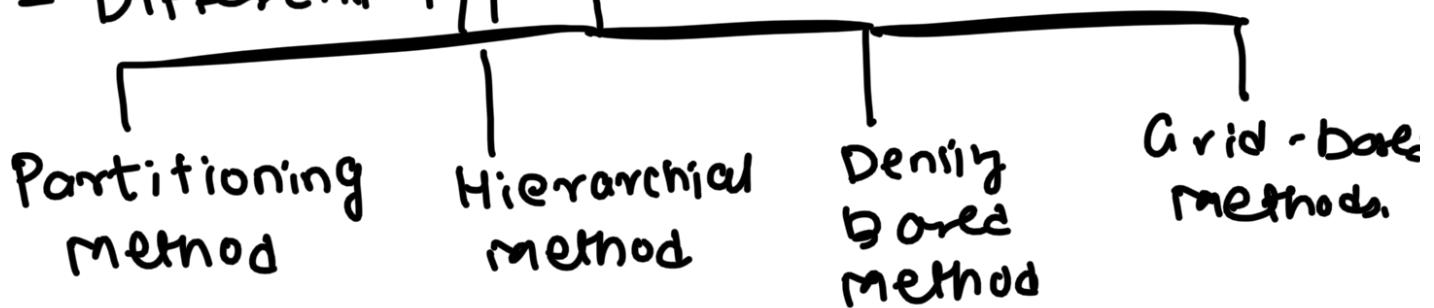


Units

- Introduction to cluster analysis.

- Different types of clustering methods:



- Partitioning Methods

Global optimal

Heuristic

k-means

k-medoid

(PAM)

- Clustering High Dim. Data

Need

Subspace

clustering

Dimension

ality reduc-

ctor

- clustering with constraints.

Constraints
on instances

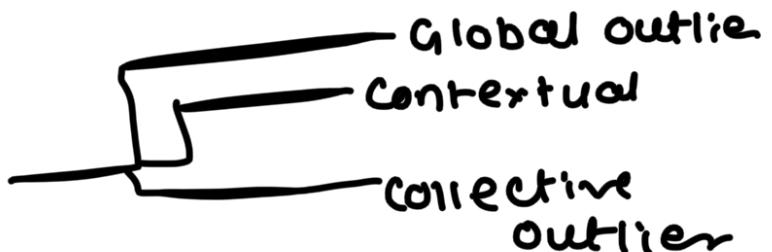
Constraints on
clusters

Constraints on
similarity measure-
ment.

