

Unit - 4

Clustering and Applications :-

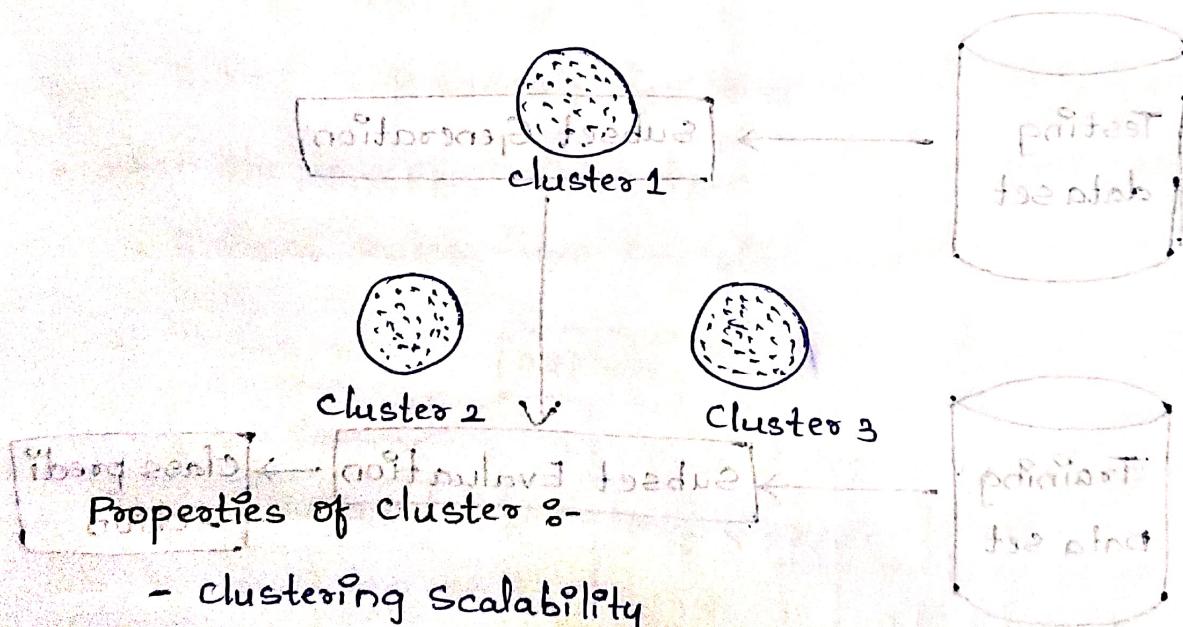
Cluster Analysis :-

Process of forming similar groups of objects together in form of a cluster.

The data items are clustered based on principles of maximising the intra class similarities and minimising the inter class similarities.

Ex:- In our class we are placing all the girls at one side and all the boys at one side. Girls are Girls cluster and boys are boys cluster.

- It is unsupervised Machine Learning algorithm.



- clustering scalability
- Algorithms usability with multiple types of data
- Dealing with Unstructured data
- Interoperability.

- clustering methods :-
- Partitioning methods
 - Hierarchical Methods
 - Density Based Methods
 - Grid based Methods.

Types of Data in cluster Analysis :-

- clustering analysis supports Data structures.
- Data structures are of two types. They are:

1. Data Matrix

2. Dissimilarity Matrix

1. Data Matrix :

Here, Data is represented as table (or) $n \times p$ matrix.

rows \rightarrow real world entities (names)

columns \rightarrow Properties of entities.

2. Dissimilarity Matrix :

- It is represented as $n \times n$ matrix.

- It identifies dissimilarities between two objects.

$$\begin{bmatrix} 0 & d(1,2) & d(1,3) \\ d(2,1) & 0 & d(2,3) \\ d(3,1) & d(3,2) & 0 \end{bmatrix}$$

Diagonal elements = 0
 Off-diagonal elements = dissimilarity

Types of Data :-

- There are four types of data. They are:

1. Interval Scaled data

2. Binary Variables

3. Categorical Variables

4. Mixed Variables

1. Interval scaled data :-

- It has continuous variables.

