

## • Ensemble learning

Ensemble learning is a machine learning paradigm that combines the predictions of multiple model to improve overall performance & generalization.

Example :- Pappu pass ho gaya (Detail Explanation)

→ Why Ensemble learning?

- 1). To overcome limitations of individual models by leveraging their collective intelligence
- 2). Reduce overfitting & enhance robustness
- 3). Improved performance on ~~a~~ imbalanced Datasets
- 5). Enhanced Generalization
  - ↳ By combining diverse model through various ensemble methods

Plus

12.26 18:33

model enables better generalization as model focuses on different aspects of data.

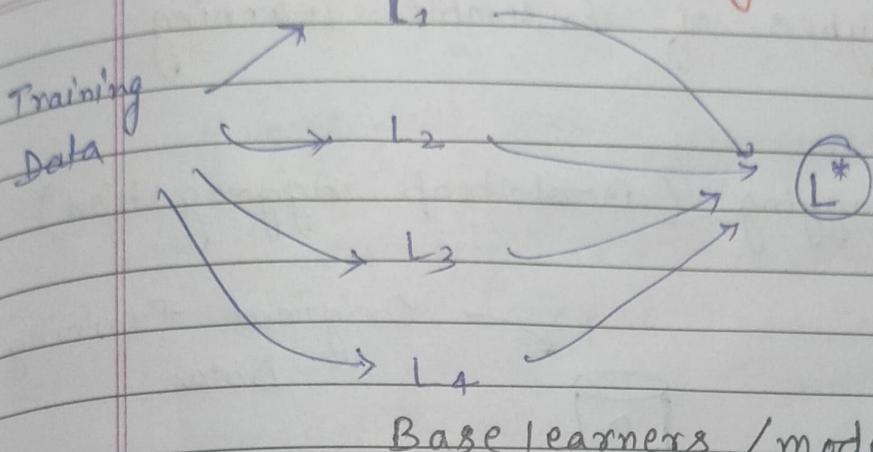
### c). Increased Model Accuracy

→ Base models: gt constituent models that form the ensemble (e.g. decision tree, neural network etc.).

→ Aggregation methods: Techniques for combining predictions (Eg voting, Averaging, stacking).

In short, Ensemble learning is a powerful strategy to enhance model performance by Combining the strength of diverse base model & Understanding the trade-offs & Selecting appropriate ensemble methods for effective predictions.

