

5/3/22

Module 5

CLUSTERING

* K-Means Clustering -

Dataset -

$$\{ 22, 9, 12, 15, 10, 27, 35, 18, 36, 11 \}$$

$K = 3$

$$C_1 = 22 \quad \{ 22, 27, 35, 18, 36 \} : \text{Iteration 1}$$

$$C_2 = 9 \quad \{ 9, 10 \}$$

$$C_3 = 12 \quad \{ 12, 15, 11 \}$$

$$\text{Mean } C_1 = 27.6$$

$$C_2 = 9.5$$

$$C_3 = 12.667$$

Iteration 2 :

$$C_1 = 27.6 \quad \{ 27, 35, 36, 22 \}$$

$$C_2 = 9.5 \quad \{ 9, 10, 11 \}$$

$$C_3 = 12.667 \quad \{ 12, 15, 18, 15 \}$$

$$\text{Mean } C_1 = 30$$

$$C_2 = 10$$

$$C_3 = 15$$

Iteration 3 :

$$C_1 = 30 \quad \{ 35, 36, 27, 18 \}$$

$$C_2 = 10 \quad \{ 9, 10, 11, 12 \}$$

$$C_3 = 15 \quad \{ 18, 15, 22 \}$$

$$\text{Mean } C_1 = 32.667$$

$$C_2 = 10.5$$

$$C_3 = 18.333$$

5 25/8/22

Iteration 4:

$$C_1 = 32.67 \{ 35, 36, 27 \}$$

$$C_2 = 10.5 \{ 9, 10, 11, 12 \}$$

$$C_3 = 18.34 \{ 22, 18, 15, 3 \}$$

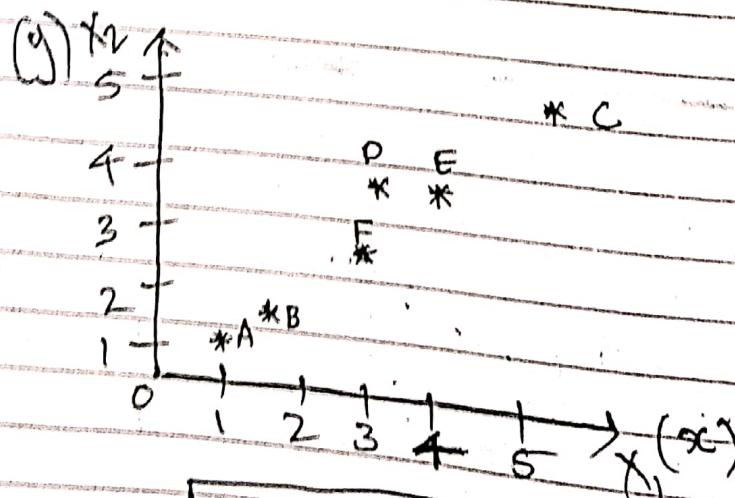
Mean $C_1 = 32.67$

$$C_2 = 10.5$$

$$C_3 = 18.34$$

* K-Means Clustering for two dimensional data

	$x_1(x)$	$x_2(y)$	
A	1		
B	1.5	1.5	$ K=2 $
C	5	5	
D	3	4	
E	4	4	
F	3	3.3	



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

5/8/22

* Distance Matrix -

	A	B	C	D	E	F
A	0	0.707				
B	0.707	0				
C	5.6568	4.95	0			
D	3.605	2.91	2.2	0		
E	4.242	3.53	1.41	1	0	
F	3.047	2.34	2.97	0.7	1.2	0

Iteration 1 -

$$G_1 = \{A, B\}$$

$$G_2 = \{C, D, E, F\}$$

	A	B	C	D	E	F
G_1	0	0.707	5.65	3.60	4.24	3.04
G_2	0.707	0	4.95	2.91	3.53	2.34
C_1	1	0	0	0	0	0
C_2	0	1	1	1	1	1

~~small~~ $G_1 = \{A, B\}$

$$G_2 = \{C, D, E, F\}$$

Calculate mean

$$G_1 = \{A, B\}$$

$$G_2 = \left(\frac{1.5 + 5 + 3 + 4 + 3}{5}, \frac{1.5 + 5 + 4 + 4 + 3}{5} \right)$$

