

28/02/23

UNIT-IV

* Types of Unsupervised learning:-

<u>Clustering</u>	<u>Association</u>
1) method of grouping the objects into clusters based on Similarities.	1) It involve discovering the patterns in data finding. Co occurrence.
2) It discover inherent grouping in data	2) used in finding the relationship b/w variables in the large database.
3) K-means clustering, hierarchical clustering, density based clustering, EM clustering [Exposition, maximization clustering]	3) Market basket Analysis, eg:- macline + Cores + Tempuglas milk + bread + egg milk + bread + cheese milk + curd + Bread .

* Challenges in Unsupervised Learning :-

- Unsupervised learning is computationally complex.
- Cannot get precise information regarding data routing
- Human intervention to validate the output variables.
- High risk of inaccurate results & lack of transparency.

* K-Means Clustering :-

- 1) It is an unsupervised learning algorithm.
- 2) given a data set of items with certain features and values for this features the algorithms will categorise the items into K group of clusters of similarities.
- 3) To calculate the similarity we can use the Euclidean, Manhattan, Hamming distances for measurement.

4) K-Means :- K indicates no. of clusters & D indicates list of data points.

5) Choose K no. of random data points as initial centroids.

6) Repeat till clusters centers stabilize.

7) Allocate each point in D to the nearest of k-th centroid.

8) Compute centroid for the cluster using all points in the clusters.

* Advantages :-

→ It is simple, easy to understand.

→ No other clustering algs performs better than K-means.

→ It is also efficient in which time taken to cluster K-means varies linear with data point.

* Disadvantages :- no.of

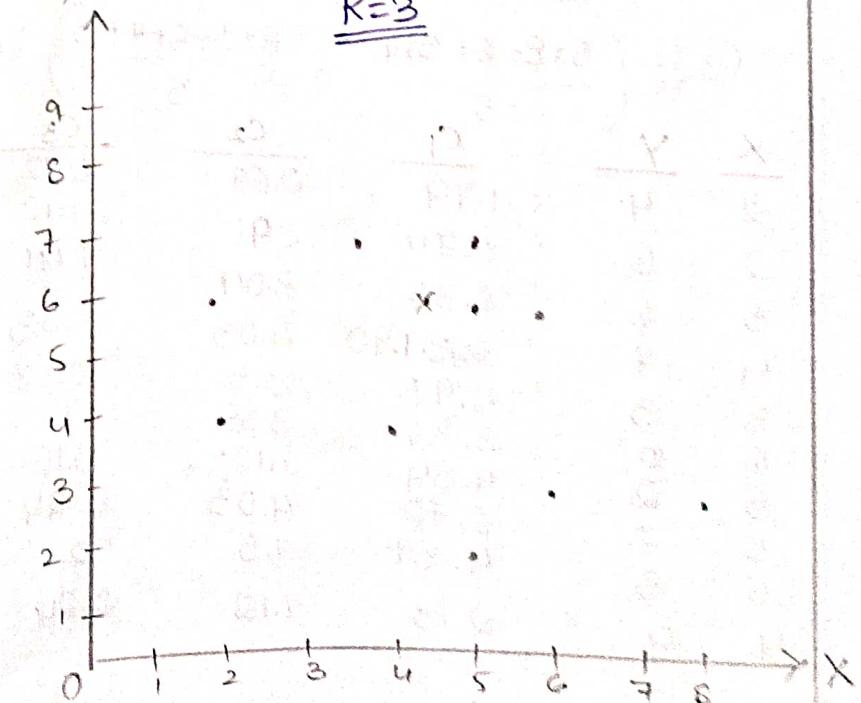
→ The user needs to specify an initial value of K.

→ The process of finding the clusters may not converge.

→ It is not suitable for discovering clusters that are ^{not} hyper ellipsoids (or) hyper spheres.

Example :-

X	Y
2	4
2	6
5	6
4	7
8	3
6	6
5	2
5	7
6	3
4	4



X	Y	(1,5)	(4,1)	(8,4)	cluster number
2	4	1.41	3.61	6	C1
2	6	1.41	5.39	6.0	C1
5	6	4.12	5.10	3.61	C3
4	7	3.61	6	5	C1
8	3	7.28	4.47	1	C3
6	6	5.10	5.39	2.83	C3
5	2	5	1.41	3.61	C2
5	7	4.47	6.08	4.24	C3
6	3	5.39	2.83	2.24	C3
4	4	3.16	3	4.	C2

$$C_1 = (1,5) \quad C_2 = (4,1) \quad C_3 = (8,4).$$

$$C_i = \left[\frac{13}{3}, \frac{10}{3} \right]$$

$$C_1 = \left(\frac{2+9+4}{3}, \frac{4+6+7}{3} \right) = \left(\frac{8}{3}, \frac{17}{3} \right) = (2.66, 5.66)$$

$$C_2 = \left(\frac{9+4}{2}, \frac{2+4}{2} \right) = (4.5, 3)$$

$$C_3 = \left(\frac{5+8+6+5+6}{5}, \frac{6+3+5+7+3}{5} \right) = (6, 5)$$

X	Y	$\frac{C_1}{5}$	$\frac{C_2}{5}$	$\frac{C_3}{5}$	cluster no.
2	4	1.79	2.69	4.12	C1
2	6	0.74	3.91	4.12	C1
5	6	2.36	3.04	1.41	C3
4	7	2.19	4.03	2.83	C1
8	3	5.97	3.5	2.83	C3
6	6	3.36	3.35	1	C3
5	2	4.34	1.12	3.16	C2
5	7	2.70	4.03	2.24	C3
6	3	4.37	1.5	2	C2
4	4	2.13	1.12	2.24	C2

