

Name - Prayag Nikul Purohit

(1)

SJSU ID - 017416737

Home-work - 2

Q1 Ch4, Ex - 3

a) Entropy of database.

$$H(\text{Annual Income, D}) = - \sum_n P(n) \log_2(P(n))$$

Annual Income $\leftarrow \begin{matrix} 5 & (25k - 50k) \\ 2 & (< 25k) \\ 1 & (> 50k) \end{matrix}$

$$\begin{aligned} H(AI, D) = & - \left(\frac{5}{8} \log_2 \left(\frac{5}{8} \right) \right) + \\ & - \left(\frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right) + \\ & - \left(\frac{1}{8} \log_2 \left(\frac{1}{8} \right) \right) \end{aligned}$$

$$H(AI, D) = - (-0.5 - 0.45 - 0.375)$$

$$= \underline{\underline{1.299 \text{ bits}}}$$

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

$$b) \text{ Gini Index (Annual Income)} =$$

$$1 - \sum p(n)^2$$

$$1 - \left[\left(\frac{5}{8} \right)^2 + \left(\frac{2}{8} \right)^2 + \left(\frac{1}{8} \right)^2 \right]$$

$$1 - (0.0625 + 0.3906 + 0.051625)$$

$$\approx 0.5312$$

$$\underline{\underline{\text{Gini Index} = 0.5312}}$$

c) So, first we need to sort the data set

ID	Age	AI
3	18	$\leq 25k$
6	24	$\leq 25k$
4	28	$25k - 50k$
5	37	$25k - 50k$
1	39	$25k - 50k$
8	40	$> 50k$
2	50	$25k - 50k$
7	52	$25k - 50k$

} 26

} 39.5

} 45

Boundary b/w ID 446 is $\frac{24+28}{2} = \underline{\underline{26}}$ (2)

Boundary b/w ID 148 is $\frac{39+40}{2} = \underline{\underline{39.5}}$

Boundary b/w ID 842 is $\frac{40+50}{2} = \underline{\underline{45}}$

H for data set = 1.299 bits

Age $\leftarrow \begin{matrix} 26 \\ 39.5 \\ 45 \end{matrix}$

$\rightarrow \text{Age} \Rightarrow 26$

$$\text{rem}(d, D) = \sum \underbrace{\frac{|Dd=l|}{|D|}}_{\text{weighting}} \times \underbrace{H(t, Dd=l)}_{\text{entropy of position}}$$

$$= \frac{2}{8} (H(<26)) + \frac{6}{8} (H(>26))$$

$$= 0 + \frac{3}{4} \left[\frac{5}{6} \log_2\left(\frac{5}{6}\right) + -\frac{1}{6} \log_2\left(\frac{1}{6}\right) \right]$$

$$= 0 + \frac{3}{4} [0.64]$$

$$= \underline{\underline{0.48}}$$

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

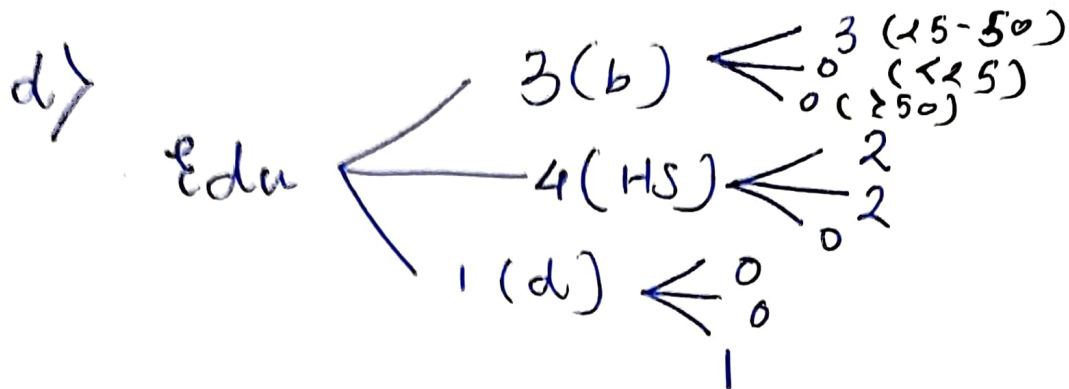
$$\begin{aligned} IG &= H(t, D) - \text{rem}(d, D) \\ &= 1.3 - 0.48 \\ &= \underline{\underline{0.82}} \end{aligned}$$

