

09-09-24.

## 22 AGE 301 Probabilistic Reasoning

Credits : 3.

### Continuous Evaluation Pattern

- \* 2 Quizzes ( $2 \times 10 = 20$ )
- \* 30 marks Assignment.
  - 15  $\rightarrow$  Assgn. 1 (Paper study) (Viva)
  - 15  $\rightarrow$  Assgn. 2 (Exam/Quiz)
- \* Mid term  $\rightarrow$  20 marks.
- \* End sem.  $\rightarrow$  30 marks.

### Sources.

- \* Christopher M. Bishop  
"Pattern Recognition & ML"

### Main Topics

- \* Sampling Methods.
- \* Probabilistic Graphical Models (PGM).

- Spam classification  
A type of prob. distribution.

- Total no. of possibilities =  $2^{n+1}$   
or parameters

- learning parameter of a distribution.  
(Task 1)

$\Downarrow$

We have some data but we don't know from which underlying prob. dist. comes from, then we have to learn those parameters ( $\theta$ ). MLE is a method used for it. (Max. Likelihood Estimation).

- Inference (Task 2)

$\Downarrow$

Once we learn these parameters, we need to know on how to use them for classification purpose / organised task.

In task 1, we model  $P_{\text{true}}$  by  $P_{\theta}$  where  $P_{\text{true}}$  is the actual data.

In task 2,

$$P_{\theta}(y=1 \mid \text{"doc"}_{x_1 \dots x_n})$$

which is the inference.

Joint distribution of  $y$  &  $x$ ;

$$P(y, x) = P(y \mid \hat{x}) \cdot p(x).$$



\* Marginal inference 2  
From joint distribution (of  $y$  &  $x$ ?)  
we find the inference / we  
find the marginal of some variable

\* Arg max  $P(x_1 \dots x_n | y=1)$   
 $x_1 \dots x_n$  < MAP inference >

for e.g. Suppose there was a certain  
no. of words what is the prob.  
that it is spam (OR)  
Given a doc. is spam what the  
max. no. of words in it.

### Modeling / Representation (Task 3)

The one in which we efficiently  
model / rep. the data (joint data).

$$p(x_1 \dots x_n, y) = p(x_1 \dots x_n | y) \cdot p(y)$$

Assume 2

$$= p(x_1 | y) \cdot p(x_2 | y) \dots$$

$$p(x_n | y) \cdot p(y)$$

$x_1 \dots x_n$  are independent given  $y$ .  
(Naive Bayes Classifier)

We are doing this modeling  
in order to make the problems  
more linear or simple so that it is  
more tractable.

$O(n)$

The Graphical repre. helps in  
forming a proper inference &  
then it will in turn help  
in creating algorithms that  
can be useful in performing  
several tasks.

In marginal inference,

$$P(y, x_1, x_2, \dots, x_n)$$

$y \rightarrow$  cost of house

$x \rightarrow$  features such as

no. of bedrooms, floors, etc.

& the above is a joint distribution.

Suppose  $x_1$  says about the bedrooms.

& we are asked some qs. on

$x_1$ .  $\therefore$  From the joint dist.,

we are only answering about  
 $x_1$  which comes under marginal  
distribution / inference.

\* If the no. of parameters ( $2^{n+1}$ )  
is very large compared to the  
training examples / data then  
there occurs the problem of  
"poor generalisation" & "computational  
issues."

$\Downarrow$

\* To overcome this we make a  
"conditional independence" assumption.



## \* Naive Bayes Assumption

→ In task 3, we look into:

- Graphical models &
- Conditional independence.

## Bayesian Networks.

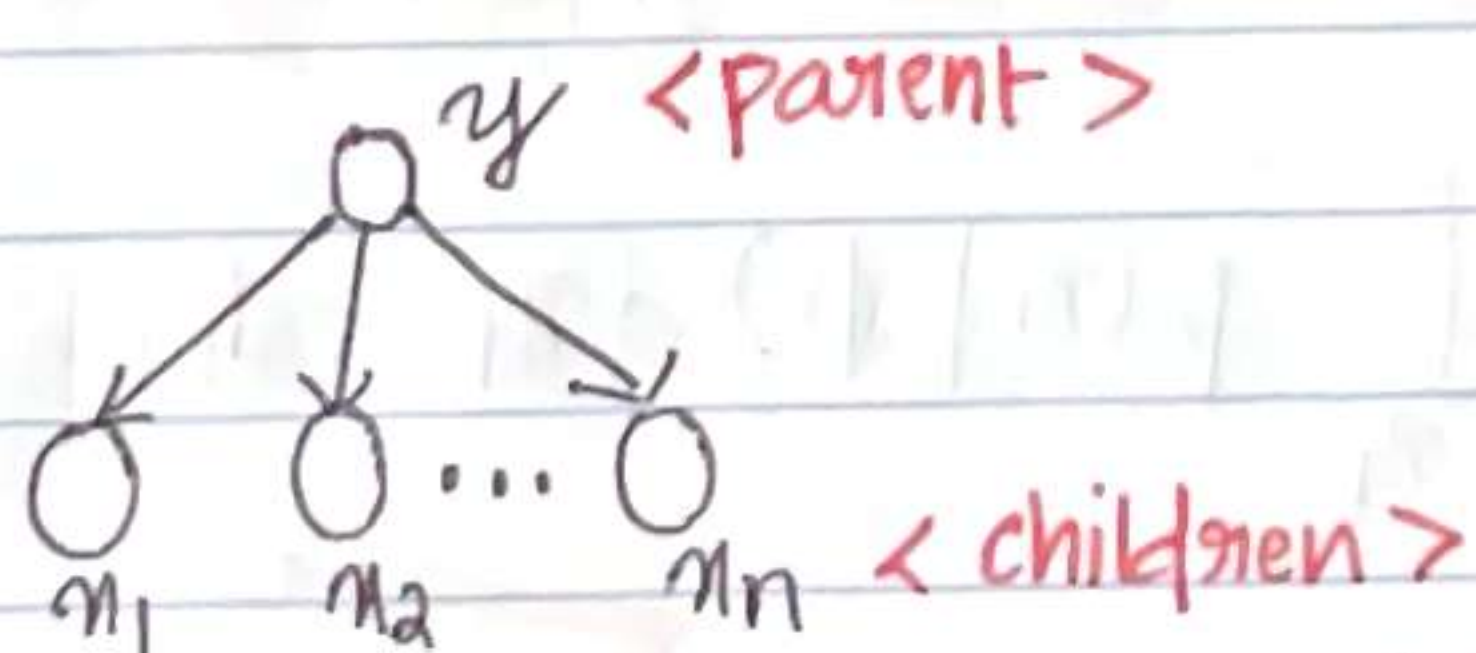
Represent  $P(\bar{x}, y)$  by directed graphs.