ThinkPad

$$C_1 = \left(\frac{0+0+4}{3}, \frac{0+5+3}{3} \right) = (0.66, 5.66)$$

$$C_2 = \left(\frac{5+6+4}{3}, \frac{0+3+4}{3} \right) = (5.33, 6.33)$$

$$C_3 = \left(\frac{5+6+6+5}{4}, \frac{6+3+6+7}{4} \right) = (5, 8.5)$$

solution - 3

X	Y	C_1 (0.66, 5.66)	C_2 (5.33)	C_3 (6.66)	cluster no
0	4	1.79	3.75	4.06	C_1
0	6	0.74	4.80	3.35	C_1
5	6	0.36	3.09	0.44	C_3
4	7	1.90	4.37	1.49	C_3
8	3	5.97	3.01	0.75	C_2
6	6	3.36	1.85	4.34	C_2
5	2	4.34	4.04	0.75	C_3
5	7	0.70	0.95	3.40	C_2
6	3	4.27	1	0.70	C_2
4	4	0.13	1.41	0.50	C_2

$$C_1 = \left(\frac{0+0+4}{3}, \frac{0+5+3}{3} \right) = \left(\frac{8}{3}, \frac{12}{3} \right) = (0.66, 5.66) \quad \boxed{C_2 = \left(\frac{5+6+4}{4}, \frac{6+3+6+7}{4} \right) = \left(\frac{23}{4}, \frac{12}{4} \right) = (5.33, 6.33)}$$

$$C_3 = \left(\frac{5+6+5}{3}, \frac{6+6+7}{3} \right) = (5.33, 6.33) \quad \boxed{= \left(\frac{03}{4}, \frac{12}{4} \right) = (5, 8.5)}$$

solution - 4

X	Y	C_1 (0.66, 5.66)	C_2 (5.33, 3)	C_3 (5.33, 6.33)	cluster no
0	4	1.79	3.75	4.06	C_1
0	6	0.74	4.80	3.35	C_1
5	6	0.36	3.09	0.44	C_3
4	7	1.90	4.37	1.49	C_3
8	3	5.97	3.01	0.75	C_2
6	6	3.36	1.85	4.34	C_2
5	2	4.34	4.04	0.75	C_3
5	7	0.70	0.95	3.40	C_2
6	3	4.27	3.02	0.68	C_2
4	4	0.13			

$$C_1 = (2.5)$$

$$C_2 = \left(\frac{03}{4}, 3 \right) = (5.75, 3)$$

$$C_3 = (5, 8.5)$$

$$C_1 = \left(\frac{2+2+4}{3}, \frac{4+6+7}{3} \right) = (2.66, 5.66)$$

$$C_2 = \left(\frac{5+6+4}{3}, \frac{2+3+4}{3} \right) = (5, 3)$$

$$C_3 = \left(\frac{5+8+6+5}{4}, \frac{6+3+6+7}{4} \right) = (6, 5.5)$$

Iteration - 3.

X	Y	C_1	C_2	C_3	<u>Cluster no</u>
2	4	(2.66, 5.66)	1.79	(5, 3)	
2	6		0.74	3.16	
5	6		3.16	4.24	
4	7		4.12	4.03	
8	3		5.97	3	
6	6		3.36	4.24	
5	2		4.34	3.16	
5	7		2.40	1	
6	3		4.27	1	
4	4		2.13	1.41	

$$C_1 = \left(\frac{2+2+4}{3}, \frac{4+6+7}{3} \right) = \left(\frac{8}{3}, \frac{17}{3} \right) = (2.66, 5.66) \quad \left| \begin{array}{l} C_2 = \left(\frac{5+6+4}{3}, \frac{2+3+4}{3} \right) \\ = \left(\frac{23}{4}, \frac{12}{4} \right) = (5.75, 3) \end{array} \right.$$

Iteration - 4.

X	Y	C_1	C_2	C_3	<u>Cluster no</u>
2	4	(2.66, 5.66)	1.79	3.66	
2	6		0.74	4.80	
5	6		2.36	3.09	
4	7		1.90	4.37	
8	3		5.97	2.25	
6	6		3.36	3.01	
5	2		4.34	1.25	
5	7		2.40	4.07	
6	3		4.27	0.95	
4	4		2.13	2.02	

$$C_1 = (2.5)$$

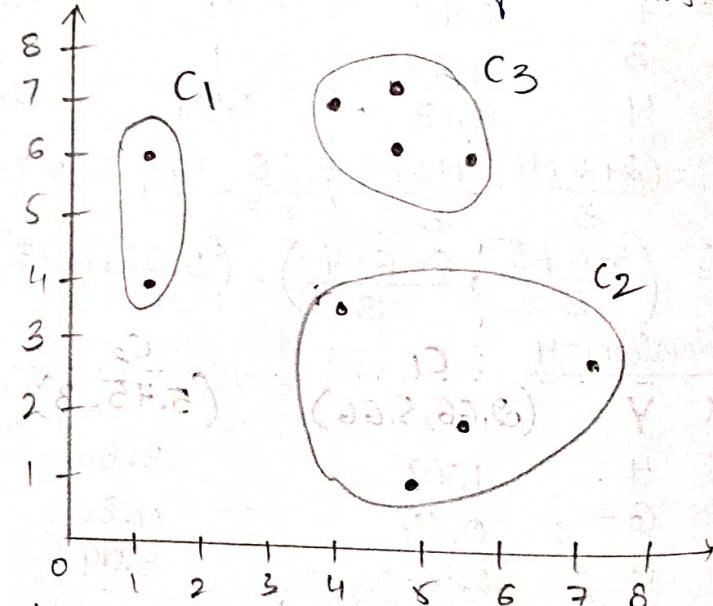
$$C_2 = \left(\frac{23}{4}, 3 \right) = (5.75, 3)$$

$$C_3 = (5, 6.5)$$

Iteration 5		c_1	c_2	c_3	Cluster number
X	Y	(2,5)	(5.75,3)	(5,6.5)	
2	4	1	3.88	3.91	C_1
2	6	1	4.80	3.04	C_1
5	6	3.16	3.09	0.50	C_3
4	7	2.83	3.37	1.12	C_3
8	3	6.32	2.25	4.61	C_2
6	6	4.17	3.01	1.12	C_3
5	2	4.24	1.25	4.50	C_2
5	7	3.61	4.07	0.50	C_2
6	3	4.47	0.25	3.64	C_3
4	4	3.24	2.02	2.69	C_2

Since there are no movements of data points, hence these are the final clusters.

$$C_1 = \left(\frac{2+2}{2}, \frac{4+6}{2} \right) = (2,5)$$



3/3/23 * Agglomerative clustering :-

for the given dataset find the clusters using a single link technique, we'll use the Euclidean distance and draw the Dendrogram.

