

Rainy = y

Yes

Sunny = c

Yes

Overcast

Yes

Rainy

No

Sunny

No

Sunny

Yes

Rainy

No

Overcast

Yes

Overcast

Yes

Step 1 Frequency Table

Weather

Yes No

Rainy

2 3

Sunny

3 2

Overcast

5 0

Total

10 4

$$\text{No of Yes} = 10$$

$$\text{No of Nos} = 4$$

$$\text{Total entries} = 10 + 4 = 14$$

Step 2 Livelihood Table

<u>Weather</u>	<u>Yes</u>	<u>No</u>	<u>prob. do this</u>
Rainy	2	2	$\frac{4}{14} = 0.29$
Sunny	3	2	$\frac{5}{14} = 0.36$
Overset	5	0	$\frac{5}{14} = 0.36$
Total	$\frac{10}{14}$	$\frac{4}{14}$	0.29
:	$\frac{10}{14}$	$\frac{4}{14}$	0.29

Steps Applying Baye's Theorem:

$$P(Y_0 | \text{Sunny}) = \frac{P(\text{Sunny} | Y_0) \times P(Y_0)}{P(\text{Sunny})}$$

$$= \frac{3 \times 0.75}{0.36} = 0.59$$

$$P(\text{no} | \text{Sunny}) = \frac{P(\text{Sunny} | \text{No}) \times P(\text{No})}{P(\text{Sunny})}$$

$$= \frac{2}{4} \times \frac{0.29}{0.36}$$

$$P(Y_0 | \text{Sunny}) > P(\text{no} | \text{Sunny})$$

Hence for the instance sunny, the class talk

Q outlook: rainy

$$P(\text{yes/rainy}) = P(\text{rainy}) \times \frac{P(\text{yes})}{P(\text{rainy})}$$

$$\therefore \frac{2}{10} \times \frac{0.71}{0.29}$$

$$\therefore 0.489 : 0.49$$

$$P(\text{no/rainy}) = \frac{P(\text{rainy})}{\text{no}} \times \frac{P(\text{no})}{P(\text{rainy})}$$

$$\therefore \frac{2}{4} \times \frac{0.29}{0.29}$$

$$\therefore \frac{1}{2} : 0.5$$

$$P(\text{no/rainy}) > P(\text{yes/rainy})$$

Hence for instance rainy the class label is No

$$\begin{matrix} \text{P yes/rainy} \\ 0.49 \end{matrix} < \begin{matrix} \text{P no/rainy} \\ 0.5 \end{matrix}$$

$$(y_{actual}, v_{act}) < (y_{predicted}, v_{pred})$$

Decision Tree

- a Consider the dataset. Identifying the root using decision tree (induction).

ID	Age	Income	Student	Credit Rating	Buy.computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle-aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle-aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle-aged	medium	no	excellent	yes
13	middle-aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Total no of Transaction - 14

No of attributes = 4 - excluding class / idle

Attributes - Age, Income, Student, Credit rating

Op Variable = Buys.computer

class labels = yes/no.

Total No of yes = 9

Total No of no = 5

Step 1

Entropy (Buys Computer)

$$= - \sum_{i=1}^2 p_i \log_2(p_i)$$

Converting $\log_2 + \log = - \left[\frac{9}{14} \log_2 \left(\frac{9}{14} \right) + \frac{5}{14} \log_2 \left(\frac{5}{14} \right) \right]$

$$\begin{aligned} &= - \left[0.64 \times \log_2(0.64) + 0.357 \log_2(0.357) \right] \\ &= -[-0.942168571] \\ &= 0.942 \end{aligned}$$