Now, for Age $\Rightarrow 39.5$
 $\text{rem}(d, D) = 0.94$

$$\begin{aligned} IG &= H(t, D) - \text{rem}(d, D) \\ &= 1.3 - 0.94 \\ &= \underline{\underline{0.35}} \end{aligned}$$

Now, for Age $\Rightarrow 45$
 $\text{rem}(d, D) = 1.095$

$$\begin{aligned} IG &= 1.3 - 1.095 \\ &= \underline{\underline{0.215}} \end{aligned}$$



3

$$IG(t, D) = H(t, d) - \text{rem}(d, D)$$

$$\text{rem}(E_d, D) =$$

$$\frac{3}{8} \left[-\left(\frac{3}{3} \log_2 \frac{3}{3}\right) - \left(\frac{0}{3} \log_2 \frac{0}{3}\right) - \left(\frac{0}{3} \log_2 \frac{0}{3}\right) \right]$$

$$+ \frac{4}{8} \left[-\left(\frac{2}{4} \log_2 \frac{2}{4}\right) - \left(\frac{2}{4} \log_2 \frac{2}{4}\right) - \left(\frac{0}{4} \log_2 \frac{0}{4}\right) \right]$$

$$+ \frac{1}{8} \left[-\left(\frac{1}{1} \log_2 \frac{1}{1}\right) - \left(\frac{0}{1} \log_2 \frac{0}{1}\right) - \left(\frac{0}{1} \log_2 \frac{0}{1}\right) \right]$$

$$= 0.5$$

$$IG = 1.3 - 0.5 = 0.8$$

Information gain for

Education = 0.8

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

non (Occupation)

$$\Rightarrow \frac{2}{8} \left[- \left(\frac{2}{2} \log_2 \frac{2}{2} \right) + 0 + 0 \right] +$$

$$\Rightarrow \frac{3}{8} \left[- \left(\frac{2}{3} \log_2 \frac{2}{3} \right) + \left(\frac{1}{3} \log_2 \frac{1}{3} \right) + 0 \right] +$$

$$\Rightarrow \frac{2}{8} \left[- \left(\frac{1}{2} \log_2 \frac{1}{2} \right) - \left(\frac{1}{2} \log_2 \frac{1}{2} \right) + 0 \right] +$$

$$\frac{1}{8} \left[- \left(\frac{1}{1} \log_2 \frac{1}{1} \right) - \left(\frac{0}{1} \log_2 \frac{0}{1} \right) - \left(\frac{0}{1} \log_2 \frac{0}{1} \right) \right]$$

$$= 0.5944$$

$$IG = 1.3 - 0.5944$$

$$= \underline{\underline{0.7056}}$$

IG for Occupation is 0.7056

$$\text{rem}(\text{Maximal Status}) =$$

(4)

$$\begin{aligned} & \frac{3}{8} \left[-\left(\frac{1}{3} \log_2 \frac{1}{3}\right) - \left(\frac{2}{3} \log_2 \frac{2}{3}\right) - 0 \right] + \\ & \frac{4}{8} \left[-\left(\frac{3}{4} \log_2 \frac{3}{4}\right) - \left(\frac{1}{4} \log_2 \frac{1}{4}\right) - 0 \right] + \\ & \frac{1}{8} \left[-\left(\frac{1}{1} \log_2 \frac{1}{1}\right) - 0 - 0 \right] \end{aligned}$$

