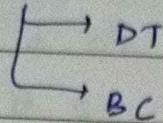Classification
- → DT
- → BC

- → BP
- → SVM
- → LR
- → KNN

⎤ Test (next week)

Clustering

Boosting algorithms

Feature selection

---

Data clustering : Unsupervised learning.
↓
↳ Class label not known

Partition based clustering
- → k-means
- → PAM
- → ~~KLARA~~ Clara
- → Clarans

⎫ Variations
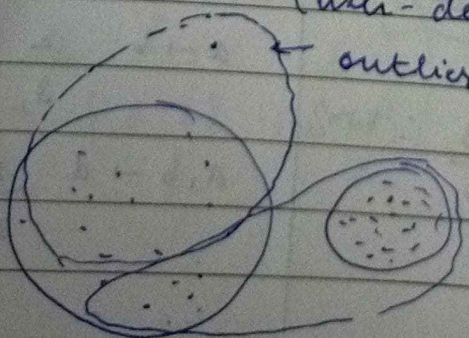⎬ in k-mean as k-mean is
⎭ not good for large set of data.
↓
- → might not give proper clusters
- → lacks accuracy & efficiency
- → Sparse matrix

---

Problem: Outliers are also part of a cluster

k: number of clusters
(user-defined parameter)

← outlier

Due to presence of outlier means converge to a value impacted largely by the outlier

Binary variables:

| | |
|---|---|
| a | 11 |
| b | 10 |
| c | 01 |
| d | 00 |

Both objects for a particular instance have 1.

Symmetric:

1 for $X_1$, 0 for $X_2$

$$d(i,j) = \frac{b+c}{a+b+c+d} = \frac{10 + 01}{11 + 01 + 10 + 00}$$

Asymmetric: $d(i,j) = \dfrac{b+c}{a+b+c}$

How to know if variable is sym/asym.

If $\neg 0 \to 1$
& $\neg 1 \to 0$ $\Big\}\Rightarrow$ Symmetric binary variable.

If $\neg 0 \neq 1$ then asymmetric binary variable.

Eg: No : N $\to$ 0
Yes : Y $\Big|$ Presence: P $\to$ 1

no instance for 10.

$$d(Jack, Mary) = \frac{10 + 01}{11 + 01 + 10} = \frac{0 + 1}{2 + 1}$$
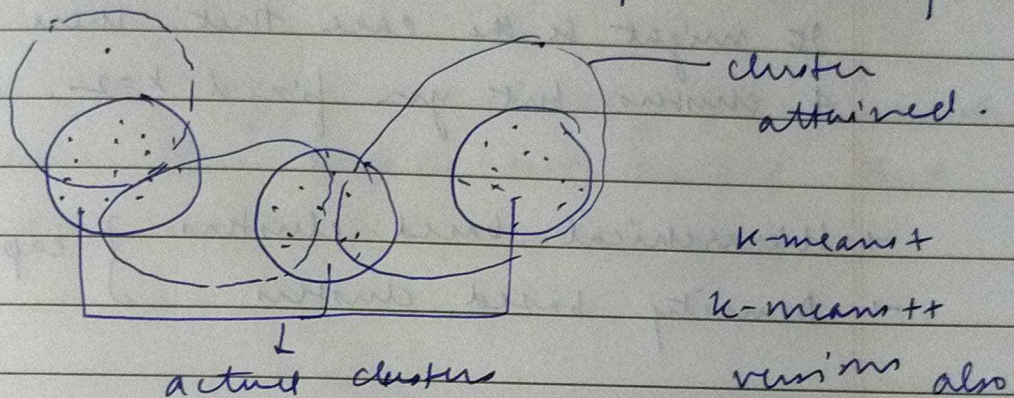
$$= \frac{1}{3} = 0.32$$

The points which were initially part of cluster 1 can now be part of cluster 2.

