4 \perp x_3 \mid x_1 \\ x_4 \perp x_2 \mid x_1 \\ x_4 \perp x_2, x_3 \mid x_1 \end{array} \right\}$$

Eric Xing

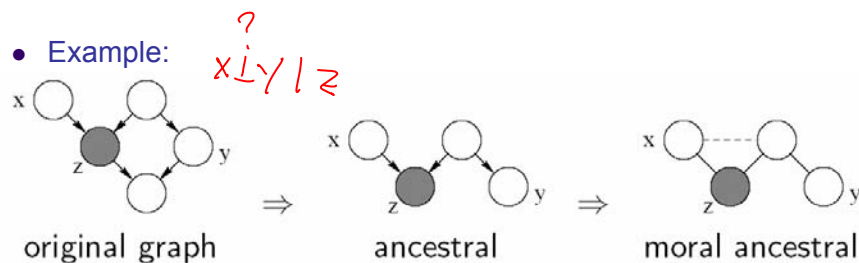
15

## Graph separation criterion

- D-separation criterion for Bayesian networks (D for Directed edges):

**Defn:** variables  $x$  and  $y$  are *D-separated* (conditionally independent) given  $z$  if they are separated in the *moralized* ancestral graph

- Example:

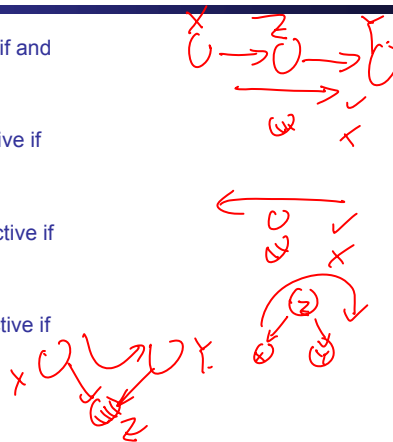


Eric Xing

16

## Active trail

- **Causal trail**  $X \rightarrow Z \rightarrow Y$  : active if and only if  $Z$  is not observed.
- **Evidential trail**  $X \leftarrow Z \leftarrow Y$  : active if and only if  $Z$  is not observed.
- **Common cause**  $X \leftarrow Z \rightarrow Y$  : active if and only if  $Z$  is not observed.
- **Common effect**  $X \rightarrow Z \leftarrow Y$  : active if and only if either  $Z$  or one of  $Z$ 's descendants is observed



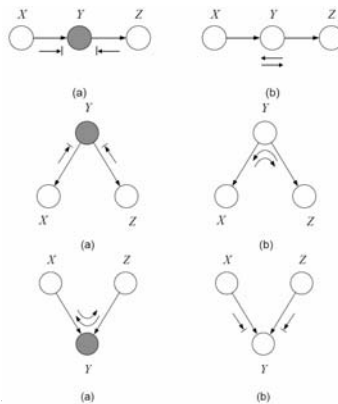
**Definition (3.3.2):** Let  $X, Y, Z$  be three **sets** of nodes in  $G$ . We say that  $X$  and  $Y$  are *d-separated* given  $Z$ , denoted  $d\text{-sep}_G(X; Y | Z)$ , if there is **no** active trail between any node  $X \in X$  and  $Y \in Y$  given  $Z$ .

Eric Xing

17

## What is $I(G)$ --- Global Markov properties of BN

- $X$  is **d-separated** (directed-separated) from  $Z$  given  $Y$  if we can't send a ball from any node in  $X$  to any node in  $Z$  using the "**Bayes-ball**" algorithm illustrated below (and plus some boundary conditions):



- Defn:  $I(G)$  = all independence properties that correspond to d-separation:

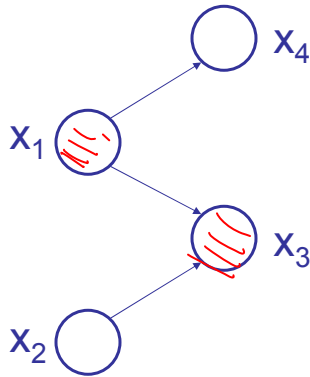
$$I(G) = \{X \perp Z | Y : d\text{-sep}_G(X; Z | Y)\}$$

