

QUESTION #02.

As following the order for splitting Age and cheese Content, here is the below 1st split.

$$\text{Age} = \frac{\text{Max} + \text{Min}}{2} = \frac{54 + 8}{2} \Rightarrow \boxed{\text{Age} = 31}$$

Then; split on the left subtree of cheese Content.

$$\text{Cheese Content} = \frac{\text{Max} + \text{Min}}{2} = \frac{9.2 + 6.2}{2} \Rightarrow \boxed{\text{cheese Content} = 7.7}$$

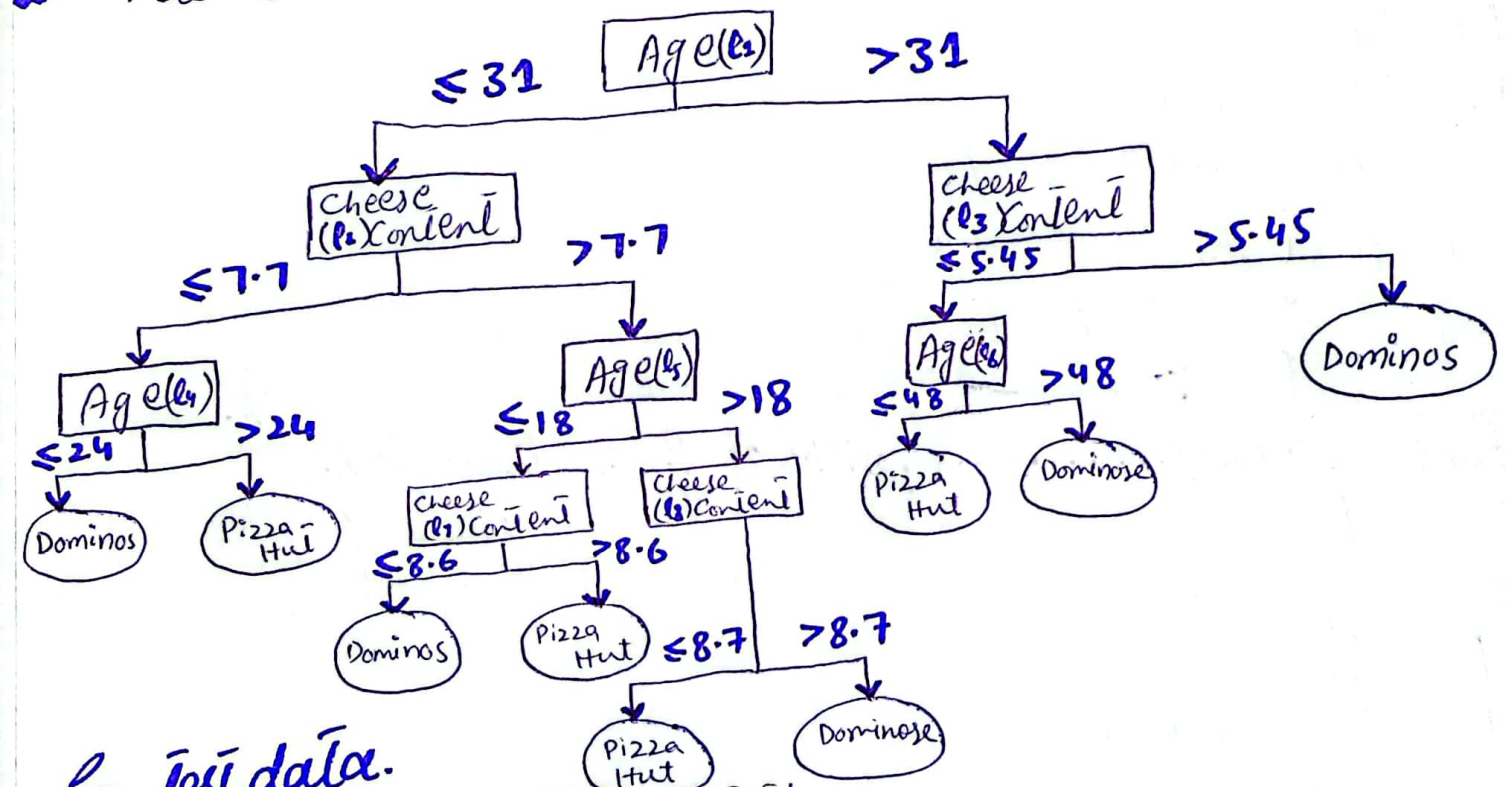
Now; split on the right subtree of cheese content.

$$\text{Cheese Content} = \frac{\text{Max} + \text{Min}}{2} = \frac{7.6 + 3.3}{2} \Rightarrow \boxed{\text{Cheese Content} = 5.45}$$

Now; Split on left subtree of cheese content.

$$\text{Age} = \frac{\text{Max} + \text{Min}}{2} = \frac{27 + 20}{2} \Rightarrow \text{Age} = 23.5 \approx 24 \Rightarrow \boxed{\text{Age} = 24}$$

So; here is below K-D Tree.



So, Test data.

Name	Age	cheeseContent	Pizza Outlet
Harry	46	7	???