Consider the dataset
Identify the root using Decision Tree gain

Step 2

Information (Age)

youth

Age middle-aged

senior

Decision Tree

Price values

Youth Yes - ?
No - 3

- 2 yes and 8 no

$$\text{Entropy (Youth)} = - \sqrt{\sum_{i=1}^3 p_i \log_2 p_i}$$

$$\begin{aligned}
 & \frac{y_0}{y_0 + n_0} = \frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \\
 & \frac{n_0}{y_0 + n_0} = -(-0.97095) \\
 & = 0.97095
 \end{aligned}$$

$$\begin{aligned}
 \text{Middle Aged} & \quad \frac{y_1}{n_1} = 4 \\
 & \quad \frac{n_1}{n_1} = 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Entropy (Middle-Aged)} & = - \left[\frac{4}{4} \log_2 \left(\frac{4}{4} \right) + \frac{0}{4} \log_2 \left(\frac{0}{4} \right) \right] \\
 & = \log_2 1 \\
 & = 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Senior} & \quad \frac{y_2}{n_2} = 3 \\
 & \quad \frac{n_2}{n_2} = 2
 \end{aligned}$$

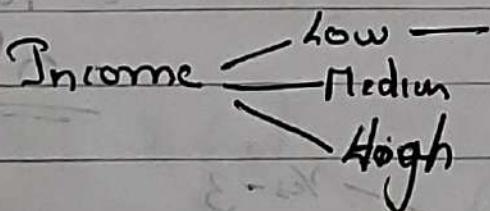
$$\begin{aligned}
 \text{Entropy (Senior)} & = - \left[\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] \\
 & = 0.97095
 \end{aligned}$$

Information gain (age)

Total entropy (age) = weighted average of entropy

$$\left[\frac{\text{young total} \times \text{entropy (young)}}{\text{total}} \right] + \frac{5}{14} \times 0.970 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.970 \\ = 0.69285 \\ - 0.693$$

$$\text{Info. gain (age)} = \text{Overall Entropy} - \text{Total entropy (age)} \\ = 0.94 - 0.693 \\ = 0.247$$

Info gain (Income)

how

```

graph LR
    how --> Yes
    how --> No
  
```

High

```

graph LR
    High --> Yes1
    High --> No1
  
```

Medium

```

graph LR
    Medium --> Yes2
    Medium --> No2
  
```

$$\text{Entropy (High)} = - \left[\frac{2}{4} \log_2 \left(\frac{1}{2} \right) + \frac{2}{4} \log_2 \left(\frac{2}{2} \right) \right]$$

$$= -(-1)$$

$$= 1$$

$$\text{Entropy (Medium)} = - \left[\frac{4}{6} \log_2 \frac{1}{3} + \frac{2}{6} \log_2 \left(\frac{2}{3} \right) \right]$$

$$= 0.916$$

$$\text{Entropy (Low)} = - \left[\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right]$$

$$= 0.811$$

Total entropy (incom) = weighted average of entropy

$$= \frac{1}{4} \times 0.811 + \frac{6}{14} \times 0.916 + \frac{4}{14} \times 1$$

Info gain (incom) = Overall entropy - Total entropy

$$= 0.94 - 0.91$$

$$= 0.029 \approx 0.03$$

0.86
0.14

classmate
Date 0.57/
Page

Info gain (Student)

Student \leftarrow Yes
No

Yes \leftarrow C
No \leftarrow I
No \leftarrow 3
No \leftarrow 9

: 3.906

$$\text{Entropy (Yes)} = -\left(\frac{3}{7} \log_2 \left(\frac{3}{7}\right) + \frac{4}{7} \log_2 \left(\frac{4}{7}\right)\right)$$

$$\text{Entropy (No)} = -\left(\frac{3}{7} \log_2 \left(\frac{3}{7}\right) + \frac{4}{7} \log_2 \left(\frac{4}{7}\right)\right)$$

: 0.985

$$\begin{aligned} \text{Total entropy (Student)} &= \text{Weighted average of entropy} \\ &= \frac{7}{14} \times 0.59 + \frac{7}{14} \times 0.985 \\ &= 0.7875 \end{aligned}$$

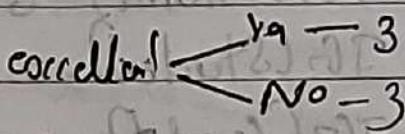
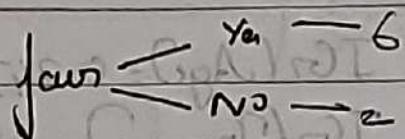
$$\begin{aligned} \text{Info gain (student)} &= \text{Overall Entropy} - \text{Total entropy (Student)} \\ &= 0.94 - 0.7875 \end{aligned}$$

$$\therefore \underline{\underline{0.1525}}$$

Info gain (credit-rating)

