

Course Number: CSCB 5380

Name: Vishnu Vardhan Kartepalli.

UNT ID: 11647518.

Homework No. 2.

Q1)

a)

$$E = \{5+, 5-\}$$

$$\therefore \text{Entropy}(E) = - \sum_{i=1}^n P_i \log_2(P_i)$$

$$\therefore \text{Entropy}(E) = - \frac{5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10}$$

$$\therefore \text{Entropy}(E) = 1$$

Attribute: Tobacco smoking.

values \rightarrow [yes, no].

$$S_{\text{yes}} \leftarrow \{+4, -1\} \quad \therefore \text{Entropy}(S_{\text{yes}}) = - \frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \\ = 0.7219$$

$$\therefore S_{\text{no}} \leftarrow \{+1, -4\} \quad \therefore \text{Entropy}(S_{\text{no}}) = - \frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} \\ = 0.7219$$

$$\therefore \text{Gain}(S) = \text{Entropy}(E) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v).$$

$$\text{Gain}(S, \text{Tobacco smoking}) = \text{Entropy}(E) - \sum_{v \in \{\text{yes}, \text{no}\}} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

$$= \text{Entropy}(E) - \frac{5}{10} \text{Entropy}(S_{\text{yes}}) - \frac{5}{10} \text{Entropy}(S_{\text{no}})$$

$$= 1 - \frac{5}{10} (0.7219) - \frac{5}{10} (0.7219)$$

$$= 0.2781$$

$$\therefore \text{Gain}(S, \text{Tobacco smoking}) = 0.2781$$

Attribute: Random Exposure

values $\rightarrow \{\text{yes}, \text{no}\}$

$$\therefore S_{\text{yes}} \leftarrow \{+2, -0\} \quad \therefore \text{Entropy}(S_{\text{yes}}) = -\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$$

$$\therefore S_{\text{no}} \leftarrow \{+3, -5\} \quad \therefore \text{Entropy}(S_{\text{no}}) = -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = 0.954$$

$$\therefore \text{Gain}(S, \text{Exposure}) = \text{Entropy}(E) - \sum_{v \in \{\text{yes}, \text{no}\}} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

$$= 1 - \frac{2}{10} \text{Entropy}(S_{\text{yes}}) - \frac{8}{10} \text{Entropy}(S_{\text{no}})$$

$$= 1 - \frac{2}{10} (0) - \frac{8}{10} (0.954)$$

$$= 0.2368$$

$$\therefore \text{Gain}(S, \text{Exposure}) = 0.2368$$

Attribute: Chronic cough

values \rightarrow {yes, no}

$$S_{\text{yes}} \leftarrow [+4, -3] \quad \therefore \text{Entropy}(S_{\text{yes}}) = -\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7}$$

$$= 0.9852$$

$$S_{\text{no}} \leftarrow [+1, -2] \quad \therefore \text{Entropy}(S_{\text{no}}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.9182$$

$$\therefore \text{Gain}(S, \text{cough}) = \text{Entropy}(E) - \sum_{v \in \{\text{yes}, \text{no}\}} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

$$= \text{Entropy}(E) - \frac{7}{10} \text{Entropy}(S_{\text{yes}}) - \frac{3}{10} \text{Entropy}(S_{\text{no}})$$

$$= 1 - \frac{7}{10} (0.9852) - \frac{3}{10} (0.9182)$$

$$= 0.0349$$

$$\therefore \text{Gain}(S, \text{cough}) = 0.0349$$

Attribute: Weight loss
values $\rightarrow \{yes, no\}$

$$\therefore S_{yes} \leftarrow \{+3, -2\} \quad \therefore \text{Entropy}(S_{yes}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\ = 0.9709.$$

$$\therefore S_{no} \leftarrow \{+2, -3\} \quad \therefore \text{Entropy}(S_{no}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\ = 0.9709.$$

$$\therefore \text{Gain}(S, \text{Weight loss}) = \text{Entropy}(E) - \sum_{v \in \{yes, no\}} \frac{|S_v|}{5} \text{Entropy}(S_v) \\ = \text{Entropy}(E) - \frac{5}{10} \text{Entropy}(S_{yes}) - \frac{5}{10} \text{Entropy}(S_{no}) \\ = 1 - \frac{5}{10} (0.9709) - \frac{5}{10} (0.9709) \\ = 0.0291$$

$$\therefore \text{Gain}(S, \text{loss}) = 0.0291$$

Here by Gain.

$$\therefore \text{Gain}(S, \text{Tobacco smoking}) = 0.2781$$

$$\text{Gain}(S, \text{Radon Exposure}) = 0.2368$$

$$\text{Gain}(S, \text{Chronic Cough}) = 0.0349$$

$$\text{Gain}(S, \text{Weight loss}) = 0.0291$$

∴ In the formation of decision Tree. We should consider the Maximum Information Gain.

