

Homework NO. 2 - Basic classification
Course: CSCE 5380.

Raghuram Nimmalapudi

OPUS ID: 11594739

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 $\rightarrow [Yes, No]$

$$S_{Yes} \leftarrow [+, -] \quad \therefore \text{Entropy}(S_{Yes}) =$$

$$\Rightarrow -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5}$$

$$\therefore \text{Entropy}(S_{Yes}) = 0.17219$$

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

$$= 0.7219$$

$$\therefore \text{Gain}(S) = \text{Entropy}(E) - \sum_{S \in Y} \frac{|S|}{|S|} \text{Entropy}(S)$$

$$\begin{aligned}
 \text{Gain}(S, \text{Tobacco smoking}) &= \text{Entropy}(E) - \sum_{v \in \{\text{Yes}, \text{No}\}} \frac{f(v)}{S} \text{Entropy}(S_v) \\
 &= \text{Entropy}(E) - \sum_{v \in \{\text{Yes}, \text{No}\}} \text{Entropy}(S_{\text{Yes}}) - \sum_{v \in \{\text{Yes}, \text{No}\}} \text{Entropy}(S_{\text{No}}) \\
 &= 1 - \frac{5}{10} (0.7219) - \frac{5}{10} (0.7219) \\
 &= 0.2981
 \end{aligned}$$

$$\therefore \text{Gain}(S, \text{Tobacco smoking}) = 0.2981$$

Attribute : Random Exposure

values $\rightarrow \{\text{Yes}, \text{No}\}$

$$\begin{aligned}
 \therefore S_{\text{Yes}} &\leftarrow [+2, -2] \quad \therefore \text{Entropy}(S_{\text{Yes}}) = -\frac{2}{5} \log_2 \frac{2}{2} - \frac{3}{5} \log_2 \frac{3}{2} \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 \therefore S_{\text{No}} &\leftarrow [-3, -5] \quad \therefore \text{Entropy}(S_{\text{No}}) = -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{2} \\
 &= 0.954
 \end{aligned}$$

$$\begin{aligned}
 \therefore \text{Gain}(S, \text{Exposure}) &= \text{Entropy}(E) - \sum_{v \in \{\text{Yes}, \text{No}\}} \frac{f(v)}{S} \text{Entropy}(S_v) \\
 &= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = 0.5
 \end{aligned}$$

$$P(\text{Yes}) = \frac{2}{10} = \frac{2}{10} \text{Entropy}(S_{\text{Yes}}) = \frac{8}{10} \text{Entropy}(S_{\text{No}})$$

$$\begin{aligned}
 P(\text{Yes}) \cdot \frac{1}{2} + P(\text{No}) \cdot \frac{1}{2} &= (1 - \frac{2}{10}) \cdot \frac{8}{10} (0.954) \\
 P(\text{Yes}) &= 0.2368
 \end{aligned}$$

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

Attribute: chronic cough
 values \rightarrow [Yes, No]

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

$$= 0.9852$$

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

$$= 0.9182$$

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

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

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

$$= 0.0349$$

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

Attribute: weight loss

values \rightarrow [Yes, No]

$$S_{\text{Yes}} \leftarrow [-3, +2] \therefore \text{Entropy}(S_{\text{Yes}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.9709$$

$$S_{\text{No}} \leftarrow [+4, -3] \therefore \text{Entropy}(S_{\text{No}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.9709$$

$$\therefore \text{Gain}(S, \text{weight loss}) = \text{Entropy}(E) - \sum_{\text{VE (Rep, No)}} \frac{S}{10} \text{Entropy}(S_{\text{Yes}})$$

$$= \text{Entropy}(E) - \frac{5}{10} \text{Entropy}(S_{\text{Yes}}) - \frac{5}{10} \text{Entropy}(S_{\text{No}})$$

$$= 1 - \sum_{\text{VE (Rep, Yes)}} \frac{S}{10} (0.9709) - \sum_{\text{VE (Rep, No)}} \frac{S}{10} (0.9709)$$

$$= 0.0291$$

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

Here by gain

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

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

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

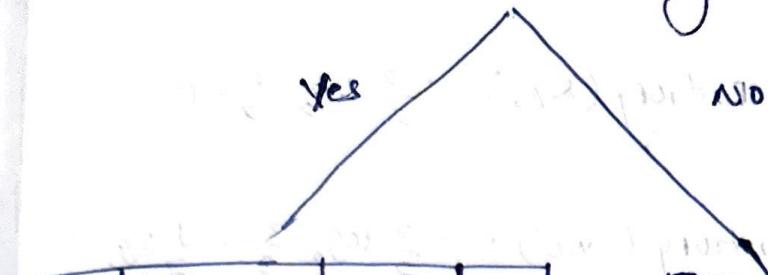
$$\text{Gain}(S, \text{weight loss}) = 0.0291$$

\therefore In the formation of decision tree, we should consider the maximum information gain.