0.75

0.25



$$\text{Entropy (fan)} = - \left(\frac{6}{8} \log_2 \frac{6}{8} + \frac{2}{8} \log_2 \frac{2}{8} \right)$$

$$= 0.811$$

$$\text{Entropy (correlation)} = - \left(\frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6} \right)$$

$$= 1$$

Total entropy = Weighted average of entropy.

$$\frac{8}{14} \times 0.811 + \frac{6}{14} \times 1$$

$$= 0.892$$

Info gain (credit-rating) = Overall - total Credit

$$= 0.94 - 0.892$$

$$= 0.048$$

The info gain with the maximum value is the root

$$\therefore \text{root} = \text{age} \quad IG(\text{Age}) = 0.347$$

$$IG(\text{Income}) = 0.029$$

$$IG(\text{Student}) = 0.152$$

$$IG(\text{Credit}) = 0.048$$

~~Grade~~ $IG(\text{Age})$ is the highest info gain.

a) Consider the dataset

Income	Student	Credit rating	Class
--------	---------	---------------	-------

medium	no	fair	yes
low	yes	fair	yes
low	yes	excellent	no
medium	yes	fair	yes
medium	no	excellent	no

Identify the root using the decision tree

Total transactions:

No. of attributes \rightarrow (Income, student, credit rating)
Op. variables: class

class labels: yes, no

Total yes, Total No: 2

$$\text{Entropy of class} = \sum_{i=1}^n P_i \log_2(P_i)$$

$$\left[\frac{3}{5} \log_2\left(\frac{3}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right) \right] = 0.971$$

Info gain (Income)

Income → medium → $y_1 - 3$
 Income → low → $y_2 - 1$
 Income → No → $y_3 - 1$

$$\text{Entropy (medium)} = - \left[\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right]$$

$$= 0.918$$

$$\text{Entropy (low)} = - \left[\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right]$$

$$= 0.918$$

$$\text{Total entropy (Income)} = \text{Total entropy} - \text{entropy of subpart}$$

$$= 0.971 - 0.918 = 0.052$$

Info gain (Student)

Student → Yes → $y_1 - 2$
 Student → No → $y_2 - 1$
 Student → No → $y_3 - 1$

$$\text{entropy (Yes)} = - \left[\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right]$$

$$= 0.918$$

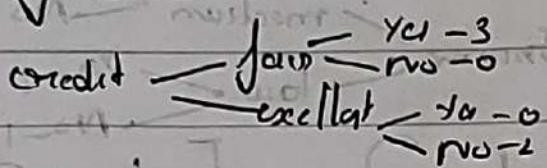
$$\text{Entropy (No)} = - \left[\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right]$$

$$= 1$$

$$\text{Total entropy (student)} = \frac{3}{5} \times 0.98 + \frac{2}{5} \times 1 = 0.951$$

$$\text{Info-gain (student)} = 0.971 - 0.951 \\ = 0.02$$

Info gain of credit ranking



$$\text{Entropy (fair)} = \left[\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] = 0$$

$$\text{Entropy (excellent)} = - \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] = 0$$

$$\text{Total entropy} = \frac{3}{5} \times 0 + \frac{2}{5} \times 0 = 0$$

$$\text{Info gain (Credit-ranking)} = 0.971 - 0 = 0.971$$

$$IG(\text{Income}) = 0.02$$

$$IG(\text{Student}) = 0.02$$

$$IG(\text{Credit-ranking}) = 0.971$$

IG with maximum = $IG(\text{Credit-ranking})$
 Root : credit-ranking

$$[(d')_{\text{good}} + (d')_{\text{bad}}] = 0.971$$

$$[(d')_{\text{good}} + (d')_{\text{bad}}] = 0.971$$

KNN Algorithm (K-Nearst Neighbour Algorithm)

- It is a classification algorithm which comes under Supervised Learning.
 - It is a lazy learner algorithm. It does not learn the training dataset.
 - It works based on similarity of data. It assumes similarity between the new data and available data and puts
- KNN works best when the dataset is small.
When dataset is big \rightarrow Decision Tree

KNN does not learn a discriminative function from the training set. \therefore it is a lazy learner.

