Till now - Probabilistic models for data ⎤ Dirichlet allocat"
⎥ Regression
Model in words ⎦ Clustering

# ORIE 4742 - Info Theory and Bayesian ML

Bayesian Networks

April 9, 2020

Sid Banerjee, ORIE, Cornell

From- Ch8, PRML
by Chris Bishop

# probabilistic graphical models

graphical representation of complex probability distributions

## types of graphical models

BayesNets: directed acyclic graphs

Markov random fields: undirected graphs
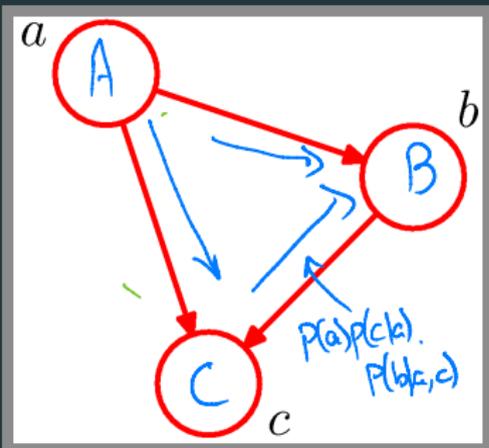
factor graphs: bipartite graphs

*(handwritten annotations: • ≡ random variable; Markov chain; ▨ ≡ functions)*

## why are they useful?

- visualizing helps in design of probabilistic models
- complex inference/learning calculations $\rightarrow$ simpler graph operations
- gives insight into properties of model: conditional independence, causal relationships
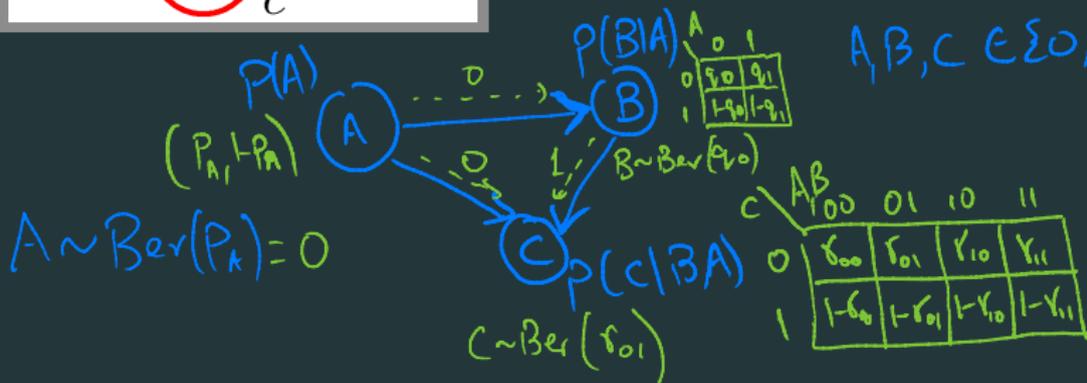
# BayesNets



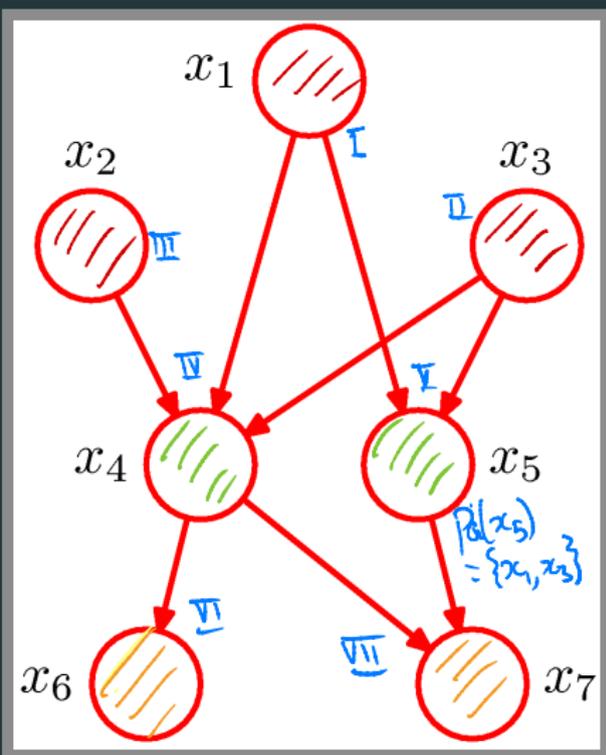directed acyclic graph (DAG) encoding conditional distributions

