Unit - 5. Unsupervised learning

clustering — grouping of data points having Simularity
→ NO. training

Soft clustering → a data point can be part of 2 cluster
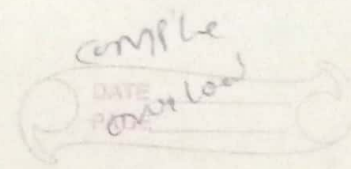or more than 2

hard clustering → only one cluster at a time.

K mean — hard clustering Algorithm

steps ① Decide K = no. of clusters
② Select K random data Points as centroids
③ find distance d b/w all and each points w'rt
each centroids.
$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$
④ Recompute centroid.
$$\left( x_c \, y_c = \frac{\sum x_i}{m} , \frac{\sum y_i}{m} \right)$$
⑤ repeat 3 & 4 until
① centroid do not change
② point remain in same cluster
③ Reached max. iteration.

Elbow Method

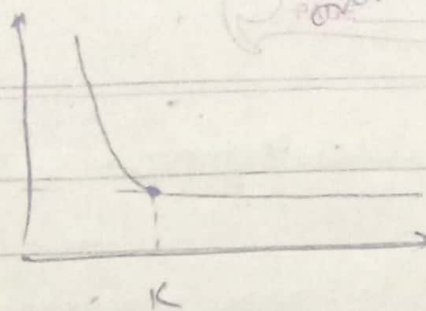graph b/w WCSS and no. of cluster K.

WCSS → Within cluster sum of square or
sum of the sq. distance b/w point in a cluster k

cluster centroid.

K medoids → partioning Algorithm.
the center of the subset is a member of the
subset called a medoid.

medoid are robust to outliers.

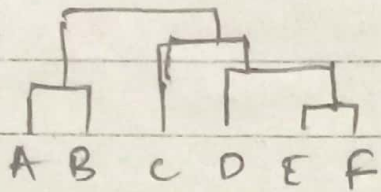PAM — partioning around medoids.

(1) Select K,
(2) Select random k points as medoids,
(3) Calculate cost    $|x_1 - x_2| + (y_1 - y_2)$ for
    remaining points for all k points.
(4) then assign that point to minimum cost.
(5) Then we Calculate total cost.
(6) Swap → Select one of non medoids O
(7) Find the same thing cost and assign points
    to cluster.
(8) find the the total cost, if totalcost is max
    then if is bad idea to swap, if found min
    then it is good idea.

# Hierarchical clustering (HC or HCA)

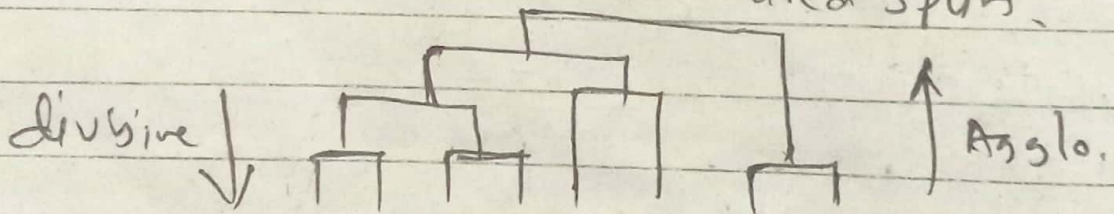→ data points are arranged in a hierarchy of cluster

dendogram → a diagram representing hierarcy

high indicate the
order

A B C D E F

## HCA Algorithm
①  Agglomerative - bottom up → Merging
②  divisive - top-down. Start with one cluster
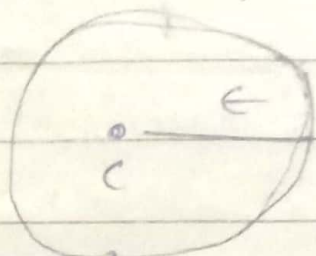        and split.

divisive ↓            ↑ Agglo.

## DBSCAN
Density Based Spatial clustering of Application
with noise.
+ unsupervised - clustering
+ based on density.

DBSCAN · (1) Epsilon (eps) ∈, measure of neighbourhood, or radius of circle



② Min_sample → How many Minimum data point required to consider an area dense
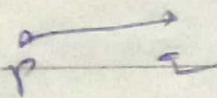
③ Density → no of data point.

④ Core point → if atleast a specified no. of neighboring points falls within specified radius ∈.

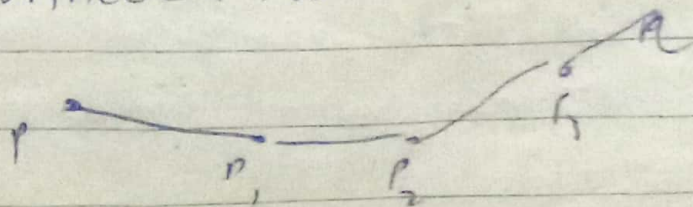⑤ Border point - which has atleast a core point

⑥ Noise point → nor core, nor Border.