1) Single link

2) Complete link

3) Average link.

(1) Single link:

Eg:

Sample no

X

Y

P ₁	0.40	0.53
P ₂	0.22	0.38
P ₃	0.35	0.32
P ₄	0.26	0.19
P ₅	0.08	0.41
P ₆	0.45	0.30

Step 1: Compute the distance matrix

	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0	0.23	0.22	0.37	0.34	0.24
P ₂	0.23	0	0.14	0.19	0.14	0.24
P ₃	0.22	0.14	0	0.15	0.28	0.10
P ₄	0.37	0.19	0.15	0	0.26	0.22
P ₅	0.34	0.14	0.28	0.28	0	0.37
P ₆	0.24	0.24	0.10	0.22	0.37	0

$$d(P_1, P_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Step 2:

	P ₁	P ₂	P _{3, P₆}	P ₄	P ₅
P ₁	0				
P ₂	0.23	0			
P _{3, P₆}	0.22	0.14	0		
P ₄	0.37	0.19	0.13	0	
P ₅	0.34	0.14	0.26	0.28	0

merge (P_{3, P₆})

merge (P_{2, P₅})

P_{3, P₆, P₄}

	P ₁	P _{2, P₅}	P _{3, P₆, P₄}
P ₁	0		
P _{2, P₅}	0.23	0	
P _{3, P₆, P₄}	0.22	0.14	0

merge (P_{2, P₅, P_{3, P₆, P₄}})

merge (P_{3, P₆, P₄})

	P ₁	P ₂	P _{3, P₆, P₄}	P ₅
P ₁	0			
P ₂	0.23	0		
P _{3, P₆, P₄}	0.22	0.14	0	
P ₅	0.34	0.14	0.26	0

	P ₁	P _{2, P₅, P_{3, P₆, P₄}}
P ₁	0	
P _{2, P₅, P_{3, P₆, P₄}}	0.23	0

4/3/23

Q1* Complete linkage:

Given a 10 dataset {1, 6, 8, 10} use the agglomerative clustering algorithm with the complete link by using the Euclidean distance to establish the hierarchical grouping relationship.

- By using the cutting threshold value of 5, how many clusters are there?
- We need to calculate the distance matrix using 1 dimensional dataset.

$$2 \text{ point } ED = \sqrt{(x_2 - x_1)^2}$$

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

(i)

1	0	4	7	9	1
5	4	0	3	5	3
8	7	3	0	2	6
10	9	5	2	0	8
2	1	3	6	8	0

1	0	4	7	9	1	①
2	4	0	3	5	3	
3	7	3	0	2	6	
4	9	5	2	0	8	
5	①	3	6	8	0	

merge (1, 5)

compute:

$$d(\{2, 6\}, \{3, 4\}) = \max \{d(2, 1), d(2, 5)\} = \max \{4, 3\} = 4$$

$$d(\{3, 6\}, \{4, 7\}) = \max \{d(3, 1), d(3, 5)\} = \max \{7, 6\} = 7$$

$$d(\{4, 7\}, \{5, 8\}) = \max \{d(4, 1), d(4, 5)\} = \max \{9, 8\} = 9.$$

merge (3, 4)

$$d(\{1, 5\}, \{3, 4\}) = \max \{d(\{1, 5\}, 3), d(\{1, 5\}, 4)\} = \max \{7, 9\} = 9$$

$$d(\{2, \{3, 4\}\}) = \max \{d(2, 3), d(2, 4)\} = \max \{3, 5\} = 5$$

(ii)

1,5	2	3	4	
1,5	0	4	7	9
2	4	0	3	5
3	7	3	0	②
4	9	5	②	0

(iii)

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

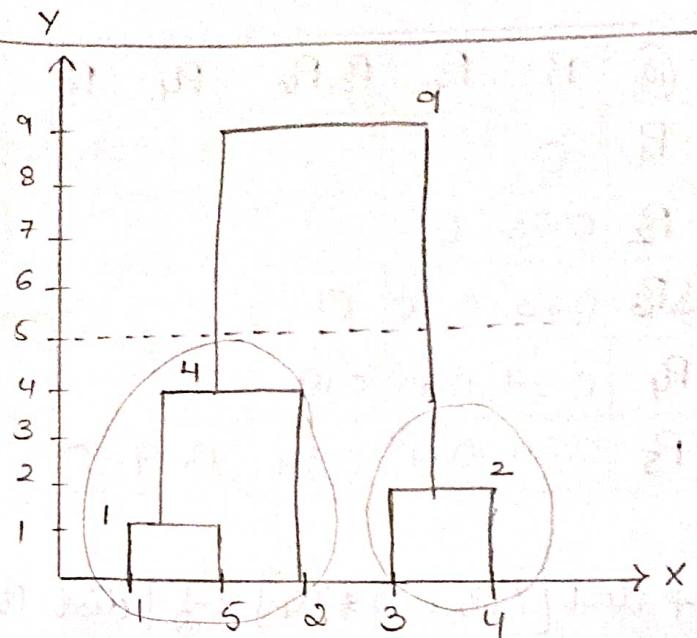
merge ({1, 5}, 2)

$$d((3, 4), (\{1, 5\}, 2)) = \max(d(\{3, 4\}, \{1, 5\}), d(\{3, 4\}, 2))$$

$$= \max(9, 5) = 9$$

(iv)	1,5,2	3,4
1,5,2	0 9	
3,4	9 0	

there are 2 clusters in the given data.



Average Link :-

Find the average link using .

	X	Y
P ₁	0.40	0.53
P ₂	0.22	0.38
P ₃	0.35	0.32
P ₄	0.26	0.19
P ₅	0.08	0.41
P ₆	0.45	0.30

①	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0					
P ₂	0.23	0				
P ₃	0.22	0.15	0			
P ₄	0.37	0.20	0.15	0		
P ₅	0.34	0.14	0.26	0.29	0	
P ₆	0.23	0.25	0.11	0.22	0.39	0

Now merge (P₃, P₆)

- distance $[(P_3, P_6), P_1] = \frac{1}{2} [\text{dist}(P_3, P_1) + (P_6, P_1)]$
 $= \frac{1}{2} [0.22 + 0.23] = 0.225$.
- distance $[(P_3, P_6), P_2] = \frac{1}{2} [\text{dist}(P_3, P_2) + (P_6, P_2)] = \frac{1}{2} [0.15 + 0.25] = 0.20$
- distance $[(P_3, P_6), P_4] = \frac{1}{2} [\text{dist}(P_3, P_4) + (P_6, P_4)] = \frac{1}{2} [0.15 + 0.22] = 0.185$.
- distance $[(P_3, P_6), P_5] = \frac{1}{2} [\text{dist}(P_3, P_5) + (P_6, P_5)] = \frac{1}{2} [0.28 + 0.39] = 0.335$.

②	P ₁	P ₂	P _{3, P₆}	P ₄	P ₅
P ₁	0				
P ₂	0.23	0			
P _{3, P₆}	0.23	0.20	0		
P ₄	0.37	0.20	0.19	0	
P ₅	0.34	0.14	0.34	0.29	0

merge(P₂, P₅)

* dist[(P₂, P₅), P₁] = $\frac{1}{2}[\text{dist}(P_2, P_1) + \text{dist}(P_5, P_1)]$

$$= \frac{1}{2}[0.23 + 0.34] = 0.29$$

* dist[(P₂, P₅), (P₃, P₆)] = $\frac{1}{2}[\text{dist}(P_2, (P_3, P_6)) + \text{dist}(P_5, (P_3, P_6))]$

$$= \frac{1}{2}[0.20 + 0.34] = 0.27$$

* dist[(P₂, P₅), P₄] = $\frac{1}{2}[\text{dist}(P_2, P_4) + \text{dist}(P_5, P_4)] = \frac{1}{2}[0.20 + 0.29]$

$$= 0.25$$

③	P ₁	P _{2, P₅}	P _{3, P₆}	P ₄
P ₁	0			
P _{2, P₅}	0.29	0		
P _{3, P₆}	0.23	0.27	0	
P ₄	0.37	0.25	0.19	0

merge[(P₃, P₆), P₁]

* dist[(P₃, P₆, P₄), P₁] = $\frac{1}{2}[\text{dist}((P_3, P_6), P_1) + \text{dist}(P_4, P_1)]$

$$= \frac{1}{2}[0.23 + 0.37] = 0.30$$

* dist[(P₃, P₆, P₄), (P₂, P₅)] = $\frac{1}{2}[\text{dist}((P_3, P_6), (P_2, P_5)) + \text{dist}(P_4, (P_2, P_5))]$

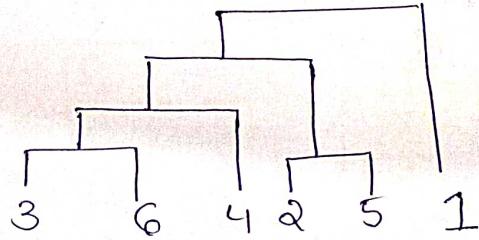
$$= \frac{1}{2}[0.27 + 0.25] = 0.26$$

④	P ₁	P _{2, P₅}	P _{3, P₆, P₄}
P ₁	0		
P _{2, P₅}	0.29	0	
P _{3, P₆, P₄}	0.30	0.26	0

merge (P_2P_5, P_3P_6, P_4)

$$* \text{ dist. } (P_2P_5P_4P_3P_6, P_1) = \frac{1}{2} [0.29 + 0.30] \\ = 0.30.$$

(5)	P_1	$P_2P_3P_4P_5P_6$
P_1	0	
$P_2P_3P_4P_5P_6$	0.30.	0

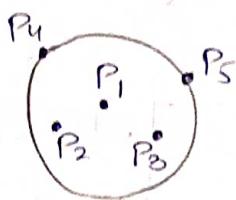


11/03/23

* DBSCAN :-

(Density based spatial clustering of applications with Noise)

- the main aim of DBSCAN is to create clusters with a minimum size & density.
- Density is defined as a minimum no. of points with a certain distance of each other.
- the concept of threshold values & minimum points and epsilon.



P₁ → Center point

ϵ = α units (EPSILON)

minimum points is 4

- DBSCAN divides the data points into 3 types.
 - 1) Core point
 - 2) Noise point
 - 3) Border point

(1) Core point :-

P₁ can form a cluster by using P₂, P₃, P₄, P₅ then P₁ is called as Core point.

If Minimum points = 6 , P₁ cannot form a cluster, thus P₁ is not a core point.

then P₁ may be Noise point or Border point.

(2) Noise point :-

P₁ is center point

ϵ = α units

minimum points = 4



As it has only 2 points in its cluster, thus, it is called as Noise point.

Example:-

Point	X	Y	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆	P ₇	P ₈	P ₉	P ₁₀	P ₁₁	P ₁₂
P ₁	3	0	P ₁	0										
(P ₂)	4	6	P ₂	1.11	0									
P ₃	5	5	P ₃	2.83	1.41	0								
P ₄	6	11	P ₄	1.24	2.83	1.11	0							
(P ₅)	7	3	P ₅	5.66	4.24	2.83	1.41	0						
P ₆	6	2	P ₆	5.83	1.11	7.31	6.00	1.41	0					
P ₇	7	0	P ₇	6.10	5.00	3.61	2.24	1.00	1.00	0				
P ₈	8	4	P ₈	5.83	4.47	3.16	2.00	1.41	2.83	2.24	0			
P ₉	3	3	P ₉	4.00	3.16	2.83	3.16	4.00	3.16	4.12	5.10	0		
P ₁₀	2	6	P ₁₀	1.11	2.00	3.16	4.11	5.63	5.66	6.10	6.32	3.16	0	
(P ₁₁)	3	15	P ₁₁	2.00	1.11	2.00	3.16	4.11	4.24	5.00	5.10	6.00	1.41	0
P ₁₂	2	4	P ₁₂	3.16	2.83	3.16	4.00	5.10	4.11	5.39	6.00	1.41	2.00	1.41

$$\text{Epsilon} = 1.9$$

$$\text{min points} = 4$$

$$P_1 : P_2, P_{10}$$

$$P_6 : P_5, P_7$$

$$P_{11} : P_2, P_{10}, P_{12}$$

$$P_8 : P_1, P_3, P_{11} \quad \checkmark$$

$$P_7 : P_5, P_6$$

$$P_{12} : P_9, P_{11}$$

$$P_3 : P_2, P_4$$

$$P_8 : P_5$$

$$\text{Core points} \geq 4 \text{ (min pts)}$$

$$\text{i.e., } P_2, P_5, P_{11}$$

$$P_4 : P_3, P_5$$

$$P_9 : P_{12}$$

$$P_5 : P_4, P_6, P_7, P_8 \quad \checkmark$$

$$P_{10} : P_1, P_{11}$$

<u>Points</u>	<u>Status</u>	
P ₁	Noise	Boulder
P ₂	Core	
P ₃	Noise	Boulder
P ₄	Noise	Boulder
P ₅	Core	
P ₆	Noise	Boulder
P ₇	Noise	Boulder
P ₈	Noise	Boulder
P ₉	Noise	
P ₁₀	Noise	Boulder
P ₁₁	Core	
P ₁₂	Noise	Boulder

$$P_1 = P_3, P_{10}$$

P₁ has a core point
so, it is borden point

$$P_3 = P_2, P_4$$

P₃ is a borden point
bcz P₂, is a core point or it has
core point.

<u>Points</u>	<u>X</u>	<u>Y</u>	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆	P ₇	P ₈
P ₁	2	10	0							
P ₂	2	5		0						
P ₃	8	4			0					
P ₄	5	8				0				
P ₅	7	5					0			
P ₆	6	4						0		
P ₇	1	2							0	
P ₈	4	9								0
$\epsilon = 2 \text{ units}$										
min points = 3										

$$P_8 | 2.2 | 4.47 | 6.40 | 4.47 | 6.40 | 4.47 | 5 | 5.39 | 0$$

$$6.40 | 4.47 | 5 | 5.39 | 0$$

$$6.40 | 4.47 | 5 | 5.39 | 0$$

$$6.40 | 4.47 | 5 | 5.39 | 0$$

P ₁ :	
P ₂ :	
P ₃ : P ₅ , P ₆	✓
P ₄ : P ₈	
P ₅ : P ₃ , P ₆	✓
P ₆ : P ₃ , P ₅	✓
P ₇ :	
P ₈ : P ₄	

Points	Status
P ₁	Noise
P ₂	Noise
P ₃	Core
P ₄	Noise
P ₅	Core
P ₆	Core
P ₇	Noise
P ₈	Noise

Evaluation of Clustering Algorithms :-

Clustering can be evaluated based on:

① Clustering tendency.

② No. of clusters

③ Clustering quality.

<https://towardsdatascience.com/clustering-evaluation-strategies-98a4006fcfc>

→ Clustering tendency :- If the data doesn't contain clustering tendency then clusters identified by any clustering algorithms may be irregular. → Non uniform distribution of points in datasets becomes important in clustering.

→ Some of the clustering algs like K-Means require 'K' as clustering parameter.

→ No. of clusters :-

→ K should not be too high nor too low.

→ There are 2 major approaches to find optimal no. of clusters.

① Domain knowledge

② Data driven approach.

* Domain knowledge :- On what data we are performing clustering

For ex: In dataset, if we have prior knowledge of species

(virginica, versicolor) then K=3.

* Data driven approach :- If the domain knowledge is not available, Mathematical methods help in finding out right no. of clusters.

→ Clustering Quality: Once clustering is done, how well the clustering is performed can be qualified by a no. of metrics. There are mainly 2 types of measures to assess the clustering performance.

- ① Extrinsic measures
- ② Intrinsic measures.

⇒ Extrinsic measures: requires ground labels.

Ex: adjusted rand index, mutual information based scores, homogeneity, completeness, F measures.

⇒ Intrinsic measures: doesn't require ground labels.

Ex: Coefficient, Davis - Bouldin Index

18/3/23 Internal measures:

① Cohesion

② Separation (coupling).

→ Cohesion: how closely related objects present in a cluster.

→ Separation: Measures how distinct or well-separated cluster is from other cluster.

* Example: Find cohesion and coupling for the data points.

cluster	F ₁	F ₂	Centroid
C ₁	1	0	$(\frac{1+1}{2}, \frac{0+1}{2}) = (1, \frac{1}{2}) = (1, 0.5)$
C ₁	1	1	
C ₂	1	2	
C ₂	2	3	$(\frac{1+2+2+1}{4}, \frac{2+3+2+1}{4}) = (\frac{3}{2}, \frac{9}{4})$
C ₂	2	2	
C ₂	1	2	
C ₃	3	1	$(\frac{3+3+2}{3}, \frac{1+3+1}{3}) = (\frac{8}{3}, \frac{5}{3})$
C ₃	3	3	
C ₃	2	1	$= (2.67, 1.67)$

Cohesion (C_1):

$$(1,0) \quad (1,0.5) \rightarrow 0.5 \quad \text{Euclidian distance.}$$

$$(1,1) \quad (1,0.5) \rightarrow 0.5 \quad \left. \right\} = \underline{\underline{1}}$$

Cohesion (C_2):

$$(1,2) \quad (1.5, 2.25) \rightarrow 0.55 \quad \left. \right\} = \underline{\underline{0.55}}$$

$$(2,3) \quad (1.5, 2.25) \rightarrow 0.90 \quad \left. \right\} = \underline{\underline{0.90}}$$

$$(2,2) \quad (1.5, 2.25) \rightarrow 0.55 \quad \left. \right\} = \underline{\underline{0.55}}$$

$$(1,2) \quad (1.5, 2.25) \rightarrow 0.55$$

Cohesion (C_3):

$$(3,1) \quad (2.67, 1.67) \rightarrow 0.75 \quad \left. \right\} = \underline{\underline{0.75}}$$

$$(3,3) \quad (2.67, 1.67) \rightarrow 1.37 \quad \left. \right\} = \underline{\underline{1.37}}$$

$$(2,1) \quad (2.67, 1.67) \rightarrow 0.95$$

Separation (C_1, C_2): $d(\text{Centroid } (C_1), \text{Centroid } (C_2)):$

$$(1,0.5), (1.5, 2.25) \rightarrow \sqrt{(0.5)^2 + (1.75)^2} = 1.82$$

Separation (C_2, C_3):

$$(1.5, 2.25), (2.67, 1.67) \rightarrow \sqrt{(1.17)^2 + (0.58)^2} = 1.36$$

Separation (C_1, C_3):

$$(1,0.5), (2.67, 1.67) \rightarrow \sqrt{(1.67)^2 + (1.17)^2} = 2.03$$

* Silhouette Coefficient :- It combines ideas of both Cohesion and Separation. But for individual points as well as clusters and clusterings. for an individual point i . Calculate a = average distance of i to the points in its cluster.

Calculate b = minimum average distance of i to points in another cluster.

Then, Silhouette Coefficient is given by :

$$S = 1 - \frac{a}{b}$$

Ex:

Point	cluster label.
P ₁	1
P ₂	1
P ₃	2
P ₄	2

Compute Silhouette Coefficient for each point, each of the 2 clusters & overall cluster.

Dissimilarity matrix:-

Points	P ₁	P ₂	P ₃	P ₄
P ₁	0	0.10	0.65	0.55
P ₂	0.10	0	0.70	0.60
P ₃	0.65	0.70	0	0.30
P ₄	0.55	0.60	0.30	0

Pointwise:-

Point P₁:

$$a = \frac{0.1}{1} = 0.1$$

$$b = \frac{0.65+0.55}{2} = 0.6$$

$$S = 1 - \frac{a}{b} = 1 - \frac{0.1}{0.6}$$

$$= \frac{5}{6}$$

$$= 0.83$$

Point P₂:

$$a = \frac{0.1}{1} = 0.1$$

$$b = 0.65$$

$$S = 0.84$$

Point P₃:

$$a = \frac{0.3}{1} = 0.3$$

$$b = \frac{0.65+0.7}{2} = \frac{1.35}{2}$$

$$= 0.67$$

$$S = 1 - \frac{a}{b}$$

$$= 1 - \frac{0.3}{0.67}$$

$$= 0.553$$

Point P₄:

$$a = 0.3$$

$$b = 0.575$$

$$S = 1 - \frac{a}{b}$$

$$= 0.48$$

cluster-wise :-

$$\text{for } C_1 = \frac{0.833 + 0.846}{2} = 0.839$$

Avg of $P_1 \& P_2$

$$\text{for } C_2 = \frac{0.55 + 0.48}{2} = 0.517$$

Avg of $P_3 \& P_4$

$$\text{overall Silhouette coefficient} = \frac{0.839 + 0.517}{2} = 0.68.$$

* External Measures of cluster Validity :- Purity.

cluster	entertain- ment	financial	foreign	metro	station	sport	Purity
1	3	5	40	506	96	27	0.7475
2	4	7	280	29	39	2	0.7756
3	1	1	1	7	4	671	0.9796
4	10	162	3	119	73	2	0.4390
5	33	22	5	70	13	23	0.7134
6	5	358	12	212	48	13	0.5525

Purity: One cluster is dominated by one class.

If a cluster has only 1 class then the purity is 1.

Purity = $\frac{\text{max val}}{\text{sum of all val.}}$

Purity for cluster 1 = $\frac{506}{3+5+40+506+96+27}$

$$= \frac{506}{677} = 0.747.$$

$$\text{Purity for cluster 2: } \frac{280}{2+4+7+280+29+39+2} = 0.7756.$$

$$\text{Purity for cluster 3: } \frac{671}{1+1+1+7+4+671} = 0.9796.$$

$$\text{Purity for cluster 4: } \frac{162}{10+162+3+119+73+2} = 0.4390.$$

$$\text{Purity for cluster 5: } \frac{70}{83+22+5+70+13+23} = 0.7134$$

$$\text{Purity for cluster 6: } \frac{358}{5+358+12+212+48+13} = 0.5525.$$

* Dimensionality Reduction :-

It refers to the techniques that are used to reduce the no. of input variables in large given data set, while reducing the dimensionality we try to preserve most of the info in the given large data set. that means we remove only unwanted/unimportant features.

Techniques:

- 1) factor analysis :- It is a variable reduction technique. reduction of set of variables into a small no. of latent factors.
- unobserved factors (latent factors) account for correlation among observed variables.

Ex:1 academic ability of a student can be done by
quantitative ability
verbal ability

→ quantitative ability covers the mathematics, physics, computer programming.

→ verbal ability covers the english, verbal reasoning.

Ex:2 restaurants are categorised on the basis of 6 variables.

1) waiting time

2) cleanliness

3) staff behaviour

4) taste of food

5) food freshness

6) food temp (hot or cool)

first 3 indicates service, next 3 indicates food quality.

latent variables are variables that are not directly observed but are inferred from other variables.

- * Principal Component Analysis (PCA):-
- A method used to project data in high dimensional space into a lower dimensional space by maximising the variance of each new dimension.
 - PCA is mostly used as a tool in exploratory data analysis & for making the predictive models.

Ex: Given following data use PCA to reduce dimensions from 2 to 1

Feature	Ex1	Ex2	Ex3	Ex4
x	4	8	13	7
y	11	4	5	14

Step 1: No. of features = 2, n = 2
No. of samples = 4, N = 4

Step 2: Computation of mean of variables

$$\bar{x} = \frac{4+6+13+7}{4} = 8$$

$$\bar{y} = \frac{11+4+5+14}{4} = 8.5$$

Step 3: Computation of co-variance matrix

→ Order pairs for xy are

→ Covariance of all order pairs

$$\text{Cov}(x,x) = \frac{1}{N-1} \sum_{k=1}^N (x_k - \bar{x})^2 = \frac{1}{4-1} [(4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2]$$

$$\boxed{\text{Cov}(x,x) = 14}$$

$$\text{Cov}(xy) = \frac{1}{4-1} [(4-8)(11-8.5) + (8-8)(4-8.5) + (13-8)(5-8.5) + (7-8)(14-8.5)]$$

$$\boxed{\text{Cov}(xy) = -11}$$

$$\text{Cov}(y,y) = \text{Cov}(y) = \frac{1}{4-1} [(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2]$$

$$\boxed{\text{Cov}(y,y) = 23}$$

Covariance matrix

$$S = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}_{2 \times 2}$$

Step 4: Eigen value, Eigen vector & Normalized Eigen vector.

(1) Eigen value: $\det(S - \lambda I) = 0$

$$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}_{2 \times 2} \Rightarrow \lambda I = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$S - \lambda I = \begin{bmatrix} 14 - \lambda & -11 \\ -11 & 23 - \lambda \end{bmatrix}$$

$$\det(S - \lambda I) = \det \begin{vmatrix} 14 - \lambda & -11 \\ -11 & 23 - \lambda \end{vmatrix} \Rightarrow (14 - \lambda)(23 - \lambda) - (-11)(-11) = 0.$$

$$\Rightarrow \lambda^2 - 37\lambda + 201 = 0 \Rightarrow \lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$\lambda_1 = 30.3849, \lambda_2 = 6.6151$$

Eigen Vector of λ_1 :

$$(S - \lambda_1 I) u_1 = 0 \quad u_1 = \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix}$$

$$\begin{bmatrix} 14 - \lambda_1 & -11 \\ -11 & 23 - \lambda_1 \end{bmatrix} \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} = 0$$

$$(14 - \lambda_1)u_{11} - 11u_{12} = 0 \Rightarrow (14 - \lambda_1)u_{11} = 11u_{12}$$

$$-11u_{11} + (23 - \lambda_1)u_{12} = 0 \quad \text{divide by } 11 \Rightarrow \frac{u_{11}}{11} = \frac{u_{12}}{23 - \lambda_1}$$

when $t = 1$

$$u_{11} = 11$$

$$u_{12} = 14 - \lambda_1$$

$$\lambda_1 = 30.3849$$

\Rightarrow Eigen vector u_1 of λ_1 =

$$= \begin{bmatrix} 11 \\ 14 - 30.3849 \end{bmatrix}$$

$$= \begin{bmatrix} 11 \\ -16.3849 \end{bmatrix}$$

\Rightarrow Normalize eigen vector u_1

$$e_1 = \begin{bmatrix} \frac{11}{\sqrt{(11)^2 + (-16.3849)^2}} \\ \frac{-16.3849}{\sqrt{(11)^2 + (-16.3849)^2}} \end{bmatrix}$$

$$e_1 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix} \quad e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

Step 5:- Derive new data set

	Ex1	Ex2	Ex3	Ex4
first PC (PC1)	P_{11}	P_{12}	P_{13}	P_{14}

$$P_{11} = e_1^T \begin{bmatrix} 4-8 \\ 11-8.5 \end{bmatrix}$$

$$P_{11} = [0.5574 - 0.8303] \begin{bmatrix} -4 \\ 2.5 \end{bmatrix}$$

$$P_{11} = [0.5574 - 0.8303] \begin{bmatrix} -4 \\ 2.5 \end{bmatrix}$$

$$P_{11} = -4.305$$

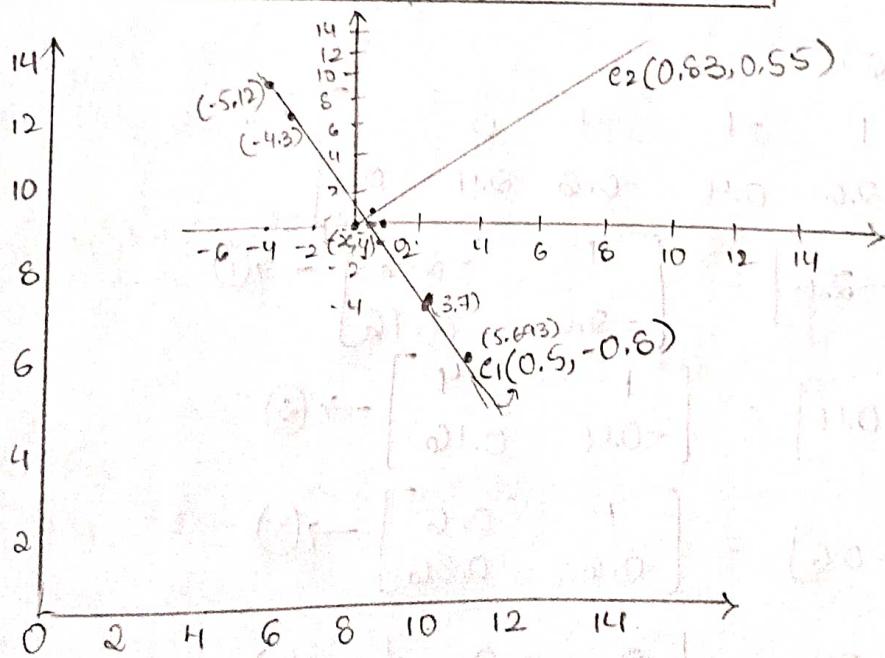
$$P_{12} = [0.5584 - 0.8303] \begin{bmatrix} 8-8 \\ 4-8.5 \end{bmatrix}$$

$$P_{12} = 3.736$$

$$P_{13} = [0.5574 - 0.8303] \begin{bmatrix} 13-8 \\ 5-8.5 \end{bmatrix} = 5.693$$

$$P_{14} = [0.5574 - 0.8303] \begin{bmatrix} 7-8 \\ 14-8.5 \end{bmatrix} = -5.124$$

	Ex1	Ex2	Ex3	Ex4
first PC (PC1)	-4.305	3.736	-5.693	-5.124



Q1/3

* Linear discriminant Analysis (LDA) :-

- LDA is a dimensionality reduction technique used as a feature processing step for pattern classification and ML applications.
- LDA is similar to PCA (Principal Component analysis) but LDA finds the axis that maximizes the separation b/w multiple classes.

Ex:-

$$C_1 = \gamma x_1 = (x_1, x_2) = \{(4, 1), (2, 4), (2, 3), (3, 6), (4, 4)\}$$

$$C_2 = \gamma x_2 = (x_1, x_2) = \{(9, 10), (6, 8), (9, 5), (8, 7), (10, 8)\}$$

Step 2: Compute within class scatter matrix (S_w)

$$S_w = S_1 + S_2$$

S_1 is the co-variance matrix for the class C_1

S_2 is the co-variance matrix for the class C_2

$$S_1 = \sum_{x \in C_1} (x - \mu_1)(x - \mu_1)^T$$

μ_1 is the mean of the class C_1

μ_2 is the mean of the class C_2

$$\mu_1 = (3, 3.6)$$

$$\mu_2 = (8, 4, 7, 6)$$

$$(x - \mu_1) = \begin{bmatrix} 1 & -1 & -1 & 0 & 1 \\ -2.6 & 0.4 & -0.6 & 2.4 & 0.4 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ -2.6 \end{bmatrix} \begin{bmatrix} 1 & -2.6 \end{bmatrix}^T = \begin{bmatrix} 1 & -2.6 \\ -2.6 & 6.76 \end{bmatrix} \rightarrow ①$$

$$\begin{bmatrix} -1 \\ 0.4 \end{bmatrix} \begin{bmatrix} -1 & 0.4 \end{bmatrix}^T = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 0.16 \end{bmatrix} \rightarrow ②$$

$$\begin{bmatrix} -1 \\ -0.6 \end{bmatrix} \begin{bmatrix} -1 & -0.6 \end{bmatrix}^T = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 0.36 \end{bmatrix} \rightarrow ③$$

$$\begin{bmatrix} 0 \\ 2.4 \end{bmatrix} \begin{bmatrix} 0 & 2.4 \end{bmatrix}^T = \begin{bmatrix} 0 & 0 \\ 0 & 5.76 \end{bmatrix} \rightarrow ④$$

$$\begin{bmatrix} 1 \\ 0.4 \end{bmatrix} \begin{bmatrix} 1 & 0.4 \end{bmatrix} = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 0.16 \end{bmatrix} \rightarrow ⑤$$

$$S_1 = \text{add } \frac{\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}}{5} = S_1 = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 0.6 \end{bmatrix}$$

$$S_2 = \sum_{x=w_2}^5 (x-\mu_2)(x-\mu_2)^T$$

$$\mu_2 = (6.4, 7.6)$$

$$(x-\mu_2) = \begin{bmatrix} 0.6 & -2.4 & 0.6 & -0.4 & 1.6 \\ 2.4 & 0.4 & -2.6 & -0.6 & 0.4 \end{bmatrix}$$

$$\begin{bmatrix} 0.6 \\ 2.4 \end{bmatrix} \begin{bmatrix} 0.6 & 2.4 \end{bmatrix} = \begin{bmatrix} 0.36 & 1.44 \\ 1.44 & 5.76 \end{bmatrix} \rightarrow ①$$

$$\begin{bmatrix} -2.4 \\ 0.4 \end{bmatrix} \begin{bmatrix} -2.4 & 0.4 \end{bmatrix} = \begin{bmatrix} 5.76 & -0.96 \\ -0.96 & 0.16 \end{bmatrix} \rightarrow ②$$

$$\begin{bmatrix} 0.6 \\ -2.6 \end{bmatrix} \begin{bmatrix} 0.6 & -2.6 \end{bmatrix} = \begin{bmatrix} 0.36 & -1.56 \\ -1.56 & 6.76 \end{bmatrix} \rightarrow ③$$

$$\begin{bmatrix} -0.4 \\ -0.6 \end{bmatrix} \begin{bmatrix} -0.4 & -0.6 \end{bmatrix} = \begin{bmatrix} 0.16 & 0.24 \\ 0.24 & 0.36 \end{bmatrix} \rightarrow ④$$

$$\begin{bmatrix} 1.6 \\ 0.4 \end{bmatrix} \begin{bmatrix} 1.6 & 0.4 \end{bmatrix} = \begin{bmatrix} 2.56 & 0.64 \\ 0.64 & 0.16 \end{bmatrix} \rightarrow ⑤$$

$$S_2 = \frac{\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}}{5} = S_2 = \begin{bmatrix} \frac{9.2}{5} & \frac{-0.20}{5} \\ \frac{-0.20}{5} & \frac{13.2}{5} \end{bmatrix} = \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix}$$

$$S_W = S_1 + S_2$$

$$= \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 0.6 \end{bmatrix} + \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix} = \begin{bmatrix} 0.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

$$S_W = \begin{bmatrix} 0.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Step 2: Compute b/w class Scatter matrix (S_B)

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$= \begin{bmatrix} -5.4 \\ 4 \end{bmatrix} \begin{bmatrix} -5.4 & 4 \end{bmatrix} = \begin{bmatrix} 59.16 & 21.6 \\ 21.6 & 16.00 \end{bmatrix}$$

Step 3: find the best LDA projection vector

Similar to PCA we find this using eigen vectors having largest eigen value.

$$S_w^{-1} S_B V = \lambda V$$

↳ projection vector

$$|S_w^{-1} S_B - \lambda I| = 0$$

$$\begin{vmatrix} 11.89 - \lambda & 8.81 \\ 5.08 & 3.76 - \lambda \end{vmatrix} = 0$$

We get $\lambda = 15.65$

$$\begin{bmatrix} 11.89 & 8.81 \\ 5.08 & 3.76 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 15.65 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0.91 \\ 0.39 \end{bmatrix}$$

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = S_w^{-1}(\mu_1 - \mu_2)$$

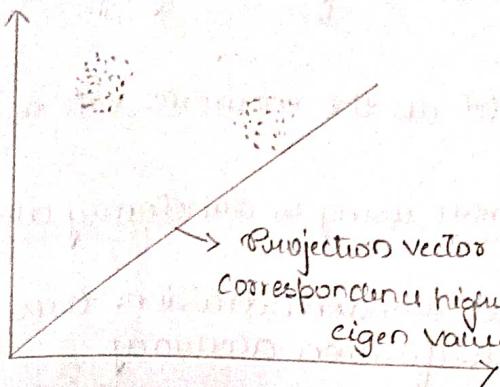
$$= \begin{bmatrix} 0.1921 & -0.031 \\ -0.032 & 0.36 \end{bmatrix} \begin{bmatrix} -5.4 \\ 4 \end{bmatrix} = \begin{bmatrix} -0.91 & -0.39 \end{bmatrix}^T$$

$$S_w^{-1} = \frac{1}{13.74} \begin{bmatrix} 5.08 & 0.44 \\ 0.44 & 2.64 \end{bmatrix}$$

$$S_w^{-1} = \begin{bmatrix} 0.384 & 0.032 \\ 0.032 & 0.192 \end{bmatrix}$$

Step 4 : Dimension reduction

$$Y = \omega^T x$$



Projection vector

Correspondence highest
eigen value

class 1 class 2