$$P(x_1, \dots, x_n | y) = P(x_1 | y) \cdot P(x_2 | y) \dots P(x_n | y)$$

Conditional probability ↘

$$P(\bar{x}, y) = P(\bar{x} | y) \cdot P(y) \\ = P(y) \cdot P(x_1 | y) \dots P(x_n | y)$$

Conditional Independence.



Arrow goes from parent to child.

Suppose,  $x_1 \rightarrow x_2$   
then

$$P(x_1, x_2) = P(x_1) \cdot P(x_2 | x_1)$$

Likewise,  $x_1 \xrightarrow{\text{no edge}} x_2$   
 $P(x_1, x_2) = P(x_1) \cdot P(x_2)$

are discrete value that

\* Suppose  $x_1$  &  $x_2$  take one of  $K$  values.  
i.e. if we take binary then  $K=2$ .  
but we don't exactly care what  
the  $K$  values are in specific.

We can get the  $x_1$  value using,  
1-of- $K$ -representation OR One-hot.  
which means if we assign a  
value to  $x_1$  at a particular position  
then all others are zero. & is  
mathematically represented as (for any  
discrete valued random variable):-

$$P(x_1) = P_1^{x_1^{(1)}} \cdot P_2^{x_1^{(2)}} \dots P_K^{x_1^{(K)}}$$

\* Independent parameter  $\Rightarrow K-1$ .

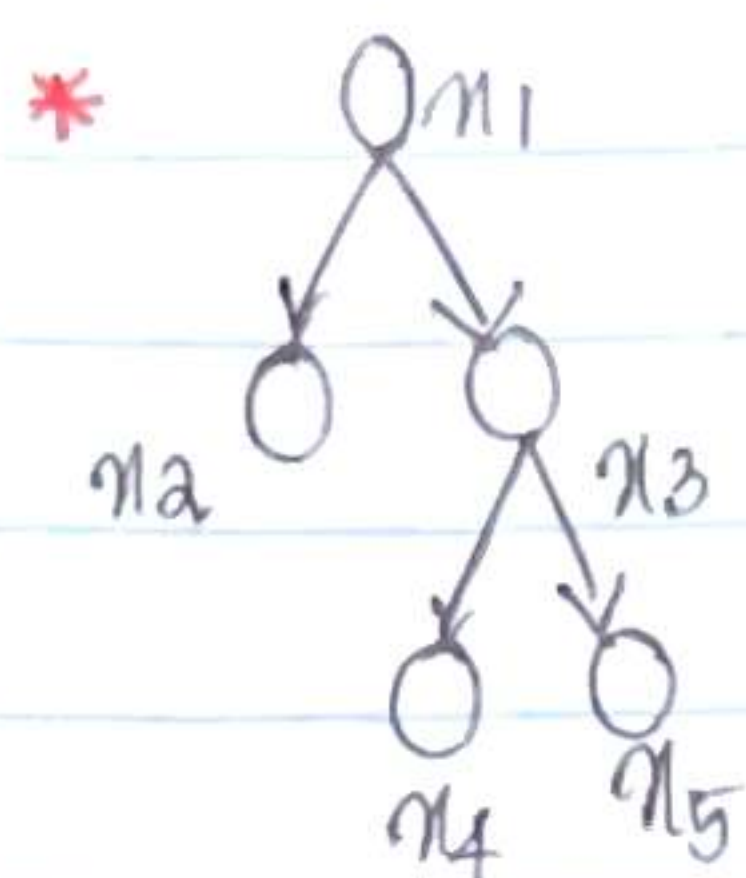
$\therefore$  for  $x_1 \rightarrow x_2$

$$\text{it is } (K-1) + K \cdot (K-1) \\ = \underline{K^2 - 1}$$

& for  $x_1, x_2$ , it is  $(K-1) + (K-1) \\ = \underline{2(K-1)}$

**Summary**  $\Rightarrow$  Making more conditional independence  
assumptions reduces no: of model  
parameters, but the class of distributions  
represented becomes restricted.





$$P(\pi_1, \pi_2 \dots \pi_5) \\ = P(\pi_1) \cdot P(\pi_2 | \pi_1) \cdot \\ P(\pi_3 | \pi_1) \cdot \\ P(\pi_4 | \pi_3) \cdot P(\pi_5 | \pi_3)$$

\* Generally,

$$P(\pi_1, \dots, \pi_k) = \prod_{i=1}^k P(\pi_i | \text{pa}(\pi_i))$$

set of parents

\* The probability distribution should not be negative & the integration should be equal to one.

$$\rightarrow P(\pi) \geq 0.$$

$$\rightarrow \int P(\pi) d\pi = 1.$$

\* Q. When will  $\prod_{i=1}^k P(\pi_i | \text{pa}(\pi_i))$  integrate to 1? Only if this happens then we can say it is a valid assumption.