So; $l_1 \rightarrow l_3 \rightarrow$ Domino's

So; Results are following

Name	Age	Cheese Content	Pizza Outlet
Harry	46	7	Domino's

QUESTION #03

a) Rule accuracy....., none of R_1 is discarded.

If the examples for R_1 are not discarded then R_2 will be chosen because, it has higher accuracy than R_3 .

Extraction from the figure are given below;

So; All positive examples = 29

All negative examples = 21.

R_1 has 12 positive, 3 negative.

R_2 has 7 positive, 3 negative.

R_3 has 8 positive, 4 negative.

$$\text{Accuracy} = \frac{\text{Positive results}}{\text{Total results}}$$

$$R_1 = 12/16 = 0.8 = 80\%$$

$$R_2 = 7/10 = 0.7 = 70\%$$

$$R_3 = 8/12 = 0.667 = 66.7\%$$

R_1 is best, than R_2 and R_3 .

So, if we discard R_1 then R_2 is best option.

$R_2 > R_3$ (because R_2 has higher accuracy than R_3).

b) Rule accuracy....., positive of R_1 is discarded.

If R_1 positive examples are discarded, then accuracy of R_2 and R_3 is affected because R_1 and R_3 are overlapped.

$$R_2 = 7/10 = 0.7 = 70\%$$

$$R_3 = 6/10 = 0.6 = 60\%$$

So, in this case there is also better to choose R_2 in this case.

$R_2 > R_3$ (because R_2 has more accuracy than R_3)

d. Rule accuracy....., when both positive and negative discard.
if the both positive and negative examples of R_1 is discarded.
So, accuracy updated:

$$R_2 = 7/10 = 0.7 = 70\% \quad R_3 = 6/8 = 0.75 = 75\%$$

So, in this case we preferred R_3 over R_2 .
 $R_3 > R_2$ (because R_3 has ^{more accuracy than} R_2).

QUESTION #04.

Step 1: Calculate distance. (L_2 Norm is Euclidian Distance).

$$d_6 = \sqrt{(120000 - 90000)^2} \Rightarrow \boxed{d_6 = 30000} \checkmark$$

Test data is given below.

(Yes, Single, 120000, No)

a) Unweighted Majority voting.

1st Neighbour: (Yes, Married, 100000, No)

2nd Neighbour: (Yes, Married, 100000, No)

3rd Neighbour: (Yes, Married, 90000, No)

As majority of the Class in Neighbour is "No" then

Our test example is also be "No." (Yes, Single, 120000, No)

b) Distance weight Majority voting.

So as discussed above our 3 nearest neighbours are same.

So, our test example is also be "No."

(Yes, Single, 120000, No)

QUESTION #05.

CONFUSION MATRIX

a). Total Number of Instances.

SUM OF ALL VALUES

Roses	100	10	5
Daisies	15	85	20
Tulips	8	18	90
Total	123	113	115

$$\text{Total} = 123 + 113 + 115$$

$$\boxed{\text{Total} = 351}$$

b). Accuracy Computation.

$$\text{Accuracy} = \frac{\text{Total correct predictions}}{\text{Total No. of predictions}} = \frac{TP + TN}{\text{Total instances}} = \frac{100 + 85 + 90}{351} = \frac{275}{351}$$

$$\text{Accuracy} = 0.7834 \Rightarrow \boxed{\text{Accuracy} = 78.34\%}$$

c). Sensitivity (Recall).

$$\frac{TP}{TP + FN} \Rightarrow \text{Roses} = \frac{100}{100 + 10 + 5} = \frac{100}{115} = 0.869 \Rightarrow \boxed{\text{Roses} = 86.9\%}$$

$$\text{Daisies} = \frac{85}{15 + 85 + 20} = \frac{85}{120} = 0.7083 \Rightarrow \boxed{\text{Daisies} = 70.83\%}$$

$$\text{Tulips} = \frac{90}{8 + 18 + 90} = \frac{90}{116} = 0.7758 \Rightarrow \boxed{\text{Tulips} = 77.58\%}$$

Roses highest sensitivity indicates that best performance is also be "correctly identifying" instance of each class.

Daisies has lower sensitivity which indicates that is "incorrectly identified".

d) Specificity.

$$TN/(TN+FP) \Rightarrow TN/N$$

$$Roses = \frac{25+20+18+90}{25+20+18+90+5} = 0.9505 \Rightarrow \boxed{Roses = 95.05\%}$$

$$Daisies = \frac{100+8+90+5}{100+8+90+5+123} = 0.8913 \Rightarrow \boxed{Daisies = 89.13\%}$$

$$Tulips = \frac{90+12+15+85}{90+12+15+85+13} = 0.9275 \Rightarrow \boxed{Tulips = 92.75\%}$$

Evaluate classified performance for individual classes

e) Precision.

$$\frac{TP}{TP+FP}$$

$$Roses = \frac{25}{25+5} = \frac{200}{223} = 0.897 \Rightarrow \boxed{Roses = 89.7\%}$$

$$Daisies = \frac{85}{85+10+18} = \frac{85}{113} = 0.752 \Rightarrow \boxed{Daisies = 75.2\%}$$

$$Tulips = \frac{90}{90+12+15} = \frac{90}{117} = 0.769 \Rightarrow \boxed{Tulips = 76.9\%}$$

Highest precision Roses \Rightarrow fewer false positives