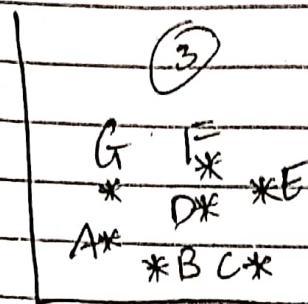
$$= (3.3, 3.56)$$

$$= (3.3, 3.56)$$

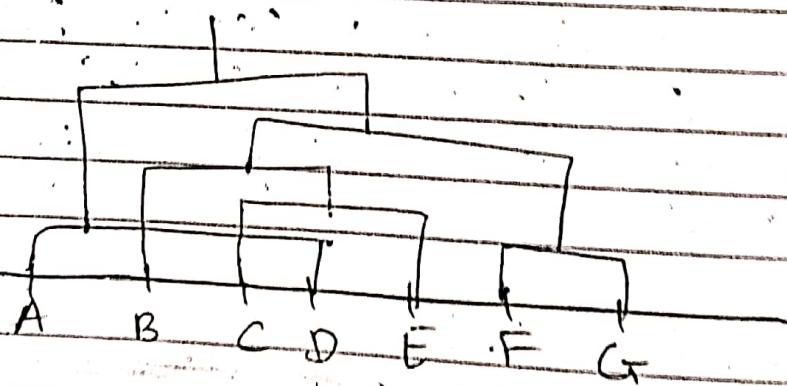
25 27/9/22

* Hierarchical Clustering

↳ Agglomerative
↳ Divisive



Dendrogram



There are 2 approaches of Hierarchical clustering

Agglomerative (Bottom up) & Divisive

(Top Down) Approaches

Following is comparison between two approaches

11/22

Agglomerative

- 1) It starts with a single data point as a cluster
- 2) Recursively add two or more appropriate clusters
- 3) The algorithm stops when k-number of clusters are achieved if k is specified otherwise it will stop after forming single large cluster

Divisive

- 1) It starts with single big cluster which consists of all the data samples
- 2) Recursively divides into smaller clusters
- 3) It will stop when k no of clusters are achieved if k is specified otherwise it will form the cluster of individual data points.

Hierarchical clustering uses distance matrix as a clustering criterion. The distance matrix or adjacency matrix can be constructed using 3 approaches which needs to 3 different techniques of forming the cluster.

1) Single linkage $\rightarrow D(C_1, C_2) = \min_{i \in C_1 \text{ & } j \in C_2} D(i, j)$

2) Complete linkage $\rightarrow D(C_1, C_2) = \max_{i \in C_1 \text{ & } j \in C_2} D(i, j)$

3) Average linkage $\rightarrow D(C_1, C_2) = \frac{\sum_{i \in C_1, j \in C_2} D(i, j)}{|C_1| \times |C_2|}$

7/9/22

where TG_{C_2} is sum of all pairwise clusters between G_1 & G_2 . NG_1 is no. of datapoints in G_1 & NG_2 is no. of datapoints in G_2

Consider the following distance matrix given & apply hierarchical clustering using single, complete & average linkages dendograms.

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

* Single linkage \rightarrow

	G_1	G_2	G_3	G_4	G_5
G_1	1	0			
G_2	2	9	0		
G_3	3	3	7	0	
G_4	4	4	5	9	0
G_5	5	5	10	2	8

Iteration 1 -

The distance between datapoint 3 & datapoint 5 is minimum.
So we can group 3 & 5 in one cluster

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

Iteration 2 -

The distance between cluster of 3, 5 & 1 is smallest. So combine 3, 5 & 1 as a single cluster.