eg. for r.v.s $A, B, C$, BN on left encodes:

$$p(A, B, C) = p(C|A, B)p(A, B)$$
$$= p(c|A, B)p(B|A)p(A)$$

$$A, B, C \in \{0, 1\}$$

$p(A)$

$p(B|A)$

| $A$ | 0 | 1 |
|---|---|---|
| 0 | $q_0$ | $q_1$ |
| 1 | $1-q_0$ | $1-q_1$ |

$B \sim Ber(q_0)$

$A \sim Ber(P_A) = 0$

$(P_A, 1-P_A)$

$p(c|BA)$

$C \sim Ber(\gamma_{01})$

| $C$ | $AB$ 00 | 01 | 10 | 11 |
|---|---|---|---|---|
| 0 | $\gamma_{00}$ | $\gamma_{01}$ | $\gamma_{10}$ | $\gamma_{11}$ |
| 1 | $1-\gamma_{00}$ | $1-\gamma_{01}$ | $1-\gamma_{10}$ | $1-\gamma_{11}$ |

$$P(x_1, x_2, \ldots, x_7) = P(x_1).$$

$$P(x_3). P(x_2) P(x_4 | x_1, x_2, x_3).$$

$$P(x_5 | x_1, x_3) P(x_6 | x_4)$$

$$P(x_7 | x_5, x_4)$$

- Any DAG has a 'topological ordering' (ie, numbering s.t. no edge from higher to lower number) · use to generate prob expansion / factorization ↓
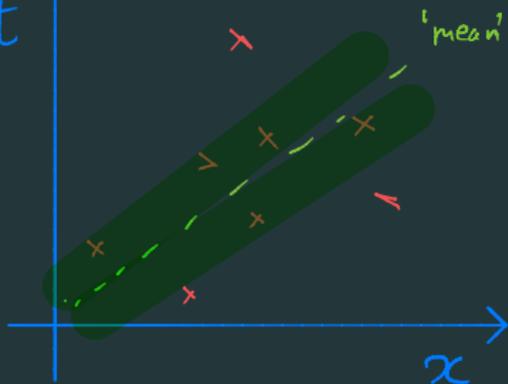
- For any $x$, $Pa(x) \equiv$ 'parents' of $x$   ·   $P(\hat{\cup}_{i=1} x_i) = \prod_{i=1}^{n} P(x_i | Pa(x_i))$

In the figure (labels):
- $x_1$ — I
- $x_3$ — II
- $x_2$ — III
- IV
- V
- $Pa(x_5) = \{x_1, x_3\}$
- VI
- VII

# example: (Bayesian) regression

Input - $(x_1, t_1), (x_2, t_2) \ldots, (x_n, t_n)$    $t$

Tasks - 1) $t_i = \sum_{j=1}^{m} w_j f_j(x_i) + w_0$

basis fns ↑    noise ↑

- $w_i \sim N(\mu_i, 1/\tau_i)$
- $\varepsilon_i \sim N(0, 1/\tau_\varepsilon)$   Want to learn $(w_1, w_2, \ldots, w_m)$
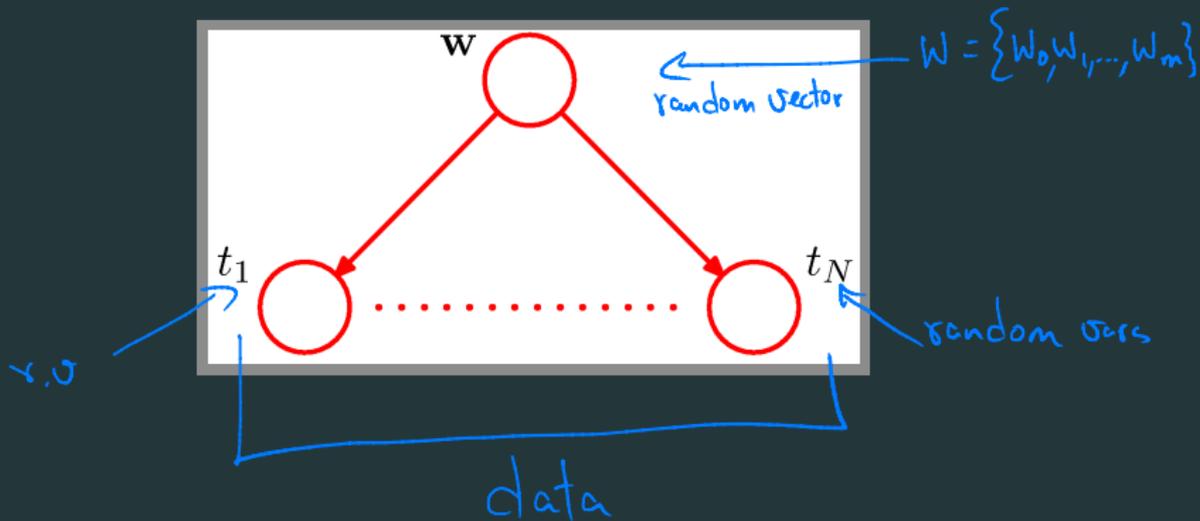  from data

