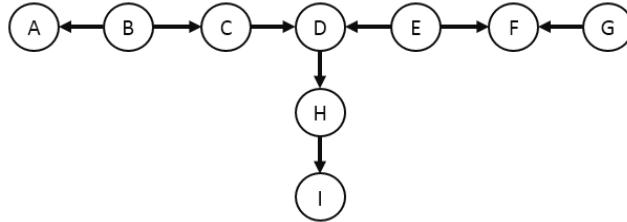


# EE 583 Probabilistic Graphical Models

## Homework 2

Due on Friday, March 25, 2016

1. [9 pts] According to the Bayesian Network shown below, list the cases, in which variables A and G become independent.



**Solution:**

A and G are independent if

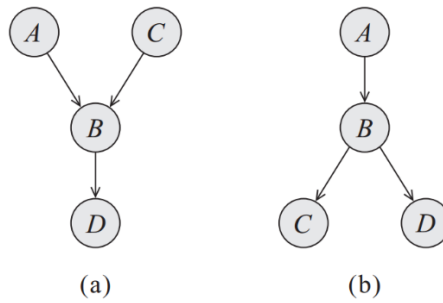
- a) B is observed
- b) C is observed
- c) E is observed
- d) F and at least one of (D,H,I) are not simultaneously observed

2. [20 pts] Let  $G$  be a BN structure and  $P$  be a distribution over random variables  $\mathbf{X} = \{X_1, \dots, X_n\}$ . Show that if  $P$  factorizes according to  $G$ , then  $(X_i \perp \mathbf{X}_{\text{NonDes}(i)} | \mathbf{X}_{\text{Pa}(i)}) \in I(P)$  for all  $X_i$ .  
Hint: Write  $P(\mathbf{X}) = P(X_i, \mathbf{X}_{\text{Pa}(i)}, \mathbf{X}_{\text{NonDes}(i)}, \mathbf{X}_{\text{Des}(i)})$  and using factorization of the RHS, show  $P(X_i | \mathbf{X}_{\text{Pa}(i)}, \mathbf{X}_{\text{NonDes}(i)}) = P(X_i | \mathbf{X}_{\text{Pa}(i)})$ .

**Solution:**

$$\begin{aligned}
 P(x_i | \mathbf{x}_{\text{Pa}(i)}, \mathbf{x}_{\text{NonDes}(i)}) &= \frac{\sum_{\mathbf{x}_{\text{Des}(i)}} P(x_i, \mathbf{x}_{\text{Pa}(i)}, \mathbf{x}_{\text{NonDes}(i)}, \mathbf{x}_{\text{Des}(i)})}{\sum_{x_i} \sum_{\mathbf{x}_{\text{Des}(i)}} P(x_i, \mathbf{x}_{\text{Pa}(i)}, \mathbf{x}_{\text{NonDes}(i)}, \mathbf{x}_{\text{Des}(i)})} \\
 &= \frac{\sum_{\mathbf{x}_{\text{Des}(i)}} P(x_i | \mathbf{x}_{\text{Pa}(i)}) \prod_{j \in \text{NonDes}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)}) \prod_{j \in \text{Des}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)})}{\sum_{x_i} \sum_{\mathbf{x}_{\text{Des}(i)}} P(x_i | \mathbf{x}_{\text{Pa}(i)}) \prod_{j \in \text{NonDes}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)}) \prod_{j \in \text{Des}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)})} \\
 &= \frac{\prod_{j \in \text{NonDes}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)}) P(x_i | \mathbf{x}_{\text{Pa}(i)}) \left[ \sum_{\mathbf{x}_{\text{Des}(i)}} \prod_{j \in \text{Des}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)}) \right]}{\prod_{j \in \text{NonDes}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)}) \sum_{x_i} P(x_i | \mathbf{x}_{\text{Pa}(i)}) \left[ \sum_{\mathbf{x}_{\text{Des}(i)}} \prod_{j \in \text{Des}(i)} P(x_j | \mathbf{x}_{\text{Pa}(j)}) \right]} \quad \begin{matrix} \rightarrow 1 \\ \rightarrow 1 \end{matrix} \\
 &= \frac{P(x_i | \mathbf{x}_{\text{Pa}(i)})}{\sum_{x_i} P(x_i | \mathbf{x}_{\text{Pa}(i)})} = P(x_i | \mathbf{x}_{\text{Pa}(i)})
 \end{aligned}$$

3. [8 pts] Consider the two networks.