Ex: 10-20, 20-30, ... etc.

- To convert individual data into continuous variables.

- do the data standardisation before that.

- Data standardisation means removal of units.
- For standardisation data we should calculate mean absolute deviation.
- Then divide the data into intervals.

2. Binary Variables :-

- It has only 2 states

0 - Variable is absent

1 - Variable is present.

- It has 2 subtypes.

- Symmetric binary → states of variables can change

- Assymmetric binary → states of variables can't change

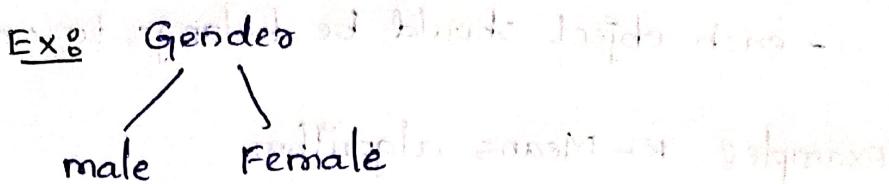
3. Categorical Variables :-

Data that can be divided into categories

- 2 types.

a) Nominal Variable :-

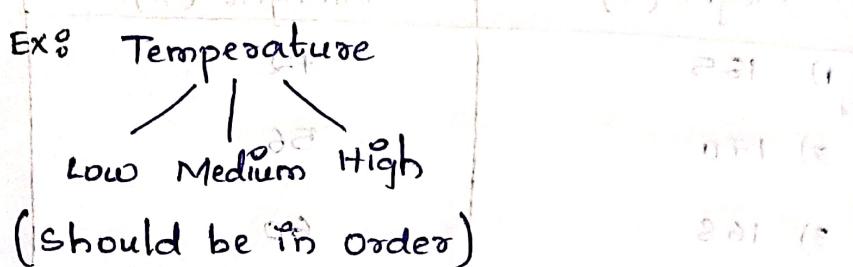
It has no particular Order to the categories.



(can be in any Order)

b) Ordinal Variable :-

It has particular Internal Order to the categories



(Should be in Order)

4. Mixed Variables :-

Combination of different types of Variables.

Categorization of Major clustering Methods :-

The clustering methods are categorized into Four types. They are :-

1. Partitioning Method
2. Hierarchical Method
3. Density-based Method
4. Grid-based Method.

1. Partitioning Methods :-

- The data is divided into partitions.
- Partition represents a cluster.
- Each cluster should be less than (\approx) equal to total data items.

$$K \leq n$$

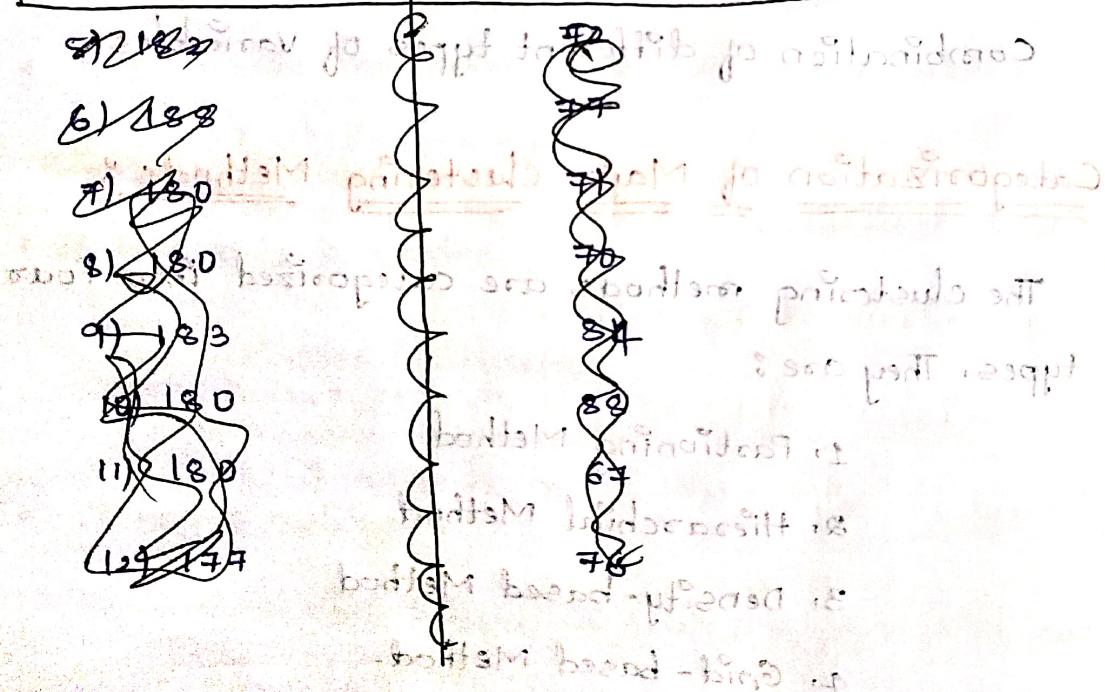
- Partition should satisfy rules.

