

CLUSTER ANALYSIS.

- Cluster Analysis.
- Types of Clustering
- Types of clusters.
- Algorithm for cluster analysis.
 - (i) K-means algorithm.
 - (ii) Agglomerative hierarchical clustering.
 - (iii) DBSCAN

Cluster Analysis

Finding groups of objects such that the objects in a group will be similar to one another and different from the objects in other groups.

Types of Cluster Analysis Clustering

1. Hierarchical versus Partitional.
2. Exclusive versus Overlapping (Non-Exclusive)
3. Complete versus Partial

Types of clusters.

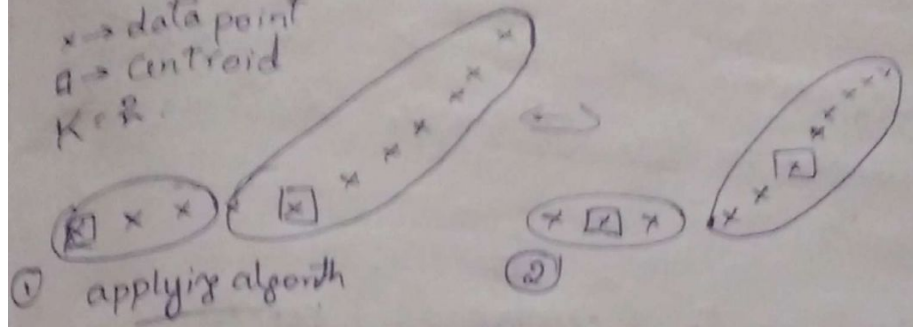
1. Well defined cluster
2. Prototype based (Centre based)
3. Graph based cluster
4. Density based
5. Shared cluster (conceptual clusters)

K-means algorithm

K - no of clusters
(family)

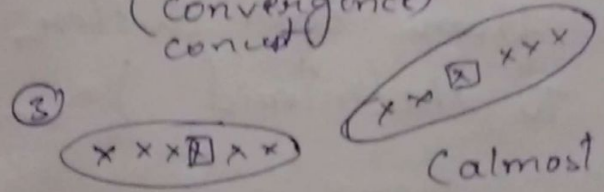
x → data point
a → centroid
K = 2

original data



1. select k points as initial centroid $\{c_1, c_2, \dots, c_k\}$
2. Repeat
3. Form k clusters by assigning each point to the closest centroid
4. Recompute the centroid for each cluster
5. Until centroids do not change

(convergence)
const



(almost all centroids are stable)

(selected random points (centroids) the distance b/w them should be maximum)

We are terminating the algorithm because all centroids stabilised and points don't move b/w the clusters.

for each point place it in the cluster whose current centroid is nearest.

After all points are assigned, update the location of centroids of k clusters.

Reassign all points to closest centroid

Repeat 2 and 3 until convergence.

(Partitional clustering)

What is right value of k_0 or
How to pick k -value

Try different k -values. By looking at the change in the average distance to centroid.

Initial k -points (centroids) picking.

approach 1:-

Sampling

- Cluster a sample of data using hierarchical clustering to obtain k clusters.
- Pick a point from each cluster (point nearest to the centroid)

approach 2:-

- Pick "dispersed" set of points.
- pick first ~~pick~~ point at random.
- pick the next point to be the one whose minimum distance from the selected points is as large as possible.
- Repeat until complete k -points

Final clustering depends on initial k -points (centroid \rightarrow how you pick)
Complexity.

1.
2.
3.
4.
5.

Complexity:

In each round we have to examine each input point exactly once to find the closest centroid.

Each round is $O(KN)$
until convergence (centroid is stabilised)

The no of rounds to reach the convergence can be very large

guls

Agglomerative Hierarchical clustering

A set of nested clusters organised as a hierarchical tree is called as hierarchical clustering.

Two types of algorithm.

1. Agglomerative (bottom-up)
2. Divisive (top down)

Agglomerative

- Initially each point in a cluster
- Repeatedly combine the two nearest clusters to one

Divisive

- Starts with one cluster and recursively split it.