## Ensemble learning

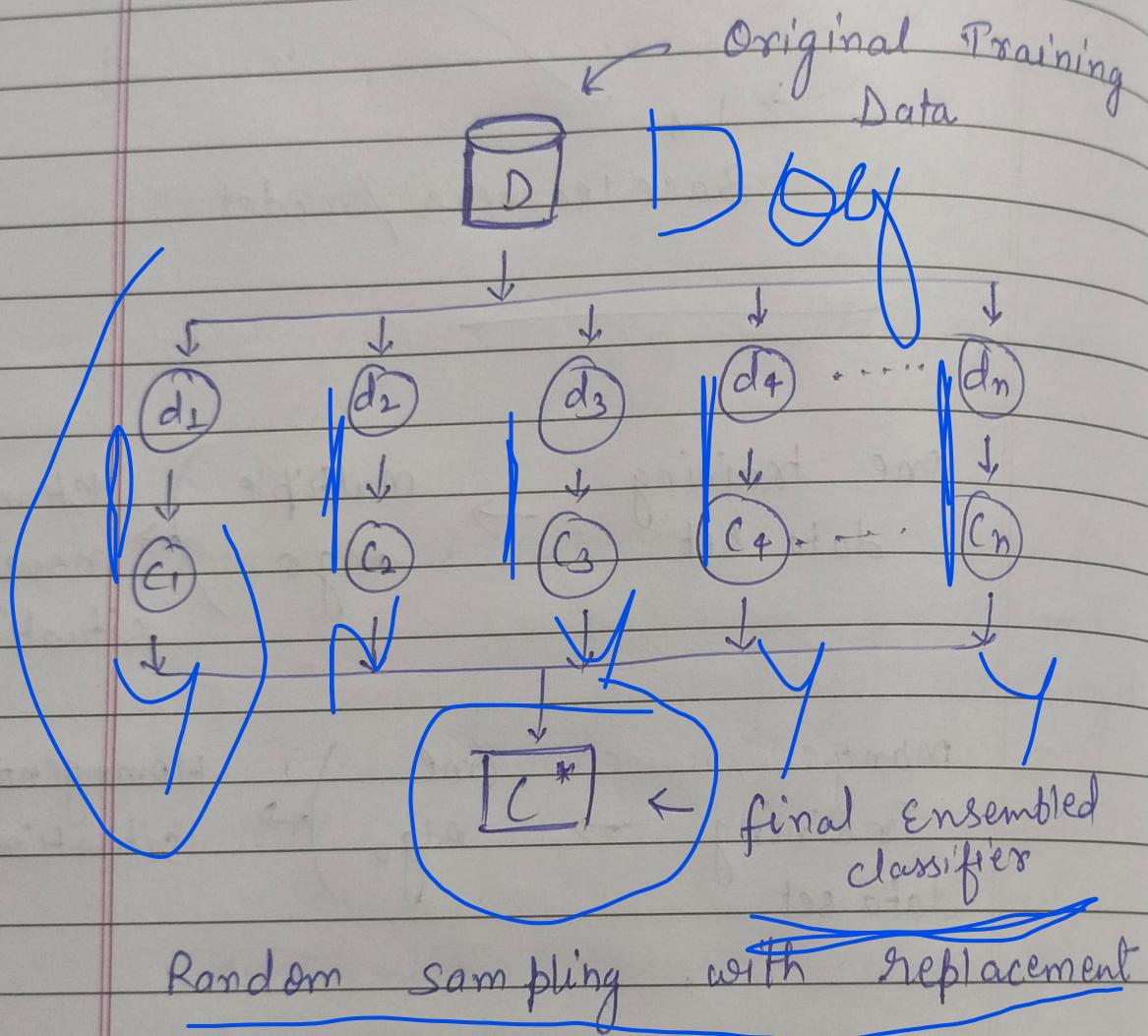


one training data set → Multiple Algo → Heterogeneous situation

Many Training data set → One Algo → Homogenous situation

## → Types of Ensemble Learning

### <1> Bagging (Bootstrap Aggregation)



multiple dataset are also known as Bootstrapped datasets.  
Sample

these samples are generated so that we can train the classifiers

Shot on OnePlus

based on these samples.

By Definition →

Bagging is an ensemble learning technique that aims to improve the stability and accuracy of machine learning algorithms by constructing multiple models in parallel on different random subsets of the training dataset

Process →

1). Bootstrap Sampling :-

Randomly draw, with replacement, multiple subsets (bags) from the original training dataset.

Same size subset, but it may contain repeated or missed data.

2). Model Training :-

Train base model independently on each bootstrap sample.

3. Prediction Aggregation:-

Aggregate prediction from all model to make final prediction.

E.g. Voting, Averaging etc.

## → Advantages

1. Reduce variance
2. Mitigate overfitting

This Bagging method is well suited for parallel processing, as base models can be trained independently on different subsets simultaneously.

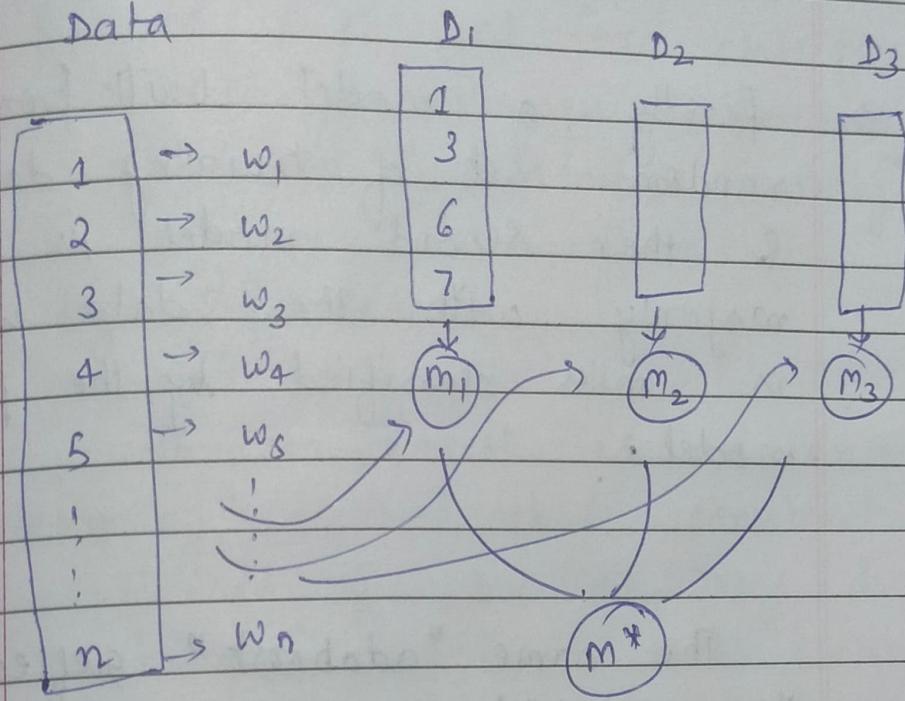
→ Boosting

Boosting is a family of ensemble ML Algos where a series of weak learners are trained sequentially with each subsequent model focusing on correcting the errors made by previous one. final prediction is weighed combination of predictions of all.

## → Boosting - AdaBoost

Training

Data



Boosting is also a homogeneous weak learner's model same as bagging but it works differently from bagging.

