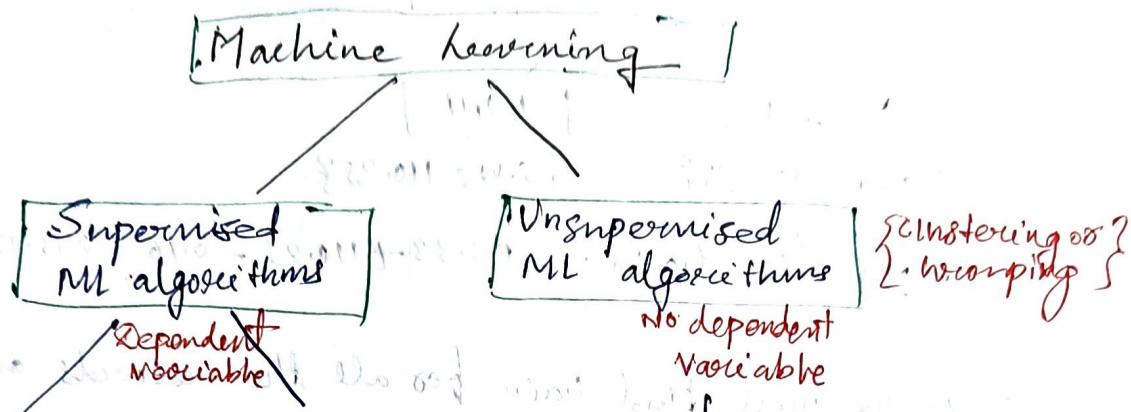


28-Nov-2022

Unsupervised Learning



Regression

- ① Linear Regression
- ② SVR
- ③ DTR
- ④ RFR
- ⑤ LBR
- ⑥ XBR
- ⑦ L1 & L2 Reg.

Classification

- ① Logistic Regression
- ② SVC
- ③ DTC
- ④ RFC
- ⑤ XBC
- ⑥ LBC/ABC

- (updated version)
- ① K-Means and K-Means++
 - ② Hierarchical
 - ③ DBSCAN

Eg: Market Segmentation

⑧ KNN (can be used for both supervised and unsupervised problem).

Dataset

Height	Weight	BMI	Country
170	60	21	IND
180	65	22	UK
180	70	20	USA
165	75	18	IND
175	85	19	USA

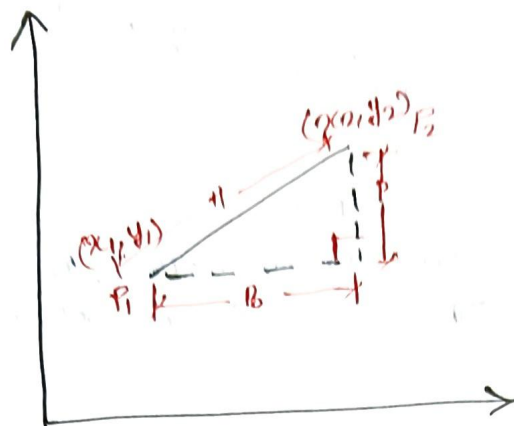
3 groups
{IND, USA, UK}

- Based on similarity, grouping happened.
(by correlation by distance)

- Euclidean distance
- Manhattan distance
- Manhattan distance
- cosine distance
- Tanimoto distance
- Squared Euclidean distance

K-Means

⇒ K-means → Data → Similarity → Distance → Euclidean distance



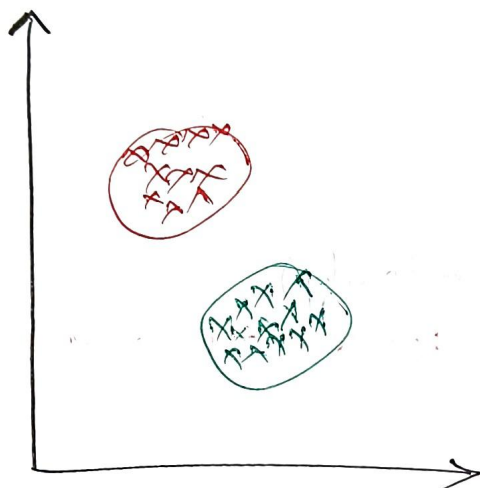
$$H^2 = P^2 + B^2$$

$$H = \sqrt{P^2 + B^2}$$

Distance b/w P_1 & P_2 ,

$$H = \sqrt{P^2 + B^2}$$

$$D(P_1, P_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



Two Group

Height

Weight

① 185	72
② 170	56
③ 168	60
④ 179	68
⑤ 182	72
⑥ 188	77
⑦ 180	71
⑧ 160	70
⑨ 183	84
⑩ 180	88
⑪ 180	67
⑫ 167	76

- ① Centroid
- ② Distance
- ③ Mean

connect these 3 to understand K-Means

> K-Means

↓
Number of centroid

→ Centroid value (I rounded it, we have to build cluster)?