KNN can identify the category to which an input belongs to.

The most preferred value for K is 5.

a. Consider the given dataset:

Maths	CS	Result
-------	----	--------

4	3	F
---	---	---

6	7	P
---	---	---

7	8	P
---	---	---

5	5	F
---	---	---

8	8	P
---	---	---

Find the result for the query $\text{maths}=6$ $\text{CS}=8$

and value of $k=3$. using KNN algorithm.

distance - Euclidean Distance

Maths	CS	Result	d	Neighbours	Rank
$d_1 = 4$	3	F	5.38	5	
$d_2 = 6$	7	P	1	1	✓
$d_3 = 7$	8	P	1	2	✓
$d_4 = 5$	5	F	3.16	4	
$d_5 = 8$	8	P	2	3	✓

$$d_1 = \sqrt{(x_1 - x_{12})^2 + (x_2 - x_{11})^2}$$

$$= \sqrt{(6-4)^2 + (8-3)^2}$$

$$= \sqrt{5^2}$$

$$= 5.38$$

→ Smallest value are ranked first

$$d_2 = \sqrt{(6-7)^2 + (8-7)^2}$$

$$d_3 = \sqrt{(6-7)^2 + (8-8)^2}$$

$$d_4 = \sqrt{(5-5)^2 + (8-8)^2}$$

$$= \sqrt{1^2 + 3^2}$$

$$= \sqrt{10}$$

$$= 3.16$$

$$d_5 = \sqrt{(6-8)^2 + (8-8)^2}$$

$$= \sqrt{2^2 + 0^2}$$

~~calculated distance b/w 6 & 8 b/w 6 & 7
2nd row - calculating~~

$\Rightarrow k=3$ means we have to select the first 3 ranks

1st, 2nd, 3rd rank ~~2nd, 3rd rank~~ ~~ranked~~

Analyze their results $1-P$ $2-P$ $3-P$

$\therefore 6$ and 8 would yield P as output

The result of the course Math: 6 and $7, 8$ will be P
 \because the first 3 neighbours rank will give the result as P . (majority)

$$(6-7)^2 + (8-7)^2$$

$$(7-7)^2 + 0^2$$

$$= 0^2 + 0^2$$

a) Consider the dataset

Brightness Saturation Class

40	20	Red
50	50	Blue
60	90	Blue
10	25	Red
70	70	Blue
60	10	Red
25	80	Blue

Find the class of the given data Brightness: 20
saturation: 35 k=5

Brightness Saturation Class d Neighbour rank

$d_1 - 40$	20	Red	25	2
$d_2 - 50$	50	Blue	33.54	3
$d_3 - 60$	90	Blue	68.007	7
$d_4 - 10$	25	Red	94.14	1
$d_5 - 70$	70	Blue	61.03	6
$d_6 - 60$	10	Red	47.16	5
$d_7 - 25$	80	Blue	45.27	4

$$d_1 = \sqrt{(60-40)^2 + (35-20)^2}$$

$$= \sqrt{20^2 + (-15)^2}$$

$$= \sqrt{400 + 225}$$

$$= \sqrt{625} = 25$$

$$d_1 = \sqrt{(\frac{20-50}{50})^2 + (\frac{35-50}{50})^2}$$

$$= \sqrt{900 + 225} = \sqrt{1125}$$

$$= 33.54$$

$$d_2 = \sqrt{(\frac{20-60}{60})^2 + (\frac{35-90}{60})^2}$$

$$= \sqrt{1600 + 3025} = \sqrt{4625}$$

$$= 68.007$$

$$d_3 = \sqrt{(\frac{14-19}{19})^2 + (\frac{14-35}{19})^2}$$

$$= \sqrt{25 + 121} = \sqrt{146}$$

$$= 12.03$$

$$d_4 = 47.16$$

$$d_5 = 45.27$$

$k = 5$

1 - Red

2 - Red

3 - Blue

4 - Blue

5 - Red

The color is Red for the query brightness
saturation = 35

Note: All Classification comes under supervised learning,
All Clustering comes under unsupervised learning.