In this model, learner learn sequentially (instead of parallel as bagging) & adaptively to improve model prediction of a learning algorithm.

Here a strong classifier is achieved by using weak models in series.

- firstly, a model built from random set of training dataset & then second model is built majorly with the data which is miss classified by the previous model.

The name "adaboost" reflects its adaptive nature; it adapts to the weaknesses of the previous models by giving more weight to miss classified data/instances.

→ Strength :

- 1) High accuracy due to sequential learning
- 2) versatility : applicable to various types of weak learners.

→ limitation :-

on OnePlus

2023.12.26 18:34

- 1) sensitive to noisy data
- 2) overfitting

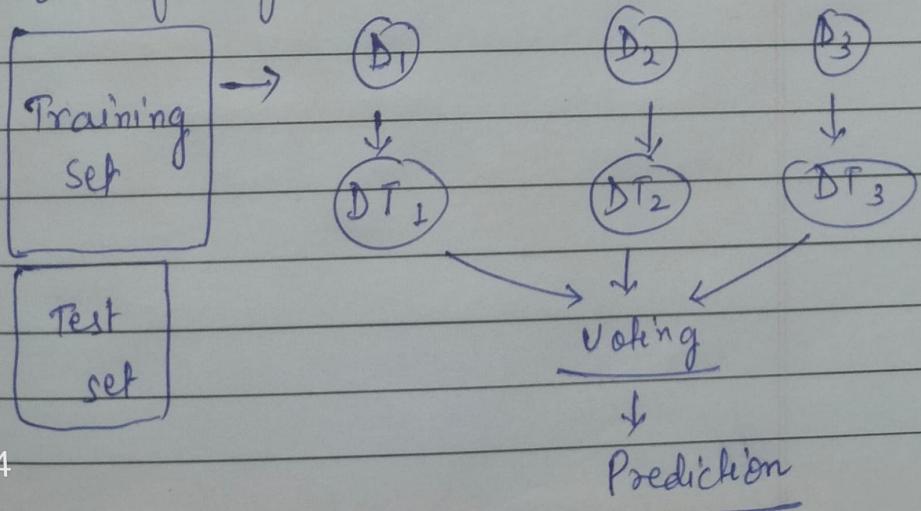
## Random forest

It belongs to supervised learning & it can be used for classification as well as Regression problems.  
It uses "Ensemble learning"

By Definition :-

Random forest is an ensemble learning method that constructs a multitude of decision trees during training & then takes average of all outputs to improve the predictive accuracy of that dataset.

- The greater number of trees in the forest leads to higher accuracy & prevents the problem of overfitting.



## Clustering

clustering is a type of Unsupervised learning in machine learning that involves grouping similar data points together based on certain criteria. The goal is to discover inherent pattern or structure within the data without prior knowledge of the class labels.

Eg:

K-means clustering

DB Scan

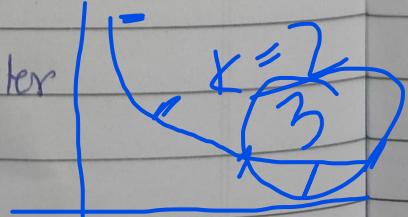
Hierarchical clustering



K-Means clustering

K-means clustering is an unsupervised machine learning algorithm, which groups the unlabeled dataset into different clusters.

K- Pre defined cluster

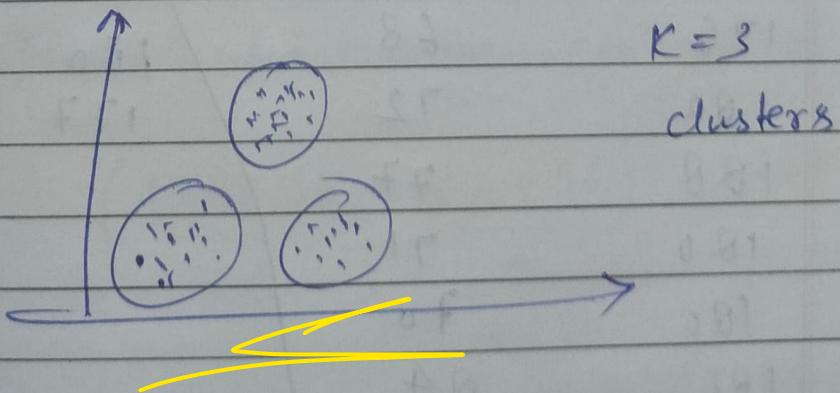


\* gt is an iterative algorithm that divides the unlabeled dataset into K clusters in such way that each instance of dataset belongs only one group that has similar properties.

→ gt a centroid based algorithm the main goal of this algorithm is to minimize the sum of distances b/w data point & their corresponding cluster.

→ first task of algorithm is - : to find best value of K center points by centroids by iterative process.

→ Assign each datapoint to its closest cluster.



→ steps involved in K-means clustering

1). Decide n clusters.

