

# Machine learning

Supervised  
ml algo

Unsupervised  
ml algo

Regression

- ① Lin reg
- ② ~~Log reg~~
- ③ SVR
- ④ DTR
- ⑤ RFR
- ⑥ GBR
- ⑦ XBR
- ⑧ KNNR

Classification

- ① Log reg
- ② SVC
- ③ DTC
- ④ RFC
- ⑤ XGB
- ⑥ GBC/ABC
- ⑦ KNNC

- ① k-means → k-means++
- ② hierarchical
- ③ DBSCAN

target/dependent / Supervisor

Supervised → ml algo

Supervisor

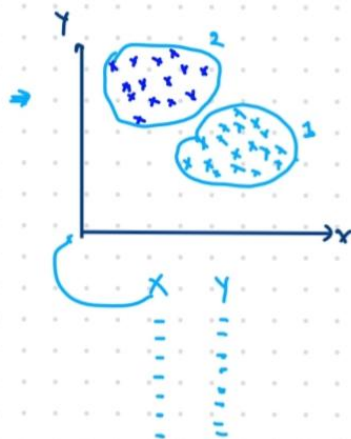
Height	Weight	BMI	Country
170	60	21	IND
180	65	22	UK
160	70	20	USA
165	75	19	IND
140	55	19	USA

3-cluster

3-group

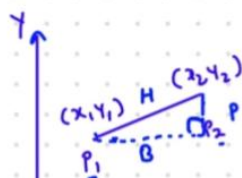
mathematically

- ① k-means
- ② hierarchical
- ③ DBSCAN



① K-means → Data → Similarity → Distance → Euclidean dist

- ① centroid.
- ② Distance.
- ③ mean.



Pythagoras

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

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

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

K-means

	Height	Weight
①	185	72
②	170	56
③	168	60
④	179	69
⑤	182	72
⑥	188	77
⑦	180	71
⑧	160	20
⑨	183	84
⑩	180	88
⑪	186	76
⑫	167	

Clustering

- ① Centroid [randomly]
- $K=2$

No of centroid = Number of cluster

$C_1 (185, 72)$   $C_2 (168, 60)$

$$= \sqrt{(168-185)^2 + (60-72)^2}$$

$$D(C_1, C_2) = 26.80$$

$$= \sqrt{\frac{(170-168)^2 + (56-60)^2}{2}} = \sqrt{\frac{2^2 + 4^2}{2}} = \sqrt{\frac{20}{2}} = \sqrt{10} \approx 3.16$$

$$C_2 = \left( \frac{170+168}{2}, \frac{56+60}{2} \right)$$

$$\approx (169, 58)$$

$$= (185, 72)$$

$C_1, 3$

$$C_1 \rightarrow 3 = 5 \quad C_2 \rightarrow 3 = 8$$

E.D. →  $d(C_1, 3)$   
 $d(C_2, 3)$

- Dunn's Index.
- Silhouette Coeff.

Inter cluster  
Intra cluster

WCSS

→ Elbow method.

2, 3, 4, 5...

no. of centroid.

K-means

K-means

$K=2$

2 centroid

dis (similarity)

min

$C_1 (185, 72)$

$$C_2 (169, 58) \quad \left( \frac{x_1+x_2}{2}, \frac{y_1+y_2}{2} \right)$$



$$\frac{179 + 155}{2} = 167$$

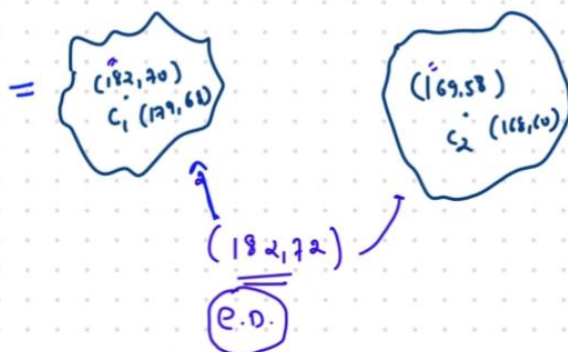
$$\frac{68 + 72}{2} = 70$$



$$D(c_1, 4) = \sqrt{(185 - 179)^2 + (72 - 68)^2} = 7.21$$

$$D(c_2, 4) = \sqrt{(165 - 179)^2 + (60 - 68)^2} = 14.14$$

Why? = low distance, min dist.



$$d(c_1, 5) = \sqrt{(182 - 182)^2 + (72 - 70)^2} = 2$$

$$d(c_2, 5) = \sqrt{(182 - 169)^2 + (72 - 58)^2} = 19.1$$

$$\begin{pmatrix} 182 \\ 70 \end{pmatrix}$$

$$\begin{pmatrix} 182 \\ 72 \end{pmatrix}$$

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

$$= (182, 71)$$



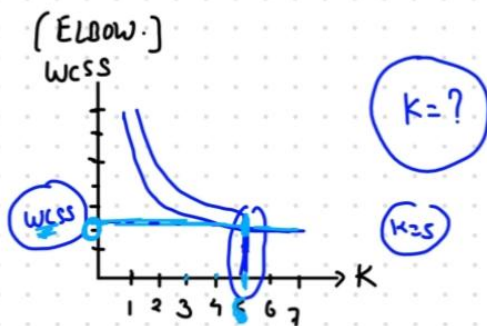
K-means

- 1 Centroid
- 2 Distance. [compare, min]
- 3 Include Point in cluster, update the centroid.

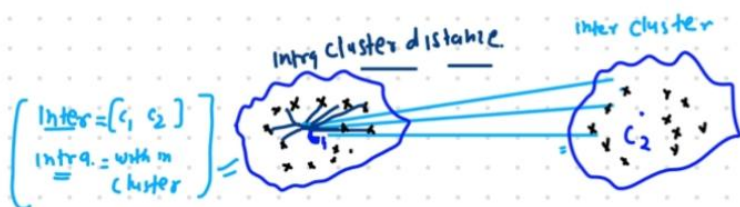
iterate while all the data points



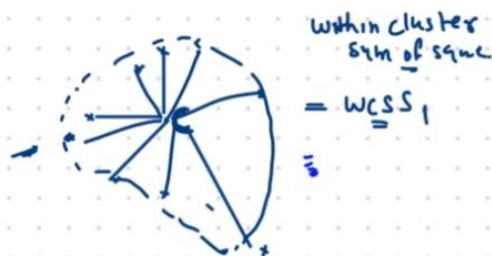
K-means.  
 $K=2, 3, 4, 5, 6$   
 [2 cluster].



$\Rightarrow$  WCSS  $\rightarrow$  With in cluster some of square



$K=1$



$K=2$



$= WCSS_1 > WCSS_2$

WCSS = within cluster sum of square

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

$K=3$

$WCSS_1 > WCSS_2 > WCSS_3$



K-means

$K=1, 2, 3, 4, 5 \dots$

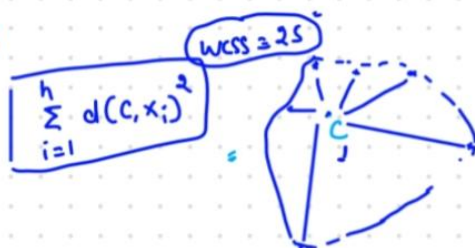
12 data point

minimum = 1

maximum = 12

WCSS

$K=1$



$K=2$

$WCSS = 20$



$K=3$

$WCSS = 15$



diff b/w K-means | K-means++ =