$$= 0.75$$

$$I_G = 1.3 - 0.75$$

$$= 0.55$$

$$I_G \text{ for } \underline{\text{Maximal Status}} = \underline{0.55}$$

$$e) \quad I_G(E_d, D) = 0.8$$

$$I_G(MS, D) = 0.55$$

$$I_G(\Theta_{cc}, D) = 0.7$$

$$\text{entropy } E_d = 1.4 \text{ bits}$$

$$\text{entropy } MS = 1.4 \text{ bits}$$

$$\text{entropy } \Theta_{cc} = 1.9 \text{ bits}$$

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

$$IGR = \frac{IG}{H}$$

$$IGR(Ed, D) = \frac{0.8}{1.4} = \underline{0.567}$$

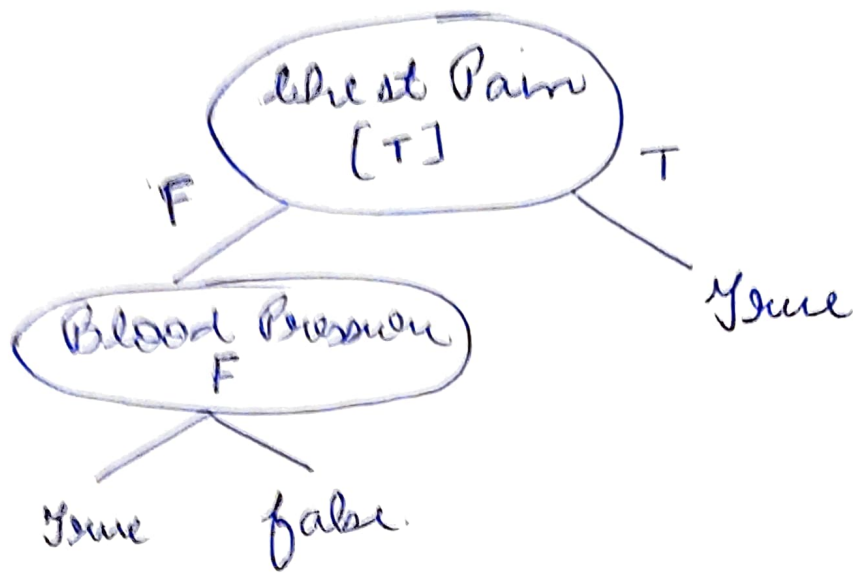
$$IGR(MS, D) = \frac{0.55}{1.4} = \underline{0.39}$$

$$IGR(Occ, D) = \frac{0.705}{1.9} = \underline{0.36}$$

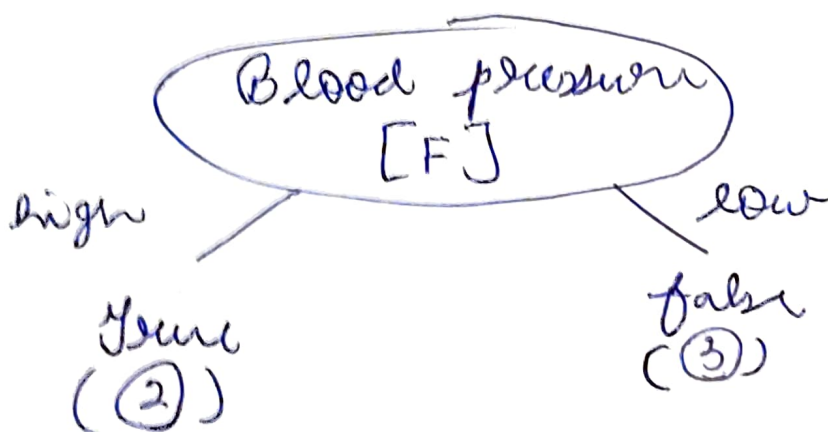
f)				
feature	Level	Partition GI	Benn	IG
Education	HS	0.5	0.25	0.2513
	bachelors	0		
	doctors	0		
Marital Status	NA	0.45	0.3542	0.1711
	M	0.375		
	D	0		
Occupation	T	0	0.2917	0.2396
	P	0.45		
	A	0.5		
	AF	0		

Q2 Ch-4, Exercise - 4

5



Given, decision tree to predict heart disease, blood pressure & chest pain are descriptive features & heart disease is the target feature.



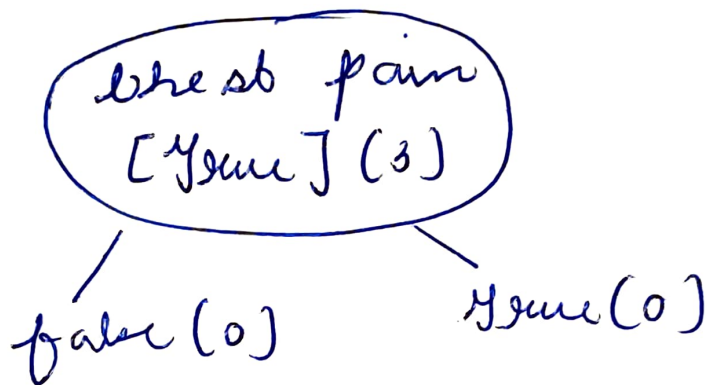
- The number in round brackets represent no. of record in pruning set of that node.

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

- The number in square bracket represent the majority target level from that node.
- The number of error coming from leaf node is greater than that of the number of error coming from the root node.



2nd iteration for pruning root node
 predicts true which causes 3
 mistake but pruning set the
 algorithm will stop here.

