

1.

a. We need 3 matrices namely Nominal, ordinal and Numeric. Another dij matrix to find distance.

b. Nominal Matrix: (dissimilarity)

	1	2	3	4	5
1	0				
2	1	0			
3	1	1	0		
4	0	1	1	0	
5	1	0	1	1	0

Condition is

$$d = \begin{cases} 0 & \text{if } x=y \\ 1 & \text{if } x \neq y \end{cases}$$

Ordinal Matrix

→ Assign values to ordinal data

Fair = 1, Good = 2

Excellent = 3

	1	2	3	4	5
1	0				
2	1	0			
3	0.5	0.5	0		
4	0	1	0.5	0	
5	1	0	0.5	1	0

Condition:

$$d = \frac{|x-y|}{n-1}$$

n=3

n-1=2

$$d = \frac{|x-y|}{2}$$

Numeric Matrix:

Min = 20

Max = 64

diff = 64 - 20 = 44

	1	2	3	4	5
1	0				
2	0.522	0			
3	0.431	0.954	0		
4	0.159	0.363	0.590	0	
5	0.568	0.045	1	0.049	0

$$d(2,1) = \frac{|22-45|}{44} = \underline{0.522}$$

$$d(2,1) = \frac{|64-45|}{44} = \underline{0.431}$$

$$d(4,1) = \frac{|38-45|}{44} = \underline{0.159}$$

$$d(5,1) = \frac{|20-45|}{44} = \underline{0.568}$$

$$d(3,2) = \frac{|64-22|}{44} = \underline{0.954}$$

$$d(4,2) = \frac{|38-22|}{44} = \underline{0.363}, \quad d(4,3) = \frac{|38-64|}{44} = \underline{0.590}$$

$$d(5,2) = \frac{|20-22|}{44} = \underline{0.045}, \quad d(5,3) = \frac{|20-64|}{44} = \underline{1}$$

$$d(5,4) = \underline{0.049}$$

c. dij

	1	2	3	4	5
1	0				
2	0.840	0			
3	0.643	0.818	0		
4	0.053	0.7876	0.6967	0	
5	0.856	0.015	0.823	0.803	0

$$d(1,1) = 0$$

$$d(2,1) = (1+1+0.522)/3 = 0.840$$

$$d(3,1) = (1+0.5+0.431)/3 = 0.643$$

$$d(3,2) = (1+0.5+0.954)/3 = 0.818$$

$$d(4,1) = (0 + 0.159) / 3 = 0.053$$

$$d(4,2) = (2 + 0.363) / 3 = 0.7876$$

$$d(4,3) = (1 + 0.5 + 0.590) / 3 = 0.6967$$

$$d(5,1) = (1 + 1 + 0.568) / 3 = 0.856$$

$$d(5,2) = (0.045) / 3 = 0.015$$

The value of $d(1,3) = 0.643$

2.

Gini Index:

$$\text{Gini Index} = 1 - \sum P_i^2$$

Age Attribute :- consider mean as 33

Two groups

 ≤ 33 and > 33

	High	Low	Total
Age: ≤ 33	3	1	4
> 33	1	1	2

$$\text{Gini Index } (\leq 33) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 6/16$$

$$\text{Gini Index } (> 33) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 1/2$$

$$\text{Gini Index (Age)} = \left(\frac{4}{6}\right)\left(\frac{6}{16}\right) + \left(\frac{1}{6}\right)\left(\frac{1}{2}\right) = 5/12$$

Car type Attribute

Car type	High	Low	Total
Family	2	1	3
Sports	2	0	2
Truck	0	1	1

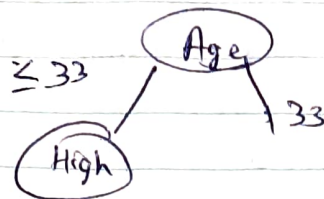
$$\text{Gini (Family)} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 4/9$$

$$\text{Gini (Sports)} = 1 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini (Truck)} = 1 - 1 = 0$$