- Outlier Analysis

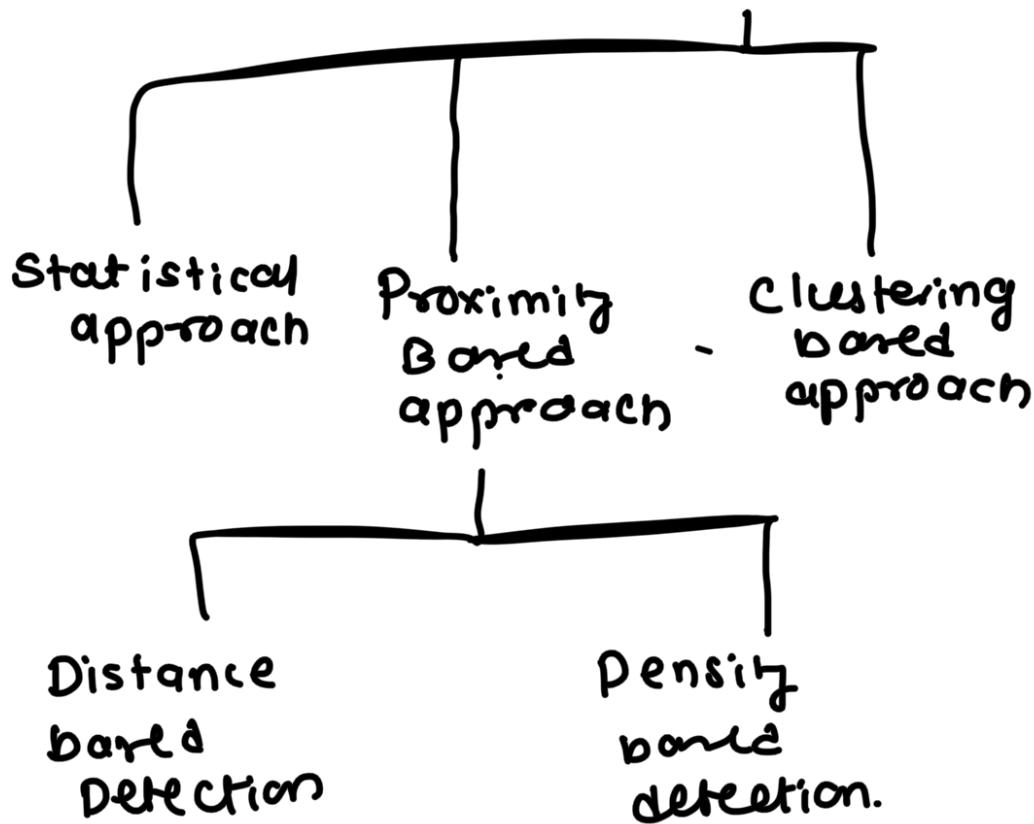
- Applications

- Types of outliers



- Need/challenges of Outlier Detection.

- Outliers Detection Methods



- Outlier Detection in High Dimensional Data.



• Clustering:

- unsupervised learning
- process of partitioning set of data objects into subsets.
- Cluster has similar data.
- Set of clusters resulting from cluster analysis is called clustering.

• Importance:

- widely used in many applications like business intelligence, image pattern recognition, web search, biology.
- Used in outlier detection.
- Business intelligence → identify customer, group them
- web search → organize the results into groups.
- helps segregate population

• TUTOR:

- **Types:**
 - 1. Partitioning methods
 - 1. Distance-based.
 - 2. Use mean or medoid to represent cluster centre
 - 3. Effective for small and medium sized data.
 - 2. Hierarchical methods
 - 1. Decomposed hierarchically
 - 2. Incorporate other tecq. like micro clustering.
 - 3. Density-based methods
 - 1. Clusters formed are dense regions
 - 2. Can find arbitrarily shaped clusters
 - 4. Grid-based methods
 - 1. Use multidimensional grid
 - 2. Fast processing Time

• **Partitioning Method:** Partition data into k clusters.



• **Clustering High Dimensional Data:**

more than 10 attributes \rightarrow High Dimensional data

Need: Example of supermarket \rightarrow customer purchase data is high dimensional \rightarrow Traditional distance methods can't be used (Euclidean or Manhattan distance).

1. **Subspace clustering approach**

- Subspace - subset of attributes in full space.
- $\dots \dots \dots$ existing subspaces of given high

- searching clusters in multidimensional data.

2. Dimensional reduction approach

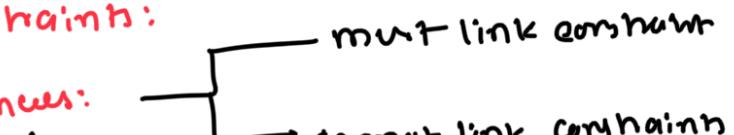
- tries to construct lower dimensional space and search for cluster

* Clustering with Constraints:

1. Constraints on instances:

- specify how pair of instances should be grouped.

must link



- constraint on a and b, then a and b in some cluster
- Transitive in nature
 $\text{must-link}(a,b)$ and $\text{must-link}(b,c)$
 then $\text{must-link}(a,c)$.

cannot link

- If cannot-link(a,b) then a and b have diff. clusters.

