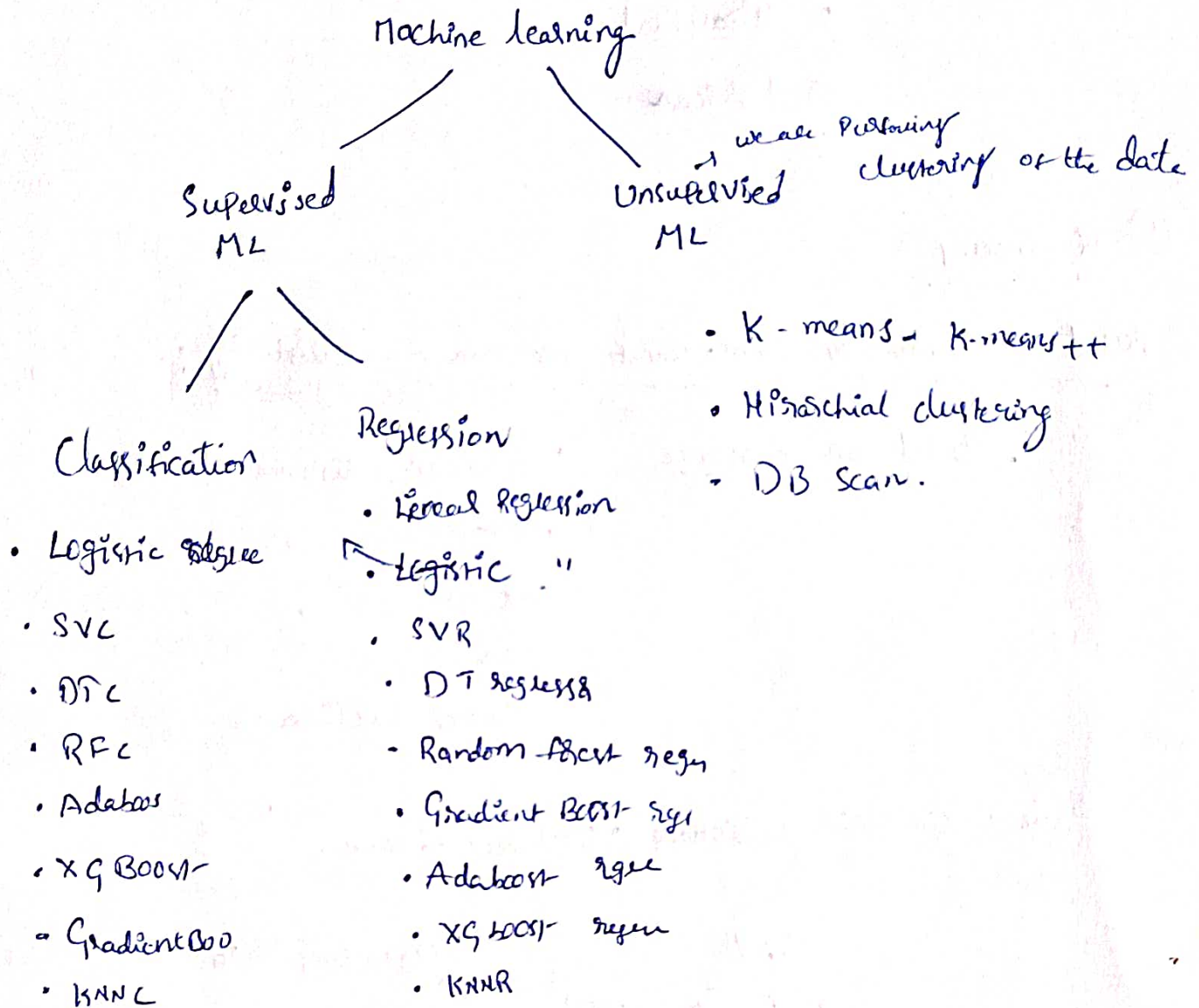


Sri  
28/11/22

# Unsupervised Machine Learning.



Target / dependent / supervisor

Supervised Algorithms

Height	Weight	BMI	Supervisor Country
170	60	21	India
180	65	22	UK
160	70	22	USA
165	75	18	India
140	55	19	USA

Unsupervised ML

↓  
Clustering → grouping

based on above dataset, we can make  
based on country  
3 clusters or 3 group.  
→ one group one group  
India, USA, UK → one group.

