

Cluster Analysis groups data objects based only on information found in the data that describes the objects & their relationships.

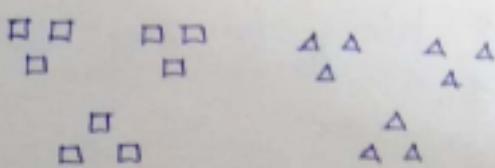
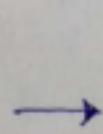
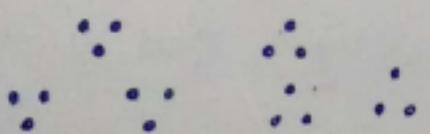
→ The goal is that the objects within a group be similar to one another & different from the objects in other groups.

The greater the similarity within a group & the greater the difference b/w groups the better are more distinct the clustering.

#### \* Different types of clustering :-

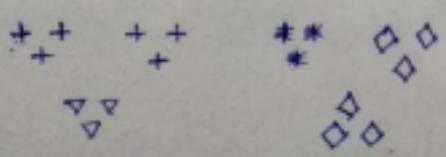
The various types of clustering are

\* Hierarchical vs partitional :- A "partitional clustering" is simply a division of the set of data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset.

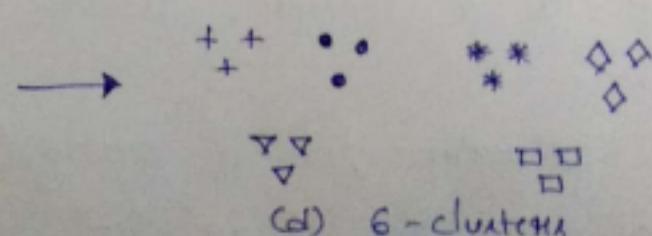


(a) Original points

(b) 2-clusters



(c) 4-clusters



(d) 6-clusters

→ "Hierarchical clustering" is a set of nested clusters that



are organized as a tree.

Each node (cluster) in the tree is a union of its children (Sub clusters) & root of the tree is the cluster containing all objects but leaf of tree are single clusters of individual data objects.

\* Exclusive vs overlapping vs Fuzzy :- In "Exclusive clustering", we assign each object to a single cluster. There are many situations in which a point could reasonably be placed in more than one cluster & these situations are better addressed by Non-Exclusive clustering.

→ "Non-Exclusive" or "Overlapping" clustering is used to reflect the fact that an object can simultaneously belong to more than one group (class).

For ex., A person at a University can be both an enrolled student & employee of University.

→ "Fuzzy" clustering, where every object belongs to every cluster with membership weight that is between "0 & 1" "0" (Absolutely doesn't belong) & "1" (Absolutely belongs). Clusters are treated as fuzzy sets in one in which an object belongs to any set with a weight that is between 0 & 1.

\* Complete vs partial :- A complete clustering assigns every object to a cluster, whereas a partial clustering does not.

→ The motivation for a partial clustering is that some objects in a data set may not belong to well-defined - classes.



→ For Example, Some Newspaper Stories May have a Common theme , while other Stories are More Generic or "One-of-a-kind".

#### \* Different Types of Clusters :-

Clustering Aim to find useful Groups of objects (Clusters), where usefulness be defined by the needs of Data Analysis.

→ There are Several different Notations of clusters that prove useful In practice

\* Well-Separated :- The data objects within a Cluster Must have Small distance & distance b/w two clusters Must be high.

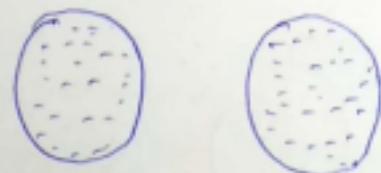


fig: Well-Separated Clusters

Sometimes a threshold is used to Specify that all the objects in a cluster Must be Sufficiently close (similar) to one Another.

\* prototype-Based Clusters (Center-Based) :- A cluster is a Set of objects in which Each object is closer (similar) to the prototype than to the prototype of Any other cluster.

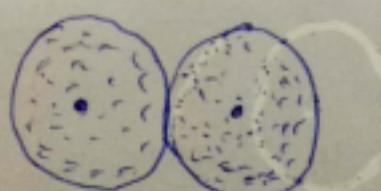


fig: proto-type-Based Clusters



\* Graph - Based Clusters (Contiguity - Based) :- If the area is represented as graph, where the nodes are objects & the links represent connection among objects then a cluster can be defined as Connected Component.