For each of them, determine whether there can be any other Bayesian network that is I-equivalent to it.

**Solution:**

- a) No, since assigning different directions to any edge in the skeleton will either make  $A \rightarrow B \leftarrow C$  no longer an immorality, or introduce new ones  $A \rightarrow B \leftarrow D$ , and  $C \rightarrow B \leftarrow D$ .
- b) Yes, since there is no immorality. Designating any node as the root in the skeleton, and directing edges outward will result in an I-equivalent Bayesian network.

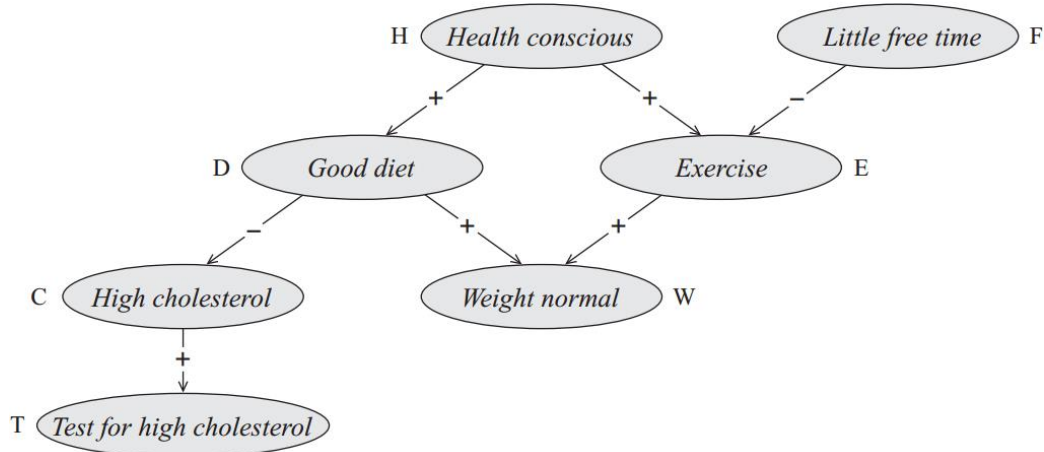
4. [20 pts] Show that global Markov property of Bayesian networks imply the local Markov property. In other words, show that in a Bayesian network, each variable is d-separated from its non-descendants by its parents.

**Solution:**

Suppose  $Y$  is a non-descendant of  $X$  and  $Y \notin \text{Pa}(X)$ . Let  $T$  be a trail between  $Y$  and  $X$ . We analyze two cases:

- a)  **$T$  contains a parent of  $X$ :** Call that parent  $Z$ . Since  $Z \rightarrow X$ ,  $Z$  cannot be a v-junction along  $T$ . Observing  $\text{Pa}(X)$  makes  $Z$  observed, and therefore  $\text{Pa}(X)$  blocks  $T$ .
- b)  **$T$  contains no parent of  $X$ :** Then  $T$  must contain at least one v-junction, because otherwise  $Y$  would be a descendant of  $X$ . To complete the proof, we have to show that the statement “all v-junctions along  $T$  contain a parent of  $X$  among their descendants” cannot be true, because if it were true observing  $\text{Pa}(X)$  would activate  $T$ . Suppose  $W$  is the first v-junction we encounter when tracing edge directions on  $T$  starting from  $X$ . Clearly, we must have a path  $X \rightarrow \dots \rightarrow W$ , because reversing the first edge from  $X$  would introduce a parent of  $X$  to  $T$ , and reversing other edge directions would introduce another v-junction before  $W$ . Then, if there is a parent  $Z \in \text{Pa}(X)$  among  $W$ ’s descendants, then we must have a cycle  $X \rightarrow \dots \rightarrow W \rightarrow \dots \rightarrow Z \rightarrow X$  violating the acyclicity of the Bayesian network, hence completing the proof.

