Karaso Roween Ajay
ID: 11646981
CSCE S380 - Homework-2

1. a)

| Attributes | Yes Count | No Count |
|---|---|---|
| 1. Tobacco Smoking | 5 | 5 |
| 2. Radon Exposure | 2 | 8 |
| 3. Chronic Cough | 7 | 3 |
| 4. Weight Loss | 5 | 5 |
| 5. Lung Cancer | 5 | 5 |

• Calculate Initial Entropy:

We have 5 +ve & 5 -ve example of lung Cancer, where $N=10$

$$E(s) = -\left[\sum\right] - \sum_{i=1}^{n} P_i \log_2 (P_i)$$

$$= -\frac{5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10}$$

$$= 1$$

Attributes:

Tobacco Smoking:

Values → [Yes, No]

$S_{Yes}$ [+4, -1]    $E(S_{Yes}) = -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5}$

$$= 0.7219$$

$S_{No}$ → [+1, -4]   ∴ $E(S_{No}) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5}$

$$= 0.7219$$

$$\therefore \text{Gain}(s) = E - \sum_{v \in r} \frac{|S_r|}{|S|} E(S_r)$$

$$= E - \frac{5}{10} E(S_{yes}) - \frac{5}{10} E(S_{No})$$

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

$$\underline{\text{Gain}(s, \text{Tobaco Smoking}) = 0.2781}$$

Attribute: Random Exposure:     Value → [Yes, No]

$$S_{yes} \leftarrow [+2, -0]  \qquad E(S_{yes}) = -\frac{2}{2}\log_2 \frac{2}{2} - \frac{0}{2}\log_2 \frac{0}{2}$$

$$= -1 \times (0) - 0$$

$$= 0 - 0$$

$$S_{No} \rightarrow [+3, -5]  \qquad = 0$$

$$E(S_{No}) = -\frac{3}{8}\log_2 \frac{3}{8} - \frac{5}{8}\log_2 \frac{5}{8}$$

$$= 0.954$$

$$\text{Gain}(S, \text{Random}) = E - \sum_{v \in (yes, no)} \frac{|S_r|}{S} E(S_r)$$

$$= 1 - \frac{2}{10} E(S_{yes}) - \frac{8}{10} E(S_{No})$$

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

$$\underline{\text{Gain}(s, \text{Random}) = 0.2368}$$

Attribute: Chronic Cough:

$S_{Yes} \rightarrow [+4, -3]$  $\quad E(S_{Yes}) = \frac{-4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7}$

$$= 6.9852$$

$S_{No} \rightarrow [+1, -2]$  $\quad E(S_{No}) = \frac{-1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$

$$= 6.9182$$

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

$$= 1 - \frac{7}{10} E(S_{yes}) - \frac{3}{10} E(S_{No})$$

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

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

Attribute:- Weight loss

$S_{Yes} [+3, -2]$  $\quad E(S_{Yes}) = \frac{-3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$

$$= 0.9709$$

$S_{No} [+2, -3]$  $\quad E(S_{No}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$
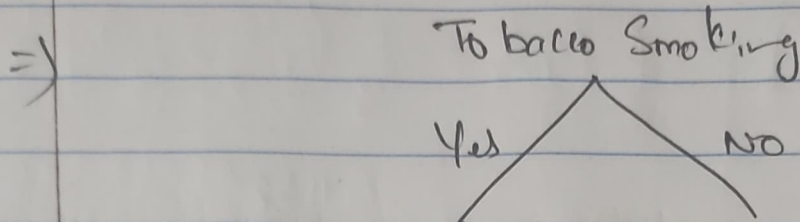
$$= 0.9709$$

$$Gain(S, weight) = E - \frac{5}{10} E(S_{Yes}) - \frac{5}{10} E(S_{No})$$

$$= 1 - \frac{5}{10} (0.9709) - \frac{5}{10} (0.9709)$$

$$Gain(S, weight) = 0.0291$$

From the information, we will consider Tobacco smoking as a root has it have max. information Gain.