Assumption A: Assume, (without any loss of generality) that all the directed edges go from lower-indexed to higher-indexed nodes.

Then,

It integrates to 1 if each factor integrates to one.

$$\int \prod_{i=1}^k P(\pi_i | \text{pa}(\pi_i)) d\pi_1 d\pi_2 \dots d\pi_k$$

$\pi_1, \pi_2, \dots, \pi_k$

need to prove this is equal to one.

Doing individually w.r.t.  $\pi_k$ :

$$\int \prod_{i=1}^{k-1} P(\pi_i | \text{pa}(\pi_i)) \left[ \int P(\pi_k | \text{pa}(\pi_k)) d\pi_k \right] d\pi_1 d\pi_2 \dots d\pi_{k-1}$$

This doesn't depend on  $k$ .

We know this integrates to 1 & to apply that we need to ensure that the other parts are not involved with  $\pi_k$ .

$\text{pa}(\pi_i)$  only involves nodes with index  $\leq k-1$ .

$\therefore$ ,

We are left with,

$$\int \prod_{i=1}^{k-1} P(\pi_i | \text{pa}(\pi_i)) d\pi_1 \dots d\pi_{k-1}$$

The same arguments keep repeating, until we end up with

$$\int P(\pi_1 | \phi) d\pi_1 = 1.$$

Q. What does the Assumption A mean in terms of the graph?

Claim  $\rightarrow$  If parents have lower index than the child, then the directed graph is "acyclic".



→ As long as the graph is a directed graph that is acyclic then the probability is valid.

→ Assumption A  $\Rightarrow$  No cycles.  
(implies)

→ There is a cycle  $\Rightarrow$  Assump. A invalidated  
< contrapositive stmt. >

### Summary 2

$\therefore$  The Bayesian n/w is always a DAG.  
A factorisation represented by a DAG is a valid probability distribution

From the above we can further say that, as long as the factors when integrated to gives '1', making it valid then their products will also be valid.

e.g.  $\begin{matrix} \bigcirc \\ \pi_1 \end{matrix} \longrightarrow \begin{matrix} \bigcirc \\ \pi_2 \end{matrix}$

$$\sum_{\pi_2} p(\pi_2 | \pi_1) = 1$$

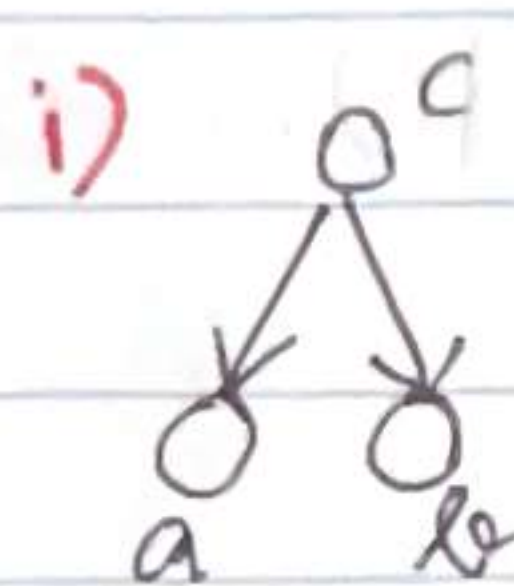
$$\sum_{\pi_1} p(\pi_1) = 1;$$

$$\text{then } \sum_{\pi_1 \pi_2} p(\pi_1) \cdot p(\pi_2 | \pi_1) = 1.$$

= This whole thing / is a comp. is an example of CG.

Goal  $\Rightarrow$  Find the set of CG (Cond. Independence) properties satisfied by a given graph.

Now to find the goal let's look into some basic graphs;



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

$$a \perp\!\!\!\perp b | c$$

means  $a$  is independent of  $b$  condition  $c$ .

Cuz  $a$  doesn't depend upon  $b$ . Thus when  $a$  &  $c$  is observed,  $a$  doesn't depend upon  $b$ , showing CG. This equation basically shows a/v CG defines. Similarly,

As we saw earlier

$$\pi_1 \longrightarrow \pi_2.$$

$$P(\pi_1, \pi_2) = P(\pi_1) \cdot P(\pi_2)$$

< Here it is conditioned on  $\phi$  >.

Thus,

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

(OR)

$$p_c(a, b) := p(a, b | c)$$

(OR)

$$p_c(a, b) = p_c(a | b) p_c(b) = p_c(a) p_c(b).$$

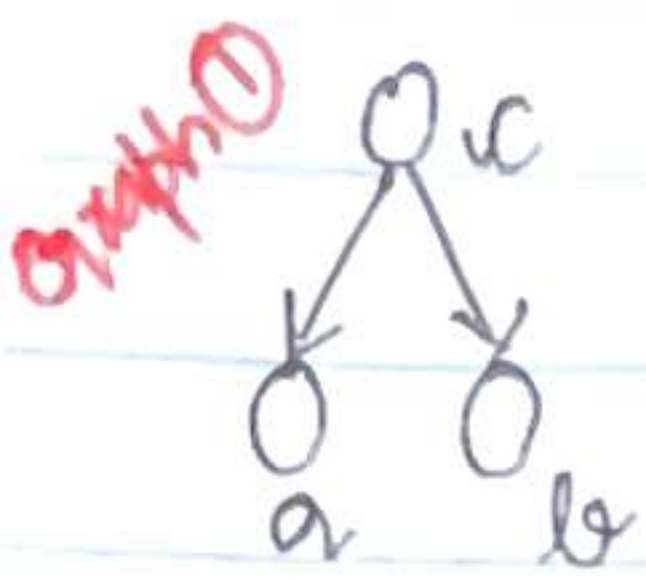
< Just use one for defining CG, everything is basically same >.