	1, 3, 5	2	4
1, 3, 5	0		
2	7	0	
4	14	5	0

Iteration 3 -

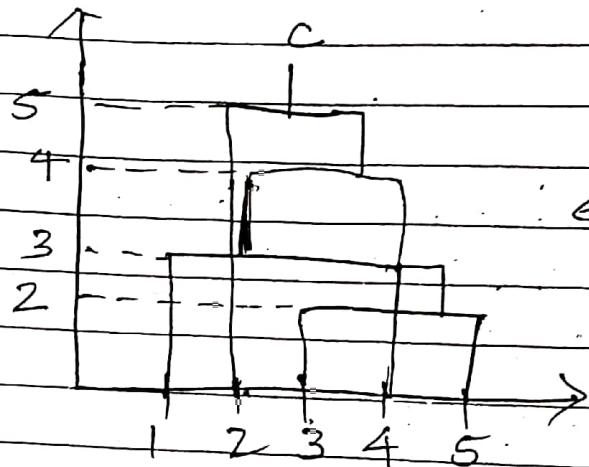
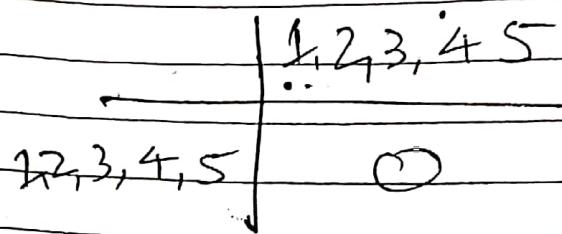
Distance between cluster 1, 3, 5 & 4 is smallest. So combine 1, 3, 5 & 4 as a single cluster.

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

Iteration 4 -

Distance between cluster 1, 3, 4, 5 & 2 is smallest. So we can club into 1, 2, 3, 4, 5 cluster.

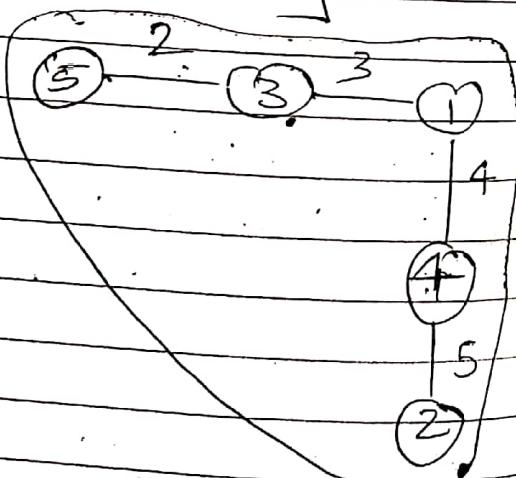
7/9/22



* Complete linkage -

<u>Edge</u>	<u>Cost</u>
5, 3	2
1, 3	3
1, 4	4
4, 2	5
4, 5	5
3, 2	7
4, 5	8
4, 3	9
1, 2	9
2, 5	10

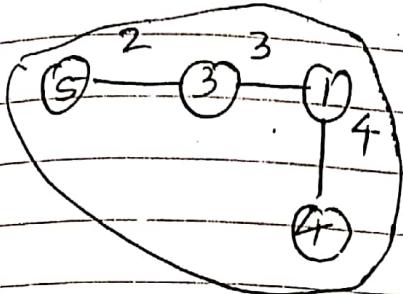
Initially



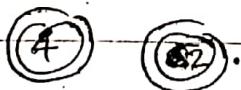
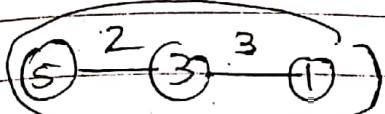
MST

Minimum Spanning Tree

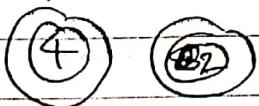
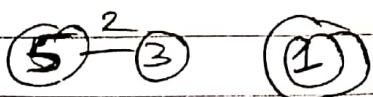
Iteration 1 -



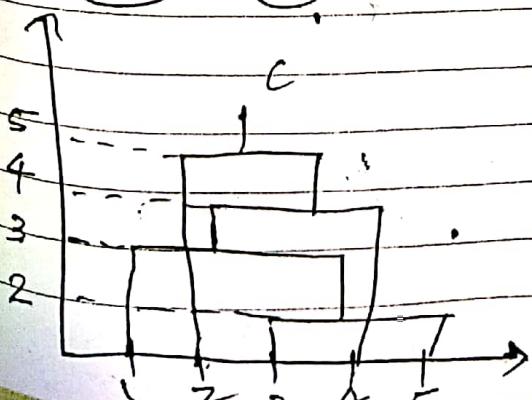
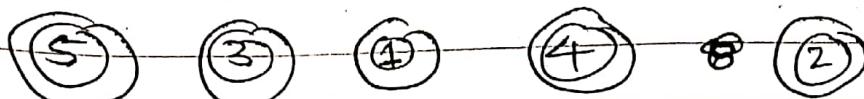
Iteration 2 -



