

## Problem 1 →

Q-1-  
Ans - The one downside of Information Gain is that it favours the predictor ~~values~~ variables with a large number of values.

Now if we split the data using these highly branching predictors, then this data might be split into subsets with low Entropy values. Hence, the disadvantages of these splits are that our model might overfit the data as the number of nodes in the tree might increase significantly.

To remedy this Information Gain Ratio introduces a normalizing term called Split Info or Intrinsic Information to reduce the bias of Information Gain. ~~Split Info is given by~~

$$\text{Split Info (F)} = - \sum \frac{|D_j|}{|D|} \cdot \log_2 \left( \frac{|D_j|}{|D|} \right)$$

$$\text{Information Gain Ratio} = \frac{\text{Information Gain}}{\text{Split Info}}$$

→ ~~We choose We split~~

→ Whichever predictor variable has the highest Gain Ratio, ~~we choose it~~, ~~we~~ we choose that predictor for splitting.

Q-2-

Ans - Entropy (Decision) =  $-\sum p(I) \log_2 p(I)$  [where  $I \in \text{Yes, No}$ ]

$$= -p(\text{Yes}) \log_2 p(\text{Yes}) - p(\text{No}) \log_2 p(\text{No})$$

$$= -\left[\frac{9}{14} \log_2 (9/14) + \frac{5}{14} \log_2 (5/14)\right] \quad \text{where } I \in \text{Yes, No}$$

$$= 0.940$$

$$\text{Gain Ratio (F)} = \text{Entropy (Decision)}$$

$$= \sum p(\text{Decision} | F) \cdot \text{Entropy (Decision} | F)$$

$$\text{Gain Ratio (F)} = \frac{\text{Gain (F)}}{\text{SplitInfo (F)}}$$

$$\text{SplitInfo (F)} = -\sum \frac{|D_j|}{|D|} \cdot \log_2 \left( \frac{|D_j|}{|D|} \right)$$

Homework →

$$\text{Gain (D, Homework)} = \text{Entropy (Decision)} - [p(\text{Decision} | \text{Homework} = \text{Much})$$

$$\cdot \text{Entropy (Decision} | \text{Homework} = \text{Much}) + p(\text{Decision} | \text{Homework} = \text{Normal})$$

$$\cdot \text{Entropy (Decision} | \text{Homework} = \text{Normal}) + p(\text{Decision} | \text{Homework} = \text{None})$$

$$\cdot \text{Entropy (Decision} | \text{Homework} = \text{None})]$$



$$\begin{aligned} \text{Entropy (Decision | Homework = Much)} &= - \left[ \frac{2}{5} \log_2 \left( \frac{2}{5} \right) + \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \right] \\ &= - [0.5288 + 0.4422] \\ &= 0.971 \end{aligned}$$

$$\text{Entropy (Decision | Homework = Normal)} = - \left[ \frac{4}{4} \log_2 \left( \frac{4}{4} \right) + 0 \right] = 0$$

$$\begin{aligned} \text{Entropy (Decision | Homework = None)} &= - \left[ \frac{3}{5} \log_2 \left( \frac{3}{5} \right) + \frac{2}{5} \log_2 \left( \frac{2}{5} \right) \right] \\ &= 0.971 \end{aligned}$$

$$\begin{aligned} \text{Gain (Decision, Homework)} &= 0.940 - \left[ \frac{5}{14} \cdot 0.971 + \frac{5}{14} \cdot 0.971 + \frac{4}{14} \cdot 0 \right] \\ &= 0.940 - \left( \frac{5}{7} \cdot 0.971 \right) \\ &= 0.2464285714 \end{aligned}$$

$$\begin{aligned} \text{SplitInfo (Decision, Homework)} &= - \left[ \frac{5}{14} \log_2 \left( \frac{5}{14} \right) + \frac{4}{14} \log_2 \left( \frac{4}{14} \right) + \frac{5}{14} \log_2 \left( \frac{5}{14} \right) \right] \\ &= - \left[ \frac{5}{7} \log_2 \left( \frac{5}{14} \right) + \frac{2}{7} \log_2 \left( \frac{4}{14} \right) \right] \\ &= - [-1.061 - 0.5164] \\ &= 1.5774 \end{aligned}$$