- D-separation is sound and complete (more details later)

Eric Xing

18

## Example:



- Complete the I(G) of this graph:

~~$X_1 \perp X_2$~~   
 $X_1 \perp X_2$   
 $X_2 \perp X_4$   
 $X_3 \perp X_4 \mid X_1$   
 $X_2 \perp X_4 \mid X_3, X_1$   
 $X_4 \perp X_3 \mid X_1$   
 $X_4 \perp X_2 \mid X_1$   
 $X_4 \perp X_2 X_3 \mid X_1$

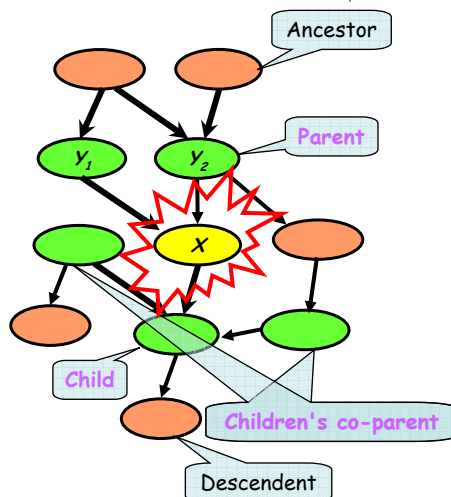
Eric Xing

19

## Summary: Conditional Independence Semantics in an BN

### Structure: **DAG**

- Meaning: a node is **conditionally independent** of every other node in the network outside its **Markov blanket**
- Local conditional distributions (**CPD**) and the **DAG** completely determine the **joint** dist.
- Give **causality** relationships, and facilitate a **generative** process



Eric Xing

20

## Toward quantitative specification of probability distribution



- Separation properties in the graph imply independence properties about the associated variables

- The Equivalence Theorem**

For a graph  $G$ ,

Let  $\mathcal{D}_1$  denote the family of **all distributions** that satisfy  $I(G)$ ,

Let  $\mathcal{D}_2$  denote the family of **all distributions** that factor according to  $G$ ,

$$P(\mathbf{X}) = \prod_{i=1:d} P(X_i | \mathbf{X}_{\pi_i})$$

Then  $\mathcal{D}_1 \equiv \mathcal{D}_2$ .

- For the graph to be useful, any conditional independence properties we can derive from the graph should hold for the probability distribution that the graph represents

*Handwritten notes:*  
 $G \rightarrow I \rightarrow P \in \mathcal{D}_1$   
 $G \rightarrow P \in \mathcal{D}_2$

Eric Xing

21

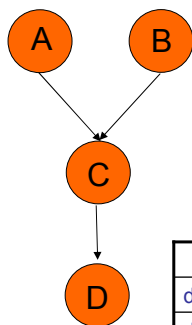
## Conditional probability tables (CPTs)



$a^0$	0.75
$a^1$	0.25

$b^0$	0.33
$b^1$	0.67

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$



	$a^0b^0$	$a^0b^1$	$a^1b^0$	$a^1b^1$
$c^0$	0.45	1	0.9	0.7
$c^1$	0.55	0	0.1	0.3

	$c^0$	$c^1$
$d^0$	0.3	0.5
$d^1$	0.7	0.5

Eric Xing

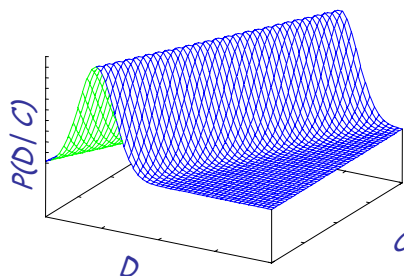
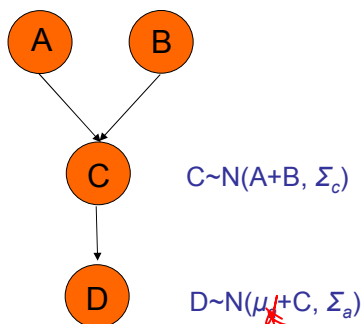
22

