

## Question Bank.

DMBI

Q:15  
O:1

Using the data for age given as 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70 answer the following. Use z-score and decimal scaling normalization to transform the values 35 for age. The standard deviation of age is 12.94 years.

Soln:

$$\text{Given: } \sigma = 12.94 ; v = 35$$

(i) Z-Score Normalization

$$v' = \frac{v - \mu}{\sigma} \quad \text{where } \mu = \text{Mean}$$

$$\mu = \frac{(13+15+16+16+19+20+20+21+22+22+25+25+25+25+30+33+33+35+35+35+35+36+40+45+46+52+70)}{27}$$

$$\mu = \frac{809}{27}$$

$$= 29.96$$

$$\therefore v' = \frac{35 - 29.96}{12.94}$$

$$= \frac{5.04}{12.94}$$

$$v' = 0.389$$

(ii) Normalization by decimal scaling

$$v' = \frac{v}{10^j}$$

where  $j$  is the smallest integer

such that  $\max(|v'|) < 1$

$$v' = \frac{v}{10^j}$$

$$v' = \frac{35}{10^2}$$

$$v' = 35$$

100

$$\boxed{v' = 0.35}$$

Q:26. Suppose a group of 12 sales price records have been  
Q:2 sorted as follows:

5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215.

Partition them into three bins by equal-frequency  
(equal-depth) partitioning.

(i) Equal frequency partitioning.

bin 1 - 5, 10, 11, 13

bin 2 - 15, 35, 50, 55

bin 3 - 72, 92, 204, 215

Q:27 Use K-Means algorithm to create 2 clusters for  
given set of values {2, 3, 8, 9, 12, 15, 18, 22}

Given K = 2

{2, 3, 8, 9, 12, 15, 18, 22}

Step 1: M1 = 3, M2 = 12

M1 = 2

M2 = 2

K1 = {2}

K1 = {2}

K1 = {2, 3}

K2 = {3, 8, 9, 12, 15, 18, 22}

Q: 27.

Q: 3.

Use K-Means algorithm to create 2 clusters for given set of values  $\{2, 3, 8, 9, 12, 15, 18, 22\}$

Given  $K = 2$

$\{2, 3, 8, 9, 12, 15, 18, 22\}$

$$P + 8 + 52 + 81 + 21 + 51 = 217$$

Step 1:  $M_1 = 2$      $M_2 = 3$

$$K_1 = \{2\}$$

$$K_2 = \{3, 8, 9, 12, 15, 18, 22\}$$

Step 2: New centroid points have to be calculated

$$M_1 = ?$$

$$M_2 = ?$$

$$M_1 = \frac{2+3}{2} = 2.5 \quad M_2 = \frac{8+9+12+15+18+22}{6} = 14$$

$$M_2 = 12.4$$

$$K_1 = \{2, 3\}$$

$$K_2 = \{8, 9, 12, 15, 18, 22\}$$

Step 3: Realign the centroid.

$$M_1 = ? \quad M_2 = ?$$

$$M_1 = \frac{2+3}{2}$$

$$M_2 = \frac{8+9+12+15+18+22}{6}$$

$$= 2.5$$

$$= 14$$

$$K_1 = \{2, 3, 8\}$$

$$K_2 = \{9, 12, 15, 18, 22\}$$

Step 4 : Realign the centroid.

$$M_1 = ?$$

$$M_2 = ?$$

$$M_1 = \frac{2+3+8}{3}$$

$$M_2 = \frac{9+12+15+18+22}{5}$$

$$= 4.3$$

$$M_2 = 15.2$$

$$K_1 = \{2, 3, 8, 9\}$$

$$K_2 = \{12, 15, 18, 22\}$$

Step 5 : Realign the centroid.

$$M_1 = \frac{2+3+8+9}{4} \quad M_2 = \frac{12+15+18+22}{4}$$

$$M_1 = 5.5 \quad M_2 = 16.75$$

$$K_1 = \{2, 3, 8, 9\}$$

$$K_2 = \{12, 15, 18, 22\}$$

Step 6 : Realign the centroid.

$$M_1 = 5.5 \quad M_2 = 16.75$$

∴ Final clusters are  $K_1 = \{2, 3, 8, 9\}$

$$K_2 = \{12, 15, 18, 22\}$$

Q: 42.  
0: 4

Use K-Means algorithm to create 3 clusters for given set of values  $\{2, 3, 8, 9, 12, 15, 18, 22\}$

$$\text{Given } k = 3 \quad \{2, 3, 8, 9, 12, 15, 18, 22\}$$

$$\text{Step 1: } M_1 = 2 \quad M_2 = 3 \quad M_3 = 8$$

$$K_1 = \{2\}$$

$$K_2 = \{3\}$$

$$K_3 = \{8, 9, 12, 15, 18, 22\}$$

Step 2: New centroid points have to be calculated.

$$M_1 = ? \quad M_2 = ? \quad M_3 = ?$$

$$M_1 = 2 \quad M_2 = \frac{3+8}{2} = 5.5 \quad M_3 = \frac{8+9+12+15+18+22}{6} = 14.2$$

$$K_1 = \{2\}$$

$$K_2 = \{3, 8\}$$

$$K_3 = \{9, 12, 15, 18, 22\}$$

Step 3: Recalculate the centroid

$$M_1 = 2$$

$$M_2 = \frac{3+8}{2} = 5.5$$

$$M_3 = \frac{9+12+15+18+22}{5} = 15.2$$

$$K_1 = \{2\}$$

$$K_2 = \{3, 8, 9\}$$

$$K_3 = \{12, 15, 18, 22\}$$

Step 4: Recalcul the centroid

$$M_1 = 2$$

$$M_2 = \frac{3+8+9}{3}$$

$$= 6.66$$

$$M_3 = \frac{12+15+18+22}{4}$$

$$= 16.75$$

$$K_1 = \{2\}$$

$$K_2 = \{3, 8, 9\}$$

$$K_3 = \{12, 15, 18, 22\}$$

Step 5: Recalcul the centroid

$$M_1 = 2$$

$$M_2 = 6.66$$

$$M_3 = 16.75$$

$\therefore$  Final clusters are  $K_1 = \{2\}$

$$K_2 = \{3, 8, 9\}$$

$$K_3 = \{12, 15, 18, 22\}$$

Q:43. Suppose that the data for analysis including the attribute age. The age value for data tuple are (in increasing order): 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 50.