2). initialize centroids

3). assign cluster

4). move centroids

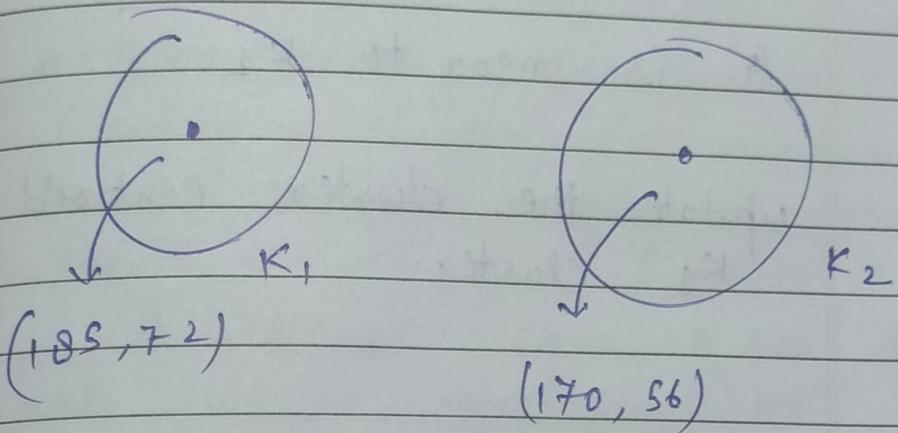
finish

→ Numerical based on K-means clustering Algo

Height	weight		
185	72		
170	56		
168	60		
179	68		
182	72	180	8867
188	77	177	76
180	71		
180	70		
183	84		
180	88		

## Euclidean Distance

$$\sqrt{(x_0 - x_c)^2 + (y_0 - y_c)^2}$$



### Euclidean Distance for ③

$$ED \text{ for } ③ \quad K_1 \rightarrow 20.00$$

$$\hookrightarrow K_2 \rightarrow 4.48$$

③ near to  $K_2$   
it will go in  $K_2$  cluster.

Now update the centroid of  $K_2$

$$\left( \frac{170+168}{2}, \frac{60+56}{2} \right) = \overbrace{\left( 169, 58 \right)}^{\uparrow}$$

New centroid for  
 $K_2$  cluster

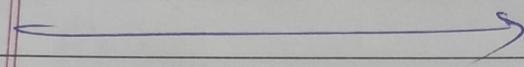
similarly for ④

$$K_1 \rightarrow 6.32$$

$$K_2 \rightarrow 14.14$$

④ is near to  $K_2$

update the cluster' centroid of  
 $K_1$  cluster.



final output

$$K_1 = \{1, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$$

$$K_2 = \{2, 3\}$$

DB Scan c  
 (Densi

DB Scan  
 machine i  
 identifies  
 based D  
 points.

Unlik  
 not he  
 number  
 & can  
 arbitra

DB  
 on thi  
 & no

each  
 neigh  
 con to  
 numb

## DB Scan Clustering (Density based clustering)\*

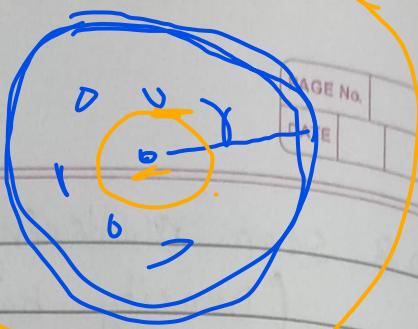
DB Scan is an unsupervised machine learning algorithm that identifies clusters in a dataset based on the density of data points.

Unlike K-means, DB Scan does not require specifying the number of clusters beforehand & can discover clusters of arbitrary shapes.

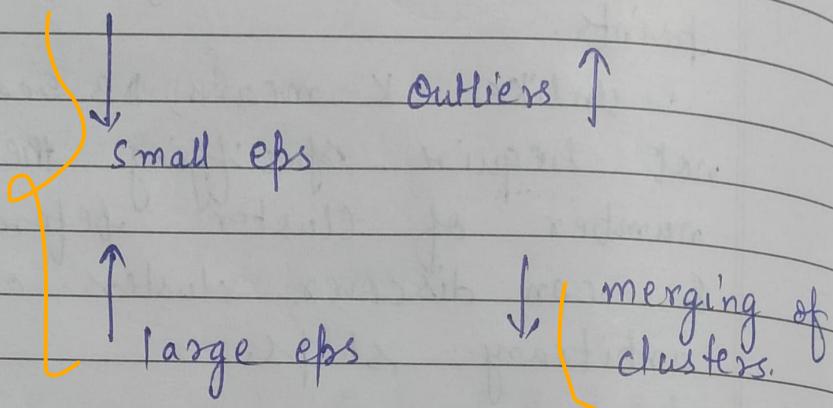
DB Scan algorithm is based on this intuitive notion of clusters & noise.