## Conditional probability density func. (CPDs)



$$A \sim N(\mu_a, \Sigma_a) \quad B \sim N(\mu_b, \Sigma_b)$$

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$



Eric Xing

23

## Summary of BN semantics



- **Defn (3.2.5):** A *Bayesian network* is a pair  $(G, P)$  where  $P$  factorizes over  $G$ , and where  $P$  is specified as set of CPDs associated with  $G$ 's nodes.
  - Conditional independencies imply factorization
  - Factorization according to  $G$  implies the associated conditional independencies.
  - Are there **other independencies** that hold for every distribution  $P$  that factorizes over  $G$ ?

Eric Xing

24

# Soundness and completeness

D-separation is sound and "complete" w.r.t. BN factorization law

**Soundness:**

**Theorem:** If a distribution  $P$  factorizes according to  $G$ , then  $I(G) \subseteq I(P)$ .

**"Completeness":**

**"Claim":** For any distribution  $P$  that factorizes over  $G$ , if  $(X \perp Y \mid Z) \in I(P)$  then  $d\text{-sep}_G(X; Y \mid Z)$ .

Contrapositive of the completeness statement

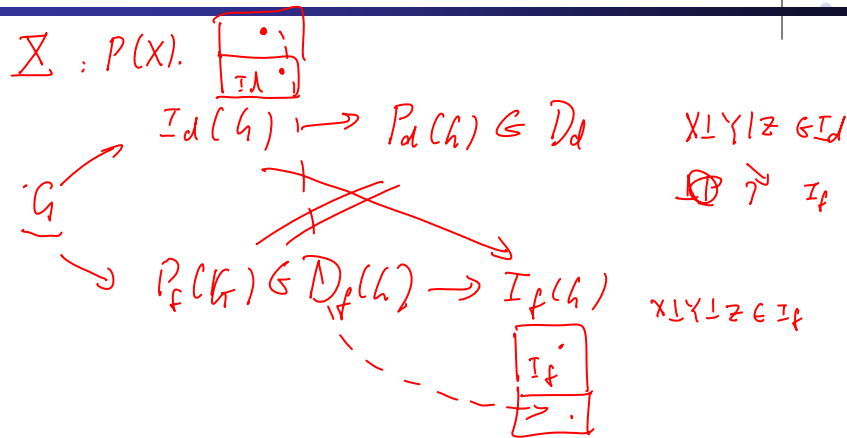
- "If  $X$  and  $Y$  are **not**  $d$ -separated given  $Z$  in  $G$ , then  $X$  and  $Y$  are **dependent** in **all** distributions  $P$  that factorize over  $G$ ."
- Is this true?

• Recitation:

- Wednesday 6-7 pm
- Thursday: 6-7pm
- Friday: 5-6pm

• Questions:

## Distributional equivalence and I-equivalence



- All independence in  $I_d(G)$  will be captured in  $I_f(G)$ , is the reverse true?
- Are "not-independence" from  $G$  all honored in  $P_f$ ?

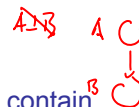
Eric Xing

27

## Soundness and completeness



- Contrapositive of the completeness statement
  - "If  $X$  and  $Y$  are **not d-separated** given  $Z$  in  $G$ , then  $X$  and  $Y$  are **dependent in all** distributions  $P$  that factorize over  $G$ ."
  - Is this true?
- No. Even if a distribution factorizes over  $G$ , it can still contain **additional independencies** that are not reflected in the structure
  - Example: graph  $A \rightarrow B$ , for actually independent  $A$  and  $B$  (the independence can be captured by some subtle way of parameterization)
- Thm: Let  $G$  be a BN graph. If  $X$  and  $Y$  are not d-separated given  $Z$  in  $G$ , then  $X$  and  $Y$  are **dependent in some** distribution  $P$  that factorizes over  $G$ .