\therefore So here tobacco smoking as a root.

As there is one "No" for "Yes" in Cancer, we have to repeat.

Tobacco smoking



smoking	Exposure	cough	weight loss	Cancer
yes	Yes	Yes	No	Yes
yes	No	Yes	No	Yes
yes	No	Yes	Yes	Yes
yes	No	Yes	Yes	Yes
yes	No	No	No	No

smoking	Exposure	cough	weight loss	Cancer
No	Yes	No	Yes	Yes
No	No	Yes	No	No
No	No	Yes	Yes	No
No	No	Yes	No	No
No	No	No	Yes	No

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

\therefore Attribute : Exposure

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

$s_{\text{No}} \leftarrow [-3, -1] \quad \therefore \text{Entropy}(s_{\text{No}}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$
 $= 0.8112$

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

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

\therefore Attribute : cough

$s_{\text{Yes}} \leftarrow (+4, -0) \quad \therefore \text{Entropy}(s_{\text{Yes}}) = -\frac{4}{4} \log_2 \frac{4}{4} - 0 = 0$

$s_{\text{No}} \leftarrow (+0, -1) \quad \therefore \text{Entropy}(s_{\text{No}}) = 0 - \frac{1}{1} \log_2 \frac{1}{1} = 0$

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

Attribute: weight loss

$$\therefore S_{\text{Yes}} \leftarrow (\text{F}_2, -0) \therefore \text{Entropy}(S_{\text{Yes}}) = -\frac{2}{2} \log_2 \frac{2}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 0$$

$$\therefore S_{\text{No}} \leftarrow (\text{F}_2, -1) : \text{Entropy}(S_{\text{No}}) = -\frac{2}{2} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= 0.9182$$

$$\therefore \text{Gain}(S, \text{weight loss}) = 0.7219 - \frac{2}{5}(0) - \frac{3}{5}(0.9182)$$
$$= 0.17096$$

Here

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

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

$$\therefore \text{Gain}(S, \text{weight loss}) = 0.17096$$

[maximum value of Gain]

$$\therefore \text{In "No" } \leftarrow (\text{F}_2, -4) \text{ is considered as } (1, 2) \rightarrow \text{No}$$

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

Attribute Exposure:

$$\therefore S_{\text{Yes}} \leftarrow (\text{F}_1, -0) \therefore \text{Entropy}(S_{\text{Yes}}) = 0$$

$$\therefore S_{\text{No}} \leftarrow (\text{F}_0, -4) \therefore \text{Entropy}(S_{\text{No}}) = 0$$

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

Attribute : cough

$$\therefore S_{yes} \leftarrow [+1, -3] \therefore \text{Entropy}(S_{yes}) = 0$$
$$\therefore S_{NO} \leftarrow [-1, -1] \therefore \text{Entropy}(S_{NO}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

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

Attribute : weight loss

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

$$\therefore S_{NO} \leftarrow [-1, -1] \therefore \text{Entropy}(S_{NO}) = 0$$

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

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

Here

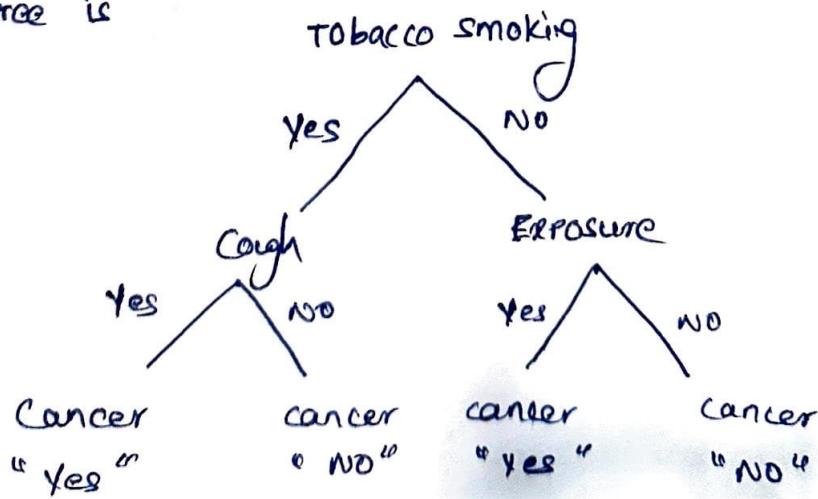
$$\therefore \text{Gain}(S, \text{Exposure}) = 0.7219$$

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

[\because maximum value of Gain]

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

\therefore The tree is



b) To calculate training error of decision tree we are considering the confusion matrix.