The key idea is that for each point of a cluster, the neighbourhood radius has to contain at least a minimum number of points.

\* Parameters



(1). eps:- It defines neighbourhood around a data point that is if the distance b/w two data points is  $\leq \text{eps}$  than they are considered neighbours.



that's why to find εps we use k-distance graph.

(2). minpts:- Minimum number of neighbours (data points) with in εps radius.

↑ large dataset      ↑ choose large minpts

in general →

$$\text{minpts} \geq D+1$$

where  $D$  is dimensions of dataset

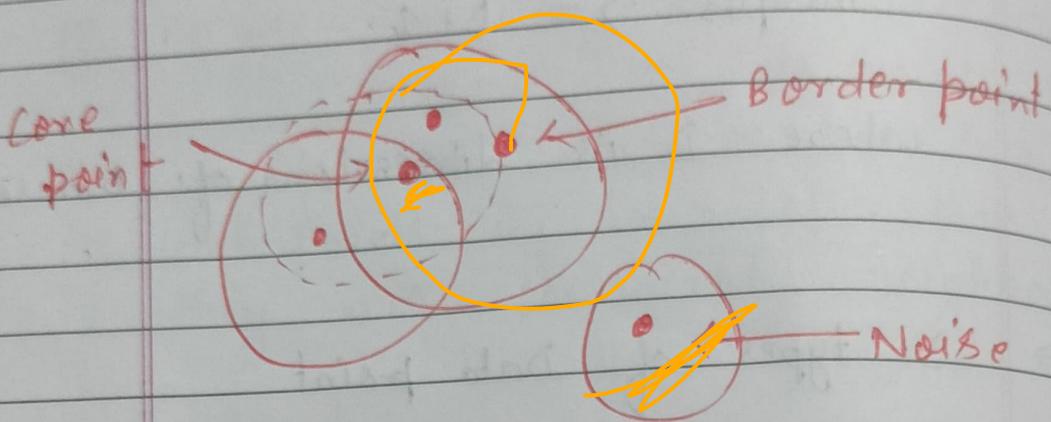
⇒ 3 types of Data point

1. Core point :- A point is core point if it has more than minpts within eps.

2. Border point :- A point which has fewer than minpts within eps but it is neighbour of core point

3. Outlier / Noise :- Not a core point nor border point.

## Numerical Based on DB Scan



is a point a  
Directly Density Reachable ? form any pair

Soln: 2 conditions (take two point p & q)

1.) q should neighbour of p

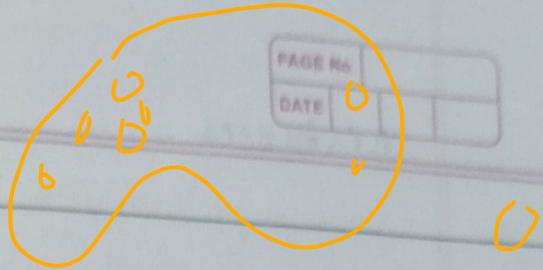
2.) p must be a core point

### Advantages of DB Scan

(1). Robust to outlier

(2). flexible in cluster shape

## Limitations



1. Difficulty with varying density
2. sensitive to parameter selection  
(choice of eps & minpts)

## Numerical

### Data Points

Ques:  $P_1 = (3, 7)$        $P_2 = (4, 6)$        $P_3 = (5, 5)$   
 $P_4 = (6, 4)$        $P_5 = (7, 3)$        $P_6 = (6, 2)$   
 $P_7 = (7, 2)$        $P_8 = (8, 1)$        $P_9 = (3, 3)$   
 $P_{10} = (2, 6)$        $P_{11} = (3, 5)$        $P_{12} = (2, 4)$

$$\text{minPts} = 4$$

$$\& \text{ epsilon } (\epsilon) / \text{eps} = 1.9$$

## Distances

DATE

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$	$P_7$
$P_1$	0	-	-	-	-	-	-
$P_2$	1.41	0	-	-	-	-	-
$P_3$	2.83	1.41	0	-	-	-	-
$P_4$	4.24	2.83	-	0	9.0	like this	-
$P_5$	5.66	-	-	-	-	-	-
$P_6$	5.83	-	-	-	-	-	-
$P_{12}$	3.16	-	-	-	-	-	-