$A$	$b^0$	$b^1$
$a^0$	0.4	0.6
$a^1$	0.4	0.6

$$P(A, B) = P(A)P(B|A)$$

Eric Xing

28



- **Theorem 3.3.6:** For **almost all** distributions  $P$  that factorize over  $G$ , i.e., for all distributions except for a set of "measure zero" in the space of CPD parameterizations, we have that  $I(P) = I(G)$

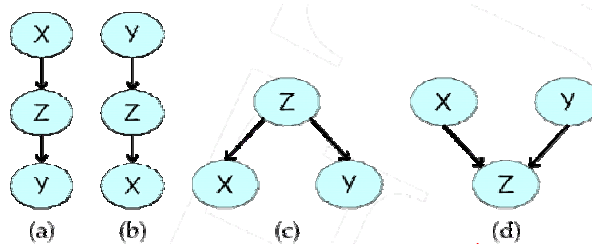
Eric Xing

29



## Uniqueness of BN

- Very different BN graphs can actually be equivalent, in that they encode precisely the same set of conditional independence assertions.



$(X \perp Y | Z).$

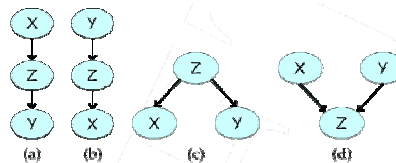
Eric Xing

30

# I-equivalence

- **Defn (3.3.9):** Two BN graphs  $G_1$  and  $G_2$  over  $X$  are *I-equivalent* if  $I(G_1) = I(G_2)$ .

- The set of all graphs over  $X$  is partitioned into a set of mutually exclusive and exhaustive *I-equivalence classes*, which are the set of equivalence classes induced by the I-equivalence relation.



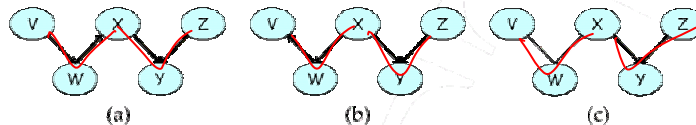
- Any distribution  $P$  that can be factorized over one of these graphs can be factorized over the other.
- Furthermore, there is no intrinsic property of  $P$  that would allow us associate it with one graph rather than an equivalent one.
- This observation has important implications with respect to our ability to determine the directionality of influence.

Eric Xing

31

# Detecting I-equivalence

- **Defn (3.3.10):** The *skeleton* of a Bayesian network graph  $G$  over  $V$  is an undirected graph over  $V$  that contains an edge  $\{X, Y\}$  for every edge  $(X, Y)$  in  $G$ .



- **Thm (3.3.11):** Let  $G_1$  and  $G_2$  be two graphs over  $V$ . If  $G_1$  and  $G_2$  have the same skeleton and the same set of v-structures then they are I-equivalent.

- graph equivalence
- Same trail
- But not necessarily active

Eric Xing

32

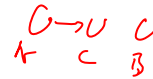
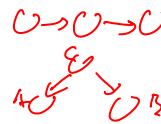
## Minimum I-MAP

$$I(\mathcal{H}_0) = \emptyset$$

- Complete graph is a (trivial) I-map for any distribution, yet it does not reveal any of the independence structure in the distribution.
  - Meaning that the graph dependence is arbitrary, thus by careful parameterization an dependencies can be captured
  - We want a graph that has the maximum possible  $I(G)$ , yet still  $\subseteq I(P)$
- **Defn 3.4.1:** A graph object  $G$  is a *minimal I-map* for a set of independencies  $I$  if it is an I-map for  $I$ , and if the removal of even a single edge from  $G$  renders it not an I-map.