- What is mean, Median, Mode, modality and Midrange of the data.
- Show a boxplot of the data and state five number summary
- Does outlier exists for given distribution?

a]

Mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$= \frac{1}{27} (809)$$

$$\bar{x} = 29.96$$

$$\text{Median} = \frac{n+1}{2}$$

$$= \frac{27+1}{2}$$

= 14<sup>th</sup> value

$$\text{Median} = 25$$

Mode

13, 15, 19, 21, 30, 40, 45, 46, 52, 70 occur once in data

16, 20, 22, 23 occur twice in data.

25, 35 occur four times in data.

Modality  $\rightarrow$  Bimodal

Midrange = Min Value + Max Value

$$= 13 + 70$$

$$= \frac{83}{2}$$

$$= 41.5$$

b)

$$Q_1 = \frac{20+20}{2}$$

$$Q_1 = 20$$

$$Q_3 = \frac{35+35}{2}$$

$$Q_3 = 35$$

b)

$$Q_1 = 20$$

$$Q_2 = 25$$

$$Q_3 = 35$$

five number summary

$$= (13, 20, 25, 35, 70)$$

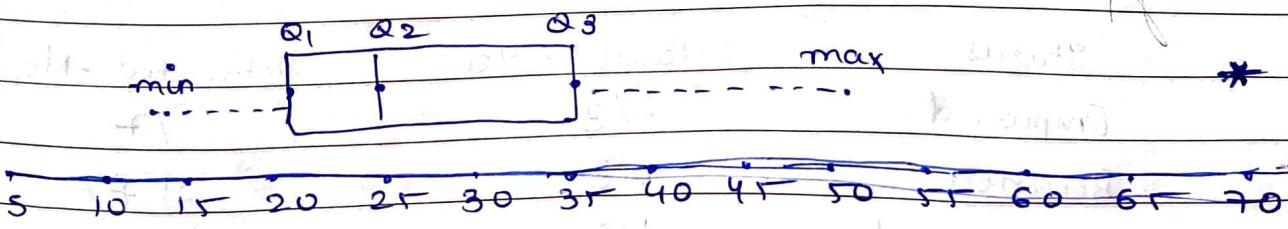
$$\text{IQR} = Q_3 - Q_1$$

$$= 35 - 20$$

$$= 15$$

$$Q_1 + (1.5 \times \text{IQR}) = Q_1 + (1.5 \times 15) \\ = 20 + 22.5 = 42.5$$

$$Q_3 + (1.5 \times \text{IQR}) = 35 + (1.5 \times 15) \\ = 35 + 22.5 = 57.5$$



c) Yes, 70'

Q-6. Illustrate Naïve Bayesian Classification technique for the given data set. Show how can classify a new tuple, with (Homeowner = yes; status = Employed; Income = Average)

ID	Homeowner	Status	Income	Defaulted
1	Yes	E	H	No
2	No	B	A	No
3	No	E	L	No
4	Yes	B	H	No
5	No	UE	A	Y
6	No	B	L	N
7	Yes	UE	H	No
8	No	E	A	Y
9	No	B	L	N
10	No	E	A	Y

$$P(\text{Yes}) = \frac{3}{10}$$

$$P(\text{No}) = \frac{7}{10}$$

HomeOwner	Defaulted - Yes	Defaulted - No
Yes	0/3	3/7
No	3/3	4/7

Status	Defaulted - Yes	Defaulted - No
Employed	2/3	2/7
Business	0/3	4/7
unemployed	1/3	1/7

Income	Defaulted - Yes	Defaulted - No
High	0/3	3/7
Average	3/3	1/7
Low	0/3	3/7

$$\begin{aligned}
 P(x' \text{ Yes}) &= P(\text{Yes} | \text{Yes}) \times P(\text{Employed} | \text{Yes}) \times P(\text{Average} | \text{Yes}) \\
 &\quad \times P(\text{Yes}) \\
 &= 0 \times \frac{2}{3} \times \frac{3}{3} \times \frac{3}{10} \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 P(x' \text{ No}) &= P(\text{Homeowner}^{\text{No}} | \text{No}) \times P(\text{Employed} | \text{No}) \times P(\text{Average} | \\
 &\quad \text{No}) \times P(\text{No}) \\
 &= \frac{4}{7} \times \frac{2}{7} \times \frac{1}{7} \times \frac{7}{10} \\
 &= 0.016 \\
 &= 0 < 0.016
 \end{aligned}$$

