

Unsupervised learning → Understand just Patterns of data ('x')

DATE

Clustering (Silhouette)

→ Explain Data

→ Filling Missing Values

'unlabelled' data into 'clusters' of similar inputs

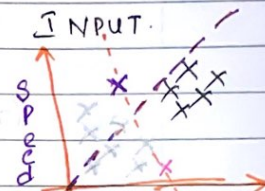
→ Describe data by Discrete 'Groups'

Group data by visual similarity

1. K-Means clustering

→ Find Similarity between points

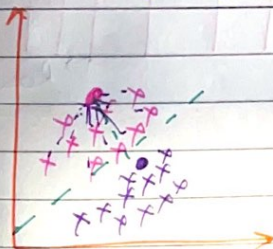
→ make them clusters



⇒ ALGORITHM Random

(Prototypes/Mean) Homopower

1. Select 'K' Value { Centroids } $K=2$ (2 clusters)
2. Initialize centroids to randomly, closest to cluster
3. Find the distance between the points
4. Select the group and find the mean
5. Update the centroids, till no more movement



⇒ Why $K=2$ / Find distance between Points

1. Euclidean distance (Pythagorean theorem)

$$P_1 (x_1, y_1) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$P_2 (x_2, y_2)$$

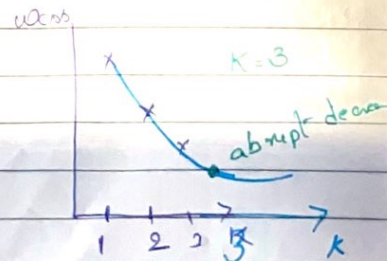
P_1 P_2

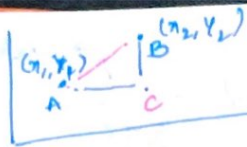
Select 'K' Value

⇒ Elbow method

$K=1$ to 20

$$W_{\text{cost}} = \sum_{i=1}^n ((c_i + x_i)^2)$$

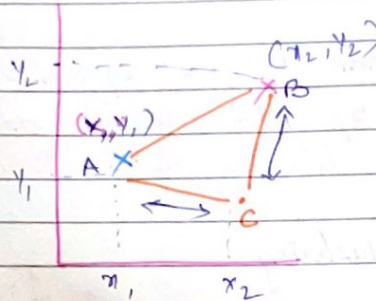




$$AB^2 = a^2 + b^2$$

DATE

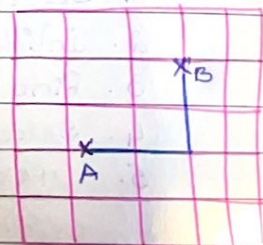
Manhattan distance



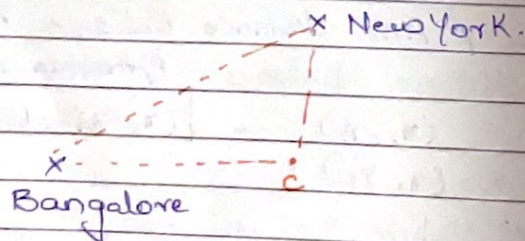
$$= |x_2 - x_1| + |y_2 - y_1|$$

di

Example:- Manhattan



Euclidean



DATE

⇒ Agglomerative Hierarchical clustering

→ recursive kmeans (Top down)

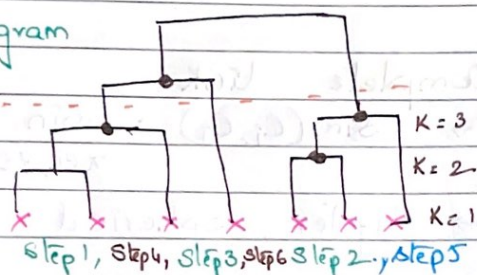
→ Build the Bottom level

→ 'm' object, every object in cluster

→ Find the closest cluster and then

m-1

→ Dendrogram



⇒ Nodes = Subsets

characteristic = Root (whole set)

leaves = individual elements

Internal nodes = union of children.

Each level = different no. of clusters

⇒ Distance = closest Pair

P_1, P_2

⇒ 1) Single Link



$$\text{Sim}(C_i, C_j) = \max_{x \in C_i, y \in C_j} \text{Sim}(x, y)$$

⇒ Most Similar Pair

2) Complete Link :-

Least Similar Pair

3) Average Link :-

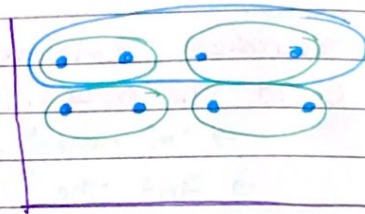
average of similarity

4) Centroid

centers of Gravity

⇒ Single link Example (Thin, long)

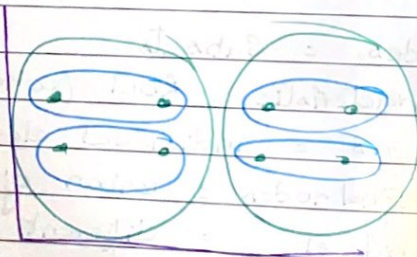
→ One Pair of link



⇒ Complete link

$$\rightarrow \text{Sim}(C_i, C_j) = \min_{x \in C_i, y \in C_j} \text{Sim}(x, y)$$

→ Uptier, Spherical



⇒ Computational complexity

HAC → N objects

N^2 similarity $O(N^2)$

$N-1$ steps

⇒ Average link

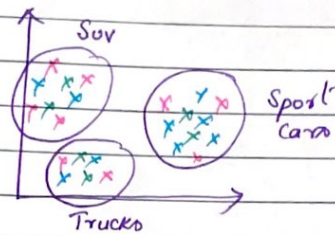
→ All ordered pairs

→ All pairs between clusters

$$\text{Sim}(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{x \in C_i} \sum_{y \in C_j} \text{Sim}(x, y)$$

DATE

Ex:- Cars between Horsepower and Speed.



- Are clusters well separated?
- Are clusters linearly separable?
- (How difficult to separate stuff)



Challenges:-

1. overlap / complicated shapes.
2. Algorithm No. of clusters
don't need no. of clusters

⇒ K-Means algorithm

⇒ How many means in K clusters?

⇒ applied on numbers



⇒ Good Centroids → distance should be minimum.