Thus the final decision tree
 is produced.

Q3 Ch-4, Ex-5

(6)

$$a) H(\text{risk}) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} \\ = 0.72 \text{ bits}$$

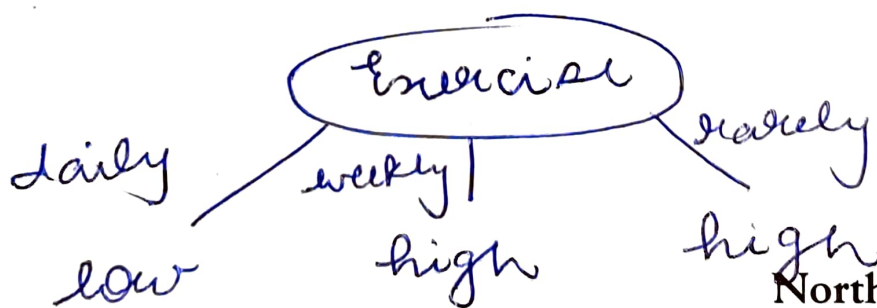
Sample A

$$\text{Rem}(E_n) = \frac{1}{5} \left[-\left(\frac{1}{1} \log_2 \frac{1}{1}\right) \right] + \\ \frac{2}{5} \left[-\frac{2}{2} \log_2 \frac{2}{2} \right] + \frac{2}{5} \left[-\frac{2}{2} \log_2 \frac{2}{2} + 0 \right] \\ = 0$$

$$\text{Rem}(\text{family}) = \frac{3}{5} \left[-\left(\frac{1}{3} \log_2 \frac{1}{3}\right) - \left(\frac{2}{3} \log_2 \frac{2}{3}\right) \right] \\ + \frac{2}{5} \left[-\left(\frac{2}{2} \log_2 \frac{2}{2}\right) + 0 \right] \\ = 0.55$$

$$IG(E_n) = 0.72 - 0 = 0.72 \checkmark$$

$$IG(\text{family}) = 0.72 - 0.55 = 0.17$$



Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

- as IG of Exercise is highest so we will choose that

for sample B

entropy = 0.72 bits

$$\text{rem}(\text{smoker}) = \frac{1}{5} \left[\left(-\frac{1}{1} \log_2 \frac{1}{1} \right) + 0 \right] + \frac{4}{4} \left[\left(-\frac{4}{4} \log_2 \frac{4}{4} \right) \right]$$

$$= 0$$

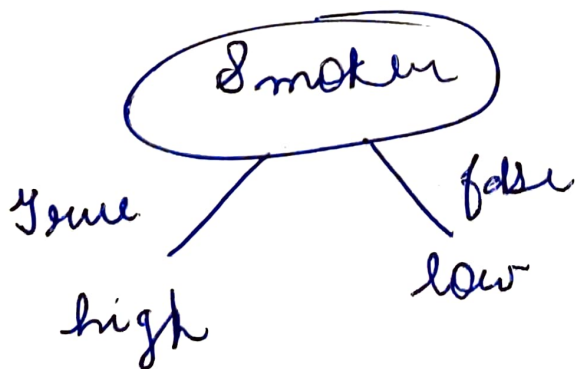
$$\text{rem}(\text{obese}) = \frac{3}{5} \left[\left(-\frac{1}{3} \log_2 \frac{1}{3} \right) + \left(-\frac{2}{3} \log_2 \frac{2}{3} \right) \right] + \frac{2}{3} \left[\left(-\frac{2}{2} \log_2 \frac{2}{2} \right) + \left(-\frac{0}{2} \log_2 \frac{0}{2} \right) \right]$$

$$= 0.55$$

$$IG(\text{smoker}) = 0.721 - 0 = 0.721 \checkmark$$

$$IG(\text{obese}) = 0.721 - 0.55 = 0.17$$

Smoker has higher IG



for sample c

7

$$\text{entropy} = \left(-\frac{2}{5} \log_2 \frac{2}{5} \right) - \left(\frac{3}{5} \log_2 \frac{3}{5} \right) \\ = 0.971 \text{ bits}$$