5. [18 pts] Consider the Bayesian network shown below. Assume all variables are binary valued. Suppose the signs at edges indicate how each random variable affects its child qualitatively, that is,  $X \xrightarrow{+} Y$  indicates  $P(Y = 1|X = 1, \mathbf{U} = \mathbf{u}) > P(Y = 1|X = 0, \mathbf{U} = \mathbf{u})$  for all values of  $\mathbf{u}$  of  $Y$ 's other parents (accordingly, the inequality reverses when  $X \xrightarrow{-} Y$ ).



For each pair of the probabilities listed in the rows of the table below, identify which one is larger than the other, or if they are equal, or if they are incomparable with the given edge-sign information alone.

(a)	$P(T = 1 D = 1)$	$P(T = 1)$
(b)	$P(D = 1 T = 0)$	$P(D = 1)$
(c)	$P(H = 1 E = 1, F = 1)$	$P(H = 1 E = 1)$
(d)	$P(C = 1 F = 0)$	$P(C = 1)$
(e)	$P(C = 1 H = 0)$	$P(C = 1)$
(f)	$P(C = 1 H = 0, F = 0)$	$P(C = 1 H = 0)$
(g)	$P(D = 1 H = 1, E = 0)$	$P(D = 1 H = 1)$
(h)	$P(D = 1 E = 1, F = 0, W = 1)$	$P(D = 1 E = 1, F = 0)$
(i)	$P(T = 1 W = 1, F = 0)$	$P(T = 1 W = 1)$

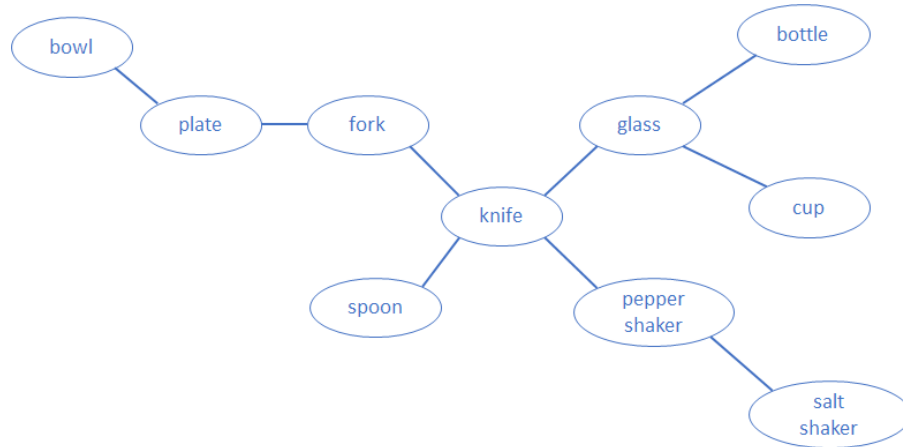
**Solution:**

(a)	$P(T = 1 D = 1)$	<	$P(T = 1)$
(b)	$P(D = 1 T = 0)$	>	$P(D = 1)$
(c)	$P(H = 1 E = 1, F = 1)$	>	$P(H = 1 E = 1)$
(d)	$P(C = 1 F = 0)$	=	$P(C = 1)$
(e)	$P(C = 1 H = 0)$	>	$P(C = 1)$
(f)	$P(C = 1 H = 0, F = 0)$	=	$P(C = 1 H = 0)$
(g)	$P(D = 1 H = 1, E = 0)$	=	$P(D = 1 H = 1)$
(h)	$P(D = 1 E = 1, F = 0, W = 1)$	>	$P(D = 1 E = 1, F = 0)$
(i)	$P(T = 1 W = 1, F = 0)$	>	$P(T = 1 W = 1)$