'mean'

$x$

2) Given new point $x_{n+1}$, predict/infer $t_{n+1}$

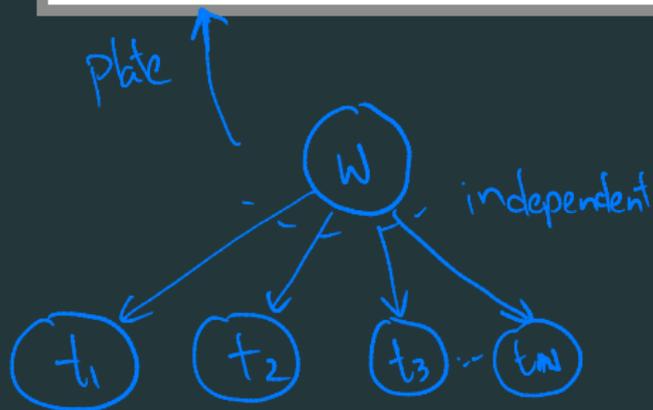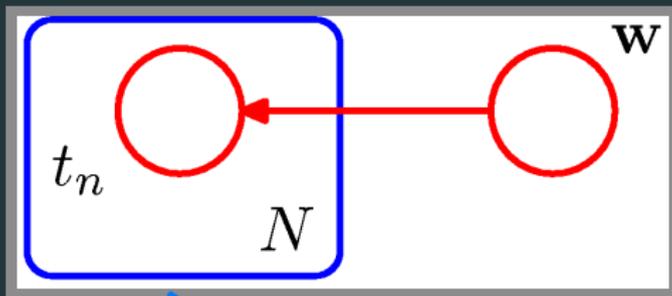Eg - $f_1(x) = 1$ (constant), $f_2(x) = x$ (linear regression)

- $f_k(x) = x^{k-1}$ (polynomial regression)

- $f_k(x) = e^{-(x-\mu)^2/2}$ (Gaussian basis fn)

# regression: basic BayesNet

# regression: inputs and hyperparameters



hyperparams are represented as solid dots
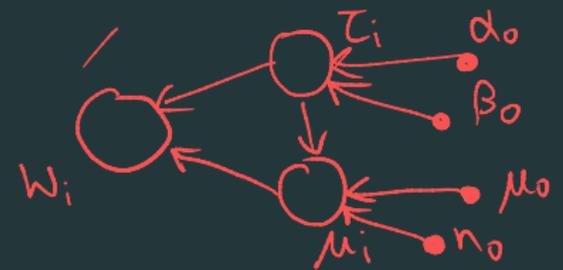
(true for any 'deterministic variable')

$x_n$

$\alpha$

$\sigma^2$

$t_n$

$N$

$\mathbf{w}$

Prior
$$\left\{ w_i \sim N\left(\mu_i, 1/\tau_i\right)\right.$$
$w \{\mu_1, \tau_1, \ldots, \mu_m, \tau_m\}$

(want to learn)
model params

$w_0 \sim N(0, \sigma^2)$

'nuisance' parameter (do not want to learn)