Step 1: Calculate Distance of each datapoint with other datapoint

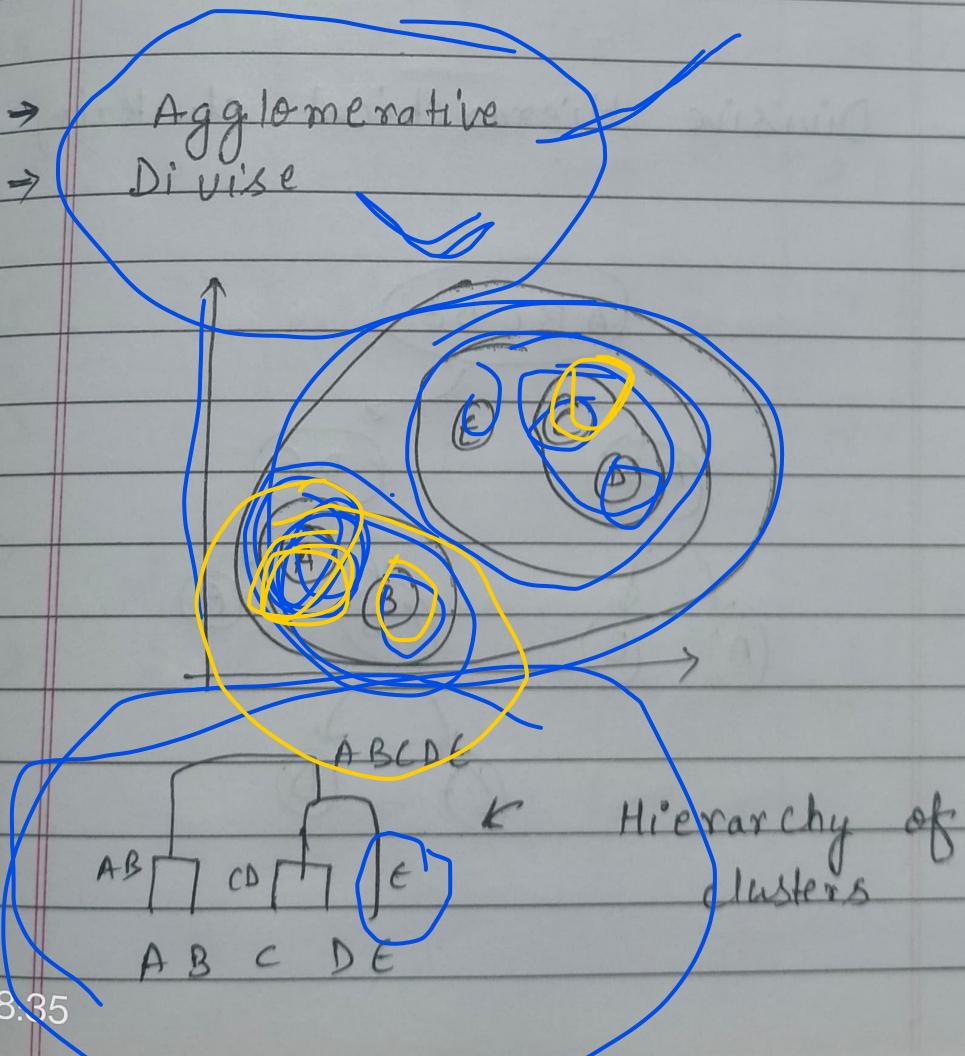
Step 2: for each point, figure out all the distances with datapoints which are  $\leq \epsilon$  - eps.

Step 3: Now Categorize all points as core / border / noise accordingly

## \* Hierarchical clustering

Hierarchical clustering is an unsupervised clustering algorithm that organizes data into hierarchical hierarchy of nested clusters.

The primary goal is to create a tree like structure (known as dendrogram) that visually represents the relationships btw data points & clusters.

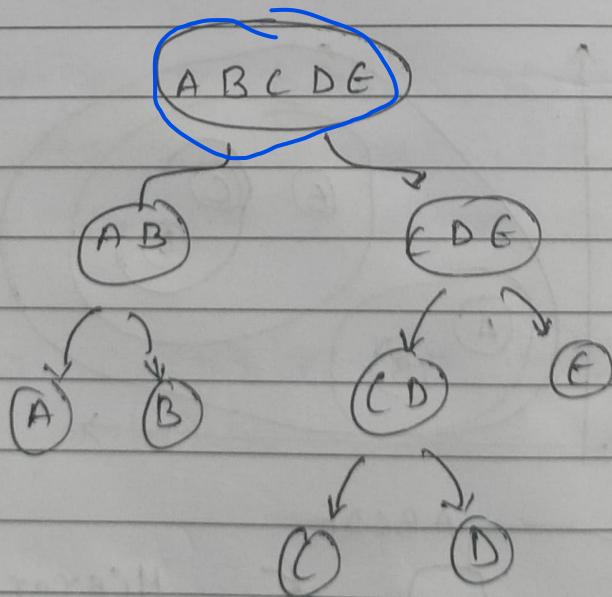


## Agglomerative Hierarchical clustering

It is bottom up approach in which the algorithm starts with taking all data points as single cluster & then merge them until they form one single cluster.

### Dendrogram

## Divisive Hierarchical clustering



Divisive Algorithm is just opposite of agglomeration clustering. It is top down approach.