Mathematically

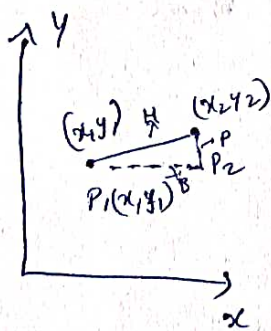
1. K-means
2. Hierarchical
3. DBscan

## ① K-means

Based on Similarity Measurement we can do clustering:-

- Based on distance : Euclidean distance measure
- Manhattan
- cosine
- Tanimoto
- Sauade Euclidean

Data  $\rightarrow$  similarity  $\rightarrow$  Distance  $\rightarrow$  Euclidean distance.

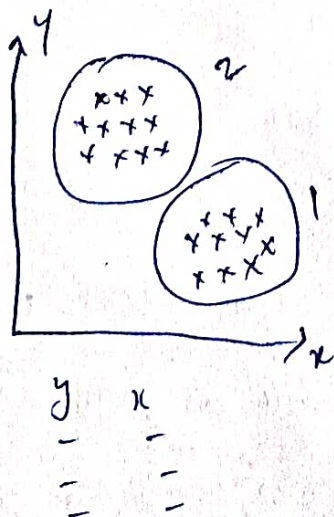


Pythagoras theorem

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

$$(P_1, P_2) H = \sqrt{P^2 + B^2}$$

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



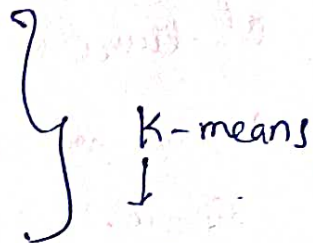
Height	Weight
185	72
170	56

$\rightarrow$  Clustering based on distance we can find the similarity

168	60
179	68
182	70
188	77
180	71
160	70
183	84
180	88
180	67



- Centroid
- Distance
- Mean



no. of Centroid



- ELBOW method

- WSS  $\rightarrow$  we can evaluate within clusters some of square

[inter cluster]  
[Intra cluster]

Height	weight
185	72
170	56
168	
179	
182	
188	

### Evaluation method:

- Dunn Index
- Silhouette Coeff

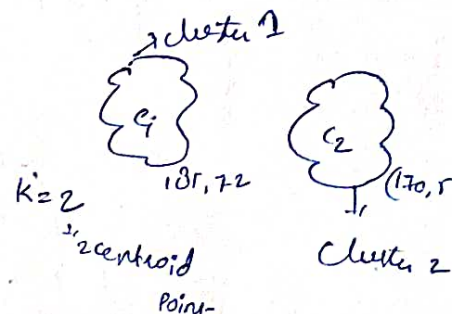
$\rightarrow$  center value around that we need 2 build the cluster

- Centroid [randomly]

here centroid = 2  
 $K=2$

(185, 72)  
 $C_1$

(170, 56)  
 $C_2$



Euclidean distance  $\rightarrow$   $d(C_1, 3)$   
 $d(C_2, 3)$

$C_1 \rightarrow 3$   
5

$C_2 \rightarrow 3$   
8

NO. of centroid = NO. of clusters

$$\begin{matrix} C_1 & 3 \\ (185, 72) & (168, 60) \\ \sqrt{(168-185)^2 + (60-72)^2} \end{matrix}$$

(170, 56)  
 $C_2$  (168, 60)

$$\begin{aligned} D(C_1, 3) &= \frac{20.6}{2} = 10.3 \\ D(C_2, 3) &= \sqrt{(170-168)^2 + (56-60)^2} \\ &= \sqrt{20} = 4.4 \end{aligned}$$

update the cluster centroid

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

$$(169, 58)$$

3rd cluster is belong to 2nd cluster.

here we are considering min distance.

$$\text{for cluster 1, } = 20.08$$

$$\text{for cluster 2, } = 4.4$$

since we can say that 3rd now belongs to cluster '2'.

$$\text{Distance 4th } D(C_1, 4) = \sqrt{(168-179)^2 + (72-68)^2} = 7.21$$

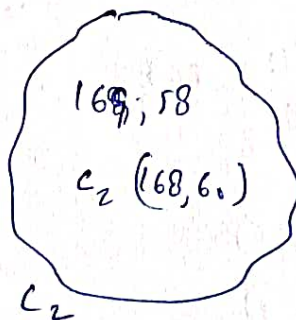
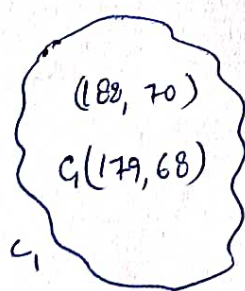
$$D(C_2, 4) = \sqrt{(169-179)^2 + (58-68)^2} = 14.42$$

4<sup>th</sup> value belongs to  $C_1$  because of less distance.

or belongs to  $C_1$  because of low distance.

update the cluster centroid

$$\left( \frac{179+181}{2} \right), \left( \frac{68+72}{2} \right) = (180, 70)$$



for (182, 72)