$w$ $\equiv$

$\tau_i$  $a_0$

$\beta_0$

$\mu_0$

$w_i$

$\mu_i$  $n_0$

Plate diagram with node $t_n$ (shaded, observed variable) inside plate $N$ with inputs $x_n$ and $\sigma^2$, and node $\mathbf{w}$ (unshaded, latent variable) with input $\alpha$.

unshaded node $\equiv$ latent variable

Shaded node $\equiv$ 'observed variable'

$(x_1, t_1), \ldots, (x_N, t_N)$

$t_i = \sum_{j=1}^{m} w_j \, f_j(x_i)$

$+ \, w_0$

learning module

- Note

$x_1, x_2, \ldots, x_N,$

and $\hat{x} \in \mathbb{R}^d$

- However $x_i$'s are not r.v.

$\hat{t} \perp\!\!\!\perp t_i \mid w$

indep $\equiv$ plate notation

hyperparam

inference

$x_n$

constants

observed vars

$t_n$

$N$

latent vars

$\mathbf{w}$

nuisance params

$\sigma^2$

$\hat{t}$

$\hat{x}$

If $x_i$ were random

$\tau_0$

$w_0$

**example: naive Bayes**

Assumption - $X_j^i \perp\!\!\!\perp X_{j'}^i \quad \forall i$ , Eg - $\left(X_1^i, X_2^i, \dots, X_d^i\right) \sim D_{i, c}\left(\alpha_{i,j}^{c}\right)$

# conditional independence

- Use a given Bayes Net to answer is. $A \perp\!\!\!\perp B \mid C$

  - $(A_1, A_2) \perp\!\!\!\perp (B_1, B_2, B_3) \mid (C_1, C_2, C_3)$

- $P(A, B \mid C) \left( \overset{?}{\underset{\neq}{=}} \right) P(A \mid C) P(B \mid C)$

- TLDR – You can answer this given a Bayes Net

  - d-separation (Pearl '88)

  - 3 building blocks



Is $A \perp\!\!\!\perp B | C$

- Question

$$P(a,b,c) = P(c) \cdot P(a|c) \cdot P(b|c)$$

without conditioning

$$P(a,b) = \sum_c P(a|c) \cdot P(b|c) \cdot P(c)$$

$$\Rightarrow \boxed{a \not\!\perp b}$$

$$P(a,b|c) = \frac{P(a,b,c)}{P(c)}$$

$$= \frac{P(c)\,P(a|c)\,P(b|c)}{P(c)}$$

$$= P(a|c) \cdot P(b|c) \Rightarrow \boxed{a \perp\!\!\!\perp b \mid c}$$

$$P(a,b,c) = P(a) \cdot P(c|a) \cdot P(b|c)$$

$$P(a,b) = P(a) \sum_c P(c|a) \, P(b|c) \neq P(a) \cdot P(b)$$

$$\Rightarrow \boxed{a \not\!\perp b} \quad \left( Eg \cdot b=c, \; c=a \right)$$

$$P(a,b|c) = \frac{P(a)\,p(c|a)\,p(b|c)}{p(c)}$$

$$= p(a|c)\,p(b|c)$$

$$\Rightarrow \boxed{a \perp\!\!\!\perp b \mid c}$$

(Markov) chain

$$P(a,b,c) = P(a) \cdot P(b) \cdot P(c|a,b)$$

$$\Rightarrow P(a,b) = P(a) \cdot P(b) \cdot \underbrace{\sum_c P(c|a,b)}_{=1}$$

$$\Rightarrow \boxed{a \perp\!\!\!\perp b}$$

('explaining away')



$$P(a,b|c) = \frac{P(a)\,P(b)\,P(c|a,b)}{P(c)}$$

$$\neq P(a|c) \cdot P(b|c)$$

$$\Rightarrow \quad a \not\perp b \mid c$$

# 'explaining away'



(from Bishop)

| AB | 00 | 01 | 10 | 11 |
|----|----|----|----|----|
| C 0 | 0.9 | 0.8 | 0.8 | 0.2 |
| 1 | 0.1 | 0.2 | 0.2 | 0.8 |

$Eg -$ $C = \mathbb{1}\{$ fever + cough burglar alarm rang $\}$

$A = \mathbb{1}\{$ allergy there is a burglar $\}$

$B = \mathbb{1}\{$ COVID-19 there is a raccoon $\}$

$\mathbb{P}[A=1] = 0.9, \quad \mathbb{P}[B=1] = 0.9$

• $\mathbb{P}[B=0 \mid C=0] = \dfrac{\mathbb{P}[C=0 \mid B=0]\,\mathbb{P}[B=0]}{\mathbb{P}[C=0]} \approx 0.25$

