

Advanced Machine LearningHW : 2Problem : 1

- (Q1) The downside of Information Gain is that it favours the predictor variables with a large number of values. Thus, if the data is split using highly branching predictors, then the data might get split into subsets with low entropy. Thus, this would lead to possible overfitting or the number of nodes in the tree might increase by a lot.

This is why Information Gain Ratio introduces a normalizing term, known as Split-Info to reduce the bias.

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

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

Finally, whichever predictor variable has the highest gain Ratio is chosen as the predictor for splitting.

(Q2) Entropy(Decision) = $-\sum p(J) \log_2 p(J)$

[where J ∈ output classes]

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

$$= - \left[\frac{9}{14} \log_2 \left(\frac{9}{14} \right) + \frac{5}{14} \log_2 \left(\frac{5}{14} \right) \right]$$

$$\approx 0.94$$

$$\text{Gain}(F) = \text{Entropy}(\text{Decision}) - \sum p(\text{Decision} | F) \times \text{Entropy}(\text{Decision} | F)$$

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

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

Homework

$$\text{Gain}(D, \text{Homework}) = \text{Entropy}(\text{Decision})$$

$$- [p(\text{Decision} | \text{Homework} = \text{Much})$$

$$\times \text{Entropy}(\text{Decision} | \text{Homework} = \text{Much}) +$$

$$p(\text{Decision} | \text{Homework} = \text{Normal}) \times \text{Entropy}(\text{Decision}, \text{Homework} = \text{Normal}) +$$

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

→ Entropy (Decision | Homework = Mech)

$$= - \left[\frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right]$$

$$= - [0.53 + 0.44]$$

$$= 0.97$$

→ Entropy (Decision | Homework = Normal)

$$= - \left[\frac{4}{4} \log_2 \left(\frac{4}{4} \right) + 0 \right] = 0$$

→ Entropy (Decision | Homework = None)

$$= - \left[\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] = 0.97$$

→ Gain (Decision, Homework) = $0.94 - \left[\frac{5}{14} \times 0.97 \right.$

$$\left. + \frac{5}{14} \times 0.97 + \frac{4}{14} \times 0 \right]$$

$$= 0.94 - \left(\frac{5}{14} \times 0.97 \right)$$

$$\approx 0.246$$

→ Shd. Info (Decision, Homework) = $- \left[\frac{5}{14} \log_2 \left(\frac{5}{14} \right) + \right.$

$$\left. \frac{4}{14} \log_2 \left(\frac{4}{14} \right) + \frac{5}{14} \log_2 \left(\frac{5}{14} \right) \right] = 1.573$$

$$\rightarrow \text{Gain Ratio (Homework)} = \frac{0.246}{1.577} = 0.156$$

Traffic

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

$$= [p(\text{Decision} | \text{Traffic} = \text{Busy}) \times \text{Entropy}(\text{Decision} | \text{Traffic} = \text{Busy}) + p(\text{Decision} | \text{Traffic} = \text{OK}) \times$$

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

$$\rightarrow \text{Entropy}(\text{Decision} | \text{Traffic} = \text{Busy})$$

$$= -\left[\frac{2}{4} \log_2\left(\frac{2}{4}\right) + \frac{2}{4} \log_2\left(\frac{2}{2}\right)\right] = 1$$

$$\rightarrow \text{Entropy}(\text{Decision} | \text{Traffic} = \text{OK}) = -\left[\frac{4}{6} \log_2\left(\frac{4}{6}\right) + \frac{2}{6} \log_2\left(\frac{2}{6}\right)\right] = 0.92$$

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

$$= 0.81$$

$$\rightarrow \text{Gain}(\text{Decision}, \text{Traffic}) = 0.94 - \left[\frac{3}{14} \times 1 + \frac{6}{14} \times 0.92 + \frac{5}{14} \times 0.81 \right] \\ = 0.029$$

$$\rightarrow \text{Split Info}(\text{Decision}, \text{Traffic}) = - \left[\frac{3}{14} \log_2 \left(\frac{3}{14} \right) + \frac{6}{14} \log_2 \left(\frac{6}{14} \right) + \frac{5}{14} \log_2 \left(\frac{5}{14} \right) \right] \\ = 1.56$$