$$\text{Gain Ratio (Homework)} = \frac{0.2464285714}{1.5774} = 0.1562245286$$

Traffic  $\rightarrow$

$$\text{Gain}(\text{Decision}, \text{Traffic}) = \text{Entropy}(\text{Decision}) - [p(\text{Decision}|\text{Traffic} = \text{Busy}) \cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Busy}) + p(\text{Decision}|\text{Traffic} = \text{Ok}) \cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Ok}) + p(\text{Decision}|\text{Traffic} = \text{Chill}) \cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Chill})]$$

$$\cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Busy}) + p(\text{Decision}|\text{Traffic} = \text{Ok}) \cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Ok}) + p(\text{Decision}|\text{Traffic} = \text{Chill}) \cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Chill})]$$

$$\text{Entropy}(\text{Decision}|\text{Traffic} = \text{Ok}) + p(\text{Decision}|\text{Traffic} = \text{Chill}) \cdot \text{Entropy}(\text{Decision}|\text{Traffic} = \text{Chill})]$$

$$\text{Entropy}(\text{Decision}|\text{Traffic} = \text{Chill})]$$

$$\text{Entropy}(\text{Decision}|\text{Traffic} = \text{Busy}) = - \left[ \frac{2}{4} \log_2 \left( \frac{2}{4} \right) + \frac{2}{4} \log_2 \left( \frac{2}{4} \right) \right] = 1$$

$$\text{Entropy}(\text{Decision}|\text{Traffic} = \text{Ok}) = - \left[ \frac{1}{6} \log_2 \left( \frac{1}{6} \right) + \frac{2}{6} \log_2 \left( \frac{2}{6} \right) \right] = 0.9183333$$

$$\text{Entropy}(\text{Decision}|\text{Traffic} = \text{Chill}) = - \left[ \frac{3}{4} \log_2 \left( \frac{3}{4} \right) + \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right] = -[-0.31125 - 0.5] = 0.81125$$

$$\text{Gain}(\text{Decision}, \text{Traffic}) = 0.940 - \left[ \frac{4}{14} \cdot 1 + \frac{6}{14} \cdot 0.9183333 + \frac{4}{14} \cdot 0.81125 \right] = 0.0289285857$$

$$\text{SplitInfo}(\text{Decision}, \text{Traffic}) = - \left[ \frac{4}{14} \log_2 \left( \frac{4}{14} \right) + \frac{6}{14} \log_2 \left( \frac{6}{14} \right) + \frac{4}{14} \log_2 \left( \frac{4}{14} \right) \right] = 1.556685793$$



$$\text{GainRatio}(\text{Traffic}) = \frac{0.0289285857}{1.5566857143} = \underline{\underline{0.0185834465}}$$

Hunger →

$$\text{Gain}(\text{Decision}, \text{Hunger}) = \text{Entropy}(\text{Decision}) -$$

$$\left[ p(\text{Decision} | \text{Hunger} = \text{A little}) \cdot \text{Entropy}(\text{Decision} | \text{Hunger} = \text{A little}) + p(\text{Decision} | \text{Hunger} = \text{A lot}) \cdot \text{Entropy}(\text{Decision} | \text{Hunger} = \text{A lot}) \right]$$

$$\text{Entropy}(\text{Decision} | \text{Hunger} = \text{A little}) = - \left[ \frac{3}{7} \log_2 \left( \frac{3}{7} \right) + \frac{4}{7} \log_2 \left( \frac{4}{7} \right) \right]$$

$$= 0.9852571429$$

$$\text{Entropy}(\text{Decision} | \text{Hunger} = \text{A lot}) = - \left[ \frac{6}{7} \log_2 \left( \frac{6}{7} \right) + \frac{1}{7} \log_2 \left( \frac{1}{7} \right) \right]$$

$$= 0.5916857143$$

$$\text{Gain}(\text{Decision}, \text{Hunger}) = 0.940 - \left[ \frac{7}{14} \cdot 0.9852571429 + \frac{7}{14} \cdot (0.5916857143) \right]$$

$$= \underline{\underline{0.1515285714}}$$

$$\text{SplitInfo}(\text{Decision}, \text{Hunger}) = - \left[ \frac{7}{14} \log_2 \left( \frac{7}{14} \right) + \frac{7}{14} \log_2 \left( \frac{7}{14} \right) \right]$$