Algorithm

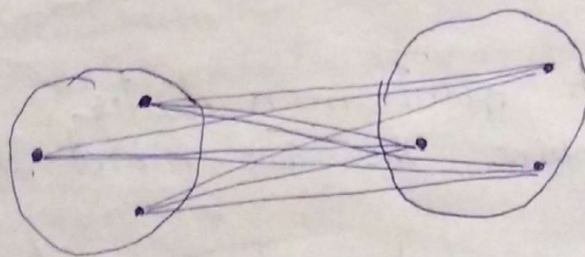
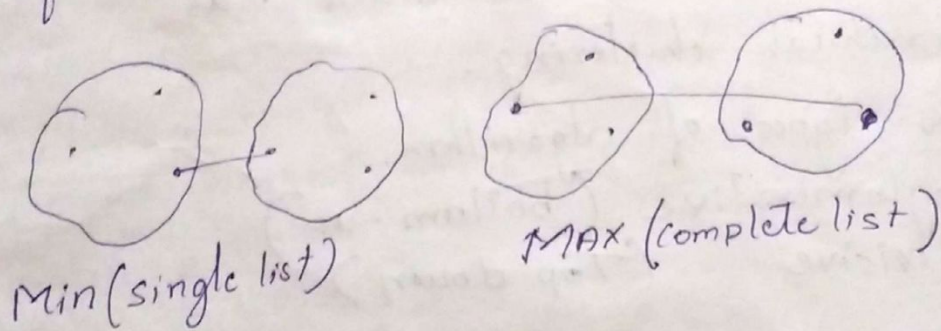
1. Compute the proximity matrix, if necessary.
2. Repeat
3. Merge the closest two clusters.
4. Update the proximity matrix to reflect the proximity b/w the new cluster and the original clusters.
5. Until only one cluster remains.

Dendrogram

The hierarchical clustering is often displayed graphically using a tree like diagram called Dendrogram, which displays both the cluster - sub cluster relationships and the order in which clusters are merged (agglomerative view) / divisive view).

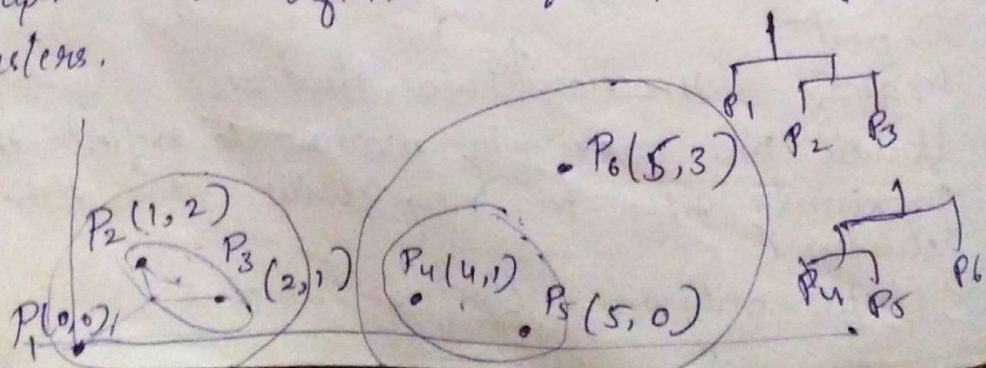
1/4/18

Define Proximity



Group average.

Graph based definitions for proximity clusters.



Agglomeration hierarchical clustering

1. How do you represent a cluster of more than one point?
2. How do you determine the measures of the clusters?
3. When to stop combining clusters?

Euclidean Space

Represent a each cluster by its centroid
= average of its point.

Fermi Condition.

Approach 1:- Pick a K number and stop when we reach K clusters.
(Naturally falls into K clusters)

Approach 2:- Stop when the next merge would create a cluster with low "cohesion".
centroid (clustroid)

Cohesion

Diameter - Maximum distance b/w the points in a cluster

Radius - Maximum distance of point from cluster.