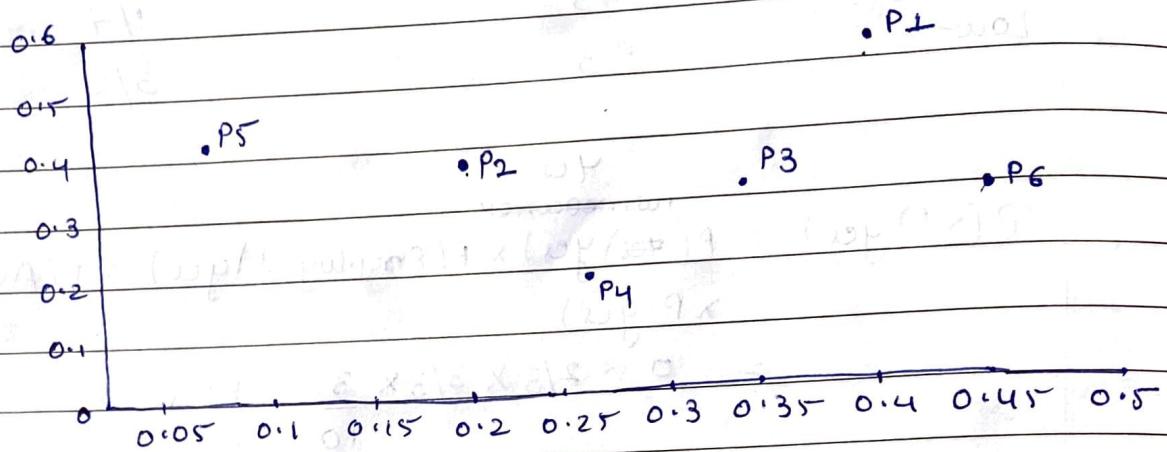
Therefore the Naive Bayesian classifier predicts  
 Defaulted = "No" for sample x.

Q.2. Assume that the database D is given by the table below. Follow a single link technique to find clusters in D. Use Euclidean distance measure.

Point	x	y
P1	0.40	0.53
P2	0.22	0.38
P3	0.35	0.32
P4	0.26	0.19
P5	0.08	0.41
P6	0.45	0.30

Step 1 -

Plot the given dataset in graph format for better understanding.



Step 2 -

Create a Distance Matrix by using the Euclidean distance formula:

$$\begin{aligned} \text{Distance } (P_1, P_2) &= \sqrt{(0.4 - 0.22)^2 + (0.53 - 0.38)^2} \\ &= \sqrt{0.0324 + 0.0225} \\ &= \sqrt{0.0549} \\ &= 0.23 \end{aligned}$$

$$\begin{aligned} \text{Distance } (P_1, P_3) &= \sqrt{(0.4 - 0.35)^2 + (0.53 - 0.32)^2} \\ &= \sqrt{0.1 + 0.412} = \sqrt{0.0025 + 0.0441} \\ &= 0.21 \end{aligned}$$

$$\begin{aligned} \text{Distance } (P_1, P_4) &= \sqrt{(0.4 - 0.26)^2 + (0.53 - 0.19)^2} \\ &= 0.36 \end{aligned}$$

$$\begin{aligned} \text{Distance } (P_1, P_5) &= \sqrt{(0.4 - 0.08)^2 + (0.53 - 0.41)^2} \\ &= 0.34 \end{aligned}$$

$$\begin{aligned} \text{Distance } (P_1, P_6) &= \sqrt{(0.4 - 0.45)^2 + (0.53 - 0.3)^2} \\ &= 0.23 \end{aligned}$$

$$\text{Distance } (P_2, P_3) = \sqrt{(0.22 - 0.35)^2 + (0.38 - 0.32)^2}$$

$$= \cancel{0.15} \quad 0.15$$

$$\text{Distance } (P_2, P_4) = \sqrt{(0.22 - 0.26)^2 + (0.38 - 0.19)^2}$$

$$= \cancel{0.20} \quad 0.19$$

$$\text{Distance } (P_2, P_5) = \sqrt{(0.22 - 0.08)^2 + (0.38 - 0.41)^2}$$

$$= \cancel{0.14} \quad 0.14$$

$$\text{Distance } (P_2, P_6) = \sqrt{(0.22 - 0.45)^2 + (0.38 - 0.30)^2}$$

$$= \cancel{0.28} \quad 0.24$$

$$\text{Distance } (P_3, P_4) = \sqrt{(0.35 - 0.26)^2 + (0.32 - 0.19)^2}$$

$$= \cancel{0.15} \quad 0.15$$

$$\text{Distance } (P_3, P_5) = \sqrt{(0.35 - 0.08)^2 + (0.32 - 0.41)^2}$$

$$= \cancel{0.28} \quad 0.28$$

$$\text{Distance } (P_3, P_6) = \sqrt{(0.35 - 0.45)^2 + (0.32 - 0.30)^2}$$

$$= \cancel{0.20} \quad 0.10$$

$$\text{Distance } (P_4, P_5) = \sqrt{(0.26 - 0.08)^2 + (0.19 - 0.41)^2}$$

$$= \cancel{0.28}$$

$$\text{Distance } (P_4, P_6) = \sqrt{(0.26 - 0.45)^2 + (0.19 - 0.30)^2}$$

$$= \cancel{0.21}$$

$$\text{Distance } (P_5, P_6) = \sqrt{(0.08 - 0.45)^2 + (0.41 - 0.30)^2}$$

$$= \cancel{0.38}$$

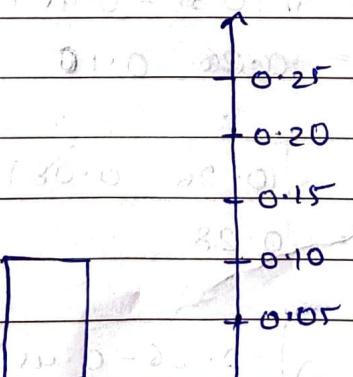
Based on the above computed values Distance Matrix can be formed as follows.

