

# ECE 175B: Probabilistic Reasoning and Graphical Models

## Lecture 12: Model Requirements and Simplifications

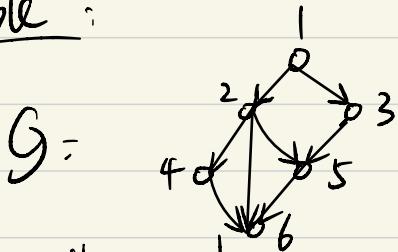
Florian Meyer & Ken Kreutz-Delgado

*Electrical and Computer Engineering Department  
University of California San Diego*

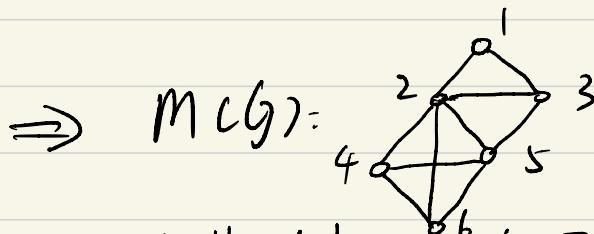
# Ancestral Graph Method for Testing CI statements in BNs.

- We simplify notation with  $\perp\!\!\!\perp_G$ :  $X \perp\!\!\!\perp_G Y | Z \triangleq \langle X | Z | Y \rangle_G$   
Then  $i \perp\!\!\!\perp_G j | k \Rightarrow X_i \perp\!\!\!\perp X_j | X_k$ .
- Testing for graph separation in MNs is easier than in BNs.
- However given a candidate CI statement  $X \perp\!\!\!\perp_G Y | Z$  for a BN  $G$ , we can't test if it is true in the moral graph  $M(G)$  (Recall that  $\perp\!\!\!\perp(M(G)) \subset \perp\!\!\!\perp(G)$ ) because CI statements involving collider nodes are lost during moralization.

Example :



" $4 \perp\!\!\!\perp_G 5 | 2$ "  $\in \perp\!\!\!\perp(G)$



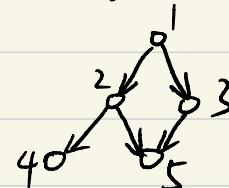
" $4 \perp\!\!\!\perp_G 5 | 2$ "  $\notin \perp\!\!\!\perp(M(G))$

- Note that it is the "downstream" collider node "6" that causes the problem because it results in the moralization link "4-5".

## Ancestral Graph Method - Cont'd.

- For  $\alpha = "X \perp\!\!\!\perp_{\text{G}} Y \mid Z"$ ; define  $U = X \cup Y \cup Z$ .  
For our example :  $U = \{2, 4, 5\}$ .
- Define  $Av(G) = U \cup \underbrace{\text{Ancestor}(U)}_{\text{nodes upstream of } U} = \text{Ancestral Graph}$
- We focus on the subgraph  $Av(G)$ , meaning that we throw away all nodes of  $G$  not in  $Av(G)$  as those nodes (including the downstream nodes) are irrelevant for ascertaining if  $\alpha \in I(G)$ .
- I.e., only the nodes in  $Av(G)$  matter for ascertaining if  $\alpha \in I(G)$ .
- For our example,

$Av(G) :$

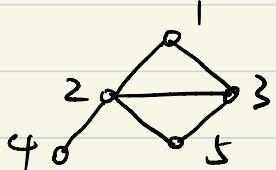


- It is the case that  $I(Av(G)) \subset I(G)$ , i.e.,  $Av(G)$  is faithful to  $G$ .

## Ancestral Graph Method - Cont'd.

- Now define  $MA_{\text{U}}(G) = M(A_{\text{U}}(G))$  the moral graph of  $A_{\text{U}}(G)$   
 then  $\mathbb{I}(MA_{\text{U}}(G)) \subset \mathbb{I}(A_{\text{U}}(G)) \subset \mathbb{I}(G)$ .  
 i.e.,  $\forall \alpha \in \mathbb{I}(MA_{\text{U}}(G)) \Rightarrow \alpha \in \mathbb{I}(G)$ .
- For our example,

$MA_{\text{U}}(G) =$



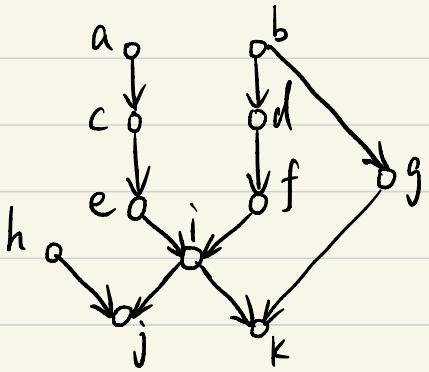
$$\left. \begin{array}{c} <4|2|5> \text{ in this MN.} \\ \downarrow \\ \alpha = "4 \perp\!\!\!\perp_{MA_{\text{U}}(G)} 5|2" \in \mathbb{I}(MA_{\text{U}}(G)) \\ \downarrow \\ \alpha \in \mathbb{I}(G), \text{i.e. } X_4 \perp\!\!\!\perp X_5 | X_2. \end{array} \right\}$$