6. [25 pts] Download the `diningData.mat` and `categoryNames.mat` files from Moodle. Each column of the `diningData` matrix corresponds to one real dining scene image from JHU table setting dataset ( <http://cis.jhu.edu/entropy-pursuit/images.php> ), and each row corresponds to a particular object category from a refined set of  $n = 10$  categories commonly seen on table settings. The  $(i, j)$ th entry of `diningData` is the binary presence indicator of category  $i$  in image  $j$ . Similarly, the  $i$ th entry of `categoryNames` array contains the name for the  $i$ th category.
- Let  $X_i \in \{0,1\}$  denote the indicator random variable for  $i = 1, \dots, n$ . Using data, estimate and provide the mutual information  $\hat{I}(X_i; X_j)$  for each pair  $i \neq j$ .
  - Setting  $\hat{I}(X_i; X_j)$  as edge weights between each possible variable pair, write a program to compute the maximally spanning tree  $T^*$  over  $\mathbf{X} = (X_i, i = 1, \dots, n)$ , and plot (or draw)  $T^*$  providing category names as node labels.
  - Let  $P(\mathbf{X})$  be the true joint distribution over  $\mathbf{X}$ , and let  $P_T(\mathbf{X})$  denote some approximation to  $P$  that uses tree structure  $T$ . Show that  $T^* = \arg \min_T D_{\text{KL}}(P || P_T)$ .

**Solution:**

b)



$$\begin{aligned}
 \text{c) } \arg \min_T D_{\text{KL}}(P || P_T) &= \arg \min_T \left( \sum_{\mathbf{x}} P(\mathbf{x}) \log \frac{P(\mathbf{x})}{P_T(\mathbf{x})} \right) \\
 &= \arg \min_T \left( -H(P) - \sum_{\mathbf{x}} P(\mathbf{x}) \sum_i \log P(x_i | x_{\text{Pa}_T(i)}) \right) \\
 &= \arg \max_T \left( \sum_i \sum_{\mathbf{x}} P(\mathbf{x}) \log \frac{P(x_i, x_{\text{Pa}_T(i)})}{P(x_{\text{Pa}_T(i)})} \right) = \arg \max_T \left( \sum_i \sum_{x_i, x_{\text{Pa}_T(i)}} P(x_i, x_{\text{Pa}_T(i)}) \log \frac{P(x_i, x_{\text{Pa}_T(i)})}{P(x_{\text{Pa}_T(i)})} \right) \\
 &= \arg \max_T \left( \sum_i \sum_{x_i, x_{\text{Pa}_T(i)}} P(x_i, x_{\text{Pa}_T(i)}) \log \frac{P(x_i)P(x_i, x_{\text{Pa}_T(i)})}{P(x_i)P(x_{\text{Pa}_T(i)})} \right) \\
 &= \arg \max_T \left( \sum_i \sum_{x_i, x_{\text{Pa}_T(i)}} P(x_i, x_{\text{Pa}_T(i)}) \log \frac{P(x_i, x_{\text{Pa}_T(i)})}{P(x_i)P(x_{\text{Pa}_T(i)})} + \sum_i \sum_{x_i} P(x_i) \log P(x_i) \right) \\
 &= \arg \max_T \left( \sum_i \sum_{x_i, x_{\text{Pa}_T(i)}} P(x_i, x_{\text{Pa}_T(i)}) \log \frac{P(x_i, x_{\text{Pa}_T(i)})}{P(x_i)P(x_{\text{Pa}_T(i)})} \right) = \arg \max_T \left( \sum_i I(X_i; X_{\text{Pa}_T(i)}) \right) = T^*
 \end{aligned}$$