2. Constraints on clusters:

specifies

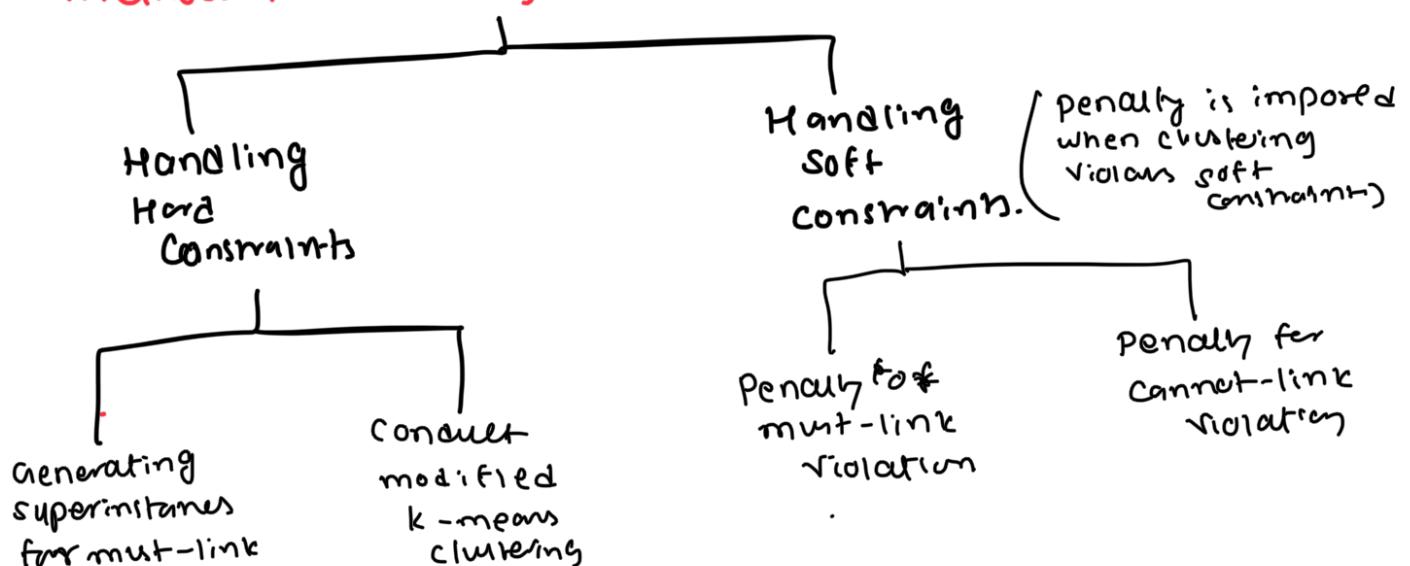
requirement on cluster, using attributes

3. Constraints on similarity measure:

measures

specifies similarity measure
 e.g. Euclidean distance must be respected.

methods for clustering with constraints:



1. Generating superinstances for must-link:

1. Transitive closure is computed on must-link constraints.
 - multiple subsets of objects, all objects

2. This closing gives multiple clusters -
in subset must be assigned in one cluster.
3. Replace all these objects by mean
4. weight of superinstance , is the no. of objects in it present
5. Must-link constraint is satisfied.

2. Conduct modified k-means clustering:

- we modify the center assignment process of k-means to nearest feasible center assignment.
- when objects are assigned to center in sequence , make sure the assignments do not violate cannot-link.
- Therefor, respects all cannot link.

3. Penalty of must-link violation: The must-link constraint (a,b) are present in diff cluster then the $\text{dist}(a,b)$ is added to objective fun. as penalty.