$\mathbb{P}[B=0 \mid C=0, A=0] = \dfrac{\mathbb{P}[C=0 \mid B=0, A=0]\,\mathbb{P}[B=0]}{\sum_{0,1}\mathbb{P}[A=i]\,\mathbb{P}[C=0 \mid A=i, B=0]} \approx 0.11$

# 'explaining away'

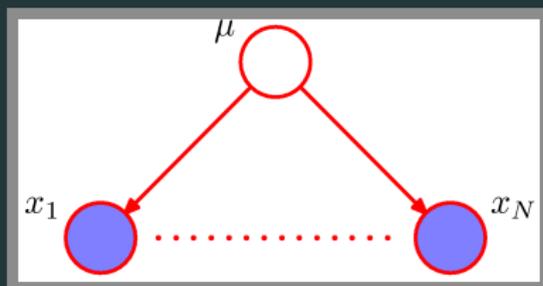- $A \perp\!\!\!\perp B \mid C$ if C is not a join or a descendant of a join



- A path from $A \rightarrow B$ is blocked by C if
i) $\rightarrow \cdot \xrightarrow{C}$ or $\leftarrow \cdot \xleftarrow{e}$
ii) C is not $\rightarrow \cdot \leftarrow$ or a descendent of $\rightarrow \cdot \leftarrow$

A,B are d-separated by C if every path $A \rightarrow B$ is blocked by C
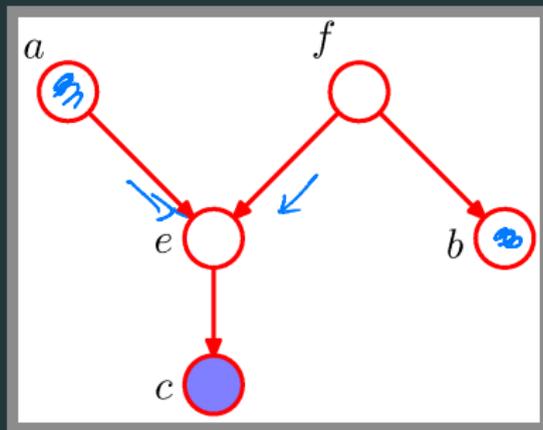
- $X_i \perp\!\!\!\perp X_j \mid \mu$

$\mu$ blocks path

- $X_i \not\perp\!\!\!\perp X_j$

$\mu$ does not block path

Q: Is
i) $A \perp\!\!\!\perp B$ ?
ii) $A \perp\!\!\!\perp B \mid C$

Ans:
- $A \not\perp\!\!\!\perp B$
- $A \not\perp\!\!\!\perp B \mid C$

Q: Is

$A \perp\!\!\!\perp B \mid F$
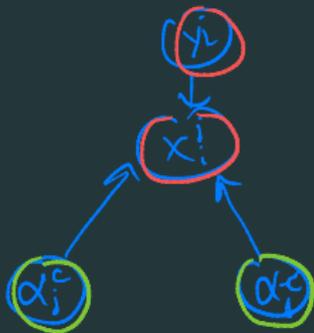
Ans     $A \perp\!\!\!\perp B \mid F$

$$\left(\alpha_1^c, \alpha_2^c \cdots \alpha_d^c\right) \perp\!\!\!\perp \left(\alpha_1^{c'}, \alpha_2^{c'} \cdots \alpha_d^{c'}\right) \Big| \left\{y^i, x^i_{1, \cdots, x^i_n}\right\}$$

# example: naive Bayes

# example: naive Bayes

# Markov random fields

- Bayes Nets encode 'local conditioning'

$$\prod_i P(x_i \mid Pa(x_i))$$

- Don't directly capture global conditional indep/dep

maximal cliques — set of nodes which form a clique, and are not subsets of a larger 'selected' clique

- Clique cover — collection of (maximal) cliques s.t. every edge is in a clique

- $P(x_1, x_2, x_3, x_4) = \dfrac{1}{Z} \cdot \prod_{\text{cliques } c} \psi_c(x_i : i \in c)$