< The \* ones are the main definitions used >.

< The above points are generalised version, not base exactly on the graph >



Let's come back to the graph,



$$p(a, b | c) = p(c) \cdot p(a | c) \cdot p(b | c). \quad \text{--- ①}$$

$$p(a, b | \phi) \stackrel{?}{=} p(a, b) \stackrel{?}{=} p(a | \phi) \cdot p(b | \phi) \stackrel{?}{=} p(a) \cdot p(b). \quad \text{--- ②}$$

From ①, how do we get ②?

Doing marginal inference,

$$\sum_c p(a, b, c) = p(a, b).$$

↓

Here since we need to eliminate  $c$  we are taking it like this. Basically we are marginalizing w.r.t.  $c$ .

$$\begin{aligned} p(a, b) &= \sum_c p(a, b, c) \\ &= \sum_c p(a | c) \cdot p(b | c) \\ &\neq p(a) \cdot p(b) \end{aligned}$$

∴ For this graph,

$$p(a, b | \phi) \neq p(a | \phi) \cdot p(b | \phi)$$

∴ In general  $\uparrow$ , is not unconditionally independent.

Now,

$$p(a, b | c) \stackrel{?}{=} p(a | c) p(b | c).$$

Now does this work?

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

Basic conditional probability.

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

Not from ①

$$= p(a | c) p(b | c).$$

Thus, we can say

$$a \perp\!\!\!\perp b | c \quad \& \quad a \not\perp\!\!\!\perp b | \phi.$$

We can also say that, the path of  $a$  &  $b$  is tail-to-tail w.r.t. node  $c$ .

saying it here coz these wordings are commonly used in graphical representations.

We will also say like, when  $c$  is observed the path from  $a$  to  $b$  is blocked.

↑  
All this in relation to this graph. ①.

OMIT



ii) Graph ②  $a \rightarrow c \rightarrow b$

Here,  $c$  is head-to-tail w.r.t path from  $a$  to  $b$ .

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

$$p(a, b | \phi) \stackrel{?}{=} p(a) \cdot p(b)$$

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

$$= \sum_c p(a) \cdot p(b|c) \cdot p(c|a)$$

↑  
Taken out as it doesn't depend on  $c$ .

$$= p(a) \cdot p(b|a)$$

$$\neq p(a) \cdot p(b)$$

Thus,  $a \not\perp b | \phi$

< Not unconditionally independent >

Now,

$$p(a, b | c) \stackrel{?}{=} p(a|c) \cdot p(b|c)$$

$$\begin{aligned} p(a, b | c) &= \frac{p(a, b, c)}{p(c)} \\ &= \frac{p(a) \cdot p(c|a) \cdot p(b|c)}{p(c)} \end{aligned}$$

$$< \frac{p(a) \cdot p(c|a)}{p(c)} = p(a|c) >$$

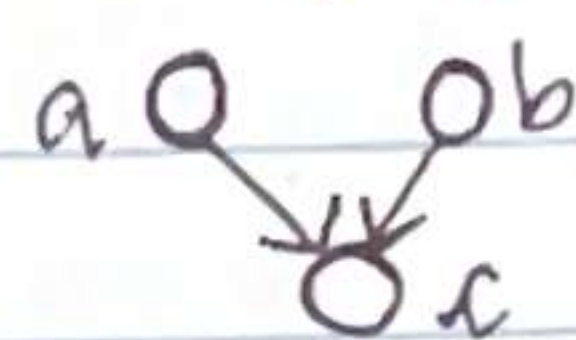
$$= p(a|c) \cdot p(b|c)$$

Thus,  $a \perp b | c$

$$\begin{aligned} < \frac{p(a) \cdot p(c|a)}{p(c)} &= \frac{p(a) \cdot p(c|a)}{p(c) \cdot p(a)} \\ &= \frac{p(c|a)}{p(c)} \\ &= p(a|c) \end{aligned}$$

This  $\perp$  is basically the ' $\perp$ ' we used before indicating the joint distribution  $\perp$  is used in case for sets. >

iii) Graph ③



$$\begin{aligned} p(a, b | \phi) &\stackrel{?}{=} p(a) \cdot p(b) \\ p(a, b | c) &\stackrel{?}{=} p(a|c) \cdot p(b|c) \end{aligned}$$

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

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

(continuing)

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

This will be equal to 1.

$$= p(a) \cdot p(b)$$

Thus,  $a \perp b | \phi$  i.e. they are unconditionally independent.

$$\begin{aligned} p(a, b | c) &= \frac{p(a, b, c)}{p(c)} = \frac{p(a) \cdot p(b) \cdot p(c|a, b)}{p(c)} \\ &\neq p(a|c) \cdot p(b|c) \end{aligned}$$

Thus,  $a \not\perp b | c$



Question  $\Rightarrow$  B (state of battery);  
 F (state of fuel tank);  
 G (state of fuel gauge).

B, F, G = 0 or 1.  
 $\downarrow$  Bad  $\downarrow$  Good < 3rd general >

B  $\begin{matrix} \nearrow 0 \\ \searrow 1 \end{matrix}$  Battery drained.  
 " " charged.