Large data $\Rightarrow$ Large noise

$\Rightarrow$ Mean affected largely

∴ The actual clusters may not be identified which impacts accuracy.



- cluster attained.
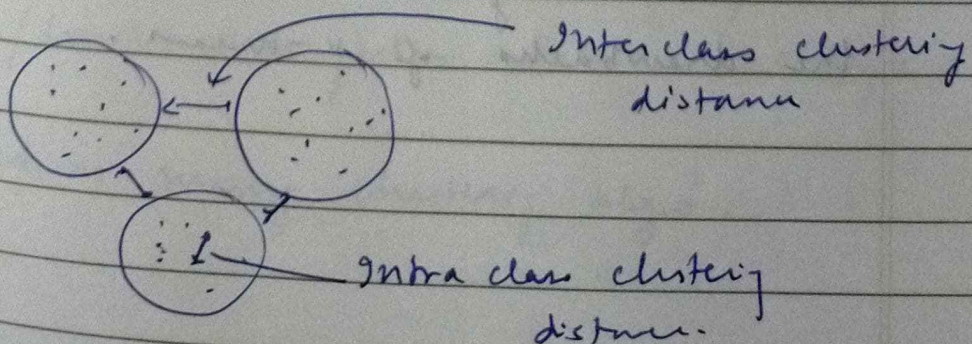- k-means +
- k-means ++ versions also available.

← actual clusters

To apply k-means on a large data, remove noise from the data.

— can be used in Text mining, etc.

Data clustering : grouped based on similarity measure (could be distance). Can be on multi-dimensional data.

→ OPTIMISATION PROBLEM



Inter class clustering distance

Intra class clustering distance

Data types

→ Interval based
→ Ratio based
→ Binary
→ Ordinal
→ Nominal

} Real data has to be transformed into a data of one of these types before applying clustering

Data from kaggle is already pre-processed & hence this step is not needed. But raw data must be pre-processed first.

Data matrix : attributes in columns
entities in rows.

$$
E \begin{bmatrix} x_{11} & \cdots & x_{1f} & \cdots & x_{1p} \\ \vdots & & \vdots & & \vdots \\ x_{i1} & \cdots & x_{if} & \cdots & x_{ip} \\ \vdots & & \vdots & & \vdots \\ x_{n1} & \cdots & x_{nf} & \cdots & x_{np} \end{bmatrix}
$$

$A_1$   $A_2$ . . . $A_n$

$E_n$

Object by variable.

Dissimilarity matrix:

Object by object : Distance b/w all objects & all other objects.

→ upper/lower triangular matrix is enough.

→ Used in many clustering algo.

Clustering requires 2 optimisations at the same time:

① Inter class distance, of clusters should be maximised

② Minimise intra class data distance ~~of clusters~~.

Problem with k-means :- Fixing k
It might be the case that there are actually 4 clusters but you fixed $k = 2$.

→ Hierarchical based clusters ⎫ capable of doing

→ Density based clusters ⎭

Within hierarchical based clustering
we have :           Single link
                    Heavy link
                    Complete link
                    Birch

        Density based :- → CURE
                         → DBSCAN
                         →

3 researchers
↳ BFR : Primarily uses k-means & overcomes the drawbacks of k-means.

Smaller unit → larger variable range.

To select basketball players based on height, we want clear distances in the differences due to closeness in heights.
∴ Lower down the unit -
Give more weightage to height.
∴ Convert Euclidean into ~~weight~~ weighted distance -

$$w(h_1 - h_2)$$

weight
assigned
to assign higher weight
to height.

Suppose 3 attributes: $A_i$, $N_i$, $B_i$
assign higher weight to the one with highest importance
$$w_1 > w_2 > w_3$$

→ Mean absolute deviation better than standard deviation due to absolute values.

Assigning weights to Minkowski distance would give weighted distance:
$$\sqrt[1/m]{w_1 |x_1 - \bar{x}|^m + \cdots}$$

Not necessary to assign weight to all attributes.