$Z$ — normalization

clique potential

$$C = \{ \{x_1, x_2\} \; \{x_1, x_4\} \; \{x_2 \, x_4 \, x_3\} \}$$

$$P(x_1, x_2, x_3, x_4) = \frac{1}{Z} \psi_{12}(x_1, x_2) \, \psi_{14}(x_1, x_4).$$

$$\psi_{234}(x_2, x_3, x_4)$$

$\Longleftrightarrow$

- Conditional indep $\Longleftrightarrow$ separation

Q: $X_3 \perp\!\!\!\perp X_1 \mid X_4, X_2$ ?

Yes as $(X_2, X_4)$ separate $X_1$ and $X_3$

disconnect

$$P(x) = p(x_1) \, p(x_2|x_1) \cdots p(x_N|x_{N-1})$$

$$P(x) = \frac{1}{z} \, \psi_{12}(x_1, x_2) \, \psi_{23}(x_2, x_3) \cdots$$
$$\cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

- Choose $\psi_{12}(x_1, x_2) = p(x_1) p(x_2|x_1)$

$$\psi_{23}(x_2, x_3) = p(x_3|x_2)$$
$$\vdots$$
$$\psi_{N-1,N}(x_{N-1}, x_N) = p(x_N|x_{N-1})$$

# BayesNet vs MRF



$P(x) = P(x_1) P(x_2) P(x_3) P(x_4 | x_1, x_2, x_3)$

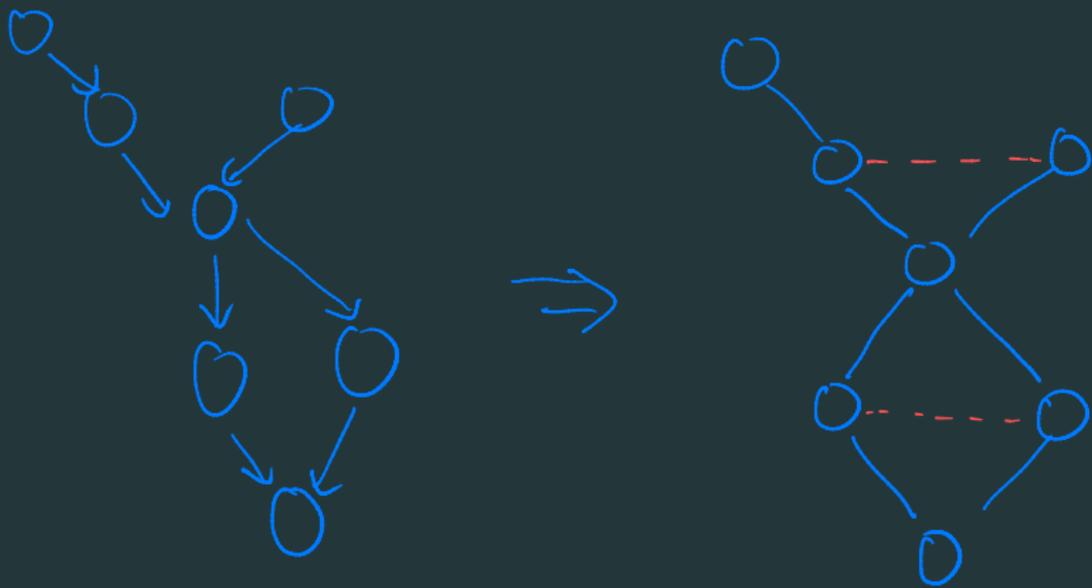$C = \{\{x_1\}, \{x_2\}, \{x_3\}, \{x_1 x_2 x_3 x_4\}\}$

$\psi_1(x_1) = P(x_1), \ldots \psi_{1234}(x_1 x_2 x_3 x_4)$

$= P(x_4 | x_1 x_2 x_3)$

$C = \{x_1 x_2 x_3\} \{x_1 x_2 x_4\} \{x_2 x_3 x_4\}$

$\Rightarrow \frac{1}{Z} \psi_{123}(x_1 x_2 x_3) \psi(x_1 x_2 x_4) \psi_{234}(x_2 x_3 x_4)$