As there is One 'No' for 'yes' in Cancer. we will the repeat process

=)

Tobacco Smoking

Yes / NO

| Exposure | Chronic | weightlou | Cancer |
|---|---|---|---|
| Yes | Yes | NO | Yes |
| No | Yes | No | Yes |
| No | Yes | Yes | Yes |
| No | Yes | Yes | Yes |
| NO | No | No | No |

| Exposure | chronic | Weightlon | Cancer. |
|---|---|---|---|
| Yes | No | Yes | Yes |
| No | Yes | No | No |
| No | Yes | Yes | No |
| No | Yes | No | No |
| No | No | Yes | No |

In 'Yes'  $E = -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.7219$

Attribute Exposure:-

$S_{Yes} \rightarrow [+1, 0]$   $E(S_{Yes}) = \frac{-1}{1} \log_2 \frac{1}{1} - 0 = 0$

$S_{No} \rightarrow [+3, -1]$   $E(S_{No}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$

$= 0.8112$

$Gain (S, Exposure) = 0.7219 - \frac{1}{5}(0) - \frac{4}{5}(0.8112) = 0.07294$

Attribute Cough:-

$S_{Yes} \rightarrow [+4, 0]$   $E(S_{Yes}) = \frac{-4}{4} \log_2 \frac{4}{4} - 0 = 0$

$S_{No} \rightarrow [0, -1]$   $E(S_{No}) = 0 - \frac{1}{1} \log_2 \frac{1}{1} = 0$

$\therefore$ Gain $(s, \text{Cough}) = 6.7219 - 0 - 0 = 0.7219$

## Attribute: weight loss:

$S_{Y_9} [+2, -0] \quad \therefore E(S_{Yes}) = \frac{-2}{2} \log_2 \frac{2}{2} - 0 = 0$

$S_{NO} [+2, -1] \quad \therefore E(S_{No}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$

$$= 0.9182$$

$Gain (s, \text{weight}) = 0.7219 - \frac{2}{5} (0) - \frac{3}{5} (0.9182) = 0.1709$

$\Rightarrow$ ~~Max Gain = Gain (S, weight) = 0.1709~~

~~from all three.~~

## In NO:-

$$E = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.7219$$

## Attribute Exposure:

$S_{Yes} [+1, -0] \qquad E(S_{Yes}) = 0$

$S_{NO} [+0, -4] \qquad E(S_{No}) = 0$

$Gain (s, \text{Radon Exposure}) = 0.7219$

## Attribute Chronic :.

$S_{Yes} [+0, -3] \qquad E(S_{Yes}) = 0$

$S_{NO} [+1, -1] \qquad E(S_{No}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$

$$= 1$$

$$\text{Gain}(S, \text{Chronic}) = 0.7219 - 0 - \frac{2}{5}(1) = 0.3219$$

**Attribute :- Weight loss :-**

$$S_{Yes}\ [+1, -2] \qquad E_4(S_{Yes}) = -\frac{1}{3}\log_2\frac{1}{3} - \frac{2}{3}\log_2\frac{2}{3}$$

$$= 0.9182$$

$$S_{NO}\ [+0, -2] \qquad E(S_{NO}) = 0$$

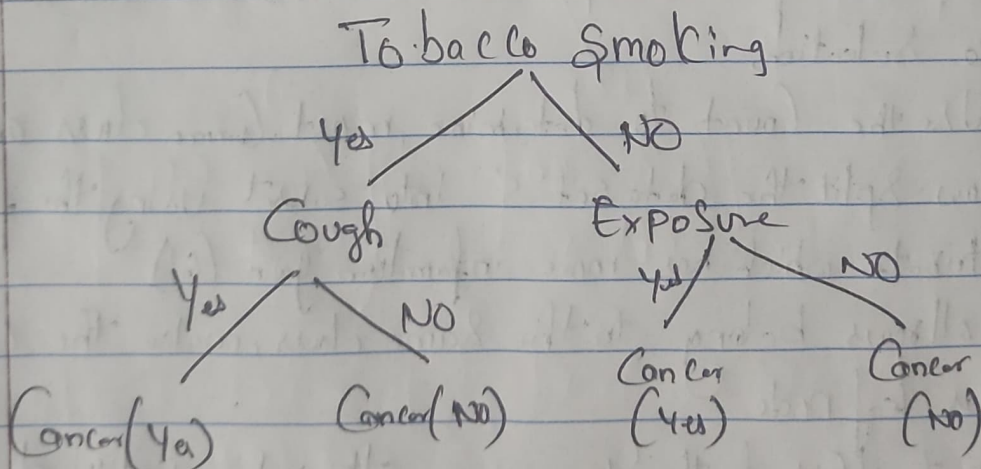$$\text{Gain}(S, \text{weight}) = 0.7219 - \frac{3}{5}(0.9182) = 0$$