$$= 1$$

$$\text{GainRatio}(\text{Hunger}) = \underline{\underline{0.1515285714}}$$

Lauren →

$$\text{Gain}(\text{Decision}, \text{Lauren}) = \text{Entropy}(\text{Decision}) - [p(\text{Decision})_{\text{Lauren=Available}} + p(\text{Decision})_{\text{Lauren=Not Available}}]$$

$$\text{Entropy}(\text{Decision} | \text{Lauren=Available}) + p(\text{Decision} | \text{Lauren=Not Available})$$

$$\text{Entropy}(\text{Decision} | \text{Lauren=Not Available})$$

$$\begin{aligned} \text{Entropy}(\text{Decision} | \text{Lauren=Available}) &= - \left[ \frac{6}{8} \log_2 \left( \frac{6}{8} \right) + \frac{2}{8} \log_2 \left( \frac{2}{8} \right) \right] \\ &= 0.81125 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Decision} | \text{Lauren=Not Available}) &= - \left[ \frac{3}{6} \log_2 \left( \frac{3}{6} \right) + \frac{3}{6} \log_2 \left( \frac{3}{6} \right) \right] \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Decision}, \text{Lauren}) &= 0.940 - \left[ \frac{8}{14} (0.81125) + \frac{6}{14} (1) \right] \\ &= 0.0478571429 \end{aligned}$$

$$\begin{aligned} \text{SplitInfo}(\text{Decision}, \text{Lauren}) &= - \left[ \frac{8}{14} \log_2 \left( \frac{8}{14} \right) + \frac{6}{14} \log_2 \left( \frac{6}{14} \right) \right] \\ &= 0.9852571429 \end{aligned}$$

$$\text{GainRatio}(\text{Lauren}) = \frac{0.0478571429}{0.9852571429} = 0.0485732514$$



If we can see that, Homework feature has the highest gain ratio, hence Homework feature will be our first split in the decision tree.

Q-3- ~~Ans~~ - The first node of our decision tree will be Homework, thus we will again make the table based on different values of Homework  $\Rightarrow$  (Homework = Much)

Homework	Traffic	Hunger	Lauren	Go Out?
Much	Busy	A little	Available	No
Much	Busy	A little	Not Available	No
Much	Ok	A little	Available	No
Much	chill	A lot	Available	Yes
Much	Ok	A lot	Not Available	Yes

From this table we can see that Going out is directly correlated with Hunger & if Homework is much.

(Homework = Normal)

Homework	Traffic	Hunger	Lauren	Go Out?
Normal	Busy	A little	Available	Yes
Normal	chill	A lot	Not Available	Yes
Normal	Ok	A little	Not Available	Yes
Normal	Busy	A lot	Available	Yes

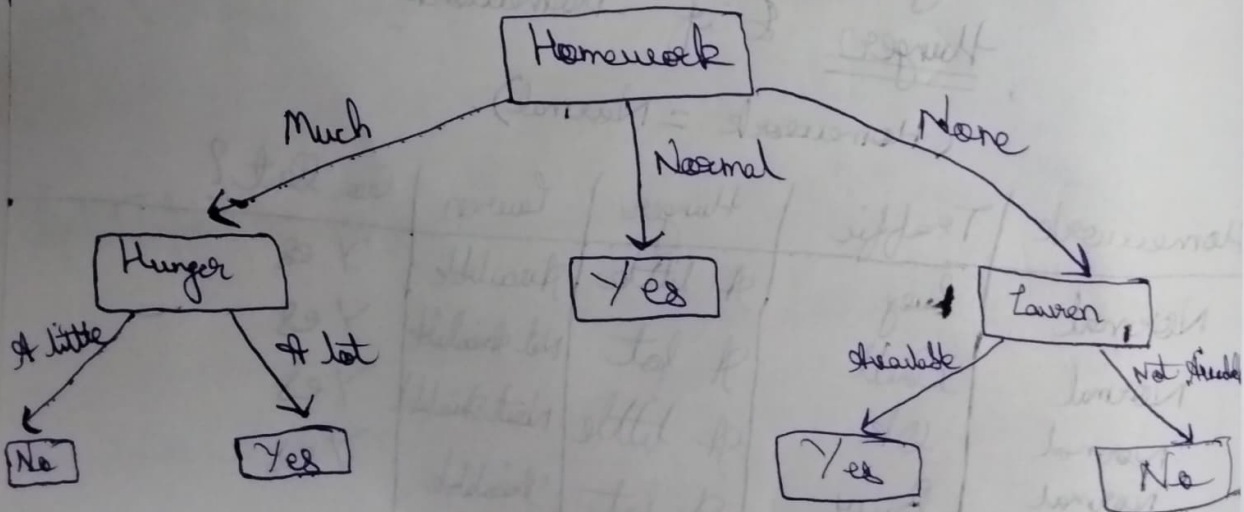
From this table we can see that if Homework is Normal then you will always go out.