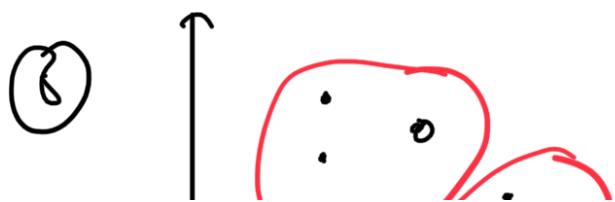
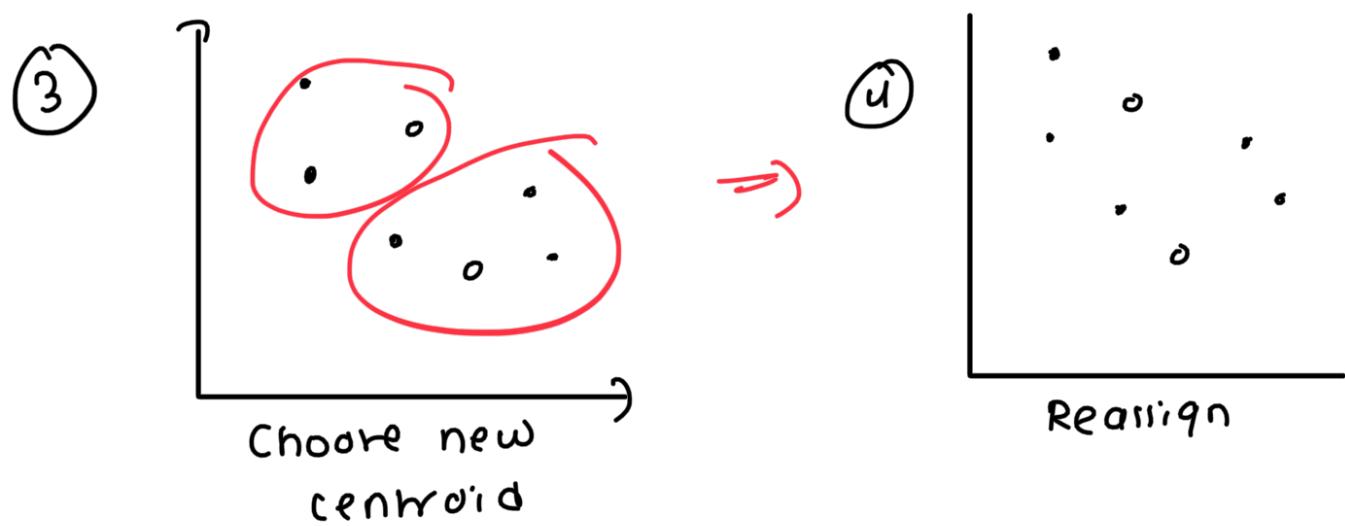
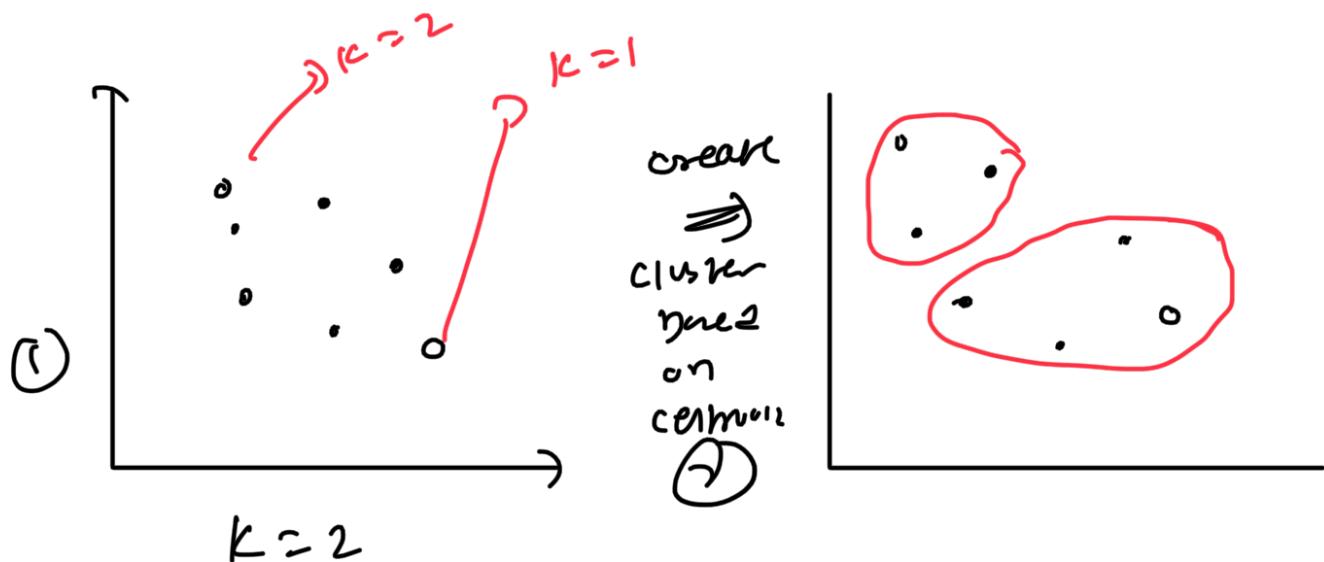
4. Penalty of cannot-link violation: The cannot-link constraint (a,b) are present in same cluster then $\text{dist}(a,b)$ is added to objective func. as penalty.