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$
$P_1$	0	$0.23$	$0.21$	$0.36$	$0.34$	$0.23$
$P_2$	$0.23$	0	$0.19$	$0.19$	$0.24$	$0.24$
$P_3$	$0.21$	$0.19$	0	$0.15$	$0.14$	$0.10$
$P_4$	$0.36$	$0.19$	$0.15$	0	$0.28$	$0.21$
$P_5$	$0.34$	$0.14$	$0.28$	$0.28$	0	$0.10$
$P_6$	$0.23$	$0.24$	$0.10$	$0.21$	$0.38$	0

Step 3 : Find the minimum value element from the distance matrix ( $P_3, P_6$ )

The minimum value element is  $(P_3, P_6)$  with a value of  $0.10 - 1^{st}$  center  $(P_3, P_6)$

Step 4 - Now Recalculate or update the distance matrix for the cluster  $(P_3, P_6)$



Step 4 - Now Recalculate or update the distance matrix for the cluster  $(P_3, P_6)$

$\min[\text{dist}(\text{point1}, \text{point2})]$

$$\min[\text{Dist}((P_3, P_6), P_1)] = \min(0.21, 0.23) = 0.21$$

$$\min[\text{Dist}((P_3, P_6), P_2)] = \min(0.14, 0.24) = 0.14$$

$$\min[\text{Dist}((P_3, P_6), P_4)] = \min(0, 0.21) = 0$$

$$\min[\text{Dist}((P_3, P_6), P_5)] = \min(0, 0.38) = 0$$

Therefore, the updated distance matrix looks as follows:

	P <sub>1</sub>	P <sub>2</sub>	P <sub>3, P<sub>6</sub></sub>	P <sub>4</sub>	P <sub>5</sub>
P <sub>1</sub>	0				
P <sub>2</sub>	0.23	0			
P <sub>3, P<sub>6</sub></sub>	0.21	0.15	0		
P <sub>4</sub>	0.36	0.19	0.15	0	
P <sub>5</sub>	0.34	0.14	0.28	0.28	0

Step 5 - Repeat steps 3 and 4

The min value element is (P<sub>2</sub>, P<sub>5</sub>) with a value is 0.14 - 2<sup>nd</sup> Cluster (P<sub>2</sub>, P<sub>5</sub>).

Recalculate or update the distance matrix for the cluster (P<sub>2</sub>, P<sub>5</sub>).

$\min[\text{dist}(\text{point1}, \text{point2})]$

$$\min[\text{Dist}((P_2, P_5), P_1)] = \min(0.23, 0.34) = 0.23$$

$$\min[\text{Dist}((P_2, P_5), P_2)] = \min(0, 0.14) = 0$$

Similarly, perform for (P<sub>3</sub>, P<sub>6</sub>) and P<sub>4</sub>.

Therefore, the updated Distance Matrix looks as follows.

		P1	P2, P5	P3, P6	P4
P1	0.0	0.23	0.21	0.36	
P2, P5	0.23	0.0	0.15	0.19	
P3, P6	0.21	0.15	0	0.15	
P4	0.36	0.19	0.15	0	

Step 6 - Repeat steps 3 & 4

The minimum value element is  $(P_2, P_5, P_3, P_6)$  with a value 0.15 - 3rd cluster  $(P_2, P_5, P_3, P_6)$ .

Here 2 values are the same then the first element is chosen as the minimum value element.

Recalculate or update the distance matrix for the cluster  $(P_2, P_5, P_3, P_6)$ .

$$\min [\text{dist}(\text{point1}, \text{point2})]$$

$$\min [\text{dist}((P_2, P_5, P_3, P_6), P_1)] = \min [0.23, 0.21] = 0.21$$

$$\min [\text{dist}((P_2, P_5, P_3, P_6), P_4)] = 0.$$

$P_1$  $P_2, P_5, P_3, P_6$  $P_4$  $P_1$ 

0

 $P_2, P_5, P_3, P_6$ 

0.21

0

 $P_4$ 

0.36

0.15

0

Step 7 -

Min [dust (point 1, point 2)]

$$\min [\text{dust} [(P_2, P_5, P_3, P_6, P_4), P_1]] = \min [\text{dust}[0.21, 0.36]] = 0.21$$

 $P_1$  $P_2, P_5, P_3, P_6, P_4$  $P_1$ 

0.21

 $P_2, P_5, P_3, P_6, P_4$ 

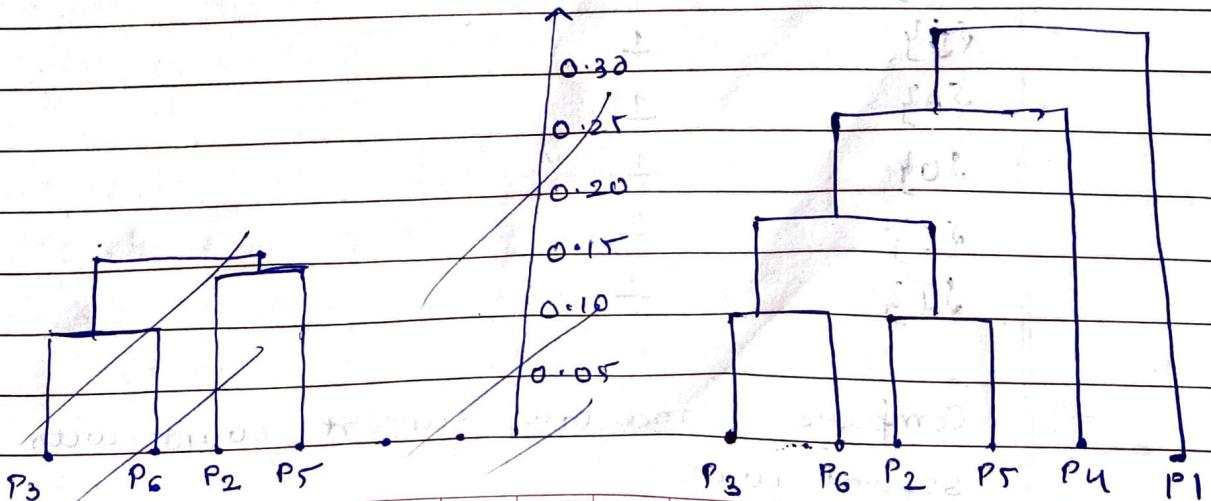
0.21

0.

Step 8 - Repeat the steps 3 &amp; 4.

The minimum value elements are  $(P_2, P_5, P_3, P_6, P_4, P_1)$  with a value is 0.21 - 5<sup>th</sup> cluster  $(P_2, P_5, P_3, P_6, P_4, P_1)$ . In this step, only 1 value is remaining so it is by default cluster.