How many centroid should we take?
→ To choose K , we use elbow method. We use WCSS (within cluster sum of square) in elbow method.

WCSS

- Inter cluster
- Intra cluster

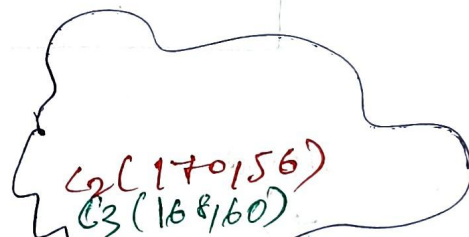
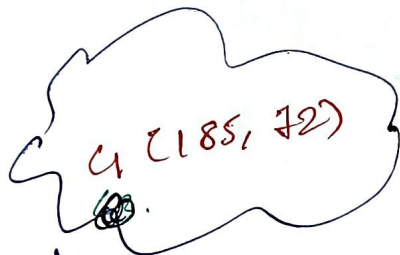
Evaluation of clustering model

- ① Dunn's Index
- ② Silhouette score

Steps for K-Means

① Initialize centroid (Randomly)

Let's take $K=2$ {Number of centroid = no. of cluster}



→ find euclidean distance, b/w C_1 and C_3 , C_2 and C_3

Let's assume, $C_3 (168, 60)$

$$d(C_1, C_3) = 3$$

$$d(C_2, C_3) = 8$$

(less distance)

↓
3rd point will be assigned to C_2 .

Actual distance

$$d(C_1, C_3) = \sqrt{(168-185)^2 + (60-72)^2} = 20.80$$

$$d(C_2, C_3) = \sqrt{170 \cdot (168-170)^2 + (60-56)^2} = 4.14 \text{ (less distance)}$$

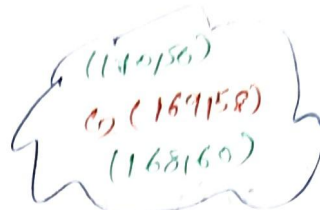
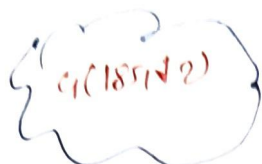
↓
 C_3 belong to 2nd cluster (C_2)

now, update the cluster 2 centroid

$$\text{updated centroid} = \left(\frac{170+168}{2}, \frac{56+60}{2} \right)$$

$$\text{new centroid} = (169, 58)$$

> Every time centroid will get updated after adding the addition of new points



> Similarly, find the distance to each point and keep updating centroid till the end.

Example



→ House 1 children has more similarity to house 1 people as compared to house 2.

4th point $c_4(179, 68)$

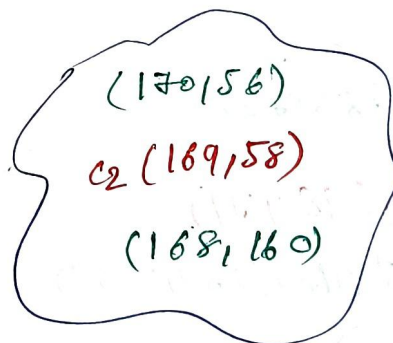
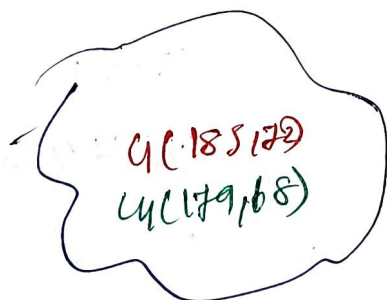
find actual distance with c_4