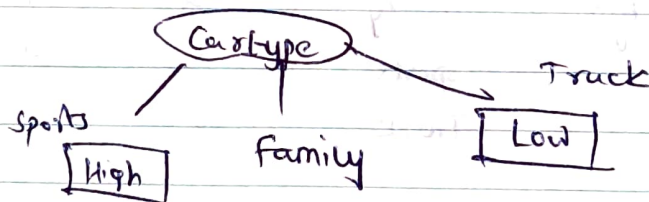
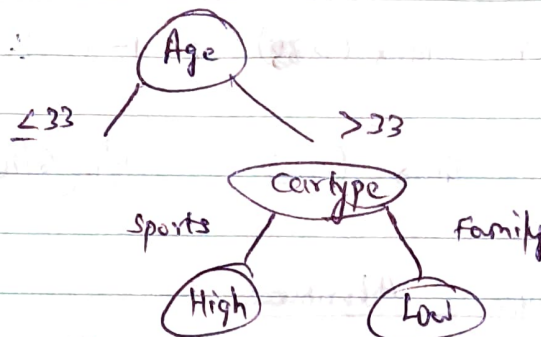
$$\text{Gini (Car Type)} = \frac{3}{6}\left(\frac{4}{9}\right) + 0 + 0 = 2/3$$

Since limit for both attributes are different
we consider Age as Root



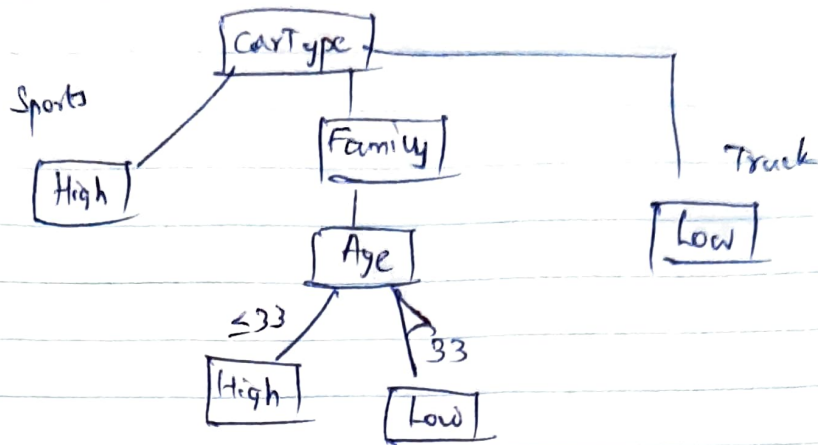
Data for > 33

Tid	Age	CarType	Class
2	43	Sports	High
3	68	family	Low



Data for family

Tid	Age	CarType	Class
0	23	family	High
3	68	family	Low
5	20	family	High



(2) Type of Attributes

Tid \rightarrow nominal

Age \rightarrow Ratio

CarType \rightarrow Ordinal

(3) Stop splitting if every node has the class attribute values

(4) age = 40 carType = family \Rightarrow class = Low

3.

Bayes predict class for x_4 Given $z = \{x_4, \text{excellent}, 4\}$

$$P(A/B) = \frac{P(B/A) P(A)}{P(B)}$$

$$(i) P(x_4|h_1) = \frac{P(h_1/x_4) \cdot P(x_4)}{P(h_1)} = \frac{1 \times 1/10}{6/10} = 1/6$$

$$P(\text{Excellent}|h_1) = \frac{P(h_1/\text{excellent}) P(\text{Excellent})}{P(h_1)}$$

$$= \frac{(1)(3/10)}{6/10} = 1/2$$

$$P(4|h_1) = \frac{P(h_1/4) P(4)}{P(h_1)} = \frac{(1) \cdot 2/10}{6/10} = 1/3$$

$$\boxed{P(z|h_1) = 1/6 \times 1/2 \times 1/3 = 1/36} \quad \text{--- (1)}$$

(ii) for class $h_2 \quad z = \{x_4, \text{excellent}, 4\}$

$$P(x_4|h_2) = \frac{(0)(1/10)}{2/10} = 0$$

$$P(\text{Excellent}|h_2) = \frac{P(h_2/\text{Excellent}) P(\text{Excellent})}{P(h_2)} = \frac{0(3/10)}{2/10} = 0$$

$$P(4|h_2) = \frac{P(h_2/4) P(4)}{P(h_2)} = \frac{(0)(2/10)}{(2/10)} = 0$$

$$\boxed{P(z|h_2) = 0} \quad \text{--- (2)}$$

(iii) for class h_3 $z = \{x_4, \text{excellent}, 4\}$

$$P(x_4 | h_3) = \frac{P(h_3 | x_4) P(x_4)}{P(h_3)} = \frac{0 \times 1/10}{1/10} = 0$$

