

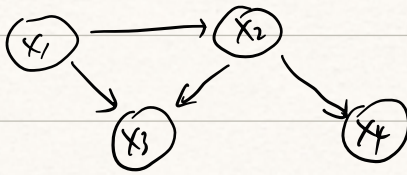
Bayesian Network, $G = (V, E)$

Directed acyclic Graph

$\begin{cases} V: \text{set of node / vertices} \\ E: \text{set of arcs (directed edges)} \end{cases}$

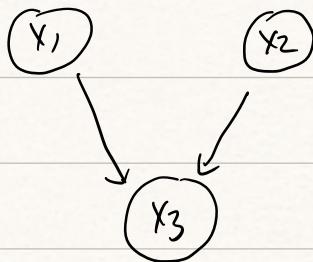
$$P(\underline{X}) = \prod_{i \in V} P(X_i | X_{pa(i)}) \quad pa(i) \rightarrow \text{direct parents of } X_i$$

[E.g.]



$$P(\underline{X}) = P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) P(X_4 | X_2)$$

[E.g.]



Is $X_1 \perp\!\!\!\perp X_2$?

① Create ancestral graph of nodes of interest



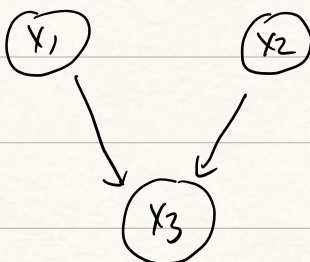
② Marry the parents

③ changes arcs to undirected edges

④ Examine separate sets

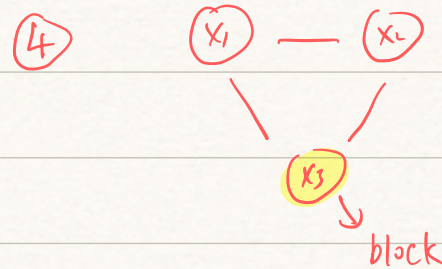
Ans: YES! $\Rightarrow X_1 \perp\!\!\!\perp X_2$

[E.g.]



Is $X_1 \perp\!\!\!\perp X_2 | X_3$?





$$\Rightarrow X_1 \not\perp X_2 \mid X_3$$

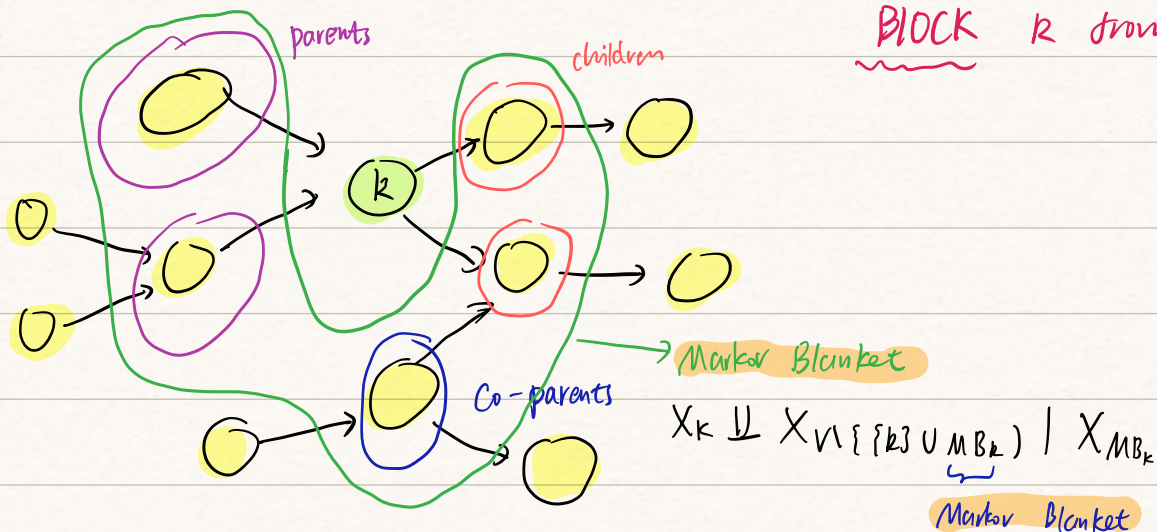
Ans: NO!