$$d(c_1, c_4) = \sqrt{(168-179)^2 + 1}$$

$$d_1(c_1, c_4) = \sqrt{(179-185)^2 + (68-72)^2}$$

$$d_2(c_2, c_4) = \sqrt{(179-169)^2 + (68-58)^2}$$

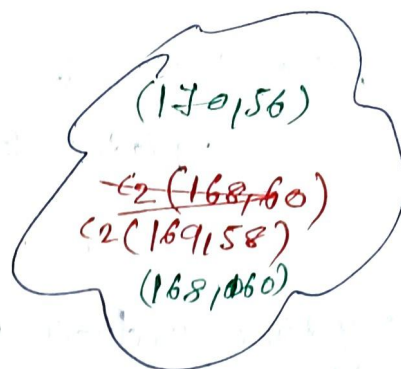
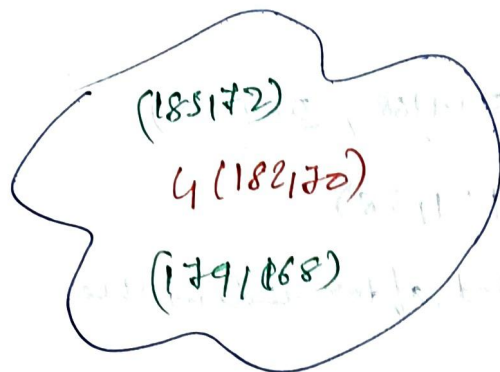
$= 7.21 \rightarrow$ (less distance as compared to d_2 .
↓
Belong to c_1 .)



now, update the cluster 1 centroid

$$\text{New centroid} = \left(\frac{185+179}{2}, \frac{72+68}{2} \right)$$

$$= (182, 70)$$



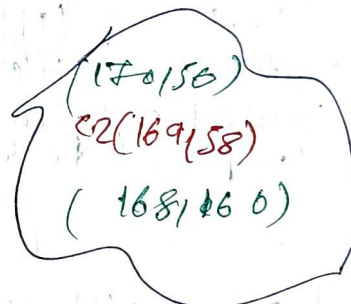
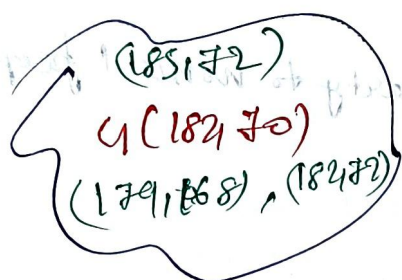
5th point, $(5(182, 72))$,

Find actual distance w.r.t c_5 ,

$$d_1(c_1, c_5) = \sqrt{(182-182)^2 + (72-70)^2} = 2 \text{ (less distance)}$$

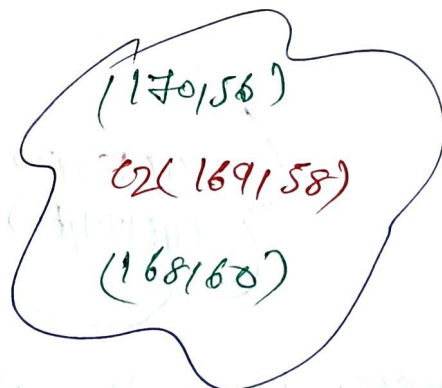
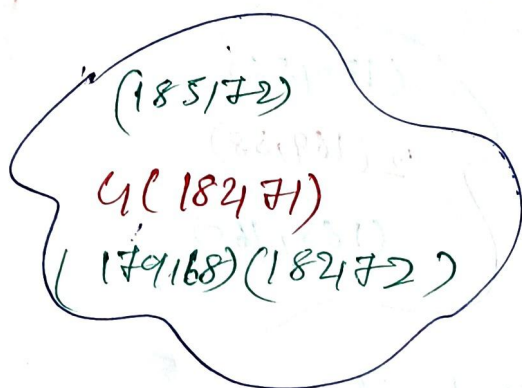
$$d_2(c_1, c_5) = \sqrt{(182-169)^2 + (72-58)^2} = 19.1$$

Belong to c_1



Now, update the centroid of cluster 1

$$\text{New centroid} = \left(\frac{182+182}{2}, \frac{70+72}{2} \right) = (182, 71)$$



> Similarly, find the distance w.r.t each point and keep updating centroid. Like this we can get final 2 clusters.

- > Keep iterating till we reach all the data points.
- > If both distance is same, any cluster can be chosen.

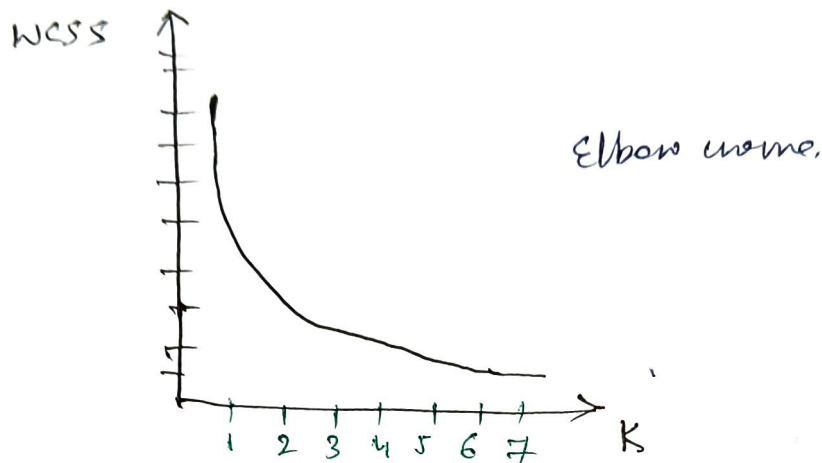
Elbow method and WCSS

Ans

$K=2$

(How can we decide??)