(Homework = None)

Homework	Traffic	Hunger	Lauren	Go Out?
None	Ok	A little	Available	Yes
None	chill	A lot	Available	Yes
None	chill	A lot	Not Available	No
None	Ok	A lot	Available	Yes
None	Ok	A little	Not Available	No

From this table we can see that the decision to go out is directly correlated to availability of Lauren if there is no Homework.

Hence from above three tables, we can construct our decision tree as →





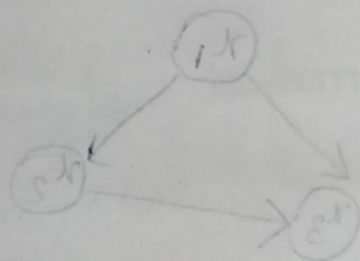
Q-4-

Ans- Since, I have a normal amount of homework,

therefore according to the decision tree, I will go out.

②

for restricted tree where we let  
 $\epsilon^x, s^x, x$  ~~area (1,2,3)~~



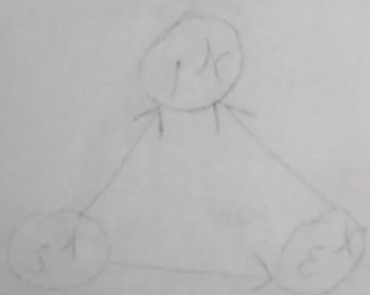
$$(x/x)q(x/sx)q(1x)q = (\epsilon^x s^x, x)q \text{ not}$$

if, given all grade we find

← if we give all the blind

$$(x/x)q(\epsilon^x)q = (\epsilon^x s^x, x)q \text{ not}$$

$$(s^x, x/1x)q$$

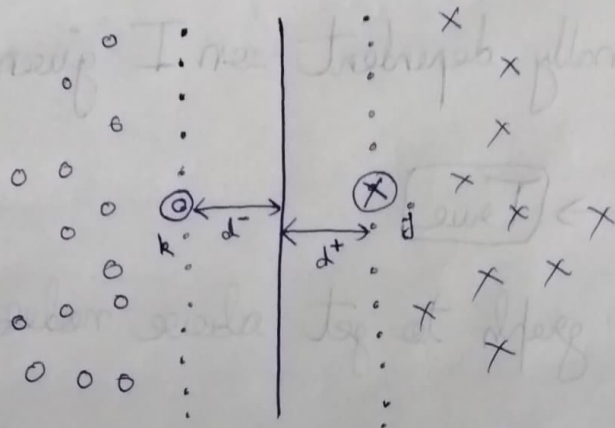


## Problem-2

Linearly separable data



Q-1-  
Ans-



$$y_i = -1$$

$$y_i = 1$$

For class  $y_i = 1 \Rightarrow$

$$w^T x_i + b \geq 1$$

$$\text{ie } w^T x_j + b = 1 \Rightarrow b = 1 - w^T x_j$$

Similarly

For class  $y_i = -1 \Rightarrow$

$$w^T x_i + b \leq -1$$

ie

$$w^T x_k + b = -1 \Rightarrow b = -1 - w^T x_k$$

$$\therefore d^+ - d^- = \frac{w^T x_i}{\|w\|_2} - \frac{w^T x_k}{\|w\|_2} = \frac{1-b}{\|w\|_2} - \frac{(-1-b)}{\|w\|_2}$$