Density - No of points per unit volume

Merits of DBSCAN

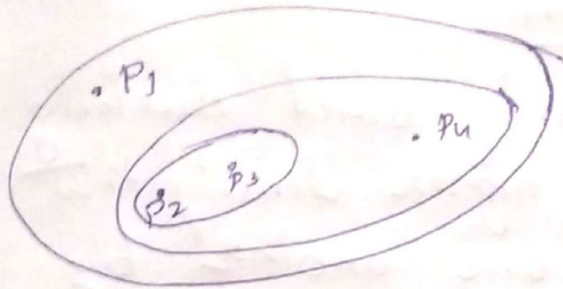
- Unlike K-mean, DBSCAN doesn't require user to specify the no of clusters to be generated.
- DBSCAN can find any shape of clusters need not be an circular in shape
- It produces efficient and accurate results.
- We can easily eliminate the noise points.
- DBSCAN algorithm is able to find arbitrary size and arbitrarily shaped clusters.

DRAWBACKS:-

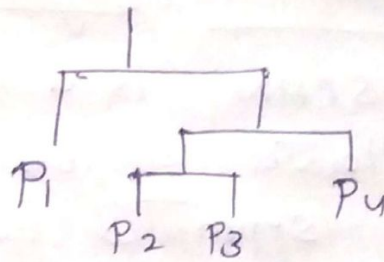
- DBSCAN algorithm fails in case of variable density clusters.
- Selecting area (EPS, ϵ) and minimum points (parameters for DBSCAN) are tricky.
- Doesnot work well in high dimensionality of data.

Strength & weakness of K-means & Agglomerative

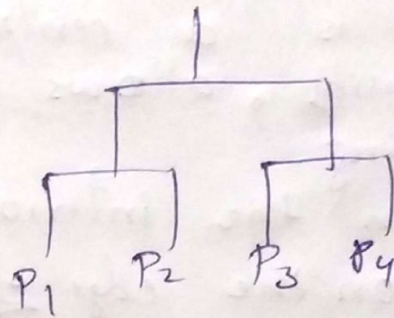
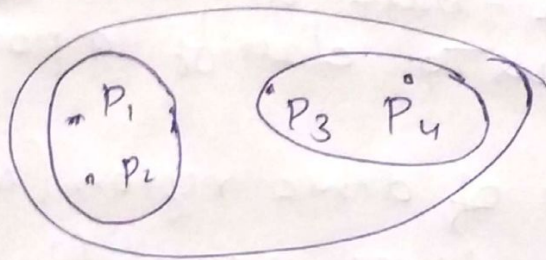
Traditional



Dendrogram



Non tradition



- Min — defines cluster proximity b/w the 2 progressive points or using (graph terms, the shortest edge b/w 2 nodes in different subset of nodes).
- Max: — takes the proximity b/w the farthest 2 points in different clusters to be the cluster proximity. longest edge b/w 2 nodes.
- Group avg — is the technique used to define the cluster proximity to be the average pair wise proximities of all the pairs of points from different clusters.

DBSCAN Algorithm

DBSCAN is a density based clustering. It locates regions of higher density that are separated from one another by regions of low density.

→ It is a centre-base approach to density, allows us to classify a point being:

1. In the interior of dense region (core)
2. On the edge of dense region (border)
3. In a sparsely occupied region (Noise)

Density = No of points per unit volume

Density based clustering that produces a partitioned clustering.

Main points Essentially given you some kind of threshold or how many you consider as being dense.
↳ EPS (area count) which you will perform the

Algorithm:-

1. Label all the points as core, border and noise points.
2. Eliminate the noise points.
3. Put an edge b/w all core points that are within ϵ of each other.
4. Make each group of connected group of core points into separate clusters.
5. Assign each border point to any no of clusters its associated core point.

DBSCAN (Density Based Spatial Clustering of Application with noise)
Graph based clustering