Clustering

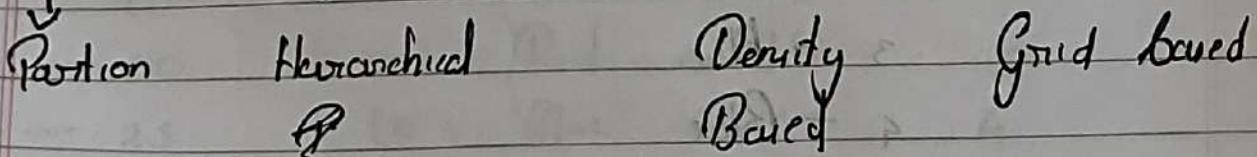
way of

- Grouping the data points into different clusters, consisting of similar data points.
- Unlabelled datasets are grouped here
- Unsupervised

Applications

- Market Segmentation
- Recommendation System
- Movie Recommendation System

Clustering



The only method where we can specify the no of partitions is the partition method.

Partition Clustering

- It is used to make partitions on the data in order to form clusters
- If n partitions are done on p objects of the collection then each partition is represented by a cluster and resp.

Two conditions to be satisfied

- ↪ One object should only belong to 1 group
- ↪ There should be no group without a purpose

Method → Hierarchical Relocation

↪ We can move data from one cluster to another

→ Centroid-based method

↪ K-Means Clustering Algorithm

Hierarchical Clustering

→ Clustering need not be specified beforehand.

Data → Cluster → Tree-like structure (cladogram)

Approach

→ Agglomerative Approach - (Bottom-up approach)

→ Divisive Approach - (Top-to-bottom)

Divisive

→ Top Down approach

Density Based Clustering :-

This method mainly focuses on density

Grid Based Clustering :-

Divided data in this technique into

Finite no of cells (grid) to form clusters

→ Faster processing time.

K-Means Algorithm

- a) Consider the data points $A_1(2, 12)$
 $A_2(2, 5)$ $A_3(8, 4)$ $B_1(5, 8)$ $B_2(7, 3)$ $B_3(6, 9)$
 $C_1(1, 2)$ $C_4(4, 9)$. The initial clusters, ^(centroids) are $C_1(2, 10)$,
 $C_2(5, 8)$ $C_3(1, 2)$. Find the clusters using K
mean algorithm.

Euclidean Distance

Data Points

		C1	C2	C3	Cluster	
A ₁	2	10	5	12	C ₁	
A ₂	2	5	10	12	C ₃	
A ₃	8	4	9	5	C ₂	
B ₁	5	8	3.61	0	C ₂	
B ₂	7	5	7.07	2.236	C ₃	
B ₃	6	4	7.21	4.123	C ₂	
C ₁	1	2	7.06	7.211	0	C ₃
C ₂	4	9	2.24	1.414	7.615	C ₂

New Centroids

$$C_1 = (2, 10)$$

$$C_2 = \left(\frac{8+5+7+6+4}{5}, \frac{1+8+5+4+9}{5} \right)$$

$$= (6, 6)$$

$$C_3 = \left(\frac{2+1}{2}, \frac{5+2}{2} \right)$$

$$\left(\frac{3}{2}, \frac{7}{2} \right)$$

$$(C_2, (1.5, 3.5))$$

$$(C_3, (2, 8, 21))$$

		c_1	c_2	c_3	cluster
A ₁	2	10	6.96	1.513.5	New cluster
A ₂	2	5	4.12	1.58	c ₃
A ₃	8	4	8.83	6.52	c ₂
B ₁	5	8	3.61	2.24	c ₂
B ₂	5	5	7.07	1.41	c ₂
B ₃	4	1	7.21	2	c ₂
C ₁	1	2	8.06	6.40	c ₃
C ₂	4	9	2.24	3.61	c ₁