- each partition should have at least one object.
- each object should belong to only one partition.

Example: K-Means algorithm:

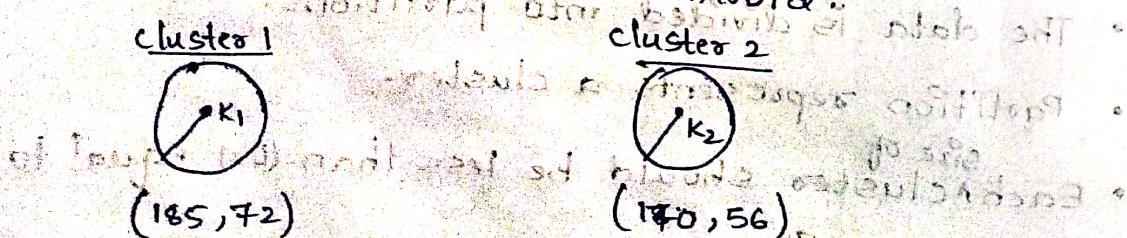
- data is divided into clusters based on distance and centroid values.

Height (x)	Weight (y)
1) 185	72
2) 170	56
3) 168	60
4) 174	68



→ Divide the data into two clusters.

Each cluster should have centroid.



We can take any order pair from example.

Now, remaining values should be divided based on

Euclidean distance (ED)

$$ED = \sqrt{(x_o - x_c)^2 + (y_o - y_c)^2}$$

For ③, $K_1 = \sqrt{(168 - 185)^2 + (60 - 72)^2} = 20.8$

$$K_2 = \sqrt{(168 - 170)^2 + (60 - 56)^2} = 4.48$$

$$K_2 < K_1$$

It belongs to cluster 2.

Now, calculate the new centroid

$$(170, 56) \text{ and } (168, 60)$$

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

For ④, $K_1 = \sqrt{(174 - 185)^2 + (68 - 72)^2} = 6.32$

$$K_2 = \sqrt{(174 - 169)^2 + (68 - 58)^2} = 14.14$$

$$K_1 < K_2 \text{ and it has } ④ \text{ than } ③ \text{ longer to}$$

∴ It belongs to cluster 1.

calculate centroid for $(185, 72)$ and $(174, 68)$

$$\left(\frac{185+174}{2}, \frac{72+68}{2} \right) = (182, 70)$$

$$K_1 = \{1, 4\}$$

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

∴ The data is divided into two different clusters.