$$= 0.1709$$

Here: $\text{Gain}(S, \text{Random}) = 0.7219$

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

$\text{Gain}(S, \text{weight}) = 0.1709 -$

$$\left( \therefore \text{ Max. Value of Gain} \right)$$

∴ The tree will be



(b) To Calculate training error of decision Tree, we are Considering the Confusion matrix

|  | NO Cancer | Cancer |
|---|---|---|
| NO Cancer | 5 | 0 |
| lung Cancer | 0 | 5 |

NO Cancel     lung Cancre

$$Accuracy = TP + TN / TP + TN + FP + FN$$

$$= 5+5 / 0 + 0 + 5 + 0$$

$$= 5 + 5 = 10$$
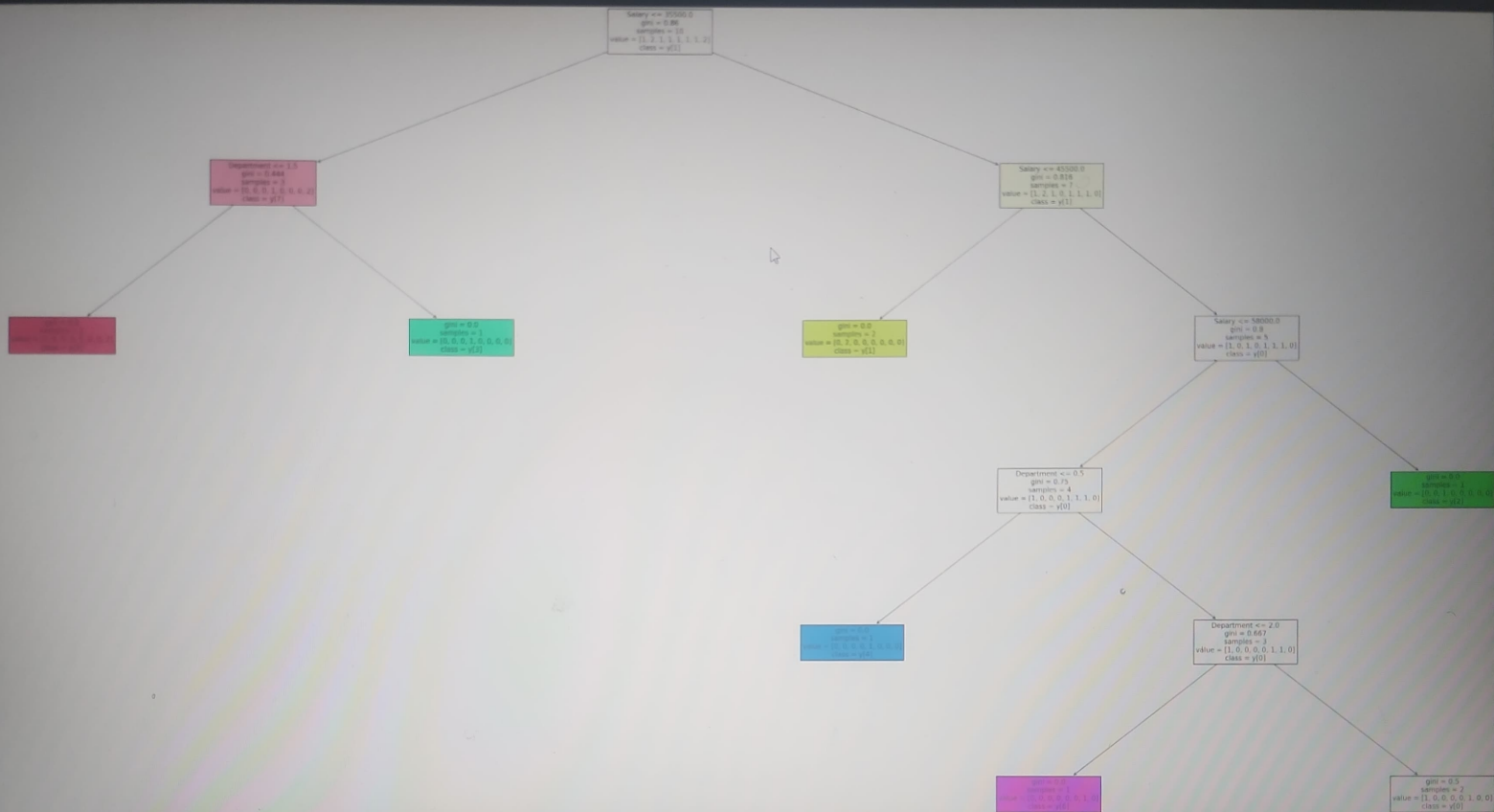
$$= \frac{5 + 5}{5 + 0 + 5 + 0}$$