Confusion matrix

No lung cancer	5	0
Lung cancer	0	5

no lung cancer lung cancer

$$\begin{aligned} \text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\ &= \frac{5+5}{5+5+0+5} \\ &= 10/10 = 1 \end{aligned}$$

Predicted values.

Training Error =

$$= \frac{FP + FN}{TP + TN + FP + FN}$$

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

$$= 0/10$$

$$= 0$$

PIER = Confusion matrix

Actual outcome

PREDICTION = Predicted class

POER = Positive predictive value

ROC curve

True positive

False negative

True negative

False positive

True positive rate = $\frac{\text{True positive}}{\text{Actual positive}}$

True negative rate = $\frac{\text{True negative}}{\text{Actual negative}}$

False positive rate = $\frac{\text{False positive}}{\text{Actual negative}}$

False negative rate = $\frac{\text{False negative}}{\text{Actual positive}}$

Q2)

- a) The basic decision tree algorithm should be modified as follows to take into consideration the count of each generalized data tuple.
- * the count of each tuple must be integrated into the calculation of the attribute selection measure. (calculate the information gain of each feature).
 - * take the count into consideration to determine the most common class among the tuples.
 - * considering that all rows don't belong to the same class, split the dataset's in to subsets using the feature for which the information gain is maximum.
 - * make a decision tree node using the feature with the maximum information gain.
 - * if all rows belong to the same class, make the current node as a leaf node with the class as its label.
 - * repeat for the remaining features until we run out of all features, or the decision tree has all leaf nodes.

(any method may be used)

Q2)

b) the resulting tree is

Salary = 26k...30k : junior

= 31k...35k : junior

= 36k...40k : senior

= 41k...45k : junior

= 46k...50k [department = secretary : junior]