⑦ 200 Density edge - two point p & q and distance b/w 2 points is less than ∈, then _____ is called density edge



⑧ Density connected points

Algorithm ① distance is calculated.

② neighbourhood. consider ∈.

③ min-Sample wit make cluster

④ repeat unti1 categorised

Advantages ① does not assume that the cluster have
Spherical Shape

② No need q finding k

③ high efficient

④ remove noise

Disadvantage.

→ It fail when. have multiple density

→ Not work well for high dimension

Spectral clustering (graph cuts)

Steps ① pre-processing

→ Construct a matrix representation q
the graph

② Decomposition

→ compute eigen value & eigen vector

→ Map each point to a lower- dimensional
representation based on one or more eigen value my

③ Grouping

→ Assign points to two or more clusters, based
on new representation.

## ~~Outlier Analysis~~

* Spectral clustering @ is unsupervised,
* we find clusters in graph
* we find subgraph and that will cluster.
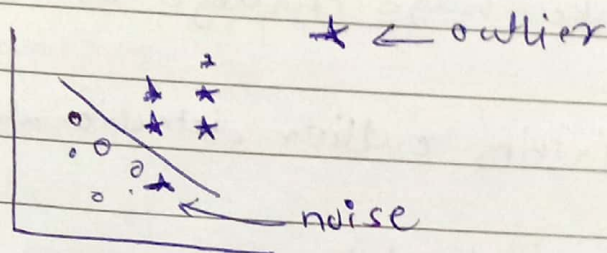
Unit - 5 (Part 2)

Outlier Analysis

outlier. Outlier are those datapoint that are significantly different from the rest of the dataset.

→ why to detect outlier

They are abnormal, and their presence can often skeow the result of statistical analyses on the dataset

- This could lead the less effective and less useful models.



noise is random error.

Types of outlier

① Global outlier
- If a data point is far away from rest of the data set.

② Contextual (conditional) outlier
- If a data point is deviates with respect to a specific context.
eg. 28°C is outlier in Shimla in winter but not in summer.

(3) Collective.

A subset of data point are deviates significantly from the entire data set.

Challenge in outlier detection.

* The boundary b/w normality & abnormality is no clear.

* differect application, may have very diff. requirement, Application Specific detection

* Noise can make huge challege for outlier detection.

* we have to justify outlier detection also.

Outlier Detection Model.

(1) Supervised. :- training + testing,
  - classifier
  - imbalanced handling.

(2) Semi supervised .- Having only some labelled daly. then we apply both Supervised & Semi supervised.

(3) Unsupervised, - clustering

(4) Statistical Method for Outlier Detection.
  Statistical model make assumption of data normality (z-score or Standard deviation, IQR (percentile

(5) proximity based Method for outlier Detection.

+ object is an outlier if the nearest of the object are far away in feature space.

+ neighbour is far away.

(6) Clustering based Method for outlier detection.

- normal data belong to large & dense cluster
- Outlier belong to Small or Sparse cluster.

Introduction to Isolation forest (iForest)

→ variation of Random forest
→ Unsupervised learning
↪ outlier are "few and different
- To make partition such that each data point is isolated

→ Isolation Tree is used (Binary tree)
(9) a regular / normal point is much harder than outlier



easy to
Isolate →
outlier

normal point

hard to
Isolate

[Partial]

Algorithm.

1) Training : Building a forest of isolation trees
→ Randomly select feature
→ Randomly partition

2) Prediction Compute outlier Score for newpoint.
let $m$ : Sample size

$$S(x,m) = 2^{\dfrac{-E\,(h(x))}{c(m)}}$$

$E(h(x))$ : Average Search height for $x$
from the tree

$c(m)$ = Average value of $h(x)$.

$E(h(x)) \ll c(m) \Rightarrow S(x,m) \cong 1$ (outlier)

$E(h(x)) \cong c\&(m) \Rightarrow S(x,m) \cong 0.5$ (Normal)

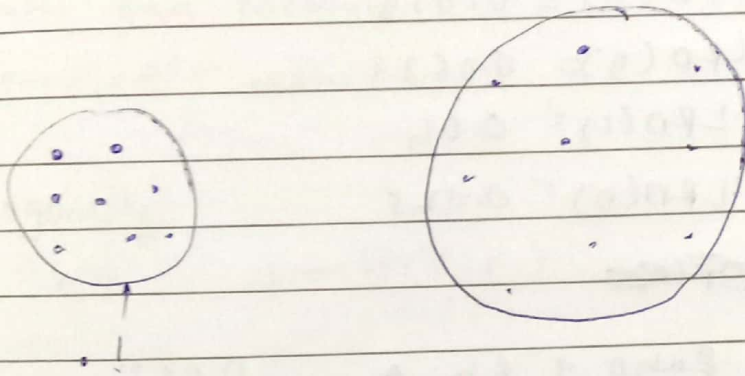hyper parameter ① No. of tree
② Sampling size
③ threshold value