∴ So, here Tobacco smoking as a root.
As there is one 'No' for 'yes' in cancer. we have to repeat.

Tobacco Smoking.

yes					No				
Smoking	Exposure	rough	weight loss	cancer	Smoking	Exposure	rough	weight loss	cancer
yes	yes	yes	No	yes	No	Yes	No	yes	yes
yes	No	yes	No	yes	No	No	yes	No	No
yes	No	yes	yes	yes	No	No	yes	yes	No
yes	No	yes	yes	yes	No	No	yes	No	No
yes	No	No	No	No	No	No	No	yes	No

For IN "YES"

$$\therefore \text{Entropy}(E) = -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.7219$$

∴ Attribute: Exposure

$$S_{yes} \leftarrow [+1, -0] \therefore \text{Entropy}(S_{yes}) = -\frac{1}{1} \log_2 \frac{1}{1} = 0$$

$$= 0$$

$$S_{No} \leftarrow [+3, -1] \therefore \text{Entropy}(S_{No}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$$

$$= 0.8112$$

$$\therefore \text{Gain}(S, \text{Exposure}) = 0.7219 - \frac{1}{5} (0) - \frac{4}{5} (0.8112)$$

$$\therefore \text{Gain}(S, \text{Exposure}) = 0.07294$$

Attribute: cough

$$\therefore S_{\text{yes}} \leftarrow \{+4, -0\} \therefore \text{Entropy}(S_{\text{yes}}) = -\frac{4}{4} \log_2 \frac{4}{4} - 0 = 0.$$

$$\therefore S_{\text{no}} \leftarrow \{+0, -1\} \therefore \text{Entropy}(S_{\text{no}}) = 0 - \frac{1}{1} \log_2 \frac{1}{1} = 0.$$

$$\therefore \text{Gain}(S, \text{cough}) = 0.7219 - 0 - 0$$

$$\therefore \text{Gain}(S, \text{cough}) = 0.7219.$$

Attribute: weight loss

$$\therefore S_{\text{yes}} \leftarrow \{+2, -0\} \therefore \text{Entropy}(S_{\text{yes}}) = -\frac{2}{2} \log_2 \frac{2}{2} - 0 = 0.$$

$$\therefore S_{\text{no}} \leftarrow \{+2, -1\} \therefore \text{Entropy}(S_{\text{no}}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ = 0.9182$$

$$\therefore \text{Gain}(S, \text{weight loss}) = 0.7219 - \frac{2}{5} (0) - \frac{3}{5} (0.9182)$$

$$\therefore \text{Gain}(S, \text{weight loss}) = 0.17098.$$

Here

$$\therefore \text{Gain}(S, \text{Exposure}) = 0.07294$$

$$\therefore \text{Gain}(S, \text{cough}) = 0.7219$$

$$\therefore \text{Gain}(S, \text{weight loss}) = 0.17098.$$

\therefore Maximum
value
of Gain

$\therefore IN$ " No "

$$\therefore Entropy(E) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.7219.$$

Attribute: Exposure :

$$\therefore S_{yes} \leftarrow [+1, -0] \quad \therefore Entropy(S_{yes}) = 0.$$

$$\therefore S_{No} \leftarrow [+0, -4] \quad \therefore Entropy(S_{No}) = 0.$$

$$\therefore Gain(S, Exposure) = 0.7219.$$

Attribute: cough

$$\therefore S_{yes} \leftarrow [+0, +3] \quad \therefore Entropy(S_{yes}) = 0.$$

$$\therefore S_{No} \leftarrow [+1, -1] \quad \therefore Entropy(S_{No}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\therefore Gain(S, cough) = 0.7219 - 0 - \frac{2}{5}(1) = 0.3219.$$

Attribute: weight loss :

$$\therefore S_{yes} \leftarrow [+1, -2] \quad \therefore Entropy(S_{yes}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \\ = 0.9182$$