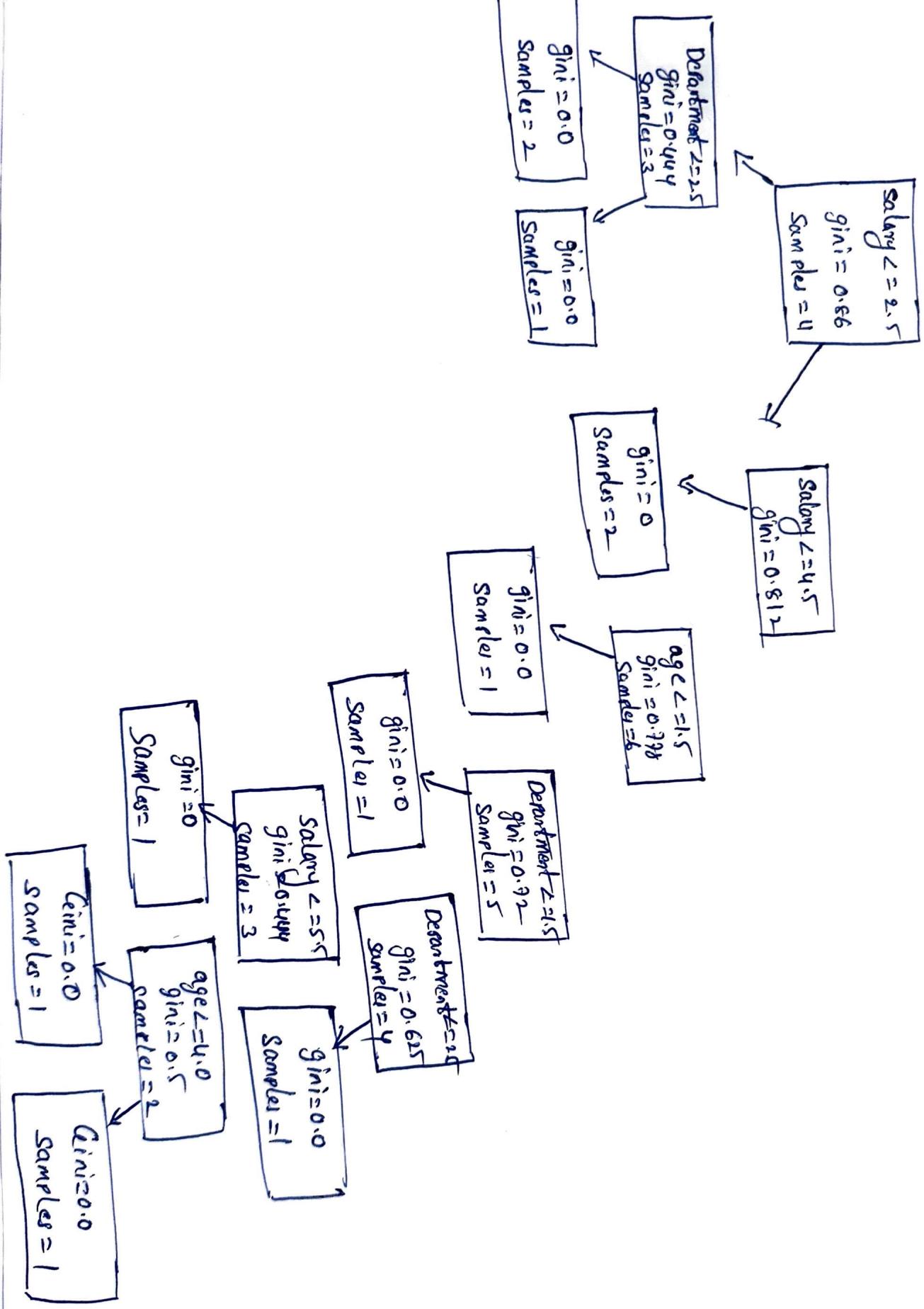
= sales : senior

= systems : junior

= marketing : senior

= 66k...70k : senior

b) Use code analysis to construct a decision tree from the given data.



c) we first estimate "prior" probabilities for "status"

class labels

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

and

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

Now conditional probabilities

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

class sales systems | marketing secretary

senior $\frac{1}{5}$ $\frac{2}{5}$ $\frac{1}{5}$ $\frac{1}{5}$

junior $\frac{2}{6} = \frac{1}{3}$ $\frac{2}{6} = \frac{1}{3}$ $\frac{1}{6}$ $\frac{1}{6}$

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

class $21 \dots 25$ $26 \dots 30$ $31 \dots 35$ $36 \dots 40$ $41 \dots 45$ $46 \dots 50$

senior $0/5 = 0$ $0/5 = 0$ $2/5 = 1/5$ $1/5$ $1/5$ $1/5$

junior $1/6$ $3/6 = 1/2$ $2/6 = 1/3$ $0/6 = 0$ $0/6 = 0$ $0/6 = 0$

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

class $26K-30K$ $31K-35K$ $36K-40K$ $41K-45K$ $46K-50K$ $51K-55K$

senior $0/5 = 0$ $0/5 = 0$ $1/5$ $(0/5 = 0)$ $p(5/5) = 1/2$ $2/5$

junior $2/6 = 1/3$ $1/6 = 0$ $0/6 = 0$ $1/6$ $2/6 = 1/3$ $0/6$

\therefore so for test instance

$\therefore v_{no} = \operatorname{argmax}_{v_j \in \{\text{yes}, \text{no}\}} p(v_j) \prod_i p(a_i | v_j)$

$\therefore (\text{systems}, "26 \dots 30", "46K-50K")$

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

$\therefore p(\text{Senior} | a) = p(\text{Senior}) * p(\text{systems/Senior}) *$

$p(26 \dots 30 | \text{Senior}) + p(46K-50K | \text{Senior})$

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

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

$$\begin{aligned}\therefore P(\text{junior} | a) &= P(\text{junior}) * P(\text{systems} / \text{junior}) * \\ &P(26 \dots 30 / \text{junior}) * P(46K - 50K / \text{junior}) \\ &= \frac{6}{11} * \frac{1}{3} * \frac{1}{2} * \frac{1}{3} = 0.030\end{aligned}$$

$$\therefore P(\text{junior} | a) = 0.030$$

\therefore by comparing the two values and by the minimum value was junior.

(corrected value)

is so

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

and since $P(\text{senior} | a) < P(\text{junior} | a)$,

so $a = 30$

incorrectly predicted

(M) 100% (V) 0% X 100% = 100%

(correct) is

(100 - 400), 0% - 40% = 60% (incorrect)

confusing model out

so 100% - 40% = 60% (incorrect)

so 100% - 40% = 60% (incorrect)