Step 9 - Dendrogram



Q:8.

Consider the following database with minimum support count = 60%. find all the frequent itemsets using apriori algorithm and also generate strong association rule if minimum confidence = 80%.

T-10      Items bought

T1      M, O, N, K, E, Y

T2      D, O, N, K, E, Y

T3      M, A, K, E

T4      M, U, C, K, Y

T5      C, O, O, K, I, E

$$\text{Min support count} = \frac{60 \times 5}{100} = 3$$

Step 1 : Generating 1-itemset frequent patterns.

C1      Itemset      Support

$\{M\}$       3

$\{O\}$       4

$\{N\}$       2

$\{K\}$       5

$\{E\}$       4

$\{Y\}$       3

$\{D\}$       1

$\{A\}$       1

$\{U\}$       1

$\{C\}$       2

$\{I\}$       1

Compare candidate support count with minimum support count.

$\square$  Step 1 : Generate frequent itemset -  $\geq 3$  occurrences

L1 Itemset Support Count

{M,Y}	3
{O,Y}	4
{K,Y}	5
{E,Y}	4
{Y,Y}	3

Step 2 : Generate C2 - Itemset frequent Pattern.

C2 Item set Support Count

{M,O,Y}	1
{M,K,Y}	3
{M,E,Y}	2
{M,Y,Y}	2
{O,K,Y}	3
{O,E,Y}	3
{O,Y,Y}	2
{K,E,Y}	4
{K,Y,Y}	3
{E,Y,Y}	2

Compare candidate support count with minimum support count

L2 Item set Support Count

{M,K,Y}	3
{O,K,Y}	3
{O,E,Y}	3
{K,E,Y}	4
{K,Y,Y}	3

## Step 3 - Generating 3-itemset Frequency Pattern

C3	Itemset	Support Count
	{MOKY}	1
	{OKEY}	3
	{KEY}	2

Compare candidate support count with minimum support count for all 3 itemsets.

L3	Itemset	Support Count
	{OKEY}	3
		2
		2

Q:9. Suppose a hospital tested the age and body fat data for 18 randomly selected adults with the following result.

age	23	23	27	27	39	41	47	49
%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2

age	52	54	54	56	57	58	58	60
%fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2

(a) Calculate the mean, median and standard deviation of age and %fat

(b) Draw the boxplot for age.

$$(a) \text{ Mean} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$= \frac{1}{18} \times 836$$

$$= 46.44$$

$$\text{Mean} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$= \frac{1}{18} \times 488.518$$

$$= 28.77$$

$$\text{Median} = \frac{50+52}{2}$$

$$= 51$$

$$\text{Median} = \frac{30+2+31.2}{2}$$

$$= 30.2$$

% fat in increasing order.

7.8, 9.5, 17.8, 25.9, 26.5, 27.2, 27.4, 28.8, 30.2, 31.2,  
31.4, 32.9, 33.4, 34.1, 34.8, 35.7, 41.2, 42.5

### Standard Deviation of age.

Age	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$
23	-23.44	549.4336
23	-23.44	549.4336
27	-19.44	377.9136
27	-19.44	377.9136
39	-7.44	55.3536
41	-5.44	29.5936
47	0.56	0.3136
49	2.56	6.5536
50	3.56	12.6736
52	5.56	30.9136
54	7.56	57.1536
54	7.56	57.1536
56	9.56	91.3936
57	10.56	111.5136
58	11.56	133.6336

58

11.56

133.6336

60

13.56

183.8736

61

14.56

211.9936

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$= \frac{2970.448}{18}$$

$$\sigma^2 = 165.02$$

Standard Deviation =  $\sqrt{165.02} = 12.84$

$$\sigma = \sqrt{165.02}$$

$$\sigma_{age} = 12.84$$

### Standard Deviation of % fat

% fat

$$(x_i - \bar{x}) \cdot P_i =$$

$$(x_i - \bar{x})^2$$

7.8

$$= 21$$

$$441$$

9.5

$$-19.3$$

$$372.49$$

17.8

$$11$$

$$121$$

25.9

$$-2.9$$

$$8.41$$

26.5

$$2.3$$

$$5.29$$

27.2

$$-1.6$$

$$2.56$$

27.4

$$-1.4$$

$$1.96$$

28.8

$$0$$

$$0$$

30.2

$$1.4$$

$$1.96$$

31.2

$$2.4$$

$$5.76$$

31.4

$$2.6$$

$$6.76$$

32.9	4.1	16.81
33.4	4.6	21.16
34.1	5.3	28.02
34.6	5.8	33.64
35.7	6.9	47.61
41.2	12.4	153.76
42.5	13.2	187.69

$$\sigma^2 = 1455.95$$

$$= \sqrt{80.886}$$

$$\sigma\% \text{ fat} = \sqrt{80.886}$$

$$= 8.99 \approx 9$$

(b) Draw boxplot for age.