Markov Blanket

Qn: Given a BN and any node $k \in V$, what is the smallest set $S \subset V \setminus \{k\}$ s.t

(+) $P(X_k \mid X_S) = P(X_k \mid X_{V \setminus \{k\}}) \Rightarrow$ find the smallest S to

BLOCK k from outside



Defn: The smallest $S \subset V \setminus \{k\}$ s.t (+) holds is called the Markov Blanket of node k .

k=1; RHS of (+)

$$P(X_i | X_{V \setminus \{i\}}) = \frac{P(X_i)}{P(X_{V \setminus \{i\}})} = \frac{\prod_{k=1}^d P(X_k | X_{pa(k)})}{\sum_{X_i} \prod_{k=1}^d P(X_k | X_{pa(k)})}$$

大部分都约掉了! (不含 X_i)

marginalization

Terms in the product in the denominator that do not depend on X_i can be brought outside the sum and cancelled with the sum in the numerator.

Qn. which terms contain X_i in $\prod_{k=1}^n P(X_k | X_{pa(k)})$

① $P(X_i | X_{pa(i)})$ [involves the set of parents of X_i]

② If $X_{ch(i)}$ is a child of X_i , then there will be terms of the form:

$$P(X_{ch(i)} | X_i, \dots)$$

[involves the children of X_i]

③ If X_1 & X_2 share a child, X_{ch} , then

there will be terms of the form:

marginalization

$$P(X_{ch} | X_1, X_2, \dots)$$

[involves co-parents of X_i]

Claim: The Markov Blanket of node i is the set

$$MB(i) = pa(i) \cup ch(i) \cup \text{co-parents}(i)$$

come from marrying!

Learning Bayesian Networks

Given ^① a DAG $G = (V, E)$ $V = \{1, 2, \dots, d\}$

↓
tells us the structure of graph

② n complete data (observations) $\underline{x}_t = (x_{t1}, \dots, x_{td})^T \quad t=1, 2, \dots, n$

$$\mathcal{D} = \{\underline{x}_t\}_{t=1}^n$$

Want to learn the conditional distribution $P(x_i | x_{pa(i)}) = \theta_{x_i | x_{pa(i)}} \quad i=1, 2, \dots, d$
(conditional probability tables)

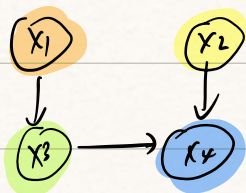
$$P(\underline{x}) = \prod_{i=1}^d P(x_i | x_{pa(i)}) = \prod_{i=1}^d \theta_{x_i | x_{pa(i)}}$$

$$\underline{\theta} = \{ \{ \theta_{x_i | x_{pa(i)}} : x_i \in [r_i], x_{pa(i)} \in [r_{pa(i)}] \}_{i=1}^d \}$$

① $\{1, \dots, r_i\}$ the set of value that x_i takes on

② $x_{pa(i)} = (x_j, x_k) \quad [r_{pa(i)}] = [r_j] \times [r_k]$

E.g.



$$[r_i] = \{1, 2\}$$

1 $\theta_{x_1=1} \quad \theta_{x_1=2}$

1 $\theta_{x_2=1} \quad \theta_{x_2=2}$

2 $\theta_{x_3=1 | x_1=1} \quad \theta_{x_3=2 | x_1=1}$
 $\theta_{x_3=1 | x_1=2} \quad \theta_{x_3=2 | x_1=2}$

4 $\theta_{x_4=1 | x_2=a, x_3=b} \quad \forall (a,b) \in [2]^2$
 $\theta_{x_4=2 | x_2=a, x_3=b}$

4
8

Maximum Likelihood parameter estimation

Given $\mathcal{D} = \{x_t\}_{t=1}^n$, want to learn $\underline{\theta}$.

$$\downarrow$$
$$x_t \sim p$$

Log-likelihood:

$$l(\mathcal{D}; \underline{\theta}, G) = \log P(\mathcal{D} | \underline{\theta})$$