$$\Rightarrow d^+ - d^- = \frac{2}{\|w\|_2}$$



∴ This SVM problem can be defined as  $\Rightarrow$

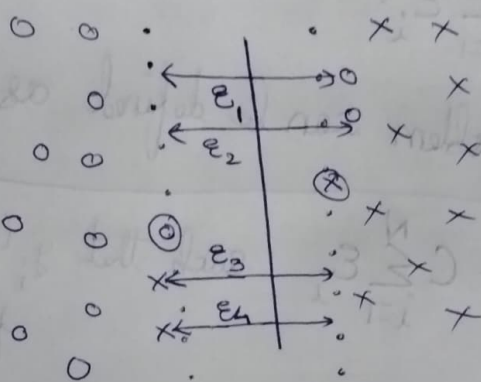
$$\begin{aligned} \max_{w, b} \quad & \frac{2}{\|w\|_2} \quad \text{such that } y_i (w^T x_i + b) \geq 1 \\ & \text{for } i \in \{1, \dots, N\} \end{aligned}$$

OR

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|_2^2 \quad \text{such that } y_i (w^T x_i + b) \geq 1 \\ & \text{for } i \in \{1, \dots, N\} \end{aligned}$$

↓ denotes no. of data points

Q-2-  
Ans-



Now if our data is not linearly separable, then we might have to deal with slack variables. The slack variables are used to minimize the misclassifications along with the objective function defined in the above problem  $\frac{1}{2} \|w\|_2^2$  (See the above figure)

In the case data points are correctly classified,  
 $y_i (w^T x_i + b) \geq 1$

In the case data is misclassified,

$$y_i(w^T x_i + b) < 1$$

The slack variable that we will be using can be defined as  $\Rightarrow$

~~$\epsilon_i$~~   $\epsilon_1, \epsilon_2, \dots, \epsilon_N$  for misclassifiers

(for  $i \in 1, \dots, N$ ) [See the above figure]

$\therefore$  Now our objective function becomes  $\Rightarrow$

$$\frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \epsilon_i$$

$\therefore$  This SVM problem can be defined as  $\Rightarrow$

$$\min_{w, b, \epsilon_i} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \epsilon_i \quad \text{such that } y_i(w^T x_i + b) \geq 1 - \epsilon_i$$

for  $i \in 1, \dots, N$

$$\epsilon_i \geq 0 \quad \text{for } i \in 1, \dots, N$$

Q-3-

Ans- SVM doesn't support multiclass classification directly. But we can still perform multiclass classification by breaking it down into smaller binary classification ~~sub~~ subproblems. There are many approaches like One vs One, One vs All and Directed Acyclic graph ~~graph~~.



- In One vs One we breakdown our multiclass classification problem into various ~~a~~ binary classification subproblems. For the final prediction we use the concept of majority voting.

- In One vs All approach we train  $N$  SVMs

SVM ①  $\rightarrow$  for class 1

SVM ②  $\rightarrow$  for class 2

$\vdots$

SVM ④  $\rightarrow$  for class  $N$

Now to predict the output of new inputs, predict output using all the SVMs & then just identify which model puts the prediction farthest into the positive region.

- In Directed Acyclic Graph approach we try to first group the classes ~~a~~ based on some logical ~~reasoning~~ grouping & then train SVMs. Thus, at the end we might need to train  $\phi$  less number of SVMs and this approach also reduces the diversity from the majority class.

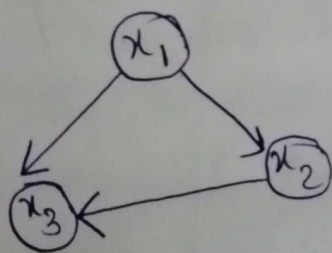
### Problem 3 →

Q-1-

Ans- (1) We are given  $N$  random

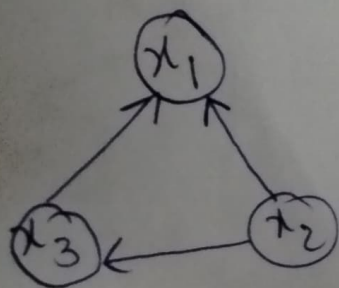
variables  $x_1, x_2, \dots, x_N$ . And now we want to graph a Bayesian Network showing the joint distribution of all  $N$  variables.