c_1 → not matching
so next iteration

New Centroids:-

$$c_1 = \frac{(2+4)}{2}, \frac{(10+9)}{2}$$

$$(3, 9.5)$$

$$c_2 = \frac{(8+5+7+6)}{4}, \frac{(4+8+5+4)}{4}$$

$$(6.5, 5.25)$$

$$c_3 = \frac{(1+2)}{2}, \frac{(5+2)}{2}$$

$$(1.5, 3.5)$$

			c_1	c_2	c_3	cluster	<u>new cluster</u>
A ₁	2	10	1.2	6.51	6.52	c ₁	c ₁
A ₂	2	5	4.61	4.5	1.58	c ₃	c ₃
A ₃	8	4	7.43	1.95	6.52	c ₂	c ₂
B ₁	5	8	2.50	3.13	5.70	c ₂	c ₁
B ₂	7	5	6.02	0.56	5.7	c ₂	c ₂
B ₃	6	4	6.26	1.85	4.53	c ₂	c ₂
C ₁	1	2	7.76	6.89	1.58	c ₃	c ₃
C ₂	4	9	7.12	4.51	6.03	c ₁	c ₂

New centroids

$$c_1 = \left(\frac{2+5}{2}, \frac{10+5}{2} \right) = (3.5, 7.5)$$

$$c_2 = \left(\frac{8+7+6+4}{4}, \frac{4+5+4+9}{4} \right) = (6.25, 5.5) \quad (7, 4.87)$$

$$c_3 = (1.5, 3.5)$$

A₁
A₂
A₃
B₁
B₂
B₃
C₁
C₂

		$\frac{c_1}{35.99}$	$\frac{c_2}{28.53}$	$\frac{c_3}{15.35}$	clusters	new cluster
A ₁	2 10	1.94	7.66	6.52	c ₁	c ₁
A ₂	2 5	4.33	5.06	1.58	C ₃	C ₃
A ₃	8 4	6.62	1.05	6.52	c ₂	C ₂
B ₁	5 8	1.67	4.18	5.70	c ₁	C _{1,2}
B ₂	7 5	5.21	0.67	5.7	c ₂	C ₂
B ₃	6 4	5.52	1.05	4.53	c ₂	C ₂
C ₁	1 2	7.49	6.44	1.58	c ₃	C _{3,1}
C ₂	4 9	0.33	0.55	6.04	q	n

New centroid

Centroids
 $(A_1, B_1), (A_3, B_2, B_3, C_2)$
 (A_2, C_1)

Consider the data points $A_1(2, 7), A_2(3, 3)$
 $A_3(6, 8), A_4(8, 8), A_5(7, 5), A_6(9, 7)$

Initial centroids - c₁(2, 0) c₂(8, 8)

Find out) The clusters using K-Means Algorithm

K-Medoids Clustering:

- It is also an unsupervised learning algorithm.
- a Consider the dataset $x_1(2,6)$, $x_2(3,4)$, $x_3(3,8)$,
 $x_4(1,2)$, $x_5(6,2)$, $x_6(6,4)$, $x_7(7,3)$, $x_8(7,6)$,
 $x_9(8,5)$, $x_{10}(7,6)$. The medoids are $c_1 = (3,4)$, $c_2 = (7,4)$. Find out the cluster using K-medoids algorithm.

			manhattan distance		
x_1	2	6	c_1	c_2	c_1
x_2	3	4	0	4	c_1
x_3	3	8	(3,8)	4	c_1
x_4	4	2	3	5	c_1
x_5	6	2	5	3	c_2
x_6	5	4	3	1	c_2
x_7	7	3	5	1	c_2
x_8	7	4	4	0	c_2
x_9	8	5	6	2	c_2
x_{10}	7	6	6	2	c_2

New medoids are

$$\{ (2,4), (3,7), (3,8), (4,2) \}$$

$$\{ (6,7), (6,9), (7,3), (7,4), (8,5), (7,6) \}$$

$$(Cost_1 + Cost(c_1) + Cost(c_2))$$

min of value

$$\frac{[3+0+4+3]}{4} = \frac{10}{4} = 2.5$$

$$\frac{[3+1+1+0+2+2]}{6} = \frac{8}{6} = 1.33$$

HW

	Data points	c_1	$\frac{c_2}{8,8}$	Centroid
A ₁	2, 3	2, 3	0	c ₁
A ₂	3, 3	1	7.07	c ₁
A ₃	6, 8	6.4	0	c ₂
A ₄	8, 8	7.81	0	c ₂
A ₅	7, 5	5.38	3.16	c ₂
A ₆	9, 7	8.06	1.41	c ₂