F  $\begin{matrix} \nearrow 0 \\ \searrow 1 \end{matrix}$  Tank empty.  
 " " full.

G  $\begin{matrix} \nearrow 0 \\ \searrow 1 \end{matrix}$  Indicator poor.  
 " " good.

B, F are independent.

$$P(B=1) = 0.9 \quad P(G=1 | B=1, F=1) = 0.8$$

$$P(F=1) = 0.9 \quad P(G=1 | B=1, F=0) = 0.2$$

$$P(G=1 | B=0, F=1) = 0.2$$

$$P(G=1 | B=0, F=0) = 0.1$$

Answer :-

i) Find  $P(F=0 | G=0)$ .

From above,

$$P(F=0) = 0.1$$

$$P(F=0 | G=0) = \frac{P(G=0 | F=0) P(F=0)}{P(G=0)}$$

$$P(G=0 | F=0) = \frac{P(G=0, F=0)}{P(F=0)}$$

$$\downarrow = \sum_{b=0,1} P(G=0 | F=0, B=b) P(B=b)$$

Just trying to add B into it coz that is the form we have in q.

$$P(G=0 | F=0, B=b) \times P(B=b)$$

$$= \frac{P(G=0, F=0, B=b)}{P(F=0, B=b)} \times P(B=b)$$

$$= \frac{P(G=0, F=0, B=b)}{P(F=0) P(B=b)} \times P(B=b)$$

given F & B are independent.

$$= \frac{P(G=0, F=0, B=b)}{P(F=0)}$$

Now we are going to do a sum over b thus due to marginalisation we will get;

$$= \sum_b \left[ \frac{P(G=0, F=0, B=b)}{P(F=0)} \right]$$

$$= \frac{P(G=0, F=0)}{P(F=0)}$$

The above was done to show that incorporating B doesn't cause any problems.

$$\therefore \textcircled{1} P(G=0 | F=0) =$$

$$P(G=0 | F=0, B=0) \cdot P(B=0) + P(G=0 | F=0, B=1) \cdot P(B=1)$$

$$= 0.9 \times 0.1 + 0.8 \times 0.9$$

$$= 0.09 + 0.72$$

$$= \underline{\underline{0.81}}$$



Similarly now we need to find,  
 $P(G=0)$ .

$$\begin{aligned} & P(G=0 | F=f, B=b) \cdot P(B=b) \cdot P(F=f) \\ &= \frac{P(G=0, F=f, B=b) \cdot P(B=b) \cdot P(F=f)}{P(F=f, B=b)} \\ &= \frac{P(G=0, F=f, B=b)}{P(F=f) \cdot P(B=b)} \cdot P(B=b) \cdot P(F=f) \end{aligned}$$

marginalizing over both  $f$  &  $b$ .

$$= \sum_b \sum_f P(G=0, B=b, F=f).$$

$$= \underline{P(G=0)}.$$

$$\Rightarrow 0.2 \times 0.9 \times 0.9 + 0.8 \times 0.9 \times 0.1 + 0.8 \times 0.1 \times 0.9 + 0.9 \times 0.1 \times 0.1$$

Applying to all the given prob.

$$= 0.162 + 0.072 + 0.072 + 0.009$$

$$= \underline{0.315}.$$

$\therefore$

$$\begin{aligned} & P(F=0 | G=0) \\ &= \frac{P(G=0 | F=0) \cdot P(F=0)}{P(G=0)} \end{aligned}$$

$$= \frac{0.81 \times 0.1}{0.315}$$

$$= \underline{0.257}$$

$< \underline{0.1} >$

ii)  $P(F=0 | G=0, B=0)$  "we give battery is drained."

$$\text{Ans: } P(F=0 | G=0)$$

$$= \frac{P_B(G=0 | F=0) \cdot P_B(F=0)}{P_B(G=0)}$$

$$\therefore P_B(G=0 | F=0)$$

$$= P(G=0 | F=0)(F=0) * (B=0) + P(G=0 | F=1)(F=1)(B=0)$$

$$= 0.9 \times 0.1 \times 0.1 + 0.8 \times 0.9 \times 0.1$$

$$= \underline{0.081}$$



\* "When will we have  $A \perp\!\!\!\perp B \mid C$ ?"

• First consider  $P = \{ \text{All paths from any node in } A \text{ to any node in } B \}$

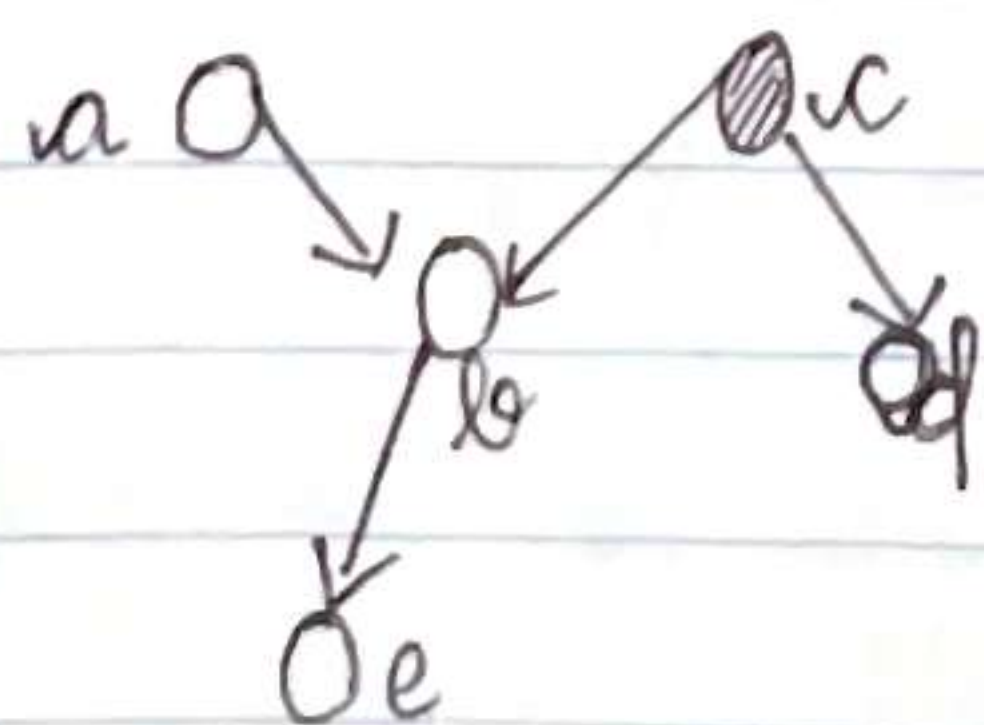
• We will say a path is blocked if either:-