$$A \perp B \mid C$$

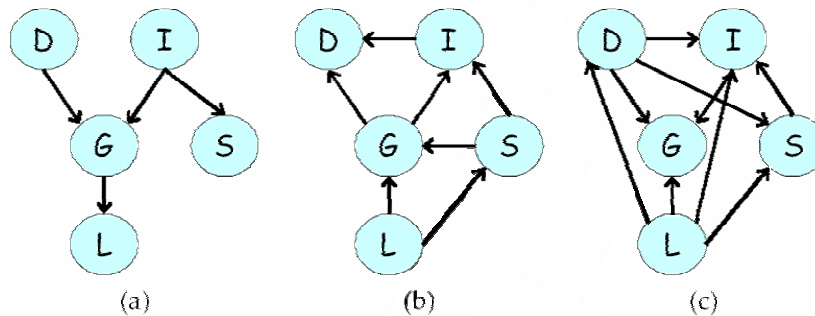
$$I(-)$$



Eric Xing

33

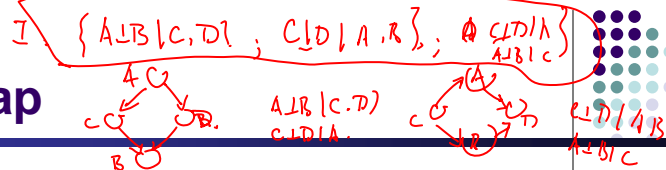
## Minimum I-MAP is not unique



Eric Xing

34

## Perfect Map

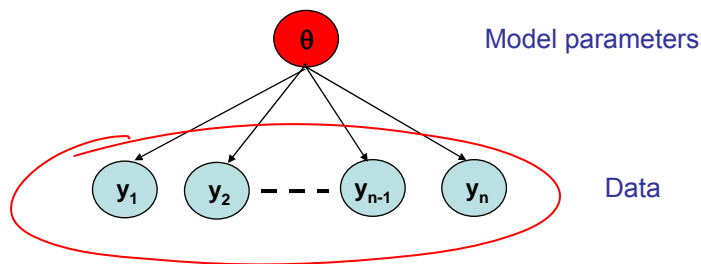


- **Defn (3.4.3):** We say that a graph object  $G$  is a *perfect map* (*P-map*) for a set of independencies  $I$  if we have that  $I(G) = I$ . We say that  $G$  is a perfect map for  $P$  if  $I(G) = I(P)$ .
- The fact that  $G$  is a minimal I-map for  $P$  is far from a guarantee that  $G$  captures the independence structure in  $P$
- Not all  $P$  has a perfect map as DAG!
- The P-map of a distribution is *unique up to I-equivalence* between networks. That is, a distribution  $P$  can have many P-maps, but all of them are I-equivalent.

Eric Xing

35

## Conditionally Independent Observations



Eric Xing

36

## “Plate” Notation

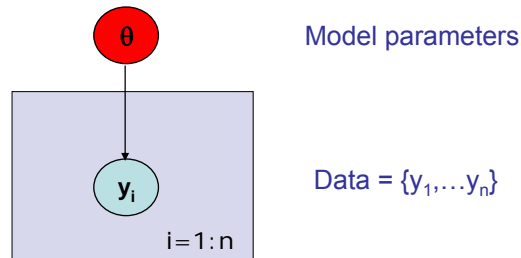


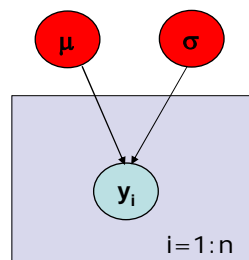
Plate = rectangle in graphical model

variables within a plate are replicated  
in a conditionally independent manner

Eric Xing

37

## Example: Gaussian Model



Generative model:

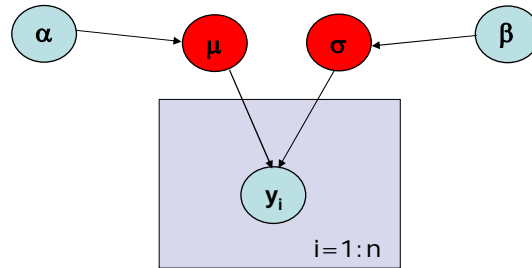
$$\begin{aligned} p(y_1, \dots, y_n \mid \mu, \sigma) &= \prod_i p(y_i \mid \mu, \sigma) \\ &= p(\text{data} \mid \text{parameters}) \\ &= p(D \mid \theta) \\ \text{where } \theta &= \{\mu, \sigma\} \end{aligned}$$