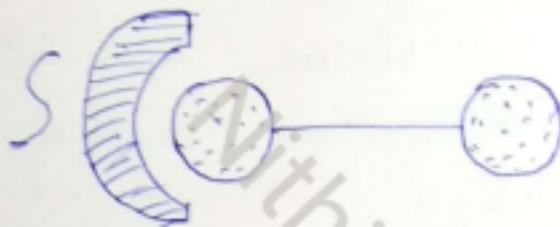


fig: Graph - Based Cluster

\* Density - Based Clusters :- A cluster is a dense region of objects that is surrounded by a region of low density. A density based defn of cluster is often employed when the clusters are irregular or intertwined & when noise & outliers are present.

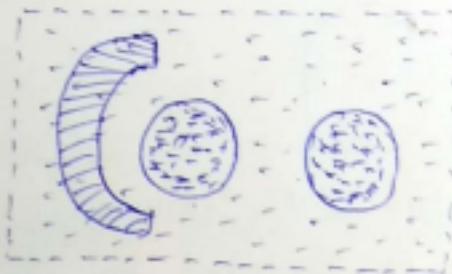


fig: Density - Based cluster

\* Shared - property clusters (Conceptual) :- More Generally, A cluster is a set of objects that share some property.

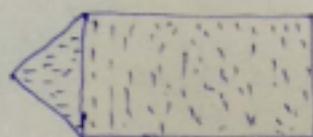


fig: Shared - property cluster

A clustering algo would need a very specific concept of a cluster to successfully detect these clusters, the process of finding such clusters is called "conceptual clustering".



### \* $k$ -Means :-

K-Means Clustering intends to partition "n" objects into  $k$ -clusters in which each object belongs to the cluster with the nearest mean.

This Method produces Exactly " $k$ " different clusters of greatest possible distinction.

→ It is a prototype-based clustering technique. Create a one-level partitioning of data objects.

There are a no of such techniques but two of most prominent are " $k$ -Means" & " $k$ -Mediod".

→ The objective of  $k$ -Means Clustering is to Minimize total "Intra-cluster Variance" or "Squared Error function".

→ In  $k$ -means, we first choose " $k$ " Initial Centroids, where " $k$ " is a user-Specified parameter, namely the no of clusters defined.

Each point is then assigned to closest Centroid & Each collection of points assigned to a centroid is a cluster.

→ The Centroid of each cluster is then updated based on the points assigned to the cluster.

We repeat the assignment & update steps until no point change clusters or equivalently until the Centroids remain the same.

### \* $k$ -Means Algorithm :-

1. Select  $k$  points as Initial Centroids

2. Repeat

3. Form  $k$  clusters by assigning each point to its closest centroid.

4. Recompute, the Centroid of each cluster



5. until centroid do not change.

k-Means Algorithm was stated generally as "recompute the Centroid of Each cluster".

→ "Data In Euclidean Space", Euclidean distance is one which measures the quality of a clustering. Using "Sum of Squared Error (SSE)" which is also known as "Scatter".

In other words, we calculate the sum of each data point i.e. Euclidean distance to the closest centroid & then compute the total sum of squared error.

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}(c_i, x)^2$$

where "dist" is the standard Euclidean distance b/w two objects in Euclidean space.

→ The Centroid (Mean) of  $i^{th}$  cluster is defined by Eqn:

$$c_i = \frac{1}{m_i} \sum_{x \in C_i} x$$

→ "Document Data", our objective in k-Means is to maximize the similarity of the documents in a cluster to cluster centroid.

The quality is known as "cohesion" of clusters.

$$\text{The total cohesion is given by} = \sum_{i=1}^k \sum_{x \in C_i} \text{similarity}(x, c_i)$$

\* k-Means Additional Point :- There are two strategies that decrease the total SSE by increasing the no. of clusters.

\* "Split a cluster" → Cluster with the largest SSE is usually chosen, but we could split the cluster with large deficit



-For for one particular attribute.

\* "Introduce a new cluster' Centroid"  $\rightarrow$  often the point that is farthest from any cluster will be chosen.

Another Approach is to choose randomly from all points or from the points with highest SSE.

There are two Strategies that decrease the no of clusters, while trying to minimize the total SSE are

\* "Divide a cluster"  $\rightarrow$  this is accomplished by removing the centroid that corresponds to cluster & reassigning the point to other clusters.

\* "Merge two clusters"  $\rightarrow$  The clusters with closest centroids are typically chosen, although another Approach to Merge the two clusters that result in small increase in total SSE.

\* Bisection k-Means - Algorithm is a straight forward extension of basic k-Means Algo that is based on simple Idea to obtain k clusters.

$\rightarrow$  Split the set of all points into two clusters, Select one of these clusters to split & so on. until k-clusters have been produced.