\* i) The path involves a tail-to-tail or head-to-tail node in  $C$  which is observed.

\* ii) The path involves a head-to-head node ~~in~~  $C$  & neither this node nor any of its descendants are observed in  $C$ .

•  $A, B, C$  - group of vertices in the graph.

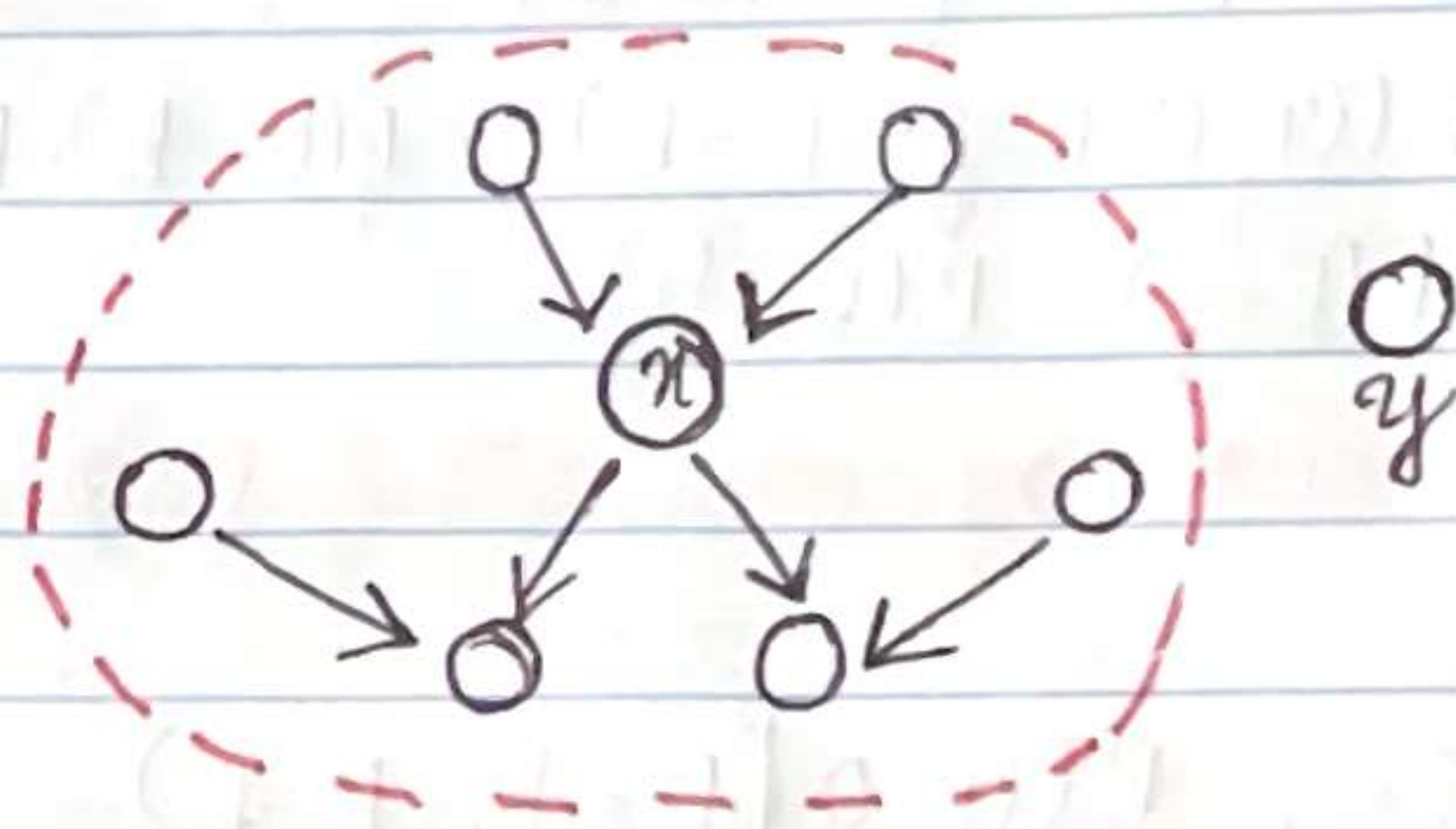
• So based on the above conditions we can draw the conclusion that  $a \perp\!\!\!\perp b \mid c$  or  $b \perp\!\!\!\perp a \mid c$ , etc. for the below graph.