Let us consider joint distribution of  ~~$p(a, b, c)$~~   $x_1, x_2, x_3$  random variables.



$$\text{Then } p(x_1, x_2, x_3) = p(x_1) \cdot p(x_2 | x_1) p(x_3 | x_1, x_2)$$

Similarly if we change the ordering, & build the following graph ⇒



$$\text{then } p(x_1, x_2, x_3) = p(x_3) \cdot p(x_3 | x_2) \cdot p(x_1 | x_3, x_2)$$



Thus, if we list out all the possibilities, then we will see that for 3 random variables we can make total  $3! = 6$  distinct graphs.

Hence, for  $N$  random variables, we can ~~say that~~ make  $N!$  distinct graphs.

② For a single discrete variable  $x$ , having  $M$  possible states, the probability distribution  $p(x|\sigma)$  is given by :-

$$p(x|\sigma) = \prod_{m=1}^M \sigma_m^{x_m}$$

Similarly for ~~2~~ discrete random variables  $x_1, x_2$ ; we have two states  $\rightarrow s_1, s_2$

$$\therefore p(x_1, x_2 | \sigma) = \prod_{a=1}^{s_1} \prod_{b=1}^{s_2} \sigma_{ab}^{x_1 a x_2 b}$$

& we can see that

total number of parameters here  $\therefore$   
are  $= (s_1 \times s_2 - 1)$

Thus, after generalizing this we can say that if we have  $N$  random variables  $x_1, x_2, x_3, \dots, x_N$  & they have  $N$  states  $s_1, s_2, s_3, \dots, s_N$  then total

(number) of parameters in the joint probability distribution will be  $\vdash [(s_1 \times s_2 \times s_3 \dots \times s_n) - 1]$

This shows that we will eventually end up with a graph that grows exponentially and ~~thus it would be very~~ so the task to calculate joint probability in this case will become very complex.

Now if we consider all ~~our~~ our random variables to be conditionally independent then,

for two random variables  $x_1, x_2$

joint probability  $p(x_1, x_2 | \sigma)$  becomes  $\rightarrow$

$$p(x_1, x_2 | \sigma) = \prod_{a=1}^{s_1} \sigma_{1a}^{x_{1a}} \prod_{b=1}^{s_2} \sigma_{2b}^{x_{2b}}$$

& no. of parameters in this joint probability will be  $(s_1 + s_2 - 1)$

After  
Generalizing this we can say that ~~our~~  $N$  random variables with  $N$  states will have  $\Rightarrow$  parameters  $(s_1 + s_2 + \dots + s_N - 1)$

& we will get a graph which grows linearly ~~as opposed~~ instead of exponentially.

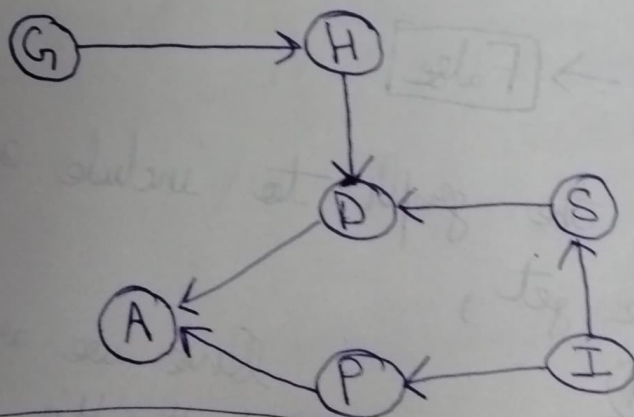


Hence, if we maintain conditional independence of all  $N$  random variables, then it will be less complex to calculate joint probability  $p(x_1, x_2, \dots, x_N)$  as the graph will grow ~~linearly~~ linearly & not exponentially.

Q-2-

- Ans-①  $G \Rightarrow$  Going out right now  
 $H \Rightarrow$  contracting Hake plague  
 $D \Rightarrow$  Drowsy  
 $S \Rightarrow$  Lack of sleep  
 $I \Rightarrow$  Insomnia  
 $P \Rightarrow$  Pay Attention  
 $A \Rightarrow$  Unable to finish the assignment.