$$P(\text{Excellent} | h_3) = \frac{0 \times 3/10}{1/10} = 0$$

$$P(4 | h_3) = \frac{0 \times 2/10}{1/10} = 0$$

$$\therefore \boxed{P(z | h_3) = 0} \quad \text{--- (3)}$$

(iv) for class h_4 $z = \{x_4, \text{excellent}, 4\}$

$$P(x_4 | h_4) = \frac{P(h_4 | x_4) P(x_4)}{P(h_4)} = \frac{0 \times 1/10}{1/10} = 0$$

$$P(\text{Excellent} | h_4) = \frac{0 \times 3/10}{1/10} = 0$$

$$P(4 | h_4) = \frac{0 \times 2/10}{1/10} = 0$$

$$\therefore \boxed{P(z | h_4) = 0} \quad \text{--- (4)}$$

from (1), (2), (3) & (4)

$$P(z | h_1) = 4/36$$

$$P(z | h_2) = 0$$

$$P(z | h_3) = 0$$

$$P(z | h_4) = 0$$

4.

KNN with $k=5$ & Euclidean Distance

From the dataset, find the euclidean distance for all

$$\text{formula} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$1. P_1 (158, 58) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-158)^2 + (61-58)^2} = \sqrt{18} = \underline{4.2426}$$

$$2. P_1 (158, 59) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-158)^2 + (61-59)^2} = \sqrt{13} = \underline{3.6055}$$

$$3. P_1 (158, 63) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-158)^2 + (61-63)^2} = \sqrt{13} = \underline{3.6055}$$

$$4. P_1 (160, 59) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-160)^2 + (61-59)^2} = \sqrt{5} = \underline{2.2360}$$

$$5. P_1 (160, 60) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-160)^2 + (61-60)^2} = \sqrt{2} = \underline{1.4142}$$

$$6. P_1 (163, 60) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-163)^2 + (61-60)^2} = \sqrt{5} = \underline{2.2360}$$

$$7. P_1 (163, 61) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-163)^2 + (61-61)^2} = \sqrt{4} = 2$$

$$8. P_1 (160, 64) \quad P_2 (161, 61)$$

$$E.D. = \sqrt{(161-160)^2 + (61-64)^2} = \sqrt{10} = \underline{3.1622}$$

$$9. \quad P_1(163, 64) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-163)^2 + (61-64)^2} = \sqrt{13} = \underline{3.6055}$$

$$10. \quad P_1(165, 61) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-165)^2 + (61-61)^2} = \sqrt{16} = 4$$

$$11. \quad P_1(165, 62) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-165)^2 + (61-62)^2} = \sqrt{17} = \underline{4.1231}$$

$$12. \quad P_1(165, 65) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-165)^2 + (61-65)^2} = \sqrt{32} = \underline{5.6568}$$

$$13. \quad P_1(168, 62) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-168)^2 + (61-62)^2} = \sqrt{50} = \underline{7.071}$$

$$14. \quad P_1(168, 63) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-168)^2 + (61-63)^2} = \sqrt{53} = \underline{7.2801}$$

$$15. \quad P_1(168, 66) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-168)^2 + (61-66)^2} = \sqrt{74} = \underline{8.6023}$$

$$16. \quad P_1(170, 63) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-170)^2 + (61-63)^2} = \sqrt{85} = \underline{9.2195}$$

$$17. \quad P_1(170, 64) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-170)^2 + (61-64)^2} = \sqrt{90} = \underline{9.4868}$$

$$18. \quad P_1(170, 68) \quad P_2(161, 61) \\ E.D. = \sqrt{(161-170)^2 + (61-68)^2} = \sqrt{133} = \underline{11.4017}$$

Here $k=5$, we take 5 points close to $(161, 61)$

Points	Distance	T-shirt Size
$P(160, 60)$	1.4142	M
$P(163, 61)$	2	M
$P(160, 59)$	2.23606	M
$P(163, 60)$	2.23606	M
$P(160, 64)$	3.1622	L

Here we have 4M and 1L

$\therefore (161, 61) \rightarrow M$ (T-shirt size)