• We should be able to tell if the conditional independencies when a graph is given & this can be done using

the given 2 conditions. This method is called "d-separation".

Q. Show that  $x \perp\!\!\!\perp y$  using d-sep.



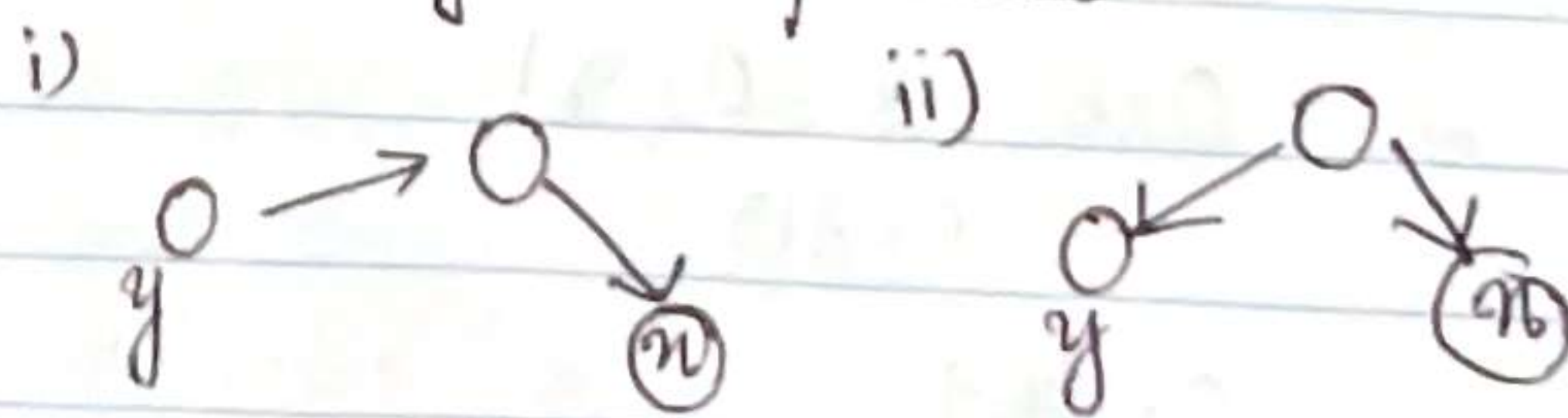
$C_x = \{ \text{parents of } x, \text{ children of } x, \text{ co-parents of } x \}$ .

$$P(x \mid \text{nodes other than } x \text{ in graph}) = P(x \mid C_x)$$

< This is like a Markov property >  
<  $C_x$  is called the Markov blanket of  $x$  >

Let  $A = \{x\}$  &  $B = \{y\}$

$y$  can be connected to  $x$  in 2 ways - via parents:-





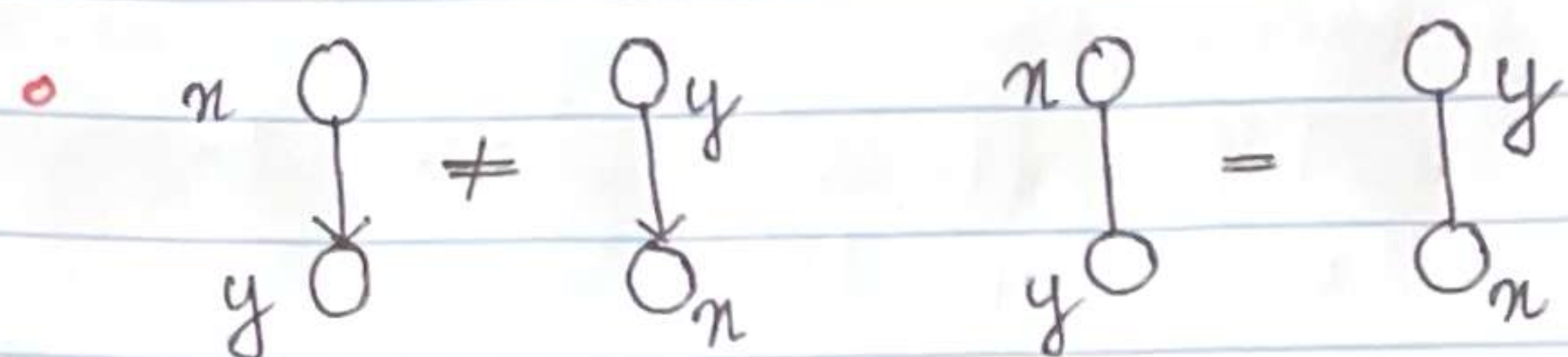
Both these ways satisfy condition 1.  
 $\therefore$  All paths via parents of  $a$  are blocked.

Now let's look into path via children of  $a$ .



## Undirected Graphical Models (Markov n/w's)

- Correlation doesn't mean causation.  
But causation means / implies correlation.