Bayesian Network  $\Rightarrow$



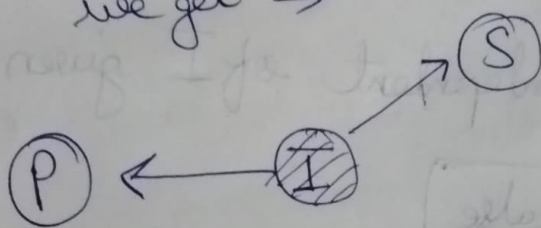
$$p(G, H, D, S, I, P, A) = p(G) \cdot p(H|G) \cdot p(D|H, S) \cdot p(S|I) \cdot p(I) \cdot p(P|I) \cdot p(A|P, D)$$

Nodes in the  $S$ 's Markov Blanket  $\Rightarrow$

$I, D \& H$  is ~~not~~ if we know the values of these nodes, then we can determine the value of  $S$ .

②  $S \perp\!\!\!\perp P | I \rightarrow \boxed{\text{True}}$

Simplifying the graph to get above nodes, we get  $\Rightarrow$



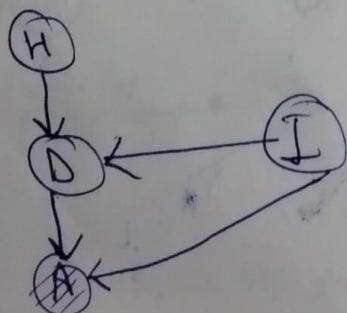
Since,  $I$  is dependent on both  $S \& P$  ~~and we~~

~~know the value of  $I$~~  this even if we know the value of  $I$  still  $P \& S$  are independent

$\therefore S$  is conditionally independent of  $P$  given  $I$

③  $H \perp\!\!\!\perp I | A \rightarrow \boxed{\text{False}}$

Simplifying the graph to get above nodes, we get  $\Rightarrow$



Since  $A$  is dependent on  $D \& I$ , &  $D$  is dependent on both  $I \& H$

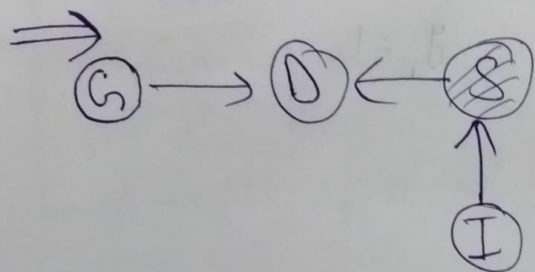


Hence,  $H$  is also dependent on  $I$ , even when  $A$  is given

$\therefore H$  is conditionally dependent on  $I$  given  $A$ .

$$\textcircled{4} G \perp\!\!\!\perp I | S \rightarrow \boxed{\text{True}}$$

Simplifying the graph to get above nodes, we get



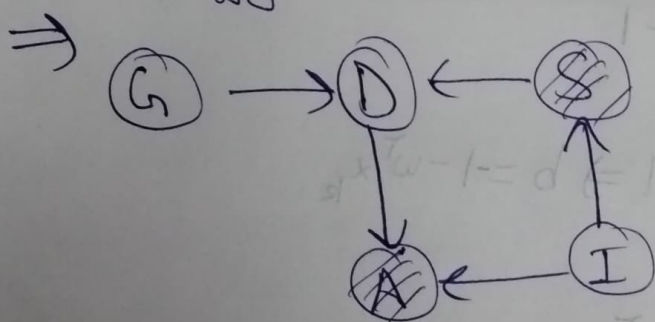
Since,  $D$  is dependent on both  $G$  &  $S$  and  $S$  is dependent on  $I$

Hence,  $G$  is independent of  $I$  even when  $S$  is given

$\therefore G$  is conditionally independent on  $I$  given  $S$ .

$$\textcircled{5} G \perp\!\!\!\perp I | S, A \rightarrow \boxed{\text{False}}$$

Simplifying the graph to get above nodes, we get



Since,  $A$  is dependent on both  $I$  &  $D$  and

$D$  is dependent on both  $S$  &  $G$  and

$S$  is dependent on  $I$

Hence,  ~~$G$~~   $G$  is dependent on  $I$  even when both  $S$  &  $A$  are known.

$\therefore G$  is dependent on  $I$  given  $S$  and  $A$ .