- Likelihood =  $p(\text{data} \mid \text{parameters})$   
 $= p(D \mid \theta)$   
 $= L(\theta)$
- Likelihood tells us how likely the observed data are conditioned on a particular setting of the parameters
  - Often easier to work with  $\log L(\theta)$

Eric Xing

38

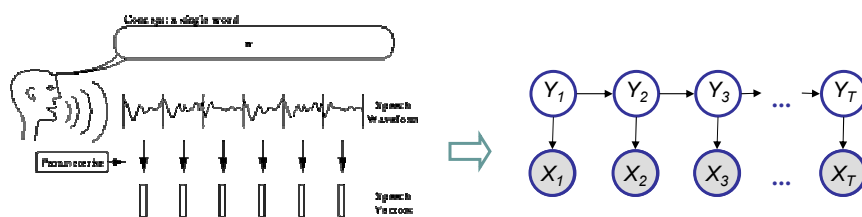
## Example: Bayesian Gaussian Model



Note: priors and parameters are assumed independent here

## Example

- Speech recognition



Hidden Markov Model

# Knowledge Engineering

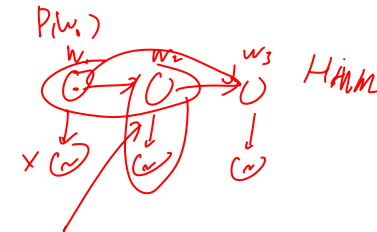
- Picking variables

- Observed
- Hidden

*Handwritten notes:*  
 $z$  (under Observed)  
 $o$  (under Hidden)  
 $am$  (under Hidden)

- Picking structure

- CAUSAL
- Generative



- Picking Probabilities

- Zero probabilities
- Orders of magnitudes
- Relative values

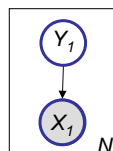
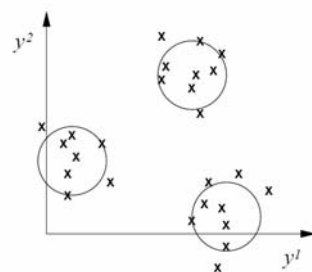
*Handwritten equations:*  
 $P(w_i)$   
 $P(x_i | w_i) \forall i$   
 $P(w_i | w_{i-1}) \forall i$

Eric Xing

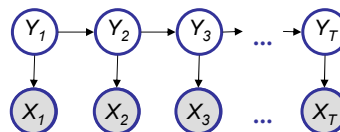
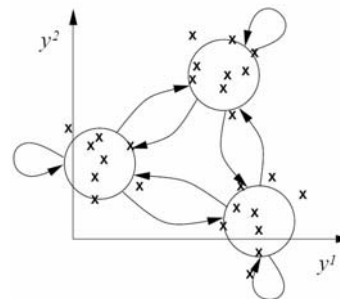
41