- Here,

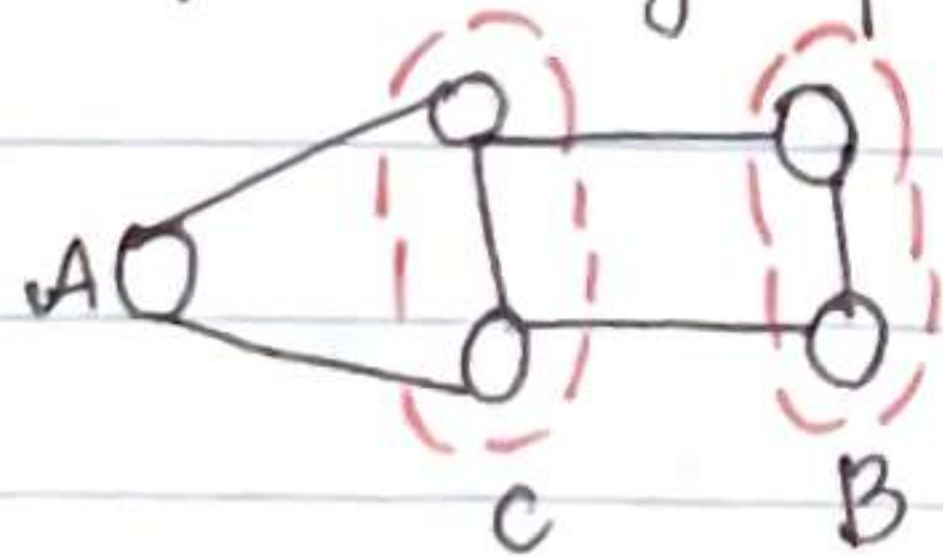
i) Conditional independence will be discussed 1<sup>st</sup> unlike directed graphs where factorization was done 1<sup>st</sup> & then conditional independence. < d-sep & all >

ii) Secondly we will do, factorization that corresponds to the definition of conditional independence.

### Conditional Indep-9

- i) Can be defined by "graph separation"

$A, B, C$  - A group of nodes.



If we remove nodes in  $C$ , & there is no path b/w a node in  $A$  to any node in  $B$ , then " $A \perp\!\!\!\perp B \mid C$ ".

< So in the example graph, now what will be the factorization? >

< When causality doesn't matter much we use undirected graphical models.

When causality is also of importance we use directed models e.g. Auto-encoders >

Let,

$\mathcal{N} \rightarrow$  set of all nodes

$n_i, n_j \rightarrow$  Nodes in graph with no direct edge b/w them.

$$\therefore P(n_i, n_j \mid \mathcal{N} - \{n_i, n_j\}) \text{ \& }$$

$$A = \{n_i\}$$

$$B = \{n_j\}$$

$$C = \mathcal{N} - \{n_i, n_j\}$$

when we remove / if we remove all the nodes in the graph except  $n_i$  &  $n_j$ , there is no direct edge b/w them, thus  $A \perp\!\!\!\perp B \mid C$ .

$$n_i \perp\!\!\!\perp n_j \mid \mathcal{N} / \text{all other nodes.}$$

OMIT



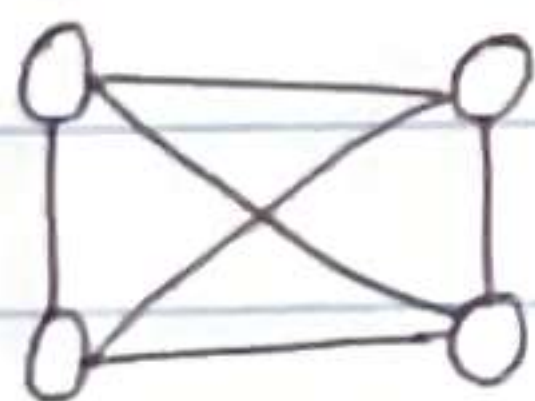
Thus,

$$P(x_i, x_j | x - \{x_i, x_j\})$$

$$= P(x_i | x - \{x_i, x_j\}) \cdot P(x_j | x - \{x_i, x_j\})$$

The factorisation should be such that the nodes that don't have a direct edge b/w them should appear in / as 2 separate factors.

Clique  $\rightarrow$  Subset of nodes where each pair is connected.



$\{x_1, x_2\}$   
 $\{x_1, x_2, x_3\}$   
 $\{x_2, x_3\}$

} All are cliques!

$\{x_2, x_3, x_4\} \rightarrow$  Not a clique.

Maximal cliques  $\rightarrow$  ones with largest size.

$\{x_1, x_2, x_3\}$  &  
 $\{x_1, x_3, x_4\}$

Let  $x_c$  be clique variables,

$$P(x_1, x_2, \dots, x_D)$$

$$= \prod_{\text{maximal clique } c} \psi(x_c)$$

$\rightarrow$  from the above we can understand that  $x_2$  &  $x_4$  needs to be 2 separate factors while writing the factorization.

Also,

$$\psi_c(x_c) \geq 0$$

$Z \rightarrow$  normalisation factor.

$$Z = \sum_{x_c} \prod_c \psi_c(x_c)$$

this will be one.

$\therefore$

$$P(x_1, x_2, \dots, x_D) = \frac{1}{Z} \prod_c \psi_c(x_c)$$

$\psi_c \rightarrow$  potential function.

$Z \rightarrow$  partition function.

$$\psi_c(x_c) = e^{-E_c(x_c)} \quad \text{+ve func.}$$

$E_c \Rightarrow$  energy functions.

If the energy func. is large, then we get a lower value as o/p.

Directed models are used in the area of generative models similarly an app. of undirected graphical modelling would be image denoising.



$X_i \Rightarrow$  Pixels in image  $\in \{-1, +1\}$

- With a probability  $p$ , we flip the sign of  $n_i$  to get a noisy pixel  $y$ .

Goal  $\rightarrow$  Given noisy img.  $y_i$  estimate  $n_i$ .

- When noise is low (i.e.  $p$  is small),  $y_i$  &  $n_i$  are highly correlated.

OR

When the noise is small that means the image is more similar to that of clean image & the not much flipping is required, then the clean & noisy img. values will be correlated.

- Neighbouring pixels in  $n$  are also highly correlated.

Some topics  $\hookrightarrow$

\* Markov blanket

\* Comparison b/w directed & undirected graphical models.