* Gain Ratio(Traffic) = $\frac{0.029}{1.56} = 0.0186$

Hunger

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

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

$$\rightarrow \text{Entropy}(\text{Decision} | \text{Hunger} = \text{A little})$$

$$= - \left[\frac{3}{7} \log_2 \left(\frac{3}{7} \right) + \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right] \\ = 0.985$$

→ Entropy (Decision | Hunger = Available)

$$= - \left[\frac{6}{7} \log_2 \left(\frac{6}{7} \right) + \frac{1}{7} \log_2 \left(\frac{1}{7} \right) \right]$$

$$= 0.59$$

→ Gain (Decision | Hunger) = $0.94 - \left[\frac{7}{14} \times 0.98 \right]$

$$+ \frac{7}{14} \times 0.59$$

$$= 0.15$$

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

$$= 1$$

* Gain Ratio (Hunger) = 0.15

Lauren

Gain (Decision, Lauren) = Entropy (Decision)

$$- \left[p(\text{Decision} | \text{Lauren} = \text{Available}) \times \text{Entropy} (\text{Decision}) \right.$$

$$\left. + p(\text{Decision} | \text{Lauren} = \text{Not Available}) \times \text{Entropy} (\text{Decision} | \text{Lauren} = \text{Not Available}) \right]$$

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

→ Entropy (Decision | Lauren = Available)

$$= - \left[\frac{6}{8} \log_2 \left(\frac{6}{8} \right) + \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right]$$

$$= 0.81$$

→ Entropy (Decision | Lauren = Not Available)

$$= - \left[\frac{3}{6} \log_2 \left(\frac{3}{6} \right) + \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right]$$

$$= 1$$

$$\rightarrow \text{Gain}(\text{Decision, Lauren}) = 0.94 - \left[\frac{8}{14} \times 0.8 + \frac{6}{14} \times 1 \right]$$

$$= 0.048$$

$$\rightarrow \text{SplitInfo}(\text{Decision, Lauren}) = - \left[\frac{8}{14} \log_2 \left(\frac{8}{14} \right) + \frac{6}{14} \log_2 \left(\frac{6}{14} \right) \right] = 0.98$$

$$\rightarrow \text{Gain Ratio}(\text{Lauren}) = \frac{0.048}{0.98} = \boxed{0.048}$$

As we can see that homework has highest gain ratio, so homework will be the first split in the decision tree.

Q3) From previous question, root node will be Homework; so the decision tree is:
 (Homework = Much)

Homework	Traffic	Hunger	Lawyer	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

thus, going out has a ~~five~~ correlation with Hunger, if Homework is much

(Homework = Normal)

Homework	Traffic	Hunger	Lawyer	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