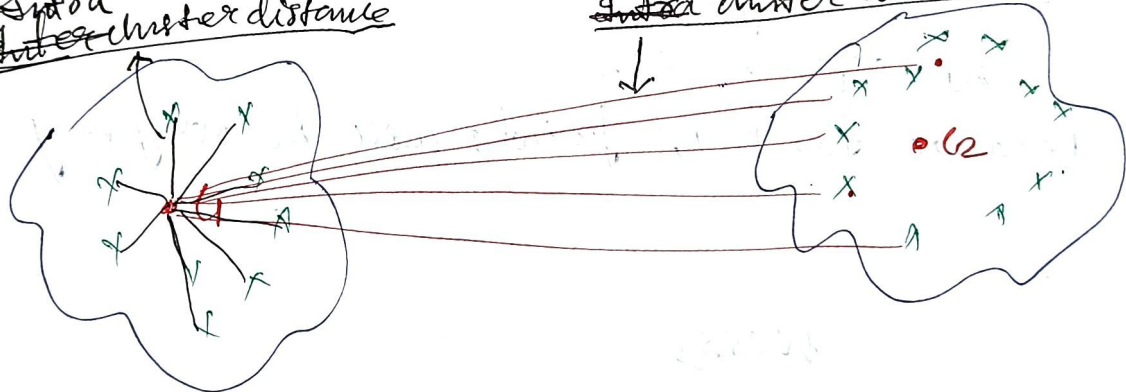
→ Using elbow method, we can decide.



WCSS → within cluster sum of squares

Intra cluster distance

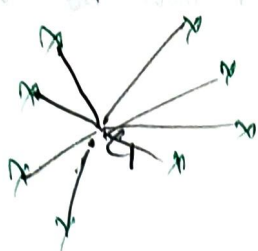
Inter cluster distance



{ Intra cluster distance - within cluster
Inter cluster distance - B/w the cluster }

WCSS

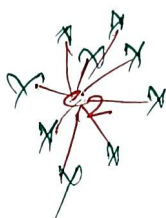
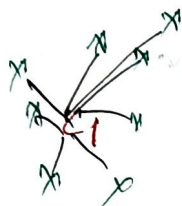
For $K=1$,



$$WCSS_1 = \sum_{i=1}^n d(c, x_i)^2$$

For $K=2$,

Two comp,



$WCSS_2$

which WCSS will be greater?

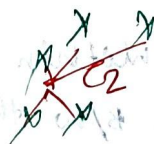
→ $WCSS_1$ will be greater.

$$WCSS_1 > WCSS_2$$

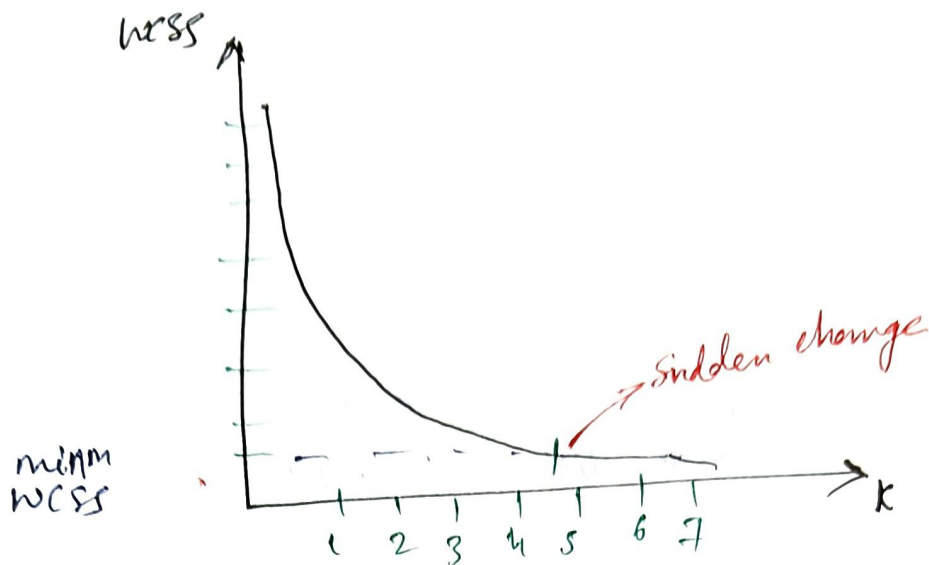
because for $K=1$, points will be very very scattered