Iteration 3 -



Iteration 4



< Dendrogram
(Bottom Up)

7/9/22



* Average Linkage -

	C_1	C_2	C_3	C_4	C_5	
C_1	0					
C_2	7	9	0			
C_3	3	3	7	0		
C_4	4	4	5	9	0	
C_5	5	5	10	12	8	0

Iteration 1 -

Distance between 3 & 5 is minimum
so we can club together

	3.5	21	2	4
3.5	0			
1	$6+5 = \frac{8}{2} = 4$			
2	$7+10 = \frac{17}{2} = 8.5$	9		
4	$9+8 = \frac{17}{2} = 8.5$	4	5	0

Iteration 2 -

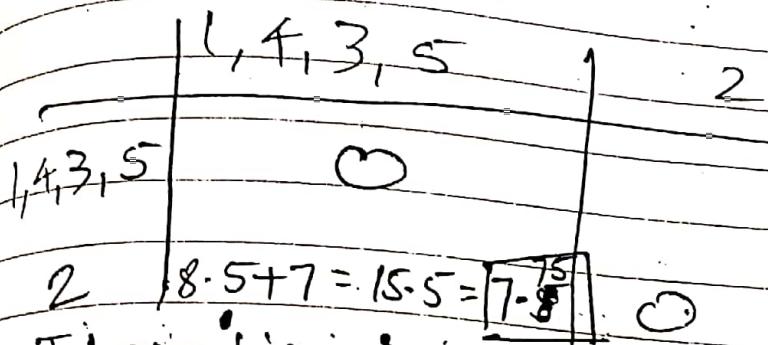
Distance between 1 & 4 is minimum
so we combine 1 & 4

	3.5	1,4	2
3.5	0		
1,4	$3+9+5+8 = \frac{25}{2} = 12.5$	0	
2	8.5	$9+5 = \frac{14}{2} = 7$	0

Iteration 3 -

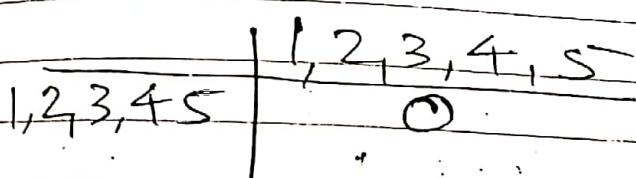
22

Distance between 3,5 & 1,4 cluster
is minimum so we can club together

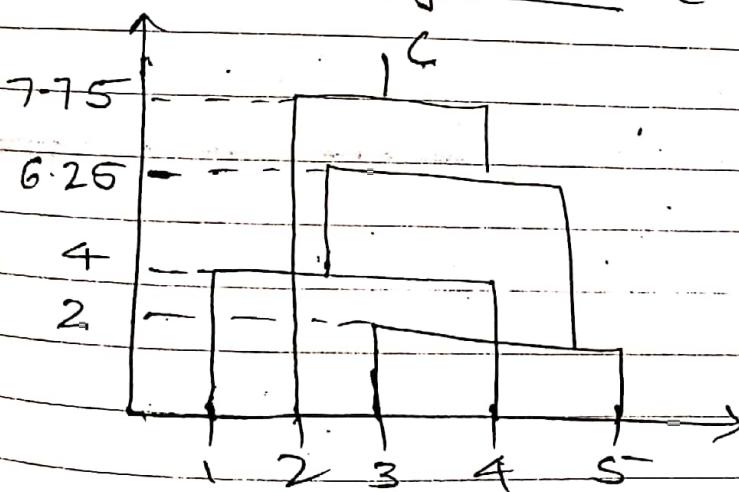


Iteration 4 :-

Distance between 1,4,3,5 & 2 is minimum so we can club together



Dendrogram (Top-Down)



Consider the following data samples with 2 features & construct the dendrogram with single, complete & average linkages.