$$= \frac{10}{10} = 1$$

Training error =

$$\frac{FP + FN}{TP + TN + TP + EN}, = \frac{0+0}{0+5+0+5} = \frac{0}{10}$$

$$= 0,$$

**2 a)** The basic decision tree algorithm should be modified as follows to taken into consideration the count of each Generalised data tuple.

- The count of each tuple must be included into calculation of attribute selection.
- Use the count to determine, most common class among tuples
- Now split the data set 'S' into subset using the feature / attribute which has more information Gain (G)
- If all rows belong to the same class, make the current node has leaf node
- Repeat this for the remaining attributes, untill decision tree has all the leaf nodes.

2 c) We first estimate prior probabilities for "status" class labels

∴ P(senior) = 5/11 & P(junior) = 6/11

⇒ Now Conditional probabilities

P(department / status)

| class | Sales | Systems | marketing | Secerdary |
|-------|-------|---------|-----------|-----------|
| Senior | 1/5 | 2/5 | 1/5 | 1/5 |
| Junior | 1/3 | 1/3 | 1/6 | 1/6 |

P(age / status)

| Class | 21-25 | 26-30 | 31-35 | 36-40 | 41-45 | 46-50 |
|-------|-------|-------|-------|-------|-------|-------|
| Senior | 0 | 0 | 2/5 | 1/5 | 1/5 | 1/5 |
| Junior | 1/6 | 1/2 | 1/3 | 0 | 0 | 0 |

P(Salary / status)

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

∴ So for test instance

$$V_{NB} = \arg\max_{z} P(v_j) \; \pi_i \; P(a_i/v_j)$$

$$v_j \; (\text{Yes, No})$$

∴ (system, "26-30", "46-50k")

Two labels $\begin{cases} \to \text{Senior} \\ \to \text{Junior} \end{cases}$

$$P\left(\text{Senior}/a\right) = P(\text{Senior}) \times P\left(\text{system} / \text{senior}\right) \times P\left(26-30 / \text{Senior}\right)$$
$$\times P\left(46k-50k / \text{senior}\right)$$

$$= \frac{5}{11} \times \frac{2}{5} \times 0 \times \frac{2}{5}$$

$$= 0$$

$$\boxed{P\left(\frac{\text{senior}}{a}\right) = 0}$$

$$P\left(\text{Junior}/a\right) = P(\text{Junior}) \times P\left(\text{system}/\text{Junior}\right) \times P\left(26-30 / \text{Junior}\right) \times$$
$$P\left(46k-50k / \text{Junior}\right)$$

$$= \frac{6}{11} \times \frac{1}{3} \times \frac{1}{2} \times \frac{1}{3} = \frac{1}{33}$$

$$\boxed{∴ P\left(\text{Junior}/a\right) = 0.030}$$

∴ By Comparing the two Values of the max value of Junior.

Hence, the label for this instance (person) is junior

$$\Rightarrow \quad P\left(Senior/a\right) < P\left(Junior/a\right)$$