} i.i.d

$$= \log \prod_{t=1}^n P(x_t | \underline{\theta})$$

$$= \sum_{t=1}^n \log P(x_t | \underline{\theta})$$

$$= \sum_{t=1}^n \log \prod_{i=1}^d P(x_{ti} | x_{t, \text{pa}(i)})$$

$$= \sum_{t=1}^n \sum_{i=1}^d \log \theta_{x_{ti} | x_{t, \text{pa}(i)}} \quad \leftarrow P(x_{ti} | x_{t, \text{pa}(i)})$$

$$= \sum_{i=1}^d \sum_{x_i, x_{\text{pa}(i)}} n(\underline{x}_i, \underline{x}_{\text{pa}(i)}) \log \theta_{x_i | x_{\text{pa}(i)}}$$

↑ 个数 n_i ↑ 个数 $[r_j] \times [r_k] \times \dots [r_{\text{pa}(i)}]$

↓

of times $(x_i, x_{\text{pa}(i)})$ occurs in
the dataset $\{(x_{ti}, x_{t, \text{pa}(i)}) : t=1, 2, \dots, n\}$

Now differentiate $l(\mathcal{D}; \underline{\theta}, G)$ w.r.t. $\underline{\theta}$.

(Notice that $\sum_{x_i} \theta_{x_i | x_{\text{pa}(i)}} = 1, \forall x_{\text{pa}(i)}$)

$$\Rightarrow \hat{\theta}_{x_i | x_{\text{pa}(i)}} = \frac{n(x_i, x_{\text{pa}(i)})}{\sum_{x_i'} n(x_i', x_{\text{pa}(i)})}$$

Empirical Frequency Counts

★

Difficult to learn in high degree! \Rightarrow have many parents

To provide one observation per configuration of parent variables require $\prod_{j \in \text{pa}(i)} r_j$ instances!

Model Selection

① Given $\hat{\theta}_{x_i | \text{pa}(i)}$ for all $(x_i, x_{\text{pa}(i)})$ and all $i \in [d]$,

we can compute

\Downarrow
suppose we have good prior models

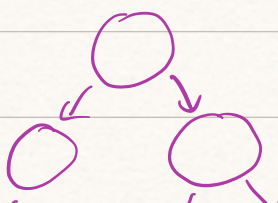
$$\text{BIC}(G) = \ell(\mathcal{D}; \hat{\theta}, G) - \frac{\dim(G)}{2} \log n$$

$\dim(G) \equiv \#$ of parameters used to describe the BN
that represents G

$$= \sum_{i=1}^d \underbrace{(r_i - 1) \prod_{j \in \text{pa}(i)} r_j}_{\substack{\theta_{x_i | x_{\text{pa}(i)}} \\ r_i - 1 \quad \prod r_j}}$$

② Don't rely on outside prior information

If $G = (V, E)$ is a tree \rightarrow connected and every nodes have in-degree ≤ 1 ,
then there's a much better way to learn G ! 入度



[d nodes] \rightarrow

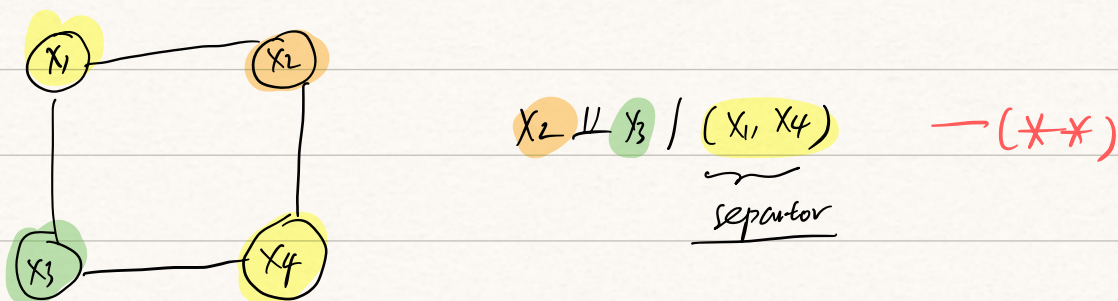
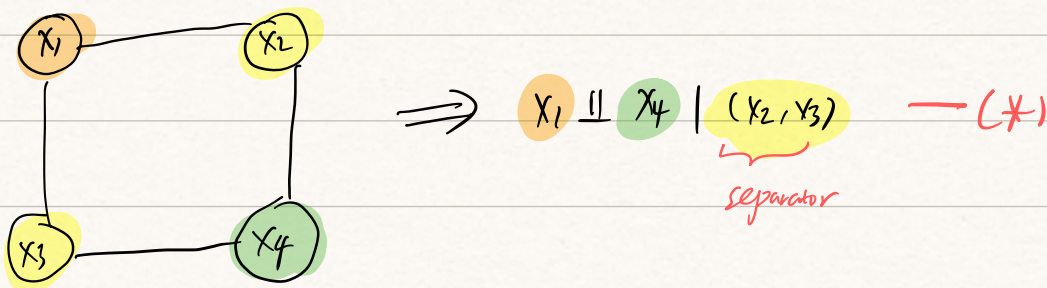
Total # of such graphs