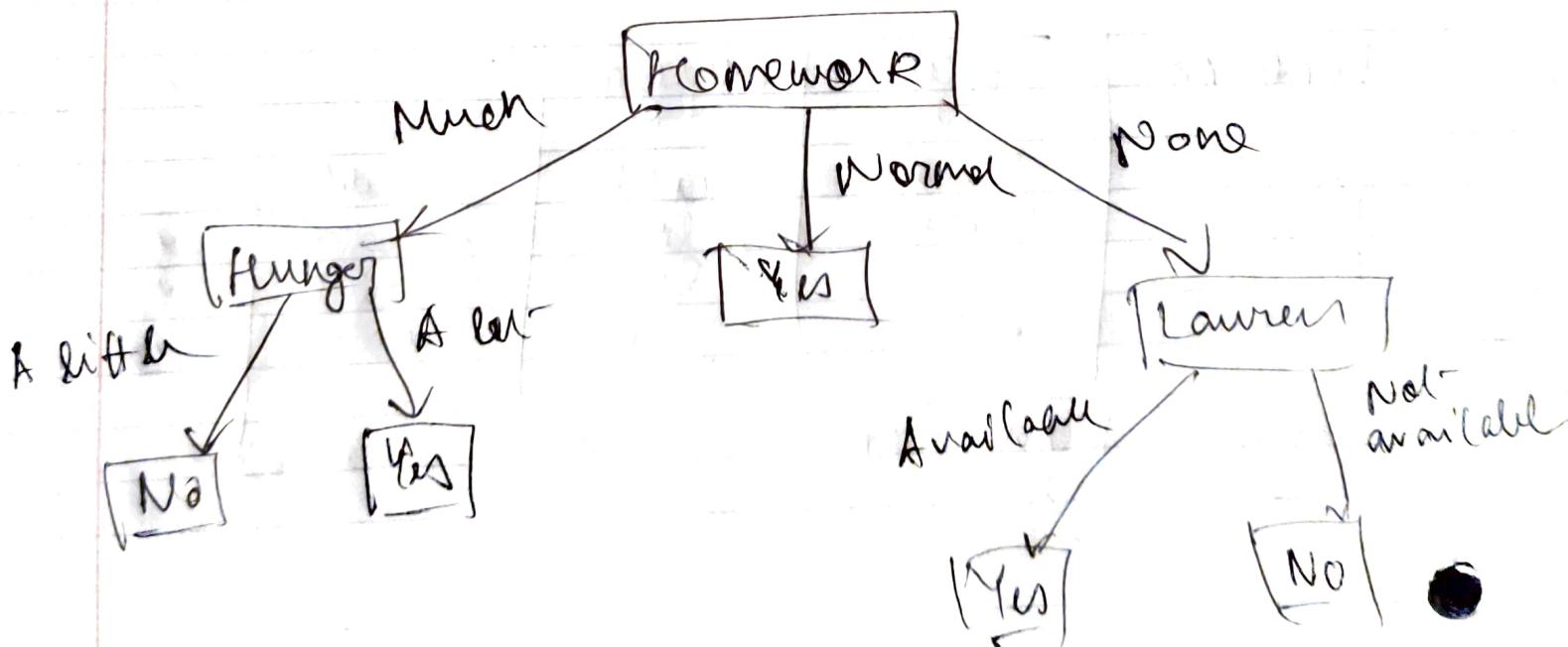
thus, if Homework is Normal, then we will always go out.

(Homework = None)

Homework	Traffic	Hunger	Lawn	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 the table, we will go out going out is directly correlated to Lauren if there is no homework.

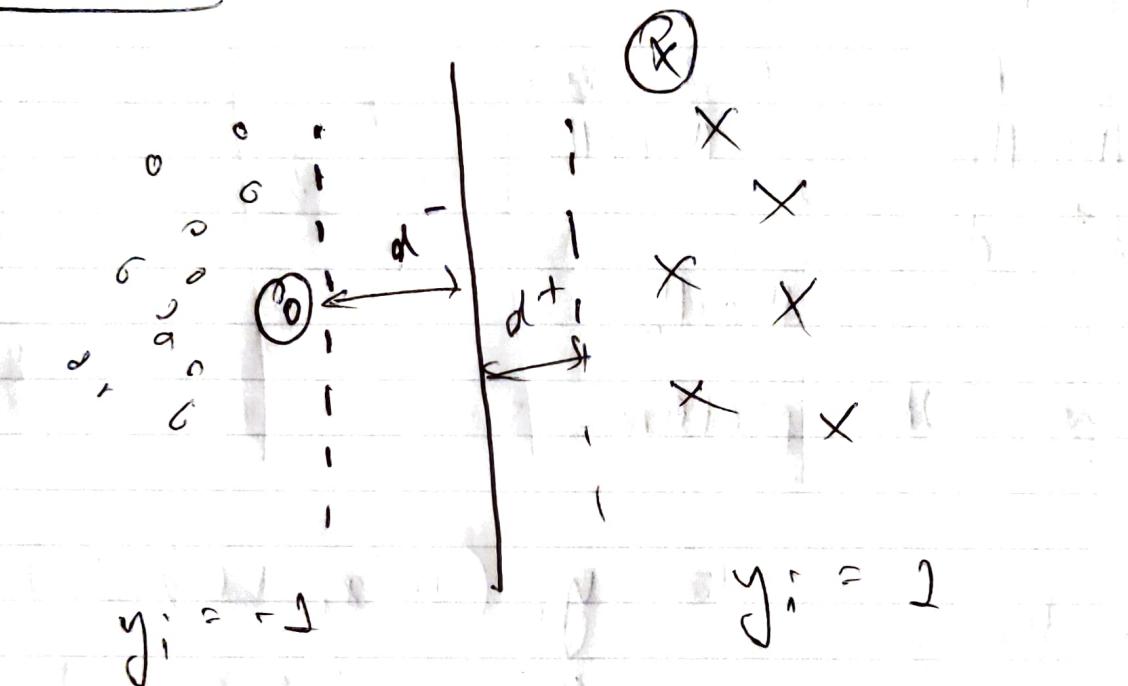
Thus, the new decision tree is:



(Q4) Since amount of homework is normal I will always go out as per the decision tree.

Problem 2

(Q5)



\Rightarrow Linearly separable data

For class $y_i = 1$

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

$$w^T x_j + b \leq 1 \Rightarrow b \leq 1 - w^T x_j$$

Similarly, for class $y_i = -1$

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

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

$$\begin{aligned}
 & \therefore d^+ - d^- \\
 &= \frac{\omega^T x_i}{\|\omega\|_2} - \frac{\omega^T x_k}{\|\omega\|_2} \\
 &= \frac{1 - b}{\|\omega\|_2} - \frac{(-1 - b)}{\|\omega\|_2} \\
 \Rightarrow & d^+ - d^- = \frac{2}{\|\omega\|_2}
 \end{aligned}$$

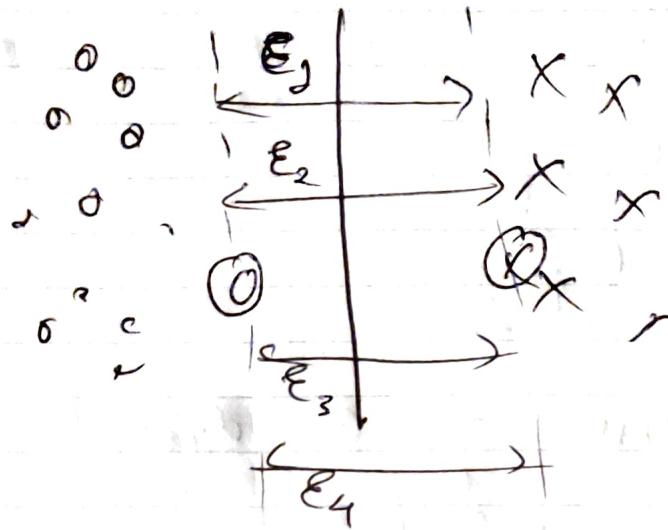
\therefore The SVM problem is \rightarrow

$$\begin{aligned}
 & \max_{\omega, b} \frac{1}{\|\omega\|_2} \quad \text{s.t. } y_i(\omega^T x_i + b) \geq 1 \\
 & \text{for } i \in \{1, \dots, N\} \\
 & \text{no of data points}
 \end{aligned}$$

or

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

Q2)



If the data is not linearly separable, then we have to further deal with slack variables. The slack variables are used to minimize the misclassifications along with the objective function defined in the above problem

$$\frac{1}{2} \|\mathbf{w}\|^2$$

In the case data points are correctly classified,

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

$$\text{else, } y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$$

The slack variables that we will use can be defined as:

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N$ for misclassified

(for $i \in \{1, \dots, N\}$)

∴ New objective function is:

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

∴ The SVM problem can be defined as:

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

$$\text{s.t. } y_i(w^\top \alpha_i + b) \geq 1 - \epsilon_i$$

$$\text{for } i \in \{1, \dots, N\}$$

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

Q3) SVM doesn't normally help with multi-class classification. However, it can be broken down into smaller binary classification subproblems. There are several approaches like One vs One, One vs All and Directed Acyclic Graph.