7/9/22

	x_1	x_2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.53

A) Distance Matrix -

(*)	A	B	C	D	E	F
A	0					
B	0.707	0				
C	5.65	4.95	0			
D	3.605	2.91	2.23	0		
E	4.24	3.53	1.41	1	0	
F	3.05	2.34	2.97	0.7	1.2	0

* Single linkage -

Distance between F & D is smallest
so combine together

Iteration 1 -

	F, D	A	B	C	E
F, D	0				
A	3.05	0			
B	2.34	0.707	0		
C	2.23	5.65	4.95	0	
E	1	4.24	3.53	1.41	0

Iteration 2 -

Distance between A & B is smallest
so combine together

	F,D	A,B	C	E
D	O			
B	2.34	O		
C	2.23	4.95	O	
E		3.53	1.41	O

Iteration 3 -

Distance between F,D & E is smallest so combine

	F,D,E	A,B	C
F,D,E	O		
A,B	2.34	O	
C	1.41	4.95	O

Iteration 4 -

Distance between F,D,E & C is smallest
so combine

	F,D,E,C	A,B
A,B,C	O	
B	2.34	O

Iteration 5 -

Distance between F,D,E,F & AB smallest
so combine

	A,B,C,D,E,F
ABCF	O

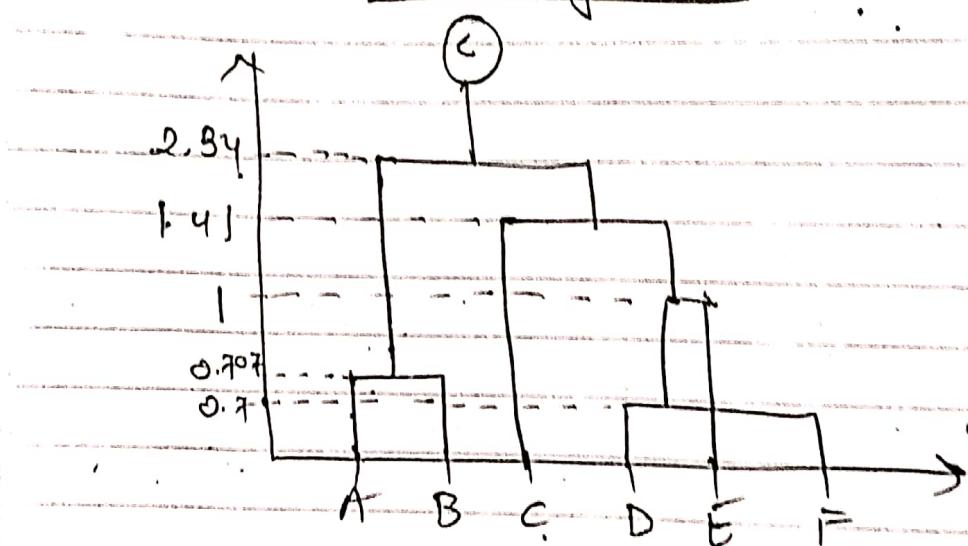
7/9/22

* Average linkage

Iteration 5 -



Dendrogram:



* Average linkage -

	A	B	C	D	E	F
A	0					
B	0.707	0				
C	5.68	4.95	0			
D	3.605	2.91	2.23	0		
E	4.24	3.53	1.41	1	0	
F	3.05	2.34	2.97	0.7	1.2	0