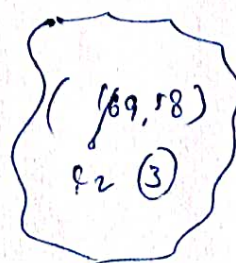
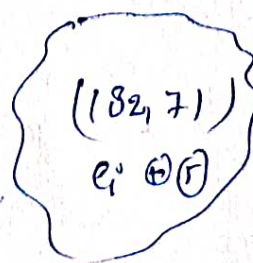
$$\text{distance b/n } d(C_1, 5) = \sqrt{(182-182)^2 + (72-70)^2} = 2$$

$$d(C_2, 5) = \sqrt{(182-169)^2 + (72-58)^2} = 19.1$$

(5 belongs to  $C_1$ )  
it belongs to 1st cluster

$$= \frac{182+182}{2}, \frac{70+72}{2}$$

$$C_1 = (182, 71)$$

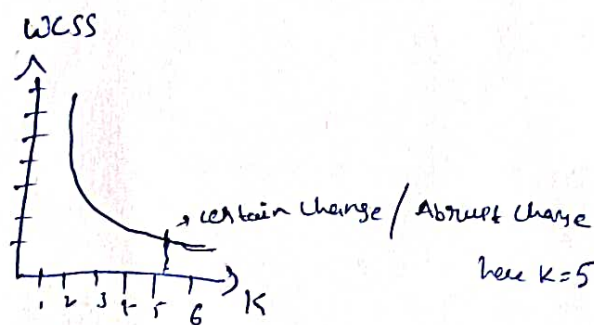




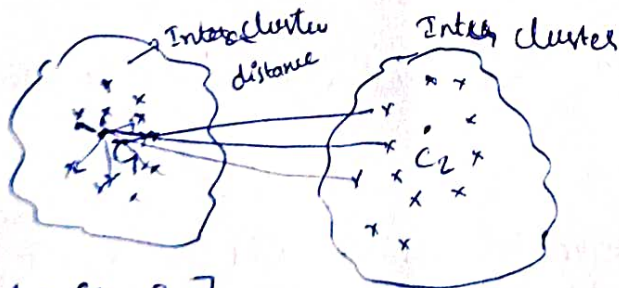
## K-means

1. Centroid
2. Distance [compute, min]
3. Include point in cluster, update the centroid

How can we consider how many clusters we need & choose, based on ELBOW METHOD.

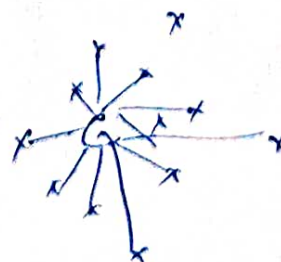


WCSS :- Within Cluster ~~Sum~~ of Squares



Inter [ $C_1, C_2$ ]  
Intra : within clusters

$K=1 \rightarrow 1 \text{ centroid} / 1 \text{ cluster}$

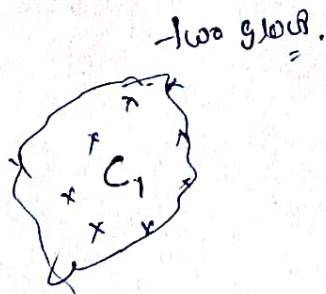


$\Rightarrow WCSS$

within cluster sum of squares

WCSS,

$K=2$



$$WCSS_2 =$$

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

$$WCSS_1 > WCSS_2$$

$K=3$



$$WCSS_1 > WCSS_2 > WCSS_3$$

Difference b/n K-means/K-means++

- The diff b/n ' ' lies in: the selection of centroids around which the clustering takes place.
- K means++ promotes the drawback of Kmeans which is it is dependend on initialization of centroid.
- Guarantee convergence :- Can warm-start the position of centroids. Easily to new examples. Generalizes to clusters of diff. shapes and sizes, such as elliptical clusters.

How to validate cluster:-  $K=5$

Regression based on  $R^2$  / adjusted  $R^2$   $R^2 = [0, 1]$  if  $R^2 = 0$  it's a worst  
Classification based on Acc/Roc / Confusion Matrix  $R^2 = 1$  it's best.

1. Dunn Index

2. Silhouette Score

1. Dunn Index =  $\frac{\max_{i \neq j} \text{dissim}(x_i, x_j)}{\max_i \text{dissim}(x_i, y_j)}$

2. Silhouette Score =  $\frac{b_i - a_i}{\max(b_i, a_i)} \Rightarrow [-1 \text{ to } 1]$  if  $-1$  worst  $+1$  best

$a_i$  = intra cluster

$b_i$  = inter cluster