Initial centroid (2, 3) (8, 8)

New centroids

$$c_1 = (2.5, 3)$$

$$c_2 = (2.5, 8)$$

	Data points	c_1^{old}	c_2^{old}	old	new
A ₁	2, 3	0.5	6.80	c ₁	c ₁
A ₂	3, 3	0.5	6.02	c ₁	c ₁
A ₃	6, 8	6.10	1.12	c ₂	c ₂
A ₄	8, 8	7.43	1.12	c ₂	c ₂
A ₅	7, 5	6.92	2.06	c ₂	c ₂
A ₆	9, 7	8.63	1.5	c ₂	c ₂

(Centroid = (A₁, A₂) , (A₃, A₄, A₅, A₆))

Remove the first medoid

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	c_1	c_2
x_1	2	6	7	8							$c_1 =$	
x_2	3	4	4	5							$c_1 =$	
x_3	3	8	8	9							$c_1 =$	
x_4	4	2	5	4							$c_2 = 12$	
x_5	6	2	3	2							$c_2 = 11$	
x_6	6	4	1	2							$c_1 =$	
x_7	7	3	1	0							$c_2 =$	
x_8	7	4	0	1							$c_1 =$	
x_9	8	5	2	3							$c_1 =$	
x_{10}	7	6	2	3							$c_1 =$	

$$\text{Cost} = \text{Cost}(c) + \text{Cost}(h)$$

$$= [7+4+8+1+0+2] + [4+2+0]$$

30

Previous Cost = 19

Current Cost = 30

Current Cost > Previous Cost

∴ The final medoids are $(3, 4), (7, 4)$

④ K-means

Exam → Explain steps in K-means algo - Part A
Problem - Part B.

Write the steps we have done in class rather than the notes

K-medoid

Improved Version of K-means algorithm.

K-medoid clustering algorithm

Exam → Difference and similarity b/w K-means and K-medoid

Density Based Clustering

- Unsupervised learning
- Clustering is based on density of data.
- The commonly used density based algo - DBSCAN
- Nesting of data is possible in density based clustering.
DBSCAN - Density Based Spatial Clustering of Applications with Noise

~~DEF~~ - Parameter of DBScan

↳ Epsilon ϵ

↳ Minimum points

Epsilon:-

Epsilon is defined as the radius of each data point around which the density is considered.

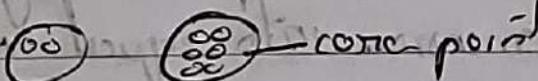
min points:

Minimum no. of data w/in the cluster

Core point

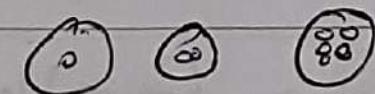
↳ If a cluster has

↳ More than minimum point within epsilon



Border points or Boundary points

↳ Fewer than minimum point,



Boundary point

None point

↳ A point which is not a core point
or boundary point

Steps in DBSCAN

* Diff b/w K-means and DBSCANs

Clique Algo

This algo uses density and grid based technique, i.e. by taking density threshold and no of grids as i/p parameter.

To handle datasets with large no of dimension
Use the Apron Property effectively

First divide data space into grids by dividing each dimension into equal intervals called bins

5th Model

Essay

{ Clustering Definition?
What is Clustering
Clustering Algorithms

a) Consider the dataset and find out the cluster using agglomerative method and draw the dendrogram.