#### Algorithm

1. Initialize the set of clusters to contain the clusters consisting of all points.

2. Repeat

3. Remove a cluster from the list of clusters

4. for  $i = 1$  to "number of trials" do

5. Bisection the Selected cluster using basic k-Means

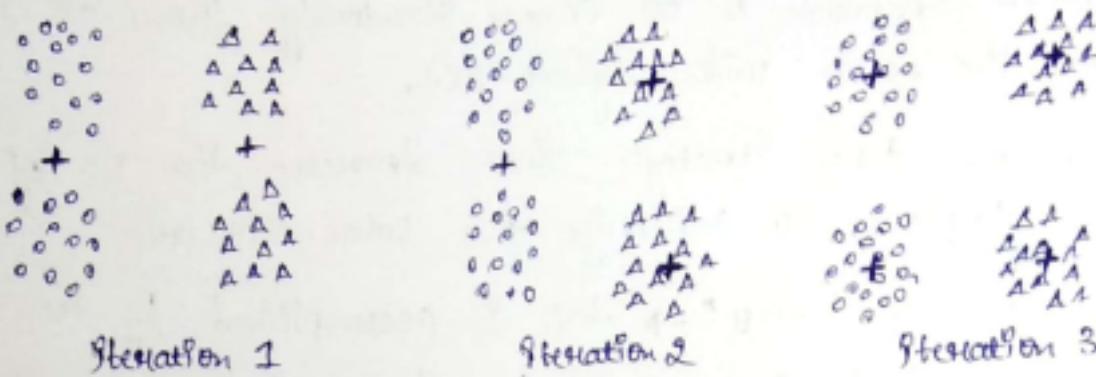
6. End for

7. Select the two clusters from Bisection with lowest total SSE



8. Add these clusters to the list of clusters  
 9. until the list of clusters contains K clusters.

→ Consider the foll example of Bisection K-Means



→ The "Time Complexity" & "Space Complexity" of k-Means Algo is

\* Space Complexity  $\rightarrow O(Cm + k)n$

\* Time Complexity  $\rightarrow O(I * k * m * n)$

$I = \text{no. of iterations}$     $m = \text{no. of points}$     $n = \text{no. of attributes}$ .

- \* Strengths & Weakness- k-Means is simple & can be used for a wide variety of data types. It is also quite efficient even though multiple runs are often performed.
- \* k-Means are even more efficient & less susceptible to "Simplification problem".
- \* k-Means is not suitable for all types of data.
- \* k-Means also have trouble in clustering data that contain outliers.
- \* k-Means is restricted to data for which there is a notion of center (centroid).

\* k-Means Example problem- Consider the foll data set

$$k = \{2, 3, 4, 10, 11, 12, 20, 25, 30\}$$

& No. of clusters to formed = 2 i.e.  $k = 2$ , Random Mean Value  $m_1$  &  $m_2 = 4.5, 12$



$$\Rightarrow k = \{2, 3, 4, 10, 11, 12, 20, 25, 30\} \quad k=2 \quad M_1 = 4 \quad M_2 = 12$$

Assign Each point to 9th cluster Centroid ( $M_1$  &  $M_2$ ) & form  $k=2$  clusters

$$\therefore k_1 = \{2, 3, 4\}$$

$$k_2 = \{10, 11, 12, 20, 25, 30\}$$

$$M_1 = \frac{2+3+4}{3}$$

$$M_2 = \frac{10+11+12+20+25+30}{6}$$

$$\boxed{\text{new } M_1 = 3}$$

$$\boxed{\text{new } M_2 = 18}$$

→ Again Assign Each point to 9th cluster new Centroid ( $M_1$  &  $M_2$ ) & form  $k=2$  clusters

$$M_1 = 3$$

$$M_2 = 18$$

$$\therefore k_1 = \{2, 3, 4, 10\}$$

$$k_2 = \{11, 12, 20, 25, 30\}$$

$$M_1 = \frac{2+3+4+10}{4} = 4.75$$

$$M_2 = \frac{11+12+20+25+30}{5} = 19.6$$

$$\boxed{\text{new } M_1 = 5}$$

$$\boxed{\text{new } M_2 = 20}$$

→ Again Assign Each point to 9th cluster new Centroid ( $M_1$  &  $M_2$ ) & form  $k=2$  clusters