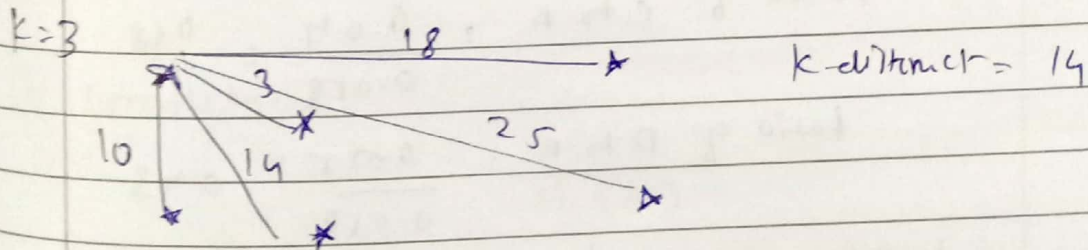| Adv. | disadv |
|---|---|
| ① faster | (1) can suffer bias |
| ② less memory. | (2) threshold value is not clear |

LOF - Local outlier factor.

→ Based on local density

→ ups. Unsupervised outlier detection.

→ It compute local density deviation of a given data point wrt to ib neighbour.



k-distance = Maximum distance b/w K neighbour.

k=3

$k$-distance = 14



reaciability distance = max (k distant, actual distance)

RD (A, B) = max (14, 10) = 14

RD (A, B) = Max (14, 3) = 14

RD (A, D) = Max (14, 14) = 14

RD (A, E) = Max (14, 18) = 18

RD (A, F) = Max (14, 25) = 25

local reacubivity density LRD = inverse of mean of the reacability between. towards other in dataset

$$LRD = \frac{14+14+14+18+25}{5} = 0.058$$

$LOF$ — Mean of the ratio of each k-neon?
to point of interest.

$LRD(A) = 0.058$

$LRD(B) = 0.062$

$LRD(C) = 0.04$

$LRD(D) = 0.025$

~~for (A)~~

Ratio of B to A $= \frac{0.062}{0.058} = 1.06$

Ratio of C to A $= \frac{0.04}{0.058} = 0.68$

Ratio of D to A $= \frac{0.025}{0.058} = 0.43$

$LOF(A) = \frac{1.06 + 0.68 + 0.43}{3} = 0.72$

~~LOF ≈ 1~~

$LOF \leq 1$ not outlier

$LOF > 1$ outlier

Here A is not outlier

Evaluation Metrics and Score.

2 method

① Extrinsic Method.

- ground truth is available, you can compare the clustering against the group truth and measure.

- Supervised method

(a) Homogeneity

$$h = 1 - \frac{H(C|K)}{H(C)}$$

It measure how the sample in a cluster are similar.

(b) Completness.

$$c = 1 - \frac{H(K|C)}{H(K)}$$

It measure how much similar sample are put together by clustering algorithm.

① Intrinsic Method

- when ground truth of a data is not available,
- It evaluate the clustering by examining how well cluster are seperated and how compact the cluster are.
- unsupervised

+ Silhoutte Coefficient

2 types Silhoutte function is used.

$$S(i) = \left\{ \frac{b(i) - a(i)}{Max(a(i), b(i))} \right\}$$

two type of data required

① Collection of all distance b/w objed

② partition obtained by Clustering

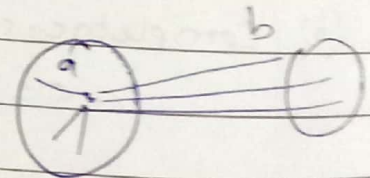Silhouette $\cancel{x}$ Value is Calculated

if  1  =  well clustered

   −1 =  'Poorley clusterd

a = Average distance of i to the points
                  in the Same cluster

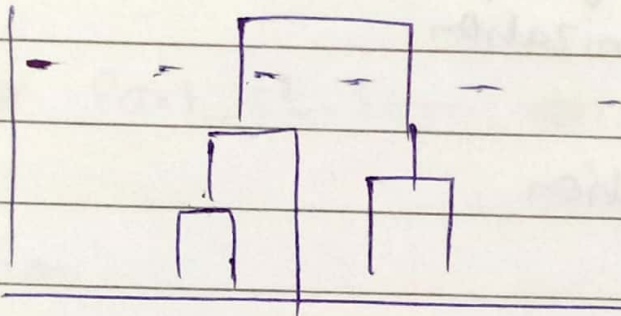b = min (average distance of i to points
                  in another cluster)

$$\underline{S = 1 - a/b \quad if \ a < b}$$

How Should we choose the Number of clusters in Hierachical clustering

→ Using dendogram.

→ Whenever two cluster are merged, we will join them in this dendogram and height of the join will be distance b/w these points.

Now we will set threshould value and draw horizental line (It should cut tallest line)



No. of cluster will be no. of intersection