$$Q_1 = 39$$

$$Q_2 = 51$$

$$Q_3 = 57$$

$$\min = 23$$

$$\max = 61$$

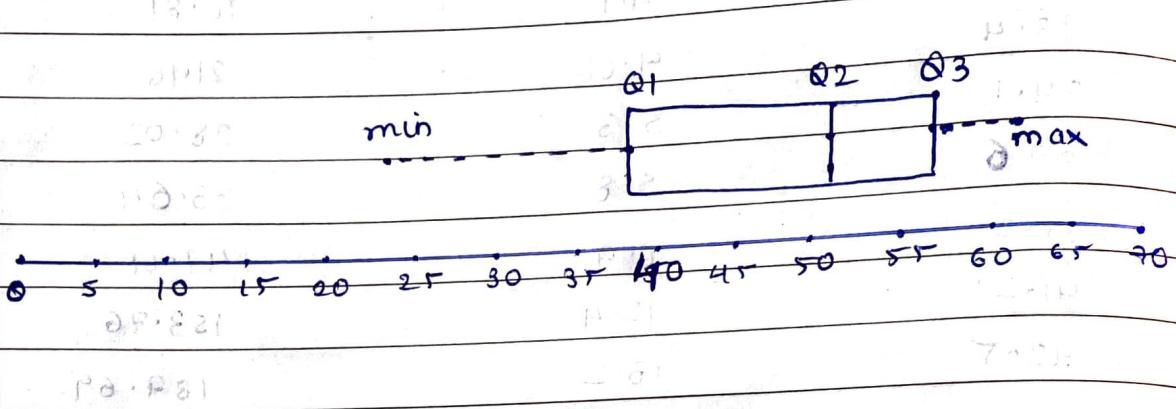
$$IQR = Q_3 - Q_1$$

$$= 57 - 39$$

$$= 18$$

$$Q_1 - (1.5 \times IQR) = 39 - (1.5 \times 18) = 12$$

$$Q_3 + (1.5 \times IQR) = 57 + (1.5 \times 18) = 84$$

Q:10.

Consider the following database with minimum support count = 50%. Find all the frequent itemset using apriori algorithm and also generate strong association rules if minimum confidence = 50%.

TID	List of ItemIDs
1	A, B, D
2	A, D
3	A, C
4	B, D, E, F

$$\begin{aligned} \text{Min Support count} &= \frac{50 \times 4}{100} \\ &= 2 \end{aligned}$$

Step 1 - Generating 1-itemset frequency pattern

CL	itemset	sup	itemset	sup
	{A}	3	{A}	3
	{B}	2	{B}	2
	{D}	3	{D}	3
	{C}	1		
	{E}	1		
	{F}	1		

Step 2 - Generating 2-itemset frequency pattern

itemset		sup	$L_2$	itemset	sup
$\{A, B\}$	1	1	$\{A, D\}$	2	2
$\{A, D\}$	2	2	$\{B, D\}$	2	2
$\{B, D\}$	2	2			

Step 3 - Generating 3-itemset frequency pattern

itemset	sup
$\{A, B, D\}$	1

Therefore frequent item sets are

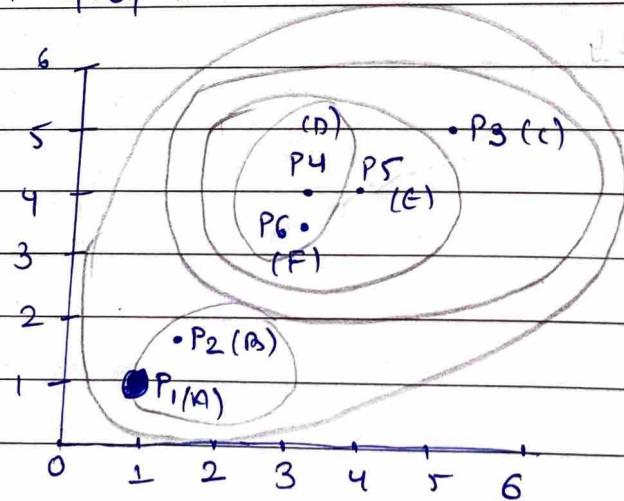
$$L_1 = \{A, D\}$$

$$L_2 = \{B, D\}$$

Q11. Suppose we have six objects (with name A, B, C, D, E and F) and each object have two numerical features [ $x_1$  and  $x_2$ ] Apply single linkage clustering and draw dendrogram.

Object	Attribute $x_1$	Attribute $x_2$
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5

Step 1 Plot Graph



Step 2.

(Creating distance matrix by using the Euclidean distance formula.)

$$\text{Distance } (P_1, P_2) = \sqrt{(1-1.5)^2 + (1-1.5)^2} \quad (\rightarrow)$$

$$\begin{aligned} \text{Distance } (A, B) &= \sqrt{(1-1.5)^2 + (1-1.5)^2} \\ &= 0.71 \end{aligned}$$

$$\text{Distance } (A, C) = \sqrt{(1-5)^2 + (1-5)^2} \\ = 5.66$$

$$\text{Distance } (A, D) = \sqrt{(1-3)^2 + (1-4)^2} \\ = 3.61$$

$$\text{Distance } (A, E) = \sqrt{(1-4)^2 + (1-4)^2} \\ = 4.24$$

$$\text{Distance } (A, F) = \sqrt{(1-3)^2 + (1-3-5)^2} \\ = 3.20$$

$$\text{Distance } (B, C) = \sqrt{(1.8-5)^2 + (1.5-5)^2} \\ = 4.95$$

$$\text{Distance } (B, D) = \sqrt{(1.5-3)^2 + (1.5-4)^2} \\ = 2.92$$

$$\text{Distance } (B, E) = \sqrt{(1.5-4)^2 + (1.5-4)^2} \\ = 3.54$$

$$\text{Distance } (B, F) = \sqrt{(1.5-3)^2 + (1.5-3-5)^2} \\ = 2.50$$

$$\text{Distance } (C, D) = \sqrt{(5-3)^2 + (5-4)^2} \\ = 2.24$$

$$\text{Distance } (C, E) = \sqrt{(5-4)^2 + (5-4)^2} \\ = 1.41$$

$$\text{Distance } (C, F) = \sqrt{(5-3)^2 + (5-3-5)^2} \\ = 2.50$$

$$\text{Distance}(D, E) = \sqrt{(3-4)^2 + (4-4)^2} \\ = 1.00$$

$$\text{Distance}(D, F) = \sqrt{(3-3)^2 + (4-3.5)^2} \\ = 0.5$$

$$\text{Distance}(E, F) = \sqrt{(4-3)^2 + (4-3.5)^2} \\ = 1.12$$

Based on the above computed values, distance matrix can be formed as follows.