## Hidden Markov Model: from static to dynamic mixture models

Static mixture



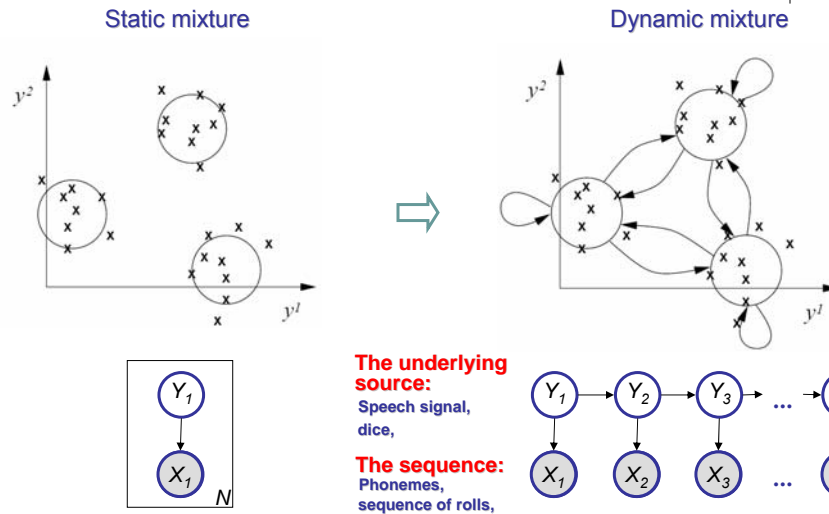
Dynamic mixture



Eric Xing

42

# Hidden Markov Model: from static to dynamic mixture models



Eric Xing

43

## The Dishonest Casino

U-U-U



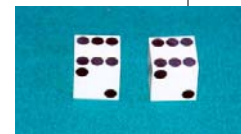
A casino has two dice:

- Fair die  
 $P(1) = P(2) = P(3) = P(5) = P(6) = 1/6$
- Loaded die  
 $P(1) = P(2) = P(3) = P(5) = 1/10$   
 $P(6) = 1/2$

Casino player switches back-&-forth  
between fair and loaded die once every  
20 turns

### Game:

1. You bet \$1
2. You roll (always with a fair die)
3. Casino player rolls (maybe with fair die,  
maybe with loaded die)
4. Highest number wins \$2

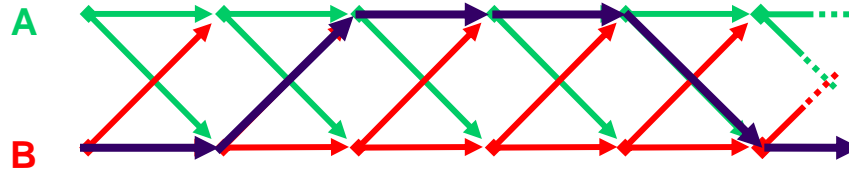
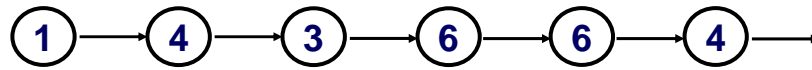


Eric Xing

44

## A stochastic generative model

- Observed sequence:



- Hidden sequence (a parse or segmentation):



Eric Xing

45

## Definition (of HMM)

- Observation space

Alphabetic set:  $C = \{c_1, c_2, \dots, c_K\}$

Euclidean space:  $\mathbb{R}^d$

- Index set of hidden states

$$I = \{1, 2, \dots, M\}$$

- 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

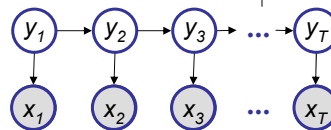
$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

- 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:

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



Eric Xing

46

# Puzzles regarding the dishonest casino



**GIVEN:** A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344  
101...

## QUESTION

- How likely is this sequence, given our model of how the casino works?
  - This is the **EVALUATION** problem in HMMs
- What portion of the sequence was generated with the fair die, and what portion with the loaded die?
  - This is the **DECODING** question in HMMs
- How “loaded” is the loaded die? How “fair” is the fair die? How often does the casino player change from fair to loaded, and back?
  - This is the **LEARNING** question in HMMs

$$P(x_i) \quad P(x_i | x_{i-1}) \\ P(\eta_i | x_i)$$

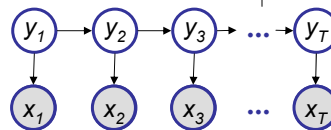
Eric Xing

47

# Probability of a parse



- Given a sequence  $\mathbf{x} = x_1, \dots, x_T$  and a parse  $\mathbf{y} = y_1, \dots, y_T$ ,
- To find how likely is the parse: (given our HMM and the sequence)