2. Hierarchical Methods :-

- It groups the data into tree of clusters.
- The structure is called dendrogram.
- Dendograms have sequences of all merges and splits.
- Hierarchical method has two sub methods

1. Agglomerative Method

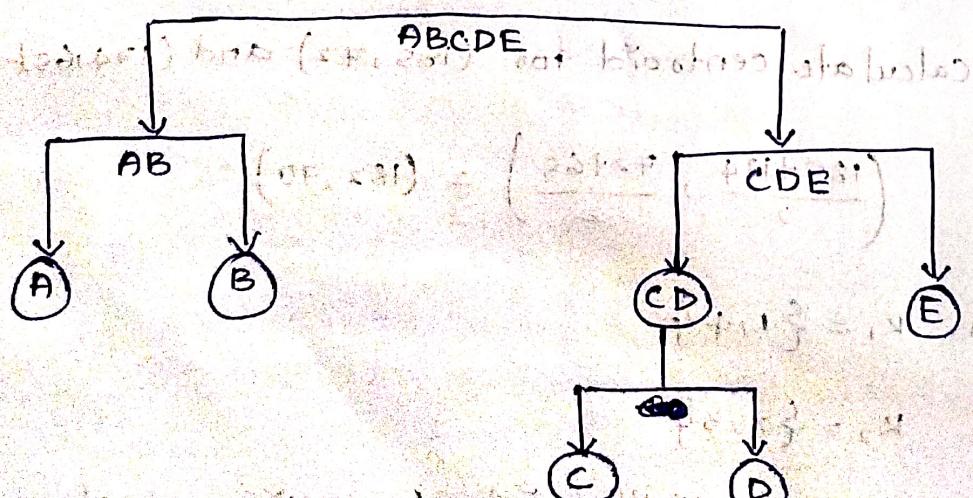
2. Divisive Method.

1. Agglomerative Method :-

- It is bottom-up method.

Steps Involved :-

- 1) calculate the similarity of one cluster with respect to all other clusters.
- 2) consider every data point as individual data.
- 3) Merge the clusters with highest similarity.
- 4) Recalculate similarity for each cluster.
- 5) Repeat ③ and ④ until the single cluster is obtained.

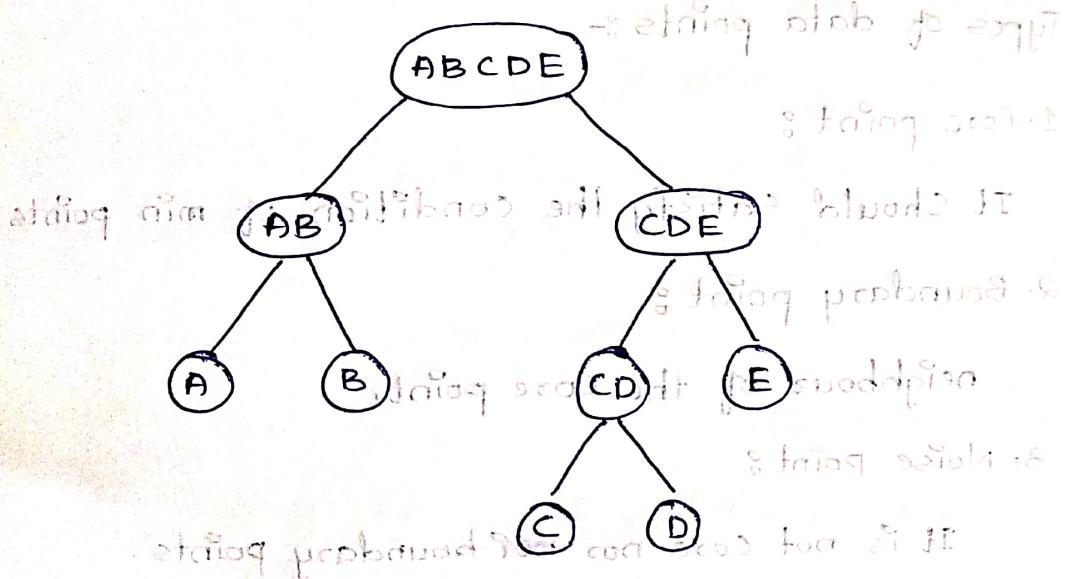


. There are 3 modes

- single linkage - min
- complete linkage - max
- Average linkage - Avg.

2. Divisive Method :-

- It is top-bottom method.
- We take all the data items into single cluster and in iterations, we split the data.
- At the end, we get N clusters.