→ In One vs One we break down the multi-class classification problem into various binary classification. Finally, majority voting is done.

- In the One vs All α , we train N SVMs
SVM(1) → for class 1
SVM(N) → for class N

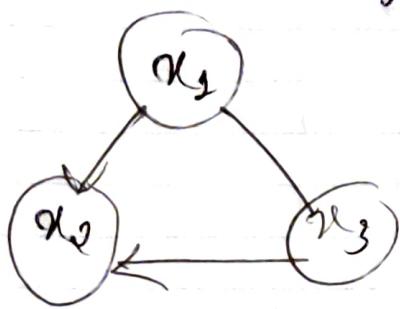
Now, in order to predict the output of new input, predict the output using all the SVMs and then just identify which model gets the prediction farther into the region of the data.

- In Directed Acyclic graphs, the classes are first grouped based on some grouping and then train the SVMs. Thus, at the end, we might need to train less number of SVMs and this approach reduces the diversity from the majority class.

Problem 3:

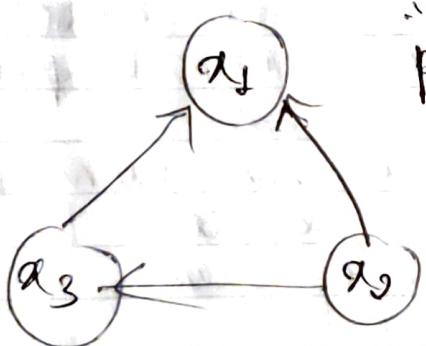
- (1)
- (2) We are given N random variables x_1, x_2, \dots, x_N . And thus, the aim is to create a Bayesian Network showing the joint distⁿ of all N variables.

Let us consider joint distⁿ of x_1, x_2, x_3



$$\therefore p(x_1, x_2, x_3) = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_1, x_2)$$

Similarly, if the ordering of the graph is changed:



$$\begin{aligned} & p(x_1, x_2, x_3) \\ &= p(x_2) \cdot p(x_3 | x_2) \cdot p(x_1 | x_3, x_2) \end{aligned}$$

Thus, if we listing out all possibilities, thus for 3 random variables there are $3! = 6$ distinct graphs.

So for N random variables, $N!$ graphs

② For a single discrete variable x having M possible states, the probability distribution is

$$P(x | \sigma) = \prod_{m=1}^M \frac{\delta_{xm}}{s_m}$$

Similarly, for discrete random variables x_1, x_2 we have two states: s_1 & s_2

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

& we can see that total parameters are $\rightarrow (s_1 \times s_2 - 1)$

Thus, after generalising this we can say that if we have N random variables x_1, x_2, \dots, x_N & they have N states s_1, s_2, \dots, s_N then total no. of parameters in the joint probability distribution will be:

$$[(s_1 \times s_2 \times \dots \times s_N) - 1]$$

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

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

for two random variables x_1, x_2

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

$$P(x_1, x_2 | \sigma) = \prod_{a=1}^{s_1} \in \underset{x_1}{\sigma_a} \prod_{b=1}^{s_2} \in \underset{x_2}{\sigma_b}$$

& no. of parameters in this joint probability will be $(s_1 + s_2 - 2)$

i. N random variables with N states will have $\rightarrow (s_1 + s_2 + \dots + s_{N-2})$ param

& thus, the graph will grow linearly instead of exponentially.