For $K=3$,

$WCSS_3$



$WCSS_1 > WCSS_2 > WCSS_3$



- ⇒ choose $k=5$, because it has abrupt changes.
- > Take average value of WCSS if no. of cluster is greater than 1.

Q what is the diff b/w K-means and K-means++?

Ans

Validation of clustering

① Dunn Index

② Silhouette Score

① Dunn Index

$$\text{Dunn Index} = \frac{\max_{i,j} \text{distance}(x_i, x_j)}{\max_i \text{distance}(x_i, y_i)}$$

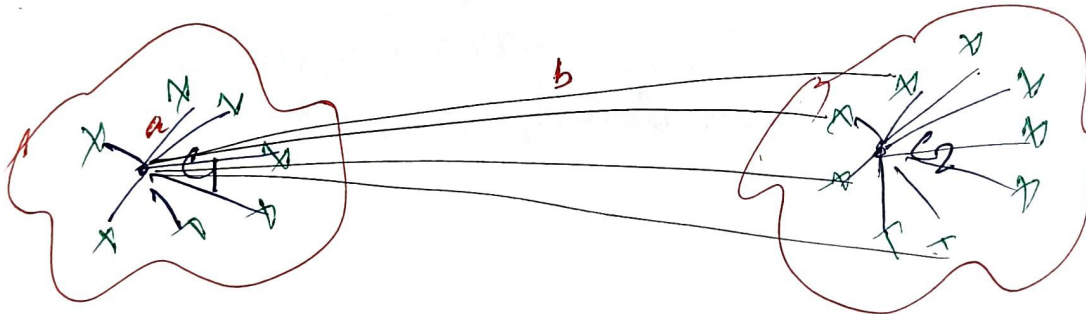
② Silhouette Score $[-1 \leq s \leq +1]$

$$\text{Silhouette score} = \frac{b_i - a_i}{\max(b_i - a_i)}$$

$$\{-1 \leq s \leq 1\}$$

a_i = Intra cluster (within cluster)

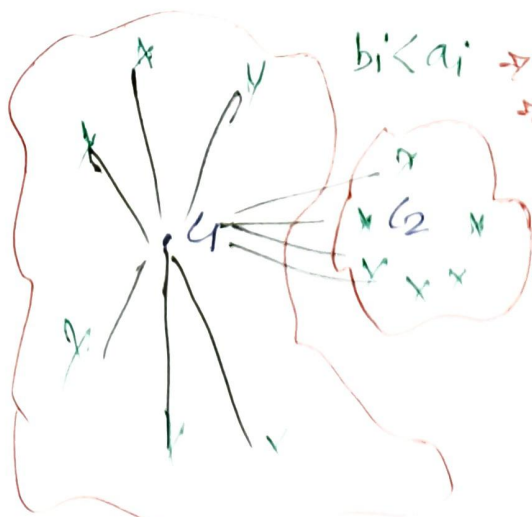
b_i = Inter cluster (B/w the clusters)



$$\left\{ \begin{array}{l} -1, (\text{worst model}) \\ +1, (\text{best model}) \end{array} \right\}$$

Example let's assume $b = 50, a = 40$
 $50 - 40 = 10$ (less scattered data, good Model)

$b = 40, a = 50$
 $40 - 50 = -10$ (more scattered data, not good)



$b_i < a_i \rightarrow$ worst model.

\rightarrow In same cluster, values are more scattered.

Q How to make a best model or optimise solution?

Ans Custom learning / Custom model / Semi-Supervised learning
 $>$ combination of supervised and unsupervised learning.

Dataset

<u>weight</u>	<u>Height</u>	<u>gender</u>
170	55	M
180	60	F
165	70	M
180	80	F
155	50	F
160	100	F

Q Build classification model and optimise result?
 Ans We can perform clustering 1st and then classification.
 Basically we are segregating the data.

LOM \rightarrow Cluster 1
 Cluster 2
 Cluster 3 \Rightarrow then any supervised learning algo algorithm.

$>$ For categorical features, we can find hamming distance.