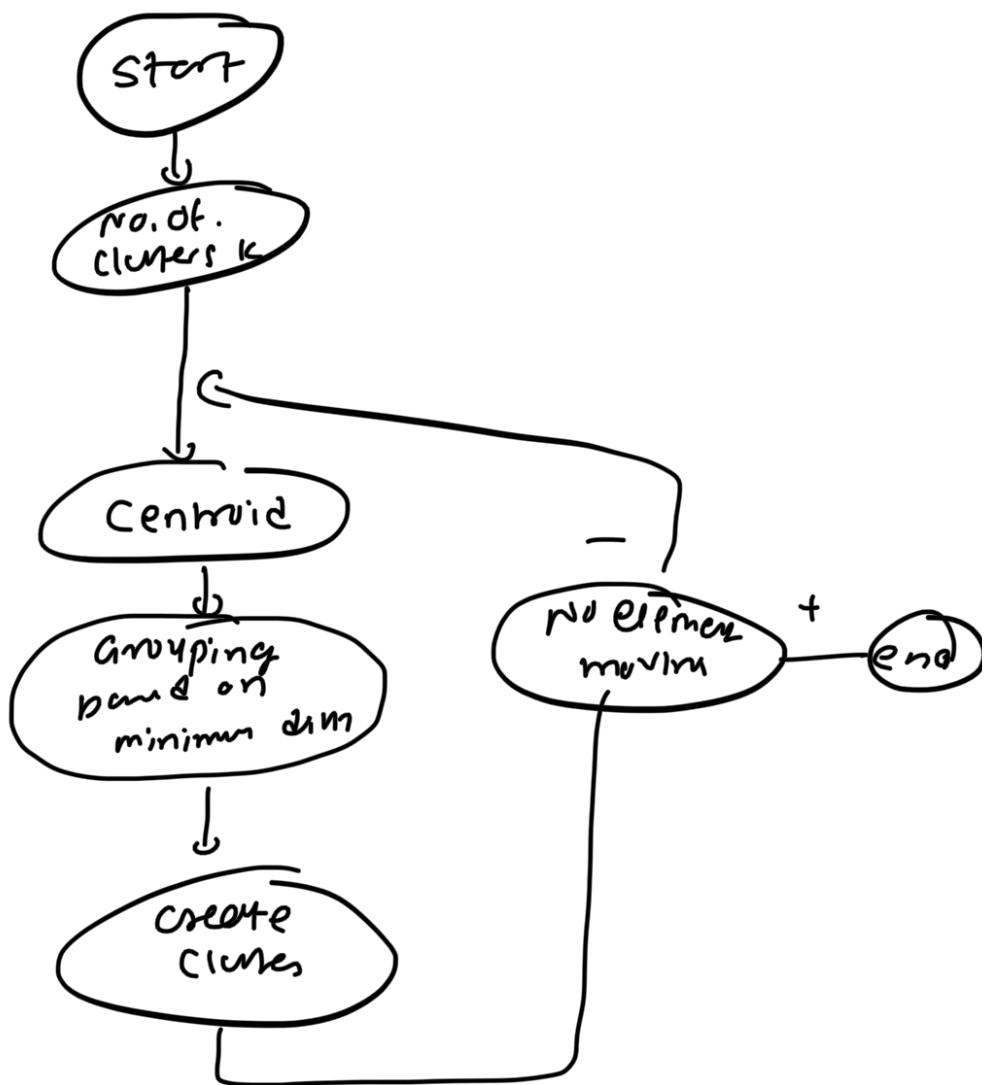
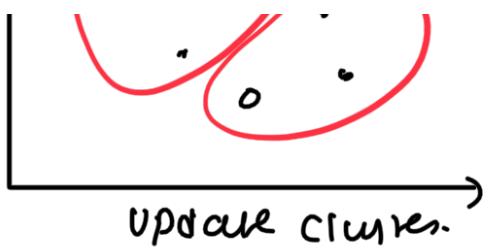
* K-mean :