Markov Random Field (Undirected Graphic Model)

Undirected graph $G=(V, E)$ $E \subset \binom{V}{2}$

$$= \{ \underbrace{\{i, j\}}_{\text{unordered}} : i, j \in V, i \neq j \}$$

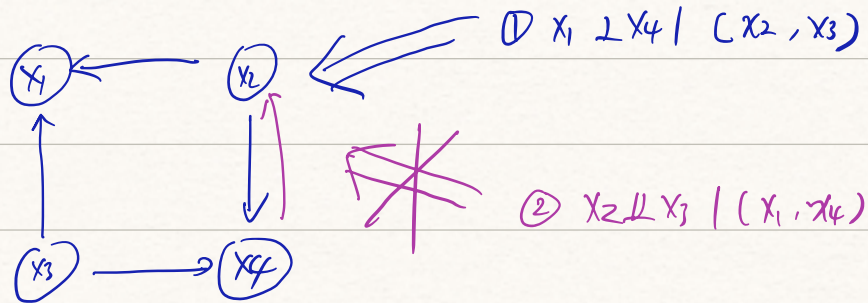
Graph encodes conditional independence relations among $X_i, i=1, 2, \dots, d$



Qn: Can you draw a Bayesian Network that encodes these two CI properties?

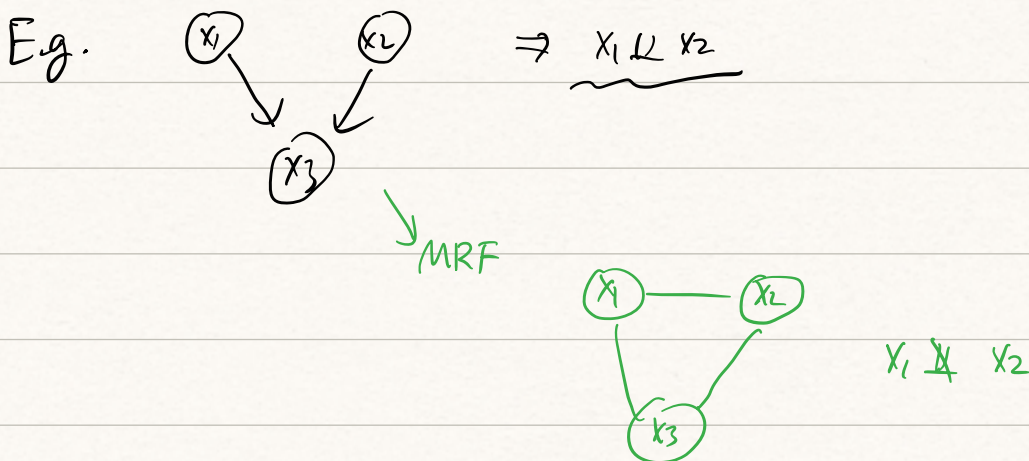
(* & **)

TRY



Impossible to draw the BN that respects (*) & (**)

In terms of ability to explicate condition independence properties
MRFs & BNs are not strict subsets to each other!