### Single Linkage

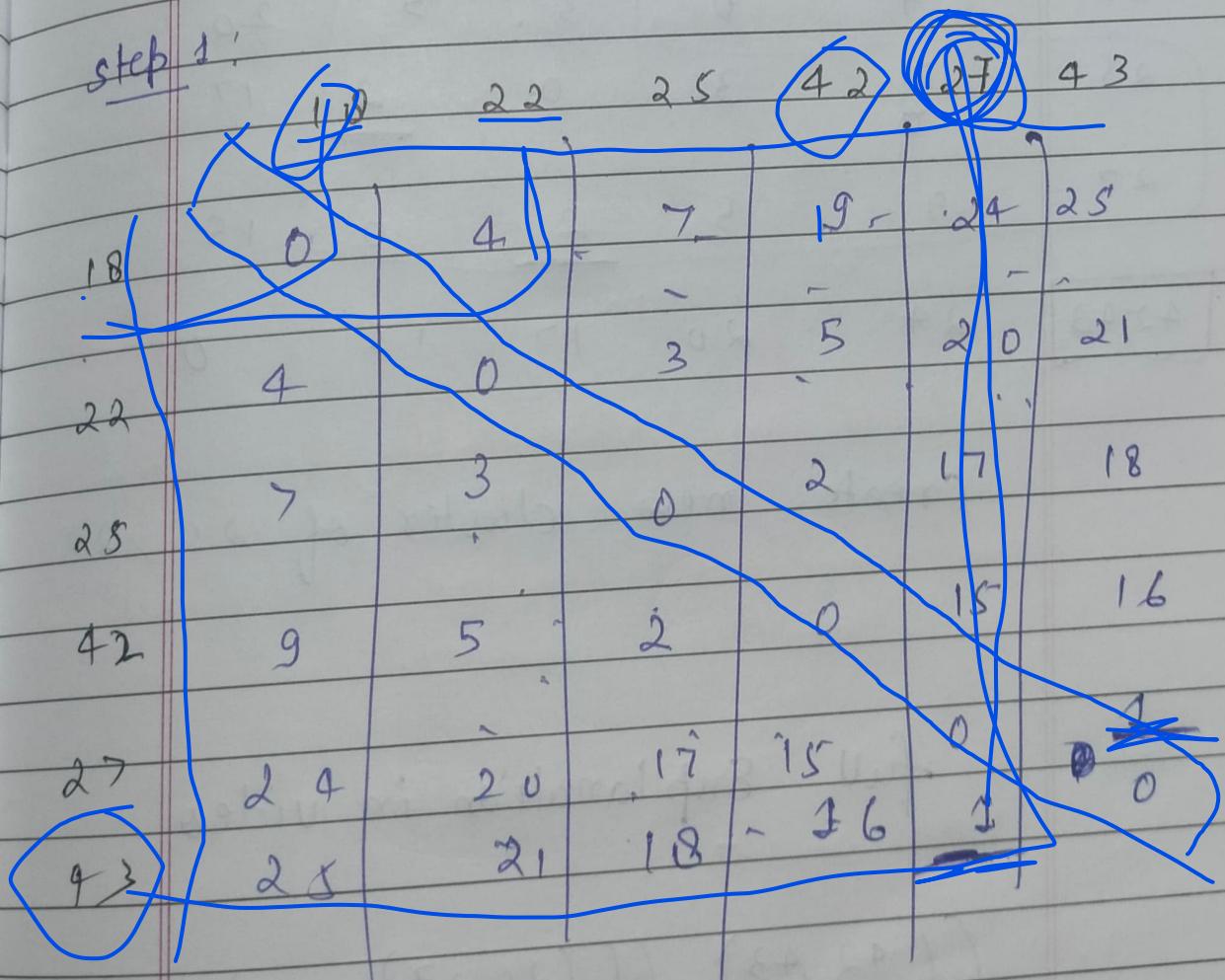
Complete Linkage 23,45 34,34

Numerical

18, 22, 25, (42, 43), 27

18, 22, 25, 42, 27, 43

step 1:



Create cluster of 42 & 43

remove the 43 now

& now will have

18, 22, 25, 27, (42, 43)

18	22	25	27	42, 43	
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	2	17
27	9	5	2	0	15
42, 43	24	20	17	15	0

Create new cluster of 25, 27

full explanation in video

$((42, 43), ((25, 27), 22), 18))$

42, 43, 25, 27, 22, 18

PAGE NO.

DATE

25, 27, 22, 18

25, 27, 22

42, 43

42

43

25, 27

25

27

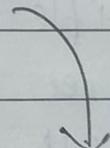
22

18

why Hierarchical clustering?

complete linkage formula

Single linkage



max Distance

minimum

Distance

## \* Apriori Algorithm

### • Fuzzy c - Mean Algorithm

Fuzzy C-mean (FCM) is a clustering algorithm that belongs to the family of fuzzy clustering methods.

It is an extension of Traditional K-means clustering, allowing data points to belong to multiple clusters with varying degree of membership rather than strictly assigning them to a single cluster.