$$\therefore S_{No} \leftarrow [+0, -2] \quad \therefore Entropy(S_{No}) = 0.$$

$$\therefore Gain(S, loss) = 0.7219 - \frac{3}{5}(0.9182) = 0.$$

$$\therefore Gain(S, weight_{loss}) = 0.1709$$

Here

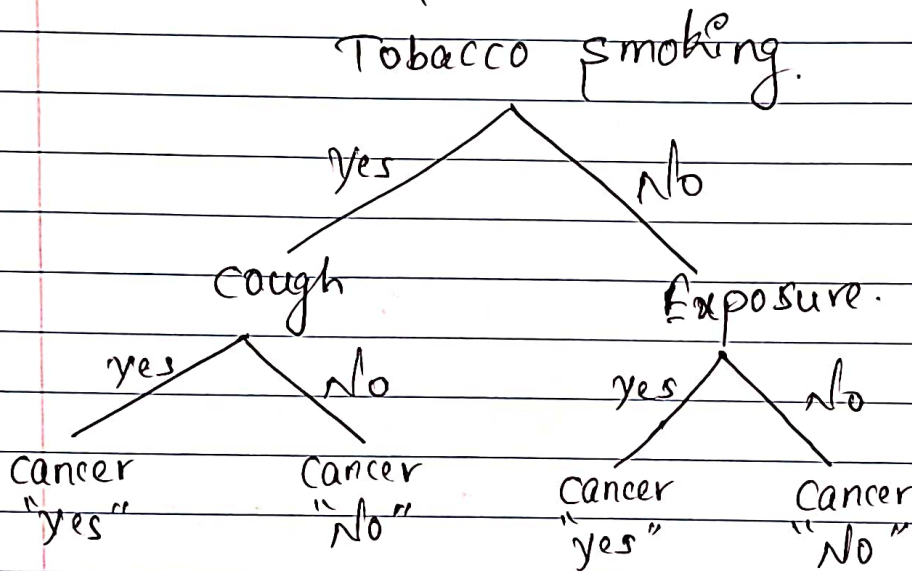
$$\therefore \text{Gain}(S, \text{Exposure}) = 0.7219$$

$$\therefore \text{Gain}(S, \text{cough}) = 0.3219$$

$$\therefore \text{Gain}(S, \text{weight loss}) = 0.1709.$$

Maximum
value of
Gain.

\therefore The tree is,



Q2

c)

We first estimate prior probabilities for "status" class labels.

$$\therefore P(\text{senior}) = 5/11$$

and

$$\therefore P(\text{junior}) = 6/11$$

Now conditional probabilities.

$$\therefore P(\text{department} / \text{status})$$

class	Sales	systems	Marketing	Secretary
Senior	$1/5$	$2/5$	$1/5$	$1/5$
Junior	$2/6 = 1/3$	$2/6 = 1/3$	$1/6$	$1/6$

$$P(\text{age} / \text{status})$$

class	21...25	26...30	31...35	36...40	41...45	46...50
Senior	$0/5 = 0$	$0/5 = 0$	$2/5$	$1/5$	$1/5$	$1/5$
Junior	$1/6$	$3/6 = 1/2$	$2/6 = 1/3$	$0/6 = 0$	$0/6 = 0$	$0/6 = 0$

$P(\text{salary/status})$

class	26k-30k	31k-35k	36k-40k	41k-45k	46k-50k	66k-70k
Senior	$0/5 = 0$	$0/5 = 0$	$1/5$	$0/5 = 0$	$2/5$	$2/5$
Junior	$2/6 = 1/3$	$1/6 = 0$	$0/6 = 0$	$1/6$	$2/6 = 1/3$	$0/6$

∴ SO for Test instance.

$$V_{NB} = \underset{V_j \in \{\text{yes}, \text{no}\}}{\text{argmax}} P(V_j) \prod_i P(a_i/V_j)$$

∴ (systems, "26...30", "46k-50k")

~~& P(senior)~~ Two labels $\begin{cases} \rightarrow \text{Senior} \\ \rightarrow \text{Junior} \end{cases}$

$$\begin{aligned} \therefore P(\text{Senior}/a) &= P(\text{senior}) * P(\text{systems/senior}) * \\ &\quad P(26...30/\text{senior}) * P(46k-50k/\text{senior}) \\ &= \frac{5}{11} * \frac{2}{5} * 0 * \frac{2}{5} \end{aligned}$$

$$\boxed{\therefore P(\text{senior}/a) = 0}$$

$$\begin{aligned} \therefore P(\text{Junior}/a) &= P(\text{junior}) * P(\text{systems/junior}) * P(26...30/\text{junior}) \\ &\quad * (P(46k-50k/\text{junior})) \\ &= \frac{6}{11} * \frac{1}{3} * \frac{1}{2} * \frac{1}{3} \end{aligned}$$

$$\boxed{\therefore P(\text{Junior}/a) = 0.030}$$

$$\therefore p(\text{Junior}/a) = 0.030$$

\therefore By comparing the two values, and by the maximum value was Junior.

\therefore So.

Hence, the label for this instance (person) is "Junior".

\therefore Since $p(\text{senior}/a) < p(\text{Junior}/a)$.