Hammersley - Clifford Thm

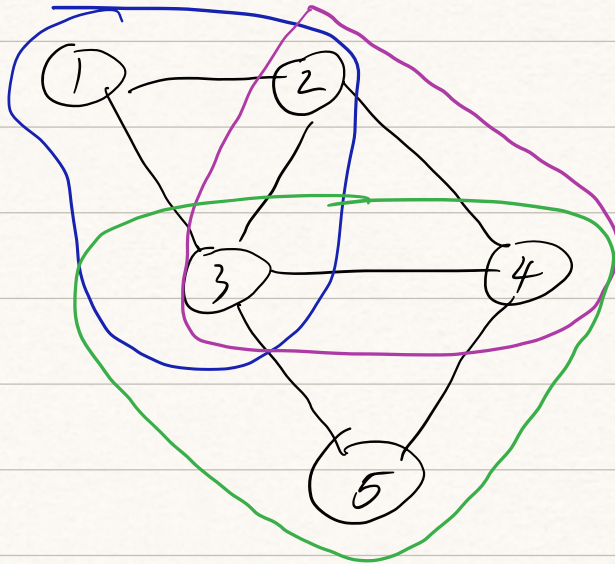
(Recall that for a BN, $p(x) = \prod_{i=1}^d p(x_i \mid x_{\text{pa}(i)})$)

If $G=(V,E)$ is an undirected graph, can we write
 $p(x)$ in the form of products of terms defined on
small subsets of nodes?

Defn: Given an undirected graph $G=(V,E)$, we say that

C is a clique if it is a fully connected set of nodes
完全图.

A maximal clique is one that cannot be extended by including one more node



Clique $C = \{1, 2\}$

is not maximal!

Clique $\tilde{C} = \{1, 2, 3\}$

is maximal

[Informal Version of the Hammersley - Clifford Thm]

For a MRF, we can write its joint distribution as follows

$$p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \phi_C(x_C)$$

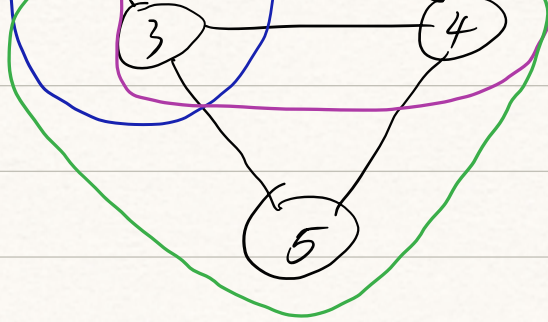
\mathcal{C} : set of all maximal cliques of G .

ϕ_C : positive f.d. defined on the domain of
 $x_C = (x_i; i \in C)$



$$p(x) = \frac{1}{Z} \phi_{1,2,3}(x_1, x_2, x_3)$$

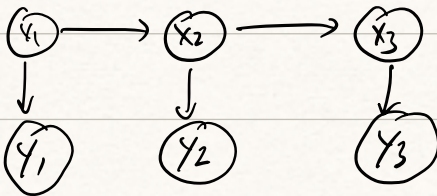
$$\phi_{2,3,4}(x_2, x_3, x_4)$$



$$\phi_{3,4,5}(x_3, x_4, x_5)$$

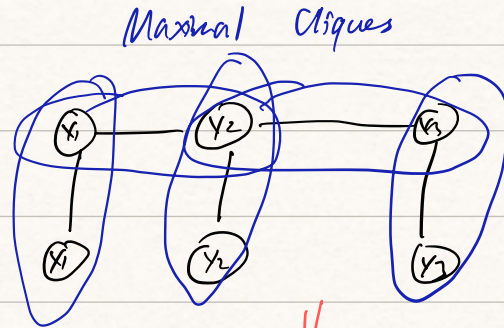
Gibbs distribution

E.g. HMM



Bayesian Network

$$p(X) = p(x_1) p(y_1 | x_1) p(x_2 | x_1) p(y_2 | x_2) p(x_3 | x_2) p(y_3 | x_3)$$



MRF

No 2 parents!

$$p(X) = \frac{1}{Z} \phi_{x_1 y_1}(x_1, y_1) \phi_{x_1 x_2}(x_1, x_2) \phi_{x_2 y_2}(x_2, y_2) \phi_{x_2 x_3}(x_2, x_3) \phi_{x_3 y_3}(x_3, y_3)$$