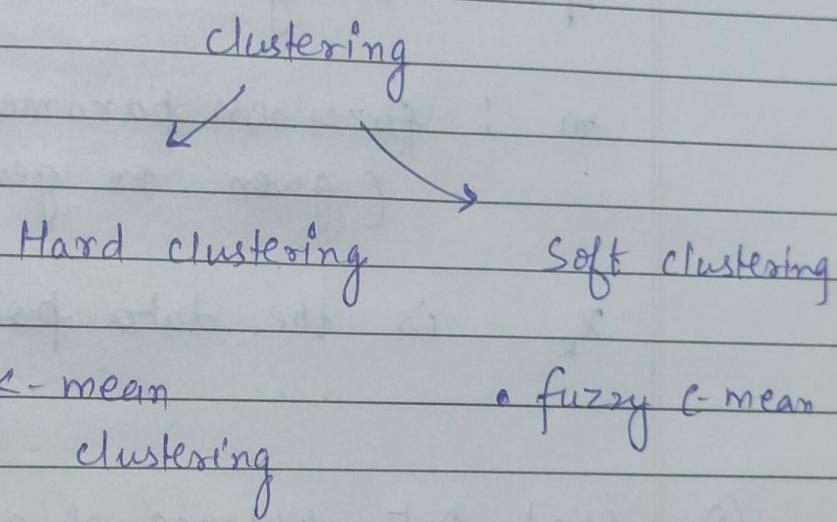
Fuzzy clustering is particularly very useful when there is uncertainty or ambiguity in the assignment of point to cluster.

#### Use cases:-

- 1). Widely used in pattern recognition, segmentation & data mining

→ One of the key parameter in Fcm is fuzziness ( $m$ ) which controls the degree to belong to multiple clusters with more equal degrees.

→ Numerical video link



- K-mean clustering
- fuzzy C-mean

### Numerical steps

- ① Given data points, based on number of clusters required initialize the membership table with random value

Note: But after all the  $\Sigma$  membership probability of each datapoint must be 1

## \* Associate Rule learning

Associate Rule learning is a machine learning technique used to discover interesting relationship or associations among a set of variables in large dataset.

The relationships are often expressed in the form of rules where one set of items implies another set of items. Some important concepts & algorithms related to association rule are

- 1. support
- 2. confidence
- 3. lift
- 4. conviction

Rules

Association Rules

1. Shoes -> socks

2. socks -> shoes

Shoes -> socks

⇒ support : Support is the proportion of transaction in the dataset that contain a particular set of item.

Confidence is the measure of reliability of a rule.

Ques: Min support = 50%.

Threshold confidence = 70%.

Eg. Tid Items

100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Itemset

Support

$$1 \quad 2/4 \rightarrow 50\%$$

$$2 \quad 3/4 \rightarrow 75\%$$

$$3 \quad 3/4 \rightarrow 75\%$$

$$4 \quad 1/4 \rightarrow 25\%$$

$$5 \quad 3/4 \rightarrow 75\%$$

X

Remove all the items from itemset having support less than min support

so new itemset is

$$\{1, 2, 3, 5\}$$

(1)

itemset

support

X

(1, 2)

 $1/4 \rightarrow 25\%$ 

X

(1, 3)

 $2/4 \rightarrow 50\%$ 

(1, 5)

 $1/4 \rightarrow 25\%$ 

(2, 3)

 $2/4 \rightarrow 50\%$ 

(2, 5)

 $3/4 \rightarrow 75\%$ 

(3, 5)

 $2/4 \rightarrow 50\%$ 

again remove all item having support less than min support

(1, 2, 3, 5) triplet

X

(1, 2, 3)

 $1/4 = 25\%$ 

(2, 3, 5)

 $2/4 = 50\%$ 

X

(1, 2, 5)

 $1/4 = 25\%$ 

{2, 3, 5}

Now building Association Rules

Rule

Support | Conf.

 $(2 \wedge 3) \rightarrow 5$ 

2

100%

 $(2 \wedge 5) \rightarrow 3$ 

2

66%

 $(3 \wedge 5) \rightarrow 2$ 

2

100%

not on OnePlus

$2 \rightarrow (3 \wedge 5)$	2	66%
$5 \rightarrow (2 \wedge 3)$	2	66%
$3 \rightarrow (2 \wedge 5)$	2	66%

$$\text{confidence} = \frac{s(A \cup B)}{s(A)}$$

Eg:  $A \quad B$   
 $(2 \wedge 3) \rightarrow 5$

$$= \frac{s(2 \wedge 3 \cup 5)}{s(2 \wedge 3)} = \frac{2}{2} = 100\%$$

Compare these Rules confidence  
with threshold confidence

those rules, with confidence  
 $\geq$  threshold confidence are  
 association rules.

heni  $(2 \wedge 3) \rightarrow 5$

$(3 \wedge 5) \rightarrow 2$

✓