	A	B	C	D	E	F
A	0					
B	0.71	0				
C	0.66	0.95				
D	3.61	2.92	2.24	0		
E	4.24	3.54	1.41	1.00	0	
F	3.20	2.50	2.10	0.50	1.12	0

Step 3  $\rightarrow$  min value = 0.5 center (D, F)

$$d(D, F) \rightarrow A = \min(3.61, 3.20) \\ = 3.20$$

$$d(D, F) \rightarrow B = \min(2.92, 2.50) \\ = 2.5$$

$$d(D, F) \rightarrow C = \min(2.24, 2.5) \\ = 2.24$$

$$d(E) \rightarrow D, F = \min(1, 1.12) \\ = 1$$

	A	B	C	D, F	E
A	0				
B	0.71	0			
C	5.66	4.95	0		
D, F	3.20	2.5	2.24	0	
E	4.24	3.54	1.41	1	0

Step 4  $\rightarrow$  min value = 0.71 center (A, B)

$$d_{C \rightarrow (A,B)} d((A,B) \rightarrow C) = \min(5.66, 4.95)$$

$$d_{(D,F) \rightarrow (A,B)} d((A,B) \rightarrow D,F) = \min(3.20, 2.5, 3.20, 2.50) \\ = 2.50$$

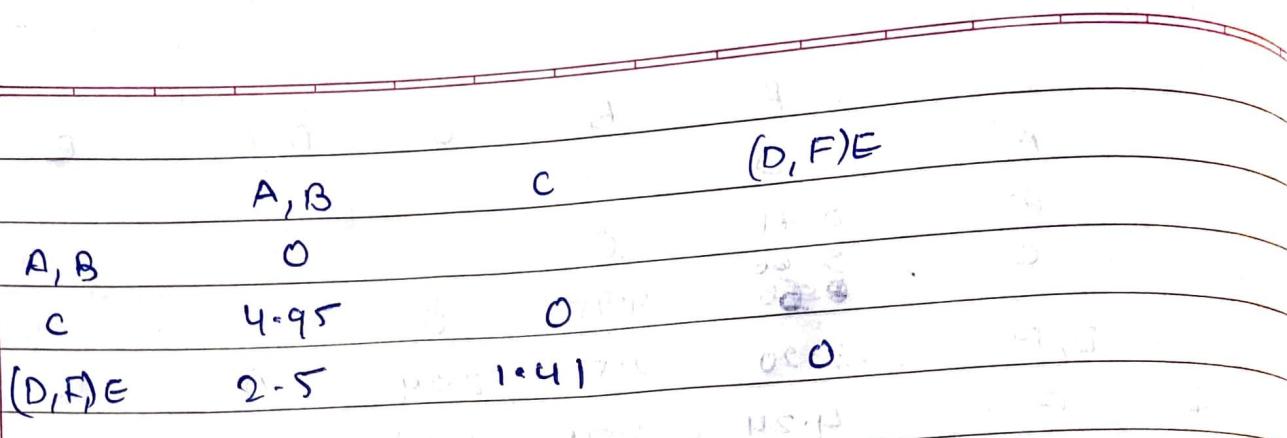
$$d_{(E) \rightarrow (A,B)} d((A,B) \rightarrow E) = \min(4.24, 3.54) \\ = 3.54$$

	A, B	C	D, F	E
A, B	(3.20, 0.71)			
C	4.95	0	0	(5.66)
D, F	2.5	2.24	0	(3.20, 2.50)
E	3.54	1.41	1.0	0

Step 5  $\rightarrow$  min value = 1 center (D, F), E

$$d \rightarrow (D,F) E \rightarrow AB = \min(2.5, 3.54) \\ = 2.5$$

$$d \rightarrow (D,F) \rightarrow C = \min(2.24, 1.41) \\ = 1.41$$

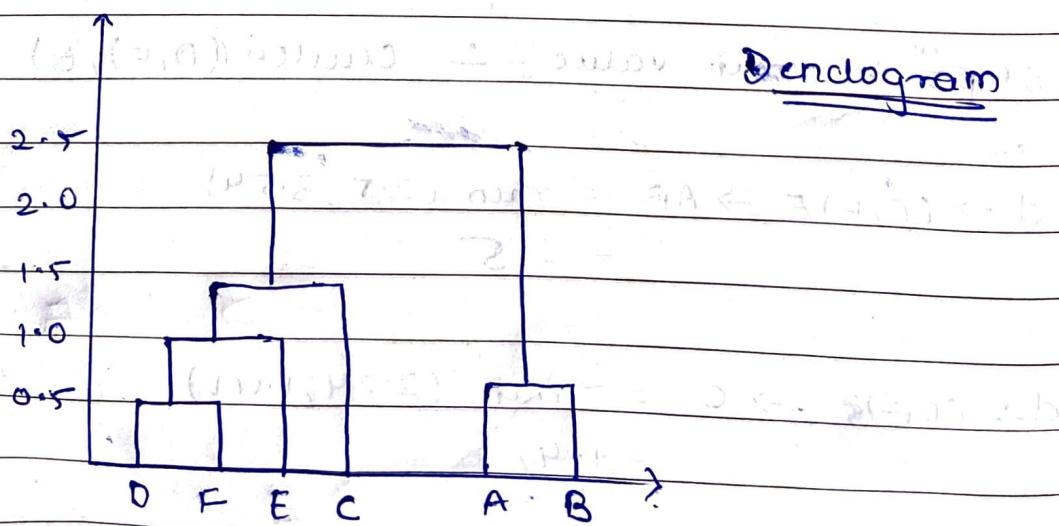


$\rightarrow (D, F) \cup E \rightarrow AB = \min(2.5, 1.41) = 1.41$

Step 6  $\rightarrow$  min value  $\rightarrow 1.41$  center  $((D, F), E), C$

$A, B$        $((D, F), E), C$   
 $(A, B)$       0      0  
 $((D, F), E), C$       2.5

Step 7  $\rightarrow$  min. value  $\rightarrow 2.5$  center  $((((D, F), E), C), (A, B))$



S.12

Find all the frequent item set and also generate association rule using Apriori algorithm for the given transaction. Assume minimum support be 50% and minimum confidence be 80%.