- unsupervised learning algo.
- In k-mean, the data object are classified based on their attributes into k number of clusters.
- Algo:
 1. Define k centroids for k clusters which are generally far away.
 2. Group elements into clusters , which are nearer to centroid of cluster
 3. Now calculate new centroid for each

cluster.

4. Again group elements based on new centroid.
5. In every step, the centroid changes the elements move from one to another cluster.
6. Do same process till no element is moving from one to another.





K-medoids:

- Take most centrally located object in a cluster.
- PAM.
- Handles outliers well.
- Ordering of input does not impact results.
- Same ...

Each cluster \Rightarrow medoid.

Algo:

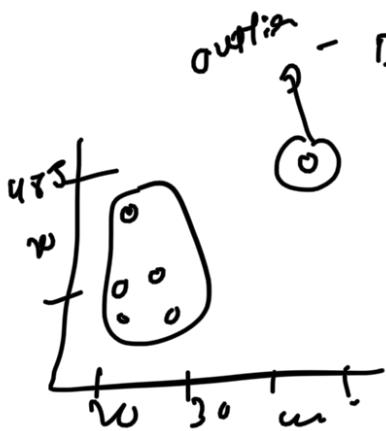
- Select k points as medoids.
- Assign all points to closest medoid.
- See for better medoids.
- Repeat until no medoids change.
- Medoids are chosen if overall cost is important.

Advantage:

- (1) effective on small datasets.
- (2) PAM robust than k-mean

* Outlier Analysis:

- observation or subset of observation diff from remaining set of data
- unusually large or small value.
- rise due to instrumental error, Human error.



outlier - Detection is imp:

1. Human error outlier needs to be deleted or corrected.
2. Outlier does not mean bad data, statistical analysis can be used to analyze.

- Applications:

1. Fraud detection.

- 2. Loan processing
- 3. Intrusion detection
- 4. Unexpected entries in database

- Types:

- 1. Global outlier : - simplest outlier to detect
- when data obj diff from rest of given data.
- 2. Contextual outlier :-
- Time series data.
- data obj anomalous within its context.
- winter season temp rises to 36°.
- 3. Collective outliers :-
- A subset of data obj outlying to entire data set.
- individual data obj is not outlier, but their collection
- Human electrocardiogram

- Detection methods.

- 1) supervised method :
 - Training data
 - normal class
 - outlier class
 - Data compare and find out which class it belongs to
 - limitation is accurately labelled training data

- 2) semi-Supervised method :

↓ ↓ ↓ ↓ ↓

- only one type of class
- If data does not fit the one model, assumed to belong to another class.