Sample x x

P ₁	1.5	1.5	1.5	1.5	1.5	1.5	1.5
P ₂	1.5	1.5					1.5
P ₃	5	5					5
P ₄	3	4					3
P ₅	4	4					4
P ₆	3	3.5	1.5	1.5	1.5	1.5	3

I Compute the distance matrix — Euclidean distance

$$d(p_1, p_2) = \sqrt{0.25 + 0.25} = 0.707$$

$$d(p_1, p_3) = \sqrt{25} = 5.05$$

$$d(p_1, p_4) = \sqrt{16} = 4.00$$

$$d(p_1, p_5) = \sqrt{16} = 4.00$$

$$d(p_1, p_6) = \sqrt{9} = 3.00$$

$$d(p_2, p_3) = \sqrt{25} = 5.00$$

$$d(p_2, p_4) = \sqrt{16} = 4.00$$

$$d(p_2, p_5) = \sqrt{16} = 4.00$$

$$d(p_2, p_6) = \sqrt{9} = 3.00$$

$$d(p_3, p_4) = \sqrt{0.25} = 0.25$$

$$d(p_3, p_5) = \sqrt{16} = 4.00$$

$$d(p_3, p_6) = \sqrt{9} = 3.00$$

(p_4, p_5) (p_4, p_6) (p_5, p_6)

0.5

1.12

 $\times \times \text{ ignore?}$

	p_1	p_2	p_3	p_4	p_5	p_6	
p_1	0	0		2.1	2.1		9
p_2	0.707	0		2	2		8
p_3	5.656	4.94	0	1	8		9
p_4	3.605	2.92	2.24	0	1		9
p_5	4.242	3.53	1.41	1.08	0		8
p_6	3.201	2.5	2.5	0.5	1.12	0	

Minimum value = 0.5

Step II Merging the two closest members in the distance matrix

	p_1	p_2	p_3	p_4, p_6	p_5	
p_1	0					
p_2	0.707	0				
p_3	5.656	4.94	0			
p_4, p_6	3.201	(2.5)	2.24	0		
p_5	4.242	3.53	1.41	1.12	0	

$$d[(p_4, p_6), p_1] = d[(p_4, p_1), (p_6, p_1)]$$

224000

classmate

Date _____

Page _____

 $m_w(3.606, 3.201)$
 3.201

$$d(p_4, p_5), p_2 = d[(p_4, p_2), (p_6, p_2)]$$

- 2.92, 2.5

$$d(p_4, p_6), p_3 = d(p_4, p_3), (d(p_6, p_3))$$

- min(2.24, 2.5)
- 2.24

$$d(p_4, p_6)_{ps} : d(p_4, p_5), d(p_6, p_5)$$

1.61, 1.12
~~1.21~~

Minimum value = 0.707

a ~~Affirmo~~

Ensemble Classifier

Combining multiple models to improve the overall performance

Type

- ↳ Bagging → Bootstrap Aggregating → Voting
- ↳ Boosting

Agglomerative and draw dm

$$d(p_1, p_2) = 1$$

$$d(p_1, p_3) = 3$$

$$d(p_1, p_4) = 5$$

$$d(p_2, p_3) = 2$$

$$d(p_2, p_4) = 6$$

$$d(p_3, p_4) = 3$$

	p_1	p_2	p_3	p_4
p_1	0			
p_2	1	0		
p_3	2	1	0	
p_4	5	6	3	0

Smallest - 1 belongs to p_2, p_1

	p_1, p_2	p_3	p_4
p_1, p_2	0		
p_3	2	0	
p_4	5	3	0

mist

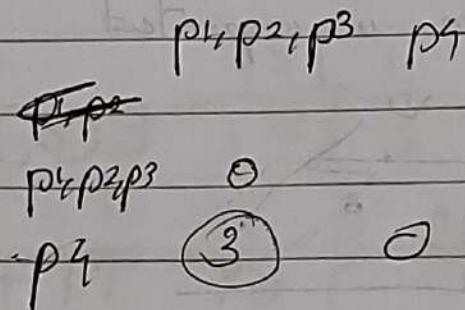
$$d(p_3, (p_1, p_2)) \therefore d[(p_3, p_1), d(p_3, p_2)]$$

$$= \min(4, 2) = 2$$

$$d[(p_4, (p_1, p_2))] = d[(p_4, p_1)], d(p_4, p_2)$$

= 5, 6
min. 5

minimum in matrix = 2



$$d(p_4, (p_1, p_2, p_3)) \cdot d[(p_4, p_1), (p_4, p_2), (p_4, p_3)]$$

= 5, 6, 3
= 3

min. p4

The Clustering core: $\{(p_1, p_2), p_3, p_4\}$