TID	item
100	A, C, D
200	B, C, E
300	A, B, C, E
400	B, E

$$\text{Min. Support count} = \frac{50 \times 4}{100}$$

$$= 2 \text{ frequent itemset}$$

TID	item	freq itemset
100	A C D	
200	B C E	<del>ABC</del>
300	A B C E	
400	B E	

Step 1  $\rightarrow$  Generating 1-item frequency pattern

1. items	sup	4	2. TID with item	sup
TID	item			
{A}	2		{A}	2
{B}	3		{B}	3
{C}	3		{C}	3
{D}	1			
{E}	3		{E}	3

$C_2$	itemset	sup	$\rightarrow L_2$	itemset	sup
	$\{A, B\}$	1		$\{A, C\}$	2
	$\{A, C\}$	2		$\{B, C\}$	2
	$\{A, E\}$	1		$\{B, E\}$	3
	$\{B, C\}$	2		$\{C, E\}$	2
	$\{B, E\}$	3			0.5
	$\{C, E\}$	2			0.5
					0.5
					0.5

$C_3$	itemset	sup	$\rightarrow L_3$	itemset	sup
	$\{A, B, C\}$	1		$\{B, C, E\}$	2
	$\{B, C, E\}$	2			0.5

Therefore frequent item sets are  
 $L = \{B, C, E\}$

Now, for strong association rule

Generate non-empty subset for  $L$

$$S = \{\{B\}, \{C\}, \{E\}, \{B, C\}, \{B, E\}, \{C, E\}\}$$

$$\text{Confidence } (X \rightarrow Y) = P(Y | X) = P(X \cup Y) / P(X)$$

Association Rule	Support	Confidence	Confidence %
$\{B\} \rightarrow \{C, E\}$	2	2/3	66.66
$\{C\} \rightarrow \{B, E\}$	2	2/3	66.66
$\{E\} \rightarrow \{B, C\}$	2	2/3	66.66
$\{B, C\} \rightarrow \{E\}$	2	2/3	66.66
$\{B, E\} \rightarrow \{C\}$	2	2/2	100
$\{C, E\} \rightarrow \{B\}$	2	2/3	66.66
		2/2	100

Therefore strong association rule are

$$\{B, C\} \rightarrow \{E\}$$

$$\{C, E\} \rightarrow \{B\}$$

Q:13. Predict a class label for an unknown example using Naive Bayesian classification on the following training dataset from all electronic customer base

$x' = \{\text{age} = " \leq 30 ", \text{income} = " \text{median} ", \text{student} = " \text{yes} ", \text{credit\_rating} = \text{fair}\}$

Age	Income	student	credit-rating	buy-computer
$\leq 30$	H	N	Fair	No
$\leq 30$	H	N	Excellent	N
$31 \dots 40$	H	N	F	Yes
$31 \dots 40$	M	N	F	Yes
$31 \dots 40$	L	Y	F	Yes
$31 \dots 40$	L	Y	E	Yes
$\leq 30$	M	N	F	N
$\leq 30$	L	Y	F	Yes
$> 40$	M	N	F	Yes
$\leq 30$	M	Y	E	Yes
$31 \dots 40$	M	N	E	Yes
$31 \dots 40$	H	Y	F	Yes
$> 40$	M	N	E	N

$$P(\text{buy-computer} = \text{yes}) = 9/14$$

$$P(\text{buys-computer} = \text{No}) = 5/14$$

age	buy-computer = Yes	buy-computer = No
$\leq 30$	2/9	3/5
31..40	4/9	0/5
$> 40$	3/9	2/5

Income Yes No

High 2/9 2/5

Low 3/9 1/5

Medium 4/9 2/5

Student Yes No

Yes 6/9 1/5

No 3/9 4/5

Credit rating Yes No

Fair 6/9 2/5

Excellent 3/9 3/5

Map Rule

$$P(x')_{\text{yes}} = P(\text{age} = \text{"}\leq 30\text{"}/\text{yes}) *$$

$$P(\text{income} > \text{"medium"}/\text{yes}) *$$

$$P(\text{student} = \text{"yes"}/\text{yes}) *$$

$$P(\text{credit\_rating} = \text{"fair"}/\text{yes}) *$$

$$P(\text{yes})$$

$$= \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9} \times \frac{9}{14}$$

$$= 0.02821$$

$$\begin{aligned}
 P(x') \text{ No} &= P(\text{age} = " \leq 30 / \text{No}) * P(\text{income} = " \text{medium} / \text{No}) * \\
 &\quad P(\text{student} = " \text{yes} / \text{No}) * \\
 &\quad P(\text{credit-rating} = " \text{fair} / \text{No}) * \\
 P(\text{No}) &= \frac{3 \times 2 \times 1 \times 2 \times 5}{5 \times 5 \times 5 \times 5 \times 5} = 0.0028 \\
 P(\text{Yes}) &= 1 - 0.0028 = 0.9972
 \end{aligned}$$

Since  $0.0028 > 0.002$ , therefore the naive Bayes classifier predicts bayes-computer = "yes" for sample X.

Q: 14. Construct a decision tree using classification algorithm and decide on which day you can play tennis.

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	S	H	H	Strong	N
D3	Overcast	H	H	W	Yes
D4	Rain	Mild	H	W	Y.
D5	R	Cool	Normal	W	Y.
D6	R	C	N	S	N
D7	O	C	N	S	Y
D8	S	M	H	W	N
D9	S	C	N	W	Y
D10	R	M	N	W	Y
D11	S	M	N	S	Y
D12	O	M	H	S	Y
D13	O	H	N	W	Y
D14	R	M	H	S	N