Iteration 1 -

Distance between D & F smallest
so combine

12

D, F	A	B	C	E
0	0			
3.05	0			
3.3275	0.707	0		
2.6	5.65	4.95	0	
1.1	4.24	3.53	1.41	0

Iteration 2 -

Distance between A & B smallest
so combine

D, F	A, B	C	E
0			
3.1887	0		
2.6	5.3	0	
1.1	3.885	1.41	0

Iteration 3 -

Distance between D, F & E smallest
so combine

D, F, E	A, B	C
0		
3.5368	0	
2.005	5.3	0

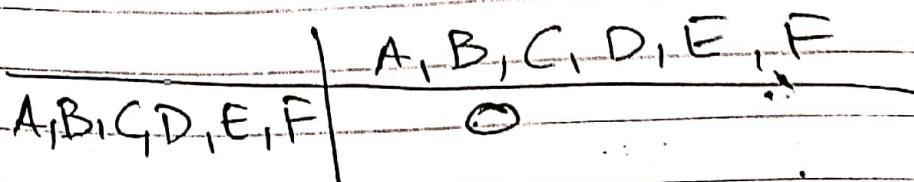
Iteration 4 -

Distance between D, F, E & C smallest
so combine

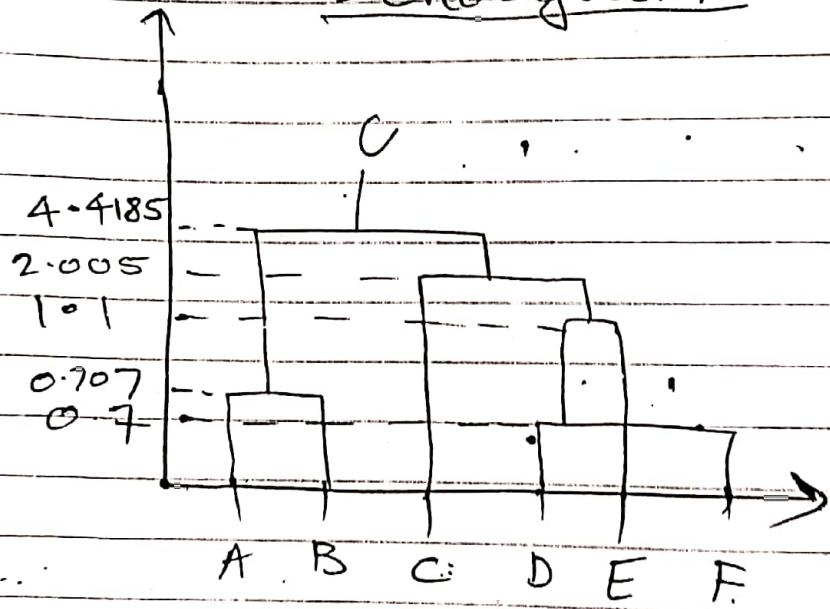
D, F, E, C	A, B
0	
4.4184	0

7/9/22

Iteration 5.
Distance between P, F, E, C &
A, B smallest so combine



Dendrogram



21/9/22

* Density Based Clustering -

- Partitioning the hierarchical methods are designed to find spherical clusters. They have difficulty in finding clusters of arbitrary shape.
- To find clusters of arbitrary shape, we can model clusters as dense regions in the object space separated by sparse regions. This is the main idea behind density based clustering.

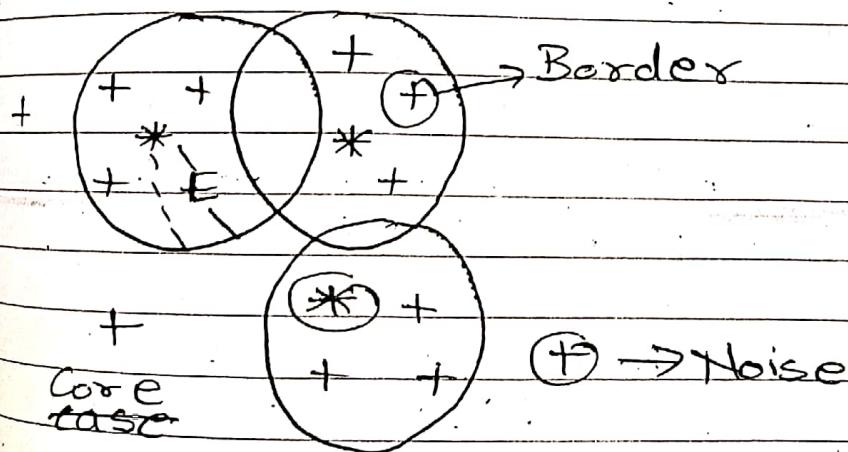
DB SCAN (Density based Spatial clustering of Application with noise) -

The density of an object O can be measured by the number of objects closed to O .