$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= p(x_1, \dots, x_T, y_1, \dots, y_T) \quad (\text{Joint probability}) \\ &= p(y_1) p(x_1 | y_1) p(y_2 | y_1) p(x_2 | y_2) \dots p(y_T | y_{T-1}) p(x_T | y_T) \\ &= p(y_1) p(y_2 | y_1) \dots p(y_T | y_{T-1}) \times p(x_1 | y_1) p(x_2 | y_2) \dots p(x_T | y_T) \\ &= p(y_1, \dots, y_T) p(x_1, \dots, x_T | y_1, \dots, y_T) \end{aligned}$$

$$\begin{aligned} \text{Let } \pi_{y_1} &\stackrel{\text{def}}{=} \prod_{i=1}^M [\pi_i^{y_1}] \quad a_{y_t, y_{t+1}} \stackrel{\text{def}}{=} \prod_{i,j=1}^M [a_{ij}]^{y_t y_{t+1}} \quad \text{and } b_{y_t, x_t} \stackrel{\text{def}}{=} \prod_{i=1}^M \prod_{k=1}^K [b_{ik}]^{y_t x_t^k} \\ &= \pi_{y_1} a_{y_1, y_2} \dots a_{y_{T-1}, y_T} b_{y_1, x_1} \dots b_{y_T, x_T} \end{aligned}$$

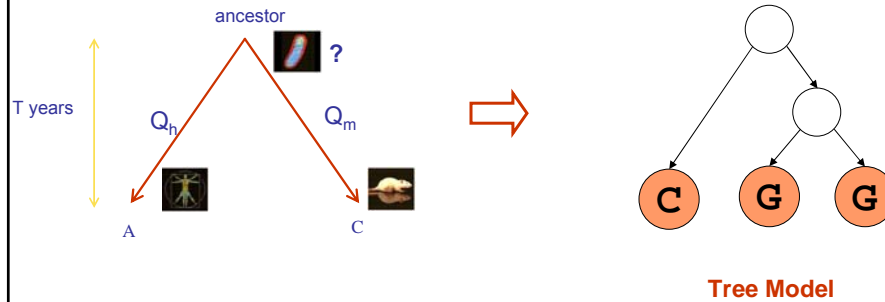
- Marginal probability:**  $p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \dots \sum_{y_T} \pi_{y_1} \prod_{t=2}^T a_{y_{t-1}, y_t} \prod_{t=1}^T p(x_t | y_t)$
- Posterior probability:**  $p(\mathbf{y} | \mathbf{x}) = p(\mathbf{x}, \mathbf{y}) / p(\mathbf{x})$

Eric Xing

48

## Example, con'd

- Evolution

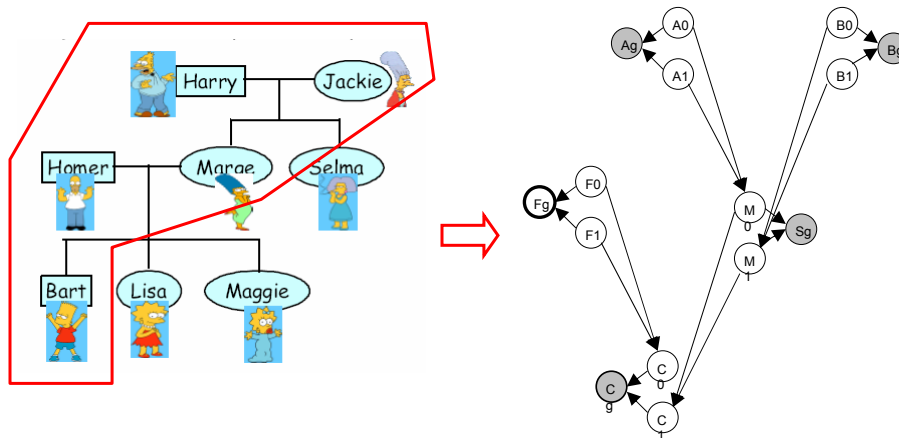


Eric Xing

49

## Example, con'd

- Genetic Pedigree



Eric Xing

50

## Summary of BN semantics



- **Defn (3.2.5):** A *Bayesian network* is a pair  $(G, P)$  where  $P$  factorizes over  $G$ , and where  $P$  is specified as set of CPDs associated with  $G$ 's nodes.