Hence if we maintain conditional independence of all N random variables, then it will be less complex to calculate joint probability

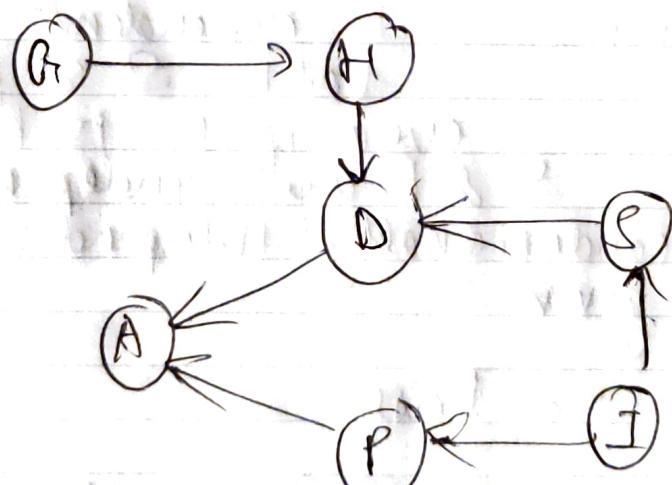
$P(x_1, x_2, \dots, x_n)$ as the graph will be linear & not exponential.

Q2)

①

- G \rightarrow Going out right now
- H \rightarrow contradicting Hokie plague
- D \rightarrow Drowsy
- S \rightarrow Lack of Sleep
- I \rightarrow Insomnia
- P \rightarrow Pay Attention
- A \rightarrow Unable to finish the assy

Bayesian Network



$$\begin{aligned} P(H, D, S, I, P, A) = & p(G) \cdot p(H|G) \cdot p(D|H, S) \\ & \cdot p(S|I) \cdot p(I) \cdot p(P|I) \\ & \cdot p(A|P, D) \end{aligned}$$

Nodes in the S 's Markov Blanket \rightarrow

$I, D \& H$ as if we knew the values of these nodes, then we can determine the value of S .

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

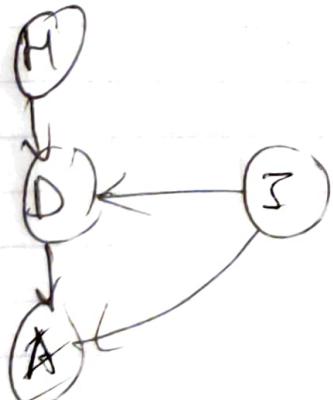
Simplifying the graph to get above nodes, we get \Rightarrow



Since, I is dependent on both S & P
even if I is known
~~so even~~ S & P are independent
 $\therefore S$ is conditionally independent
of P given I

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

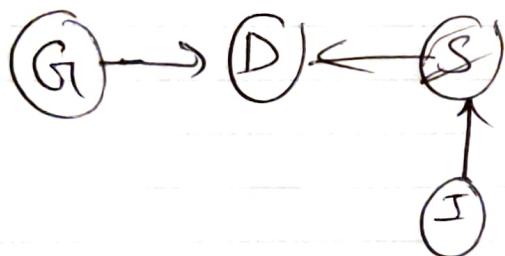
Simplifying the graph to get above nodes we get \Rightarrow



Since, A is dependent on D & I & D is dependent on both I & H .

- H is also dependent on I when A is known
- H is conditionally dependent on I given A.

④ $G \perp\!\!\!\perp I | S \rightarrow$ True

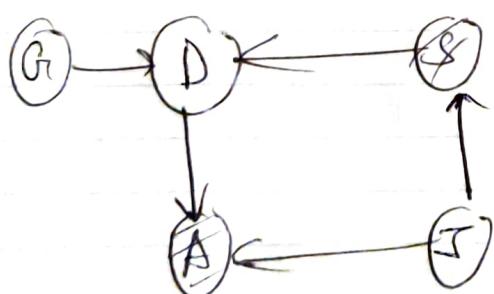


Since, D is dependent on both G & S and S is dependent on I

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

- G is conditionally independent on I given S.

⑤ $G \perp\!\!\!\perp I | S, A \rightarrow$ False



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

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

- G is dependent on I given S and A.