3) Unsupervised - method:

- uses notion of outliers and then makes use of same to detect outliers.
- Three methods
 - Statistical Based
 - Proximity Based
 - Clustering Based

1. Statistical Based:



Z-Score Interquartile range



calculate



- divide data into quartiles

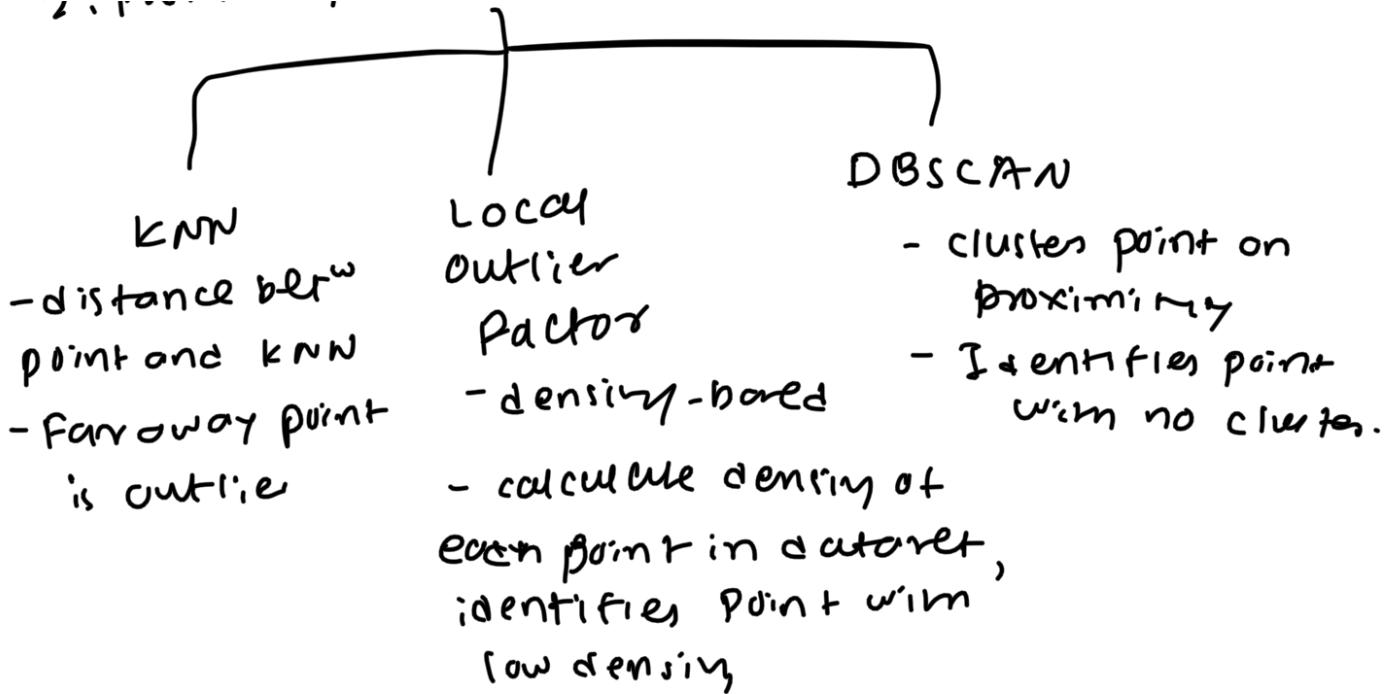
- Z-score of each point

- identifies points that are more than a certain number of IQR,

- Z-score: no. of standard deviations

that a point is from mean

~ Proximity Based:



3. clustering Based



* Outlier Detection in High Dimensional Data:

- Dimensionality increase, noise increases.
- Conventional methods cannot be used.

1) Data sparsity. — {
 - high-D often sparse
 - nearly dominated by noise.

2) Data subspace —
 - Adaptive to subspaces
 signifying the data
 - capturing local behavior of data