- Procedure:
  - Step 1: Given  $G$  and  $\alpha = "X \perp\!\!\!\perp Y | Z"$ , form  $\mathcal{U} = X \cup Y \cup Z$ ,  
 and  $A_{\mathcal{U}}(G) = \mathcal{U} \cup \text{ancestor}(\mathcal{U})$
  - Step 2: Form  $MA_{\text{U}}(G) = M(A_{\text{U}}(G))$
  - Step 3: Test if  $X \perp\!\!\!\perp_{MA_{\text{U}}(G)} Y | Z$  is true.

## Ancestral Graph Method - Cont'd.

Example : (See. BRML Fig. 4.4). Is  $a \perp\!\!\!\perp b | \{d, i\} \Rightarrow x_a \perp\!\!\!\perp x_b | \{x_d, x_i\}$ ?

$G =$

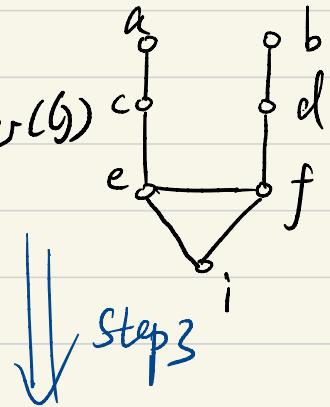
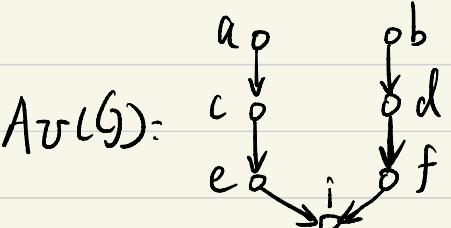


Step 1

$$U = \{a, b, d, i\}$$

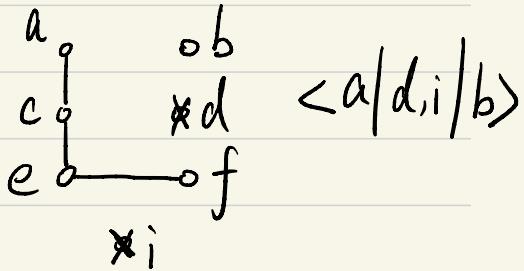
Step 2

$MA_U(G)$



$AU(G) :$

$$\begin{array}{c} a \\ \perp\!\!\!\perp_{MA_U(G)} b | d, i \\ \perp\!\!\!\perp_G b | d, i \\ \perp\!\!\!\perp x_a x_b | \{x_d, x_i\} \end{array}$$



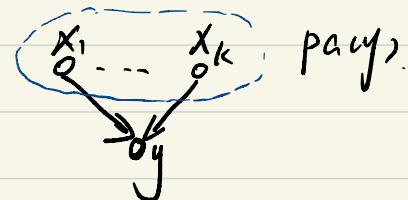
$\langle a | d, i | b \rangle$

## Model Simplification.

- Graphical Models can quickly become intractable complex.
- Model approximations can often be made that provide significant reductions in complexity which provides reasonably accurate performance.
- We discuss three.
  - Logistic Regression Model
  - "Naive Bayes" assumption
  - Gaussian Linear Model / Gaussian BNs

# Logistic Regression Model

- Parents-to-node is a "fan-in" situation.



$$P(y, x_1, \dots, x_k) = P(y | x_1, \dots, x_k) P(x_1) \dots P(x_k)$$

If all binary,  $2^k$  values must be specified or learned (e.g.  $2^{10} = 1024$ )

- Assume a Logistic Regression probability model. This is an "abstraction step" in the modeling process.

- A). Regression step: Learn or specify a linear regression model.

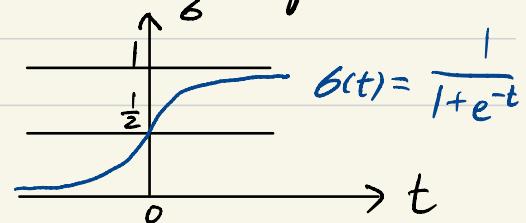
$$t(pa(y)) = b_0 + b_1 x_1 + \dots + b_k x_k$$

$k+1$  parameters are needed, (e.g.  $k=10 \Rightarrow 11$  values)

- B). Logistic step: Apply the non-linear "sigmoid" logistic function.