- Add edges between unconnected parents of each child node

- Is $A \perp\!\!\!\perp B \mid C$?
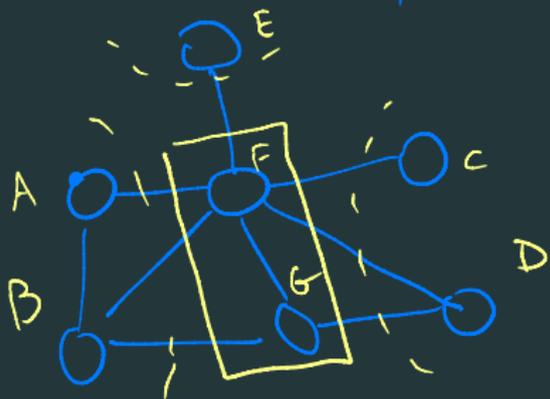
  - Convert Bayes Net of 'ancestors of C' into MRF (via moralization)

  - Check for conditional independence

**Eg:-** Markov chain

indep ← → indep

$X_1$ — $X_2$ — $X_3$ — $X_4$

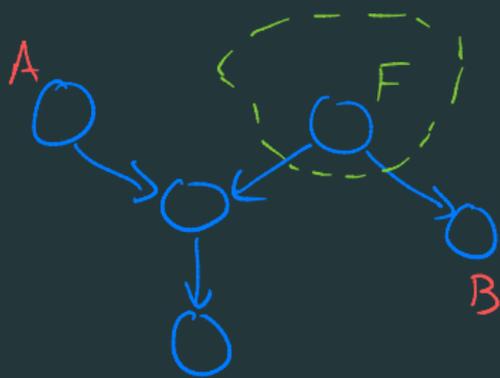- $X_i \perp\!\!\!\perp X_j \mid X_k$ if $i < k < j$
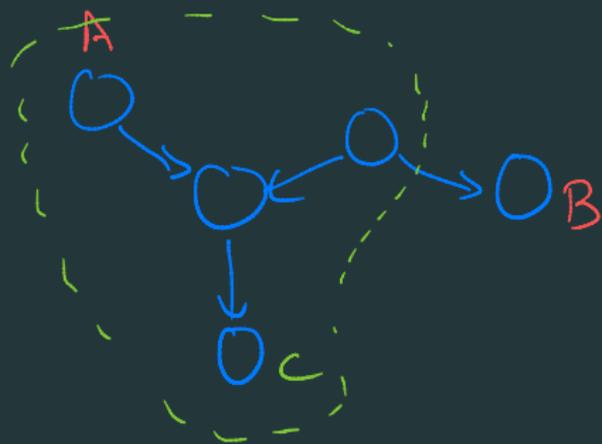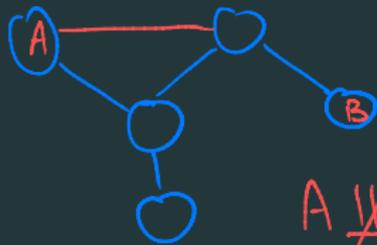  
  or $i > k > j$

**Eg**



$A, B \perp\!\!\!\perp C, D \mid F, G$

$C \perp\!\!\!\perp D \mid F, G$

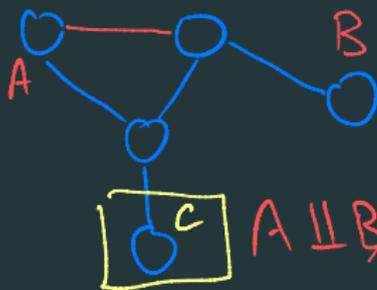$A \not\!\perp\!\!\!\perp B \mid F, G$
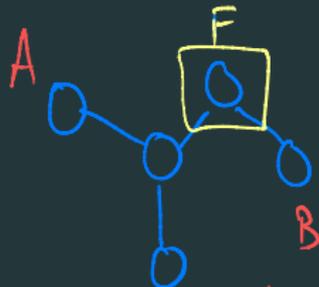
$E \perp\!\!\!\perp C, D \mid F, G$

- Unconditional

- Condn on C

- Condn on F

$A \not\!\perp\!\!\!\perp B$

$A \perp\!\!\!\perp B / C$

$A \perp\!\!\!\perp B / F$