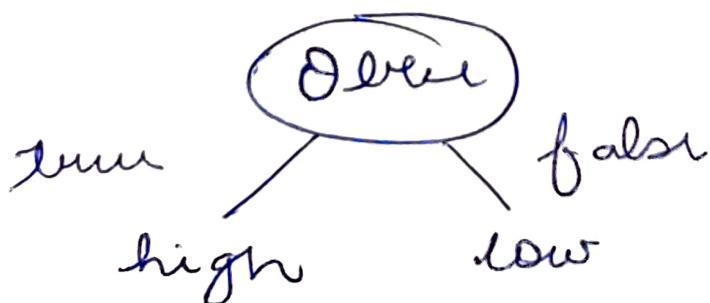
$$\text{seen (obses)} = \frac{3}{5} \left(-\frac{1}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) \\ + \frac{2}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right) \\ = 0.55$$

$$\text{seen (family)} = \frac{4}{5} \left[\left(-\frac{2}{4} \log_2 \frac{2}{4} \right) + \left(\frac{1}{4} \log_2 \frac{1}{4} \right) \right] \\ + \frac{1}{5} \left[\left(-\frac{1}{1} \log_2 \frac{1}{1} \right) \right] \\ = 0.8$$

$$IG(\text{obses}) = 0.971 - 0.55 = 0.42$$

$$IG(\text{family}) = 0.971 - 0.8 = 0.17$$

here, obses has higher IG



Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

b) Exercise = rarely
Smoker = false
Obese = true
family = yes

Tree 1

Exercise → rarely → risk → high

Tree 2

Smoker → false → risk → low

Tree 3

Obese → true → risk → high

So, majority value for risk
is high so it will predict
high as the output.

84 Cn-4, En-6

(8)

$$a) H(\text{nick}) = \left(-\frac{3}{6} \log_2 \frac{3}{6}\right) + \left(-\frac{3}{6} \log_2 \frac{3}{6}\right) \\ = 1 \text{ bit.}$$

$$\text{rem}(\text{nick}) = 2 \times \frac{3}{6} \left[\left(-\frac{1}{3} \log_2 \frac{1}{3}\right) + \left(-\frac{2}{3} \log_2 \frac{2}{3}\right) \right] \\ = 0.9183$$

$$\text{rem}(\text{romkue}) = \frac{2}{6} \left[\left(-\frac{2}{2} \log_2 \frac{2}{2}\right) + 0 \right] \\ + \frac{4}{6} \left[\left(-\frac{3}{4} \log_2 \frac{3}{4}\right) + \left(-\frac{1}{4} \log_2 \frac{1}{4}\right) \right] \\ = 0.54$$

$$\text{rem}(\text{Drinks Alcohol}) = \frac{5}{6} \left[\left(-\frac{2}{5} \log_2 \frac{2}{5}\right) + \left(-\frac{3}{5} \log_2 \frac{3}{5}\right) \right] \\ + \frac{1}{6} \left[\left(-\frac{1}{1} \log_2 \frac{1}{1}\right) \right] \\ = 0.809$$

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

$$IG(\text{obese}) = 1 - 0.9183 = 0.0817$$

$$IG(\text{smoker}) = 1 - 0.54 = \underline{\underline{0.46}} \checkmark$$

$$IG(\text{Drink alcohol}) = 1 - 0.809 = 0.191$$

Smoker has the highest IG so it will be the root node according to ID3 algorithm.

b) Other features which can be more descriptive features such as workout, skin care & medication, can also be added to the dataset that can possibly indicate the low target level.

85 Ch - 4, En - 7

9

$$\begin{aligned} a) H(\text{Buys}, D) &= \left[-\frac{9}{14} \log_2 \frac{9}{14} \right] + \\ &\quad \left[-\frac{5}{14} \log_2 \frac{5}{14} \right] \\ &= \underline{0.9403} \end{aligned}$$

$$\begin{aligned} \text{Hem}(\text{student}) &= \frac{7}{14} \left[-\left(\frac{4}{7} \log_2 \frac{4}{7} \right) \right. \\ &\quad \left. - \left(\frac{3}{7} \log_2 \frac{3}{7} \right) \right] \\ &= 0.788 \end{aligned}$$

$$\begin{aligned} IG(\text{student}) &= 0.9397 - 0.788 \\ &= 0.1517 \end{aligned}$$

$$IG(\text{student}) < IG(\text{age}),$$

so is not a better option.

Northeastern University
San Francisco

Northeastern University
Silicon Valley

NORTHEASTERN.EDU/BAYAREA/

b) for ID

all the values in ID column is
unique & different

$$\Rightarrow \text{sum(ID)} = 0$$

hence, it is not a good choice to
use ID as a ~~var~~ as it will
result in over fitting of the
model