3. Density-Based Methods :-

- Data objects are clustered "based on density".
- Density means mass/volume

Example :- DBSCAN

(Density Based spatial clustering of Applications with Noise).

- It has 2 inputs (~~2~~ and ~~1~~)

→ ϵ - Epsilon

→ minimum points.

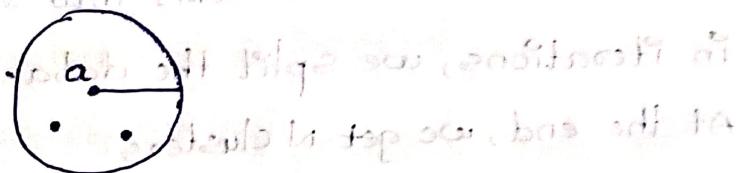
Where,

ϵ - radius of circle formed with data objects as centre.

min points - minimum no. of datapoints inside the circle

Ex: min pts = 3

Two objects inside circle with radius ϵ .



Types of data points:-

1. Core point :-

It should satisfy the condition of min points.

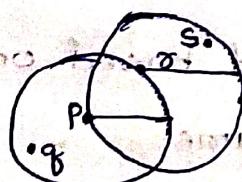
2. Boundary point :-

neighbour of the core point.

3. Noise point :-

It is not core nor boundary points.

Ex:-



$p, r \rightarrow$ core

$q, s \rightarrow$ boundary

$t \rightarrow$ noise.

Object = 3.

4. Grid-Based Methods :-

- It uses a multi resolution of grid data structure.
- It divides the objects into finite no. of cells that form a grid like structure.

- then density is calculated for these cells
- Sort the Cells according to density.
- Identify cluster centers

→ Update neighbour cells.

- Grid Based clustering is quick processing time.

Example: STING

(statistical Information Grid clustering Algorithm).

- Spatial data divided into rectangular cells at different levels of resolution, these cells forms a tree structure.
- cells at higher levels contains small cells composed to lower levels
- clustering done based on parameters
↓
mean, count, min, max.
- calculations of these parameters should start from root and go down till bottom layer.

Outlier Analysis :-

- Among the data items in Data Base, there may be some items which do not follow the general behavior of data.
- Those data items are called outliers.
- Analysis of outliers is known as Outlier Analysis.

Outlier Detection :-

The process of identifying outliers and subsequently removing them.

These are two methods. They are :-

1. Statistical Approach.
(and finding probability distribution)
2. Proximity Approach.

1. Statistical Approach :-

- It is based on probability of data points.

- Lowest probability is considered as Outlier.

- parametric method.

- Non Parametric method.

2. Proximity Approach :-

- It is based on location of data points :-

- Density based Approach

- Distance based Approach

- Grid based Approach

- Deviation based Approach.

Types of Outliers :-

There are three types of Outlier. They are :

1. Global / point Outlier.

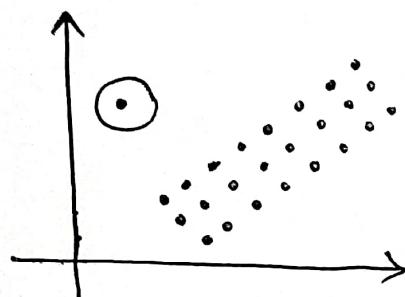
2. Collective outliers

3. Conditional outliers.

1. Global / point Outlier :

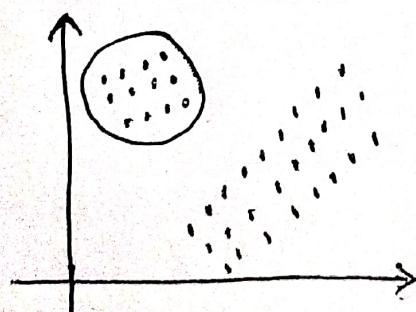
When a single data objects deviates from the rest of data.

points \rightarrow Global / point outliers.



2. collective outliers :

When a group of data objects deviates from the rest of data.



3. conditional outliers :

Data Objects deviates from others because of the specific condition.