$$M_1 = 5$$

$$M_2 = 20$$

$$\therefore k_1 = \{2, 3, 4, 10, 11, 12\}$$

$$k_2 = \{20, 25, 30\}$$

$$M_1 = \frac{2+3+4+10+11+12}{6}$$

$$M_2 = \frac{20+25+30}{3}$$

$$\boxed{\text{new } M_1 = 7}$$

$$\boxed{\text{new } M_2 = 25}$$

→ Again Assign Each point to 9th cluster new Centroid ( $M_1$  &  $M_2$ ) & form  $k=2$  clusters

$$M_1 = 7$$

$$M_2 = 25$$

$$\therefore k_1 = \{2, 3, 4, 10, 11, 12\}$$

$$k_2 = \{20, 25, 30\}$$

$$M_1 = \frac{2+3+4+10+11+12}{6} = 7$$

$$M_2 = \frac{20+25+30}{3} = 25 //$$



Thus we are getting the same mean value we have to stop.

∴ Two well Separated Clusters of 'k' are

$$K_1 = \{2, 3, 4, 10, 11, 12\} \quad K_2 = \{20, 25, 30\}$$

#### \* Agglomerative Hierarchical Clustering :-

Hierarchical Clustering Tech

- These are a Second Important Category of Clustering Methods.

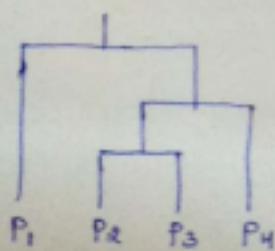
These Approaches are relatively old compared to Many Clustering Algorithms but they Still Enjoy Widespread Use.

→ There are two basic Approaches for generating a Hierarchical Clustering

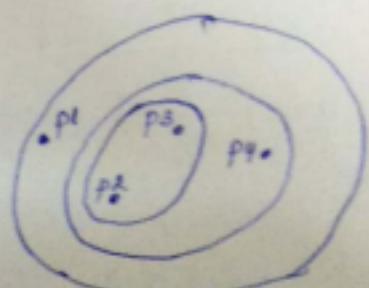
\* Agglomerative → Start with point as Individual Clusters & at Each Step, "Merge" the closest pair of Clusters.

\* Divisive → Start with one, all Individual cluster & at Each Step "Split" a Cluster until only Singleton cluster of Individual points remain.

→ A Hierarchical Clustering is often displayed graphically using a tree-like diagram called "dendrogram" which displays both Cluster-Subcluster Relationships & Order in which the Clusters were Merged or Split.



(Dendrogram)



(Nested-cluster Diagram)

A Hierarchical Clustering can also be Graphically represented



Using a "Nested cluster diagram".

- \* Agglomerative Hierarchical Clustering Algorithms: Starting with individual points as clusters, Successively Merge the two closest clusters until only one cluster remains.

The Algorithm is as follows:-

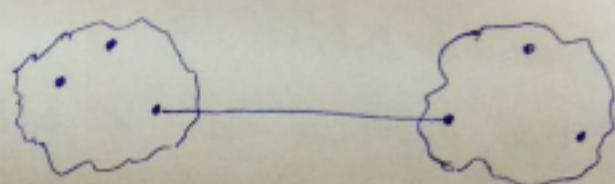
1. Compute the proximity Matrix, if necessary
2. Repeat
3. Merge the two closest clusters
4. Update the proximity Matrix to reflect the proximity between New cluster & the original cluster
5. Repeat until only one cluster remains.

- \* Defining proximity between clusters :- Cluster proximity is typically defined with a particular type of cluster in mind.

- The different proximity Measures in Agglomerative clustering are

- \* MIN
- \* MAX
- \* Group Average

- \* MIN → Defines cluster proximity as the proximity between the closest two points that are in different clusters or the shortest edge between two nodes in different subset of nodes.

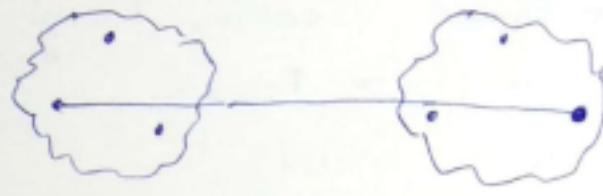


(MIN)

- \* MAX → Defines cluster proximity as the proximity between



the farthest two points that are in different clusters or  
the longest edge between two nodes in different subset of  
nodes.



(MAX)

- \* Group Average → Define cluster proximity to be the Average pairwise proximities of all pairs of points from different clusters.



(Group Average)

→ An alternative technique "Ward Method" also assumes that a cluster is represented by its centroid, but it measures proximity between two clusters in terms of "SSE" that result from merging the two clusters like "K-Means".

- \* Strength & Weakness :- of Agglomerative clustering are
  - \* Agglomerative Algorithm can produce Better-Quality Clusters.
  - \* Agglomerative Algorithm are expensive in terms of their computation & storage requirements.
  - \* All Merge can cause trouble of Noise in high-Dimension Data Such as Document Data.