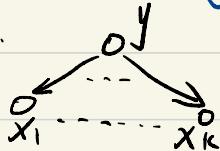
$$\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{e^t + 1} \text{ to } t = t(pa(y))$$

$$\Rightarrow P(y | pa(y)) = \sigma(t(pa(y)))$$

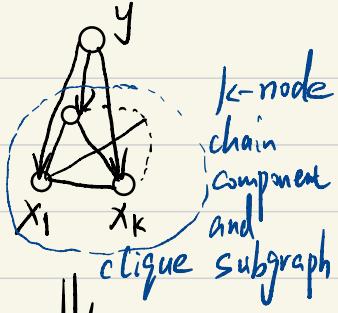


# Naive Bayes Model

- The "fan-out" model.



as an significant simplification of



The naive Bayes Assumption

$$\begin{aligned} & P(x_1, \dots, x_k | y) \\ &= P(x_1 | y) \dots P(x_k | y) \\ &\text{If binary, need specify } 2^k \text{ values (e.g. } k=10 \rightarrow 2^{10} \text{ values)} \end{aligned}$$

$P([x_1, \dots, x_k] | y)$   
if binary, we need specify  
 $2^{k+1} - 1$  values (e.g.  $k=10 \rightarrow 2^{10} - 1$  values)

- The naive Bayes assumption can work quite well if  $y$  is a "hidden factor" that explains a lot of the covariance of  $x_1, \dots, x_k$
- The theory of finding such (often hidden) factors is called "Factor Analysis" (FA).
- It is closely related to Principal Component Analysis (PCA), which attempts to explain the variance of  $x_1, \dots, x_k$ .

# The Gaussian Model.

- As we know, a fully connected graph is highly complex, and this giant clique is hard to unravel. In probability,  $P(X) = P(x_1, \dots, x_N)$  need  $2^N - 1$  specific values even binary. (e.g.  $N=100, \rightarrow 2^{100} / 0^{30}$  values)
- Suppose  $X_i \in \mathbb{R} = \mathbb{R}$  are jointly Gaussian,  $i = 1, \dots, N$ .  
i.e., that  $X \in \mathbb{R}^N$  is a  $N$ -dimensional, Gaussian random vector, denoted as  $X \sim \mathcal{N}(\mu, \Sigma)$ , where the mean  $\mu$  needs  $N$  parameters, and the covariance matrix  $\Sigma$  needs  $\frac{N(N+1)}{2}$  parameters. Then a total of  $\frac{N(N+3)}{2}$  parameters.
- The drop from  $(2^N - 1)$  parameters to  $\frac{N(N+3)}{2}$  parameters is a stunning reduction of complexity arising from the "abstracting step" of using an abstract probability model.
- For  $N=100$ , we need 5150 parameters for the Gaussian model.
- We gain additional efficiency by modularizing (divide and conquer) the world using a Bayesian Network (BN).

# Gaussian BN

- Let  $X = \{X_l, l=1, \dots, N\}$ :  $P(X) = \prod_{l=1}^N P(X_l | \text{pa}(X_l))$
- Gaussian Linear Model (see eq. 10.15 Murphy) (\*)  

$$X_l = \mu_l + \sum_{x_s \in \text{pa}(X_l)} w_{l,s} (X_s - \mu_s) + e_l$$
, where  $e_l \sim N(0, \sigma_l^2)$ , for  $l=1, \dots, N$ .
- Note that this is equivalent to  $P(X_l | \text{pa}(X_l)) \sim N(\mu_l + \sum_{x_s \in \text{pa}(X_l)} w_{l,s} (X_s - \mu_s), \sigma_l^2)$
- Stacking  $X_l, \mu_l, e_l$  into vectors as  $X, \mu, e$ , and  $w_{l,s}$  into matrix  $W$ .  
 we have  $(I - W)(X - \mu) = e$  (\*) elements of  $W$  are either 0 or  $w_{l,s}$   
 Assume ancestral ordering.  $W$  is lower-triangular (see eq 10.19 Murphy)
- Since one can empirically estimate  $\mu = E(X)$  and  $\Sigma = \text{cov}(X)$  from (\*)  
 one can obtain estimates of  $w_{l,s}$  and  $\sigma_l^2$  as discussed in eq. 10.22 Murphy
- Let  $|\text{pa}(X_l)| = p_l$ , take the value  $p=5$  as the average. The number of  
 parameters is  $N$  for means  $\mu_l$ ,  $N$  for variance  $\sigma_l^2$ , and  $SN$  for weights  $w_{l,s}$ .  
 (e.g.,  $N=100 \rightarrow 700$  parameters). We have reduced the complexity  $10^3 \rightarrow 5/150 \rightarrow 700$ .