The algorithm finds core objects, the objects whose neighbourhoods are connected to form dense regions as clusters.

The algo has 2 user defined parameters
E-neighbourhood - to specify neighbourhood of an object O with radius absolute (E) centered at O .

Mean point - To decide the density threshold as dense region. In this concept, we come across 3 different data points



Core - This is a point that has at least 'm' points within distance n from itself

Border - This is a point that has at least one core point at distance n

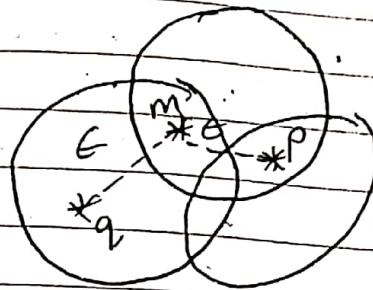
Noise - This is a point that is neither a core nor a border & it has less than 'm' points within distance n from itself

21/9/22

* Concept of DBSCAN -

- The Algo proceeds by arbitrarily picking up a point in the dataset (until all points have been visited).
- If there are atleast mean points, points within radius of ϵ to the point then we consider all these points to be part of the same clusters.
- The clusters are then expanded by recursively repeating the neighbourhood calculation for each neighbourhood.
- Given a Set D of objects, the algo first identifies the core objects using ϵ and mean points, these ~~core~~ objects and their neighbourhoods are used to form a denser region.
- For a core object q & an object p it is said that p is directly density reachable from q if p is ~~written~~ ϵ neighbourhood of q .
- The object p is density reachable from a core obj q , if there is chain of obj P_1, P_2, P_n such that $P_1 = q$ & so on upto $P_n = p$ & P_{i+1} is directly density reachable from P_i where $i \leq n$ wrt ϵ & mean points.

8/9/22



- M is directly density reachable from q & p is density directly reachable from M. So P is density reachable from q
- If p & q are core, then they both are density reachable from each other if q is a core obj & p is not then p is density reachable from q but q is not from p
- Two points p_1, p_2 are density connected if there is an obj q such that both p_1 & p_2 are density reachable from q w.r.t E & mean p
- A subset C of Dataset D is a dense region & so a cluster if -
 - i) for any 2 obj $O_1, O_2 \in C$, O_1 & O_2 are density connected
 - ii) There exists no obj $O \in C$ & another obj $O' \in (D-C)$ such that O & O' are density connected

* The Algo -

Step 1 - Initially all the objects in the given set are unvisited $D = \{ \} \rightarrow$

21/9/22

Step 2 - An Unvisited Object p is selected randomly $D = \{A, B, C\}$

3 - p is marked as visited

4 - If p is not a core obj then it is marked as noise point

5 - If p is a core obj then a new cluster C is created for p & all the objects in ϵ neighbourhood of p are added to a candidate set N , $N = \{ \dots \}$

6 - For each object p' in N which is unvisited, the algo labels it visited & checks if it is a core obj.

A : $N = \{ D, E, F, G \}$

(B)

C :

$p \in N$

D :

E F G

7 - If p' is a core obj, then all the obj in ϵ neighbourhood of p' are added to N .

8 - Algo adds the objects to clusters from C the objects of N that are not already added to clusters & such objects are removed from N .

9 - The process is repeated until N is empty & the cluster C is completed

10 - To find the next cluster, algo randomly selects an unvisited obj from the remaining objects ($D - C$)

21/9/22

H - The clustering process is continued until all the points are visited.

Module 6

21/9/22

* EM Algorithm -

It is an approach for maximum likelihood estimation (MLE) in the presence of latent variables. The most discussed application of EM Algorithm is for clustering with mixture model. The EM Algorithm is an iterative approach that cycles between 2 steps i.e. estimation step(E) and maximization step(M)

E: Classify the data using current theory

M: Generate the best theory using the current classification of data.

Step E generates the expected classification for each example and Step M generates the most likely theory given the classified data.

Eg) Imagine there are two coins A and B. One is more likely to get heads, the other more likely to get tails.

You pick one at random and toss it and then find out which one was it.

Suppose we do it five times i.e.