Above Mentioned two problems can be addressed to some degree by first partially clustering the data using Another Technique Such as K-Means.



### \* DBSCAN :-

Density - Based clustering Locates regions of High Density that are Separated from one another by regions of Low - Density

→ DBSCAN Is a Simple & Effective Density - Based Clustering Algorithm that Illustrates a no. of Important Concepts that are Important for any Density - Based Clustering Approach.

\* Center - Based Approach :- There are Several distinct Methods but DBSCAN Is based on Center Based Approach

→ In the "Center - Based Approach", density Is Estimated for a particular point P in data set by counting the no. of points within Specified Radius "Eps" of that point (this includes point P itself).

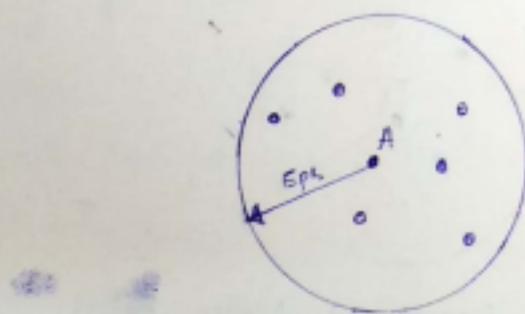


Figure Center - Based Approach

The no. of points within a Radius of Eps of point A Is "#".

\* Classification of points According to Center - Based Approach :- The Center - Based Approach to density allows us to classify a point into following three categories

- \* Core point
- \* Border point
- \* Noise or Background point

\* Core points → These points are In the Interior of density Based Cluster.



"A point  $P$  is a core point, if the no. of points within a given neighborhood around the point as determined by user specified parameter (Distance Function ( $Epsilon$ )) exceeds a certain threshold "Minpt" which is also user-specified parameter"

Ex:- point A is core point for radius  $Epsilon$  if  $Minpt \leq 7$

\* Border points  $\rightarrow$  A border point is not core point, but fall within the neighborhood of core point.

A border point can fall within neighborhoods of several core points.

\* Noise points  $\rightarrow$  A noise point is any point that is neither a core point nor a border point.

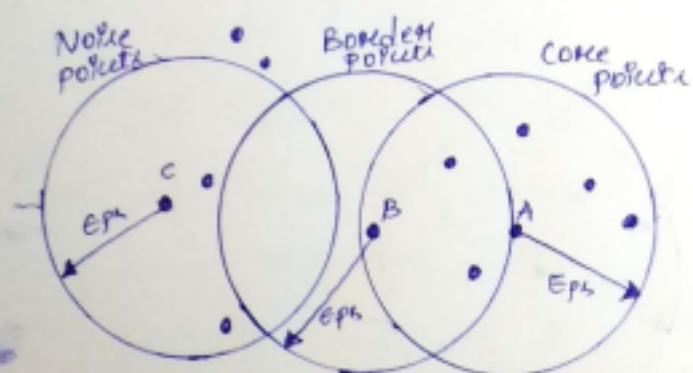


Fig:- Core, Border & Noise points

\* DBSCAN Algorithm - can be informally described as follows.  
Any two core points that are close enough within a distance  $Epsilon$  of one another are put in same cluster.

$\rightarrow$  Any border point that is close enough to a core point is put in same cluster as core point.  
Noise points are discarded.

$\rightarrow$  "DBSCAN Algorithm"

1. Label all points as core, border & noise
2. Eliminate noise points.



3. put an Edge between all core points that are within  $E_{pi}$
  4. Make Each group of connected core points to Separate Clusters.
  5. Assign Each Border point to one of clusters of its Associated core points.
- \* Strength & Weakness- DBSCAN uses a density-based defn of a cluster
  - \* It is relatively resistant to Noise & can handle clusters of arbitrary shapes & sizes
  - \* DBSCAN has trouble when the clusters have widely varying densities.
  - \* It also has trouble with high-dimensional data because density is more difficult to define for such data.
  - \* DBSCAN can be expensive when the computation of Nearest Neighbours requires computing all pairwise proximity numbers.
  - \* Cluster Evaluation - (Cluster Validation)

The Evaluation of the resulting classification Model is an integral part of the process of developing a classification Model & there are well-separated Measures & procedures ex "Accuracy & Cross-validation" etc.

→ "Cluster Evaluation" should be a part of any cluster Analysis, A key motivation is that almost every clustering Algo will find clusters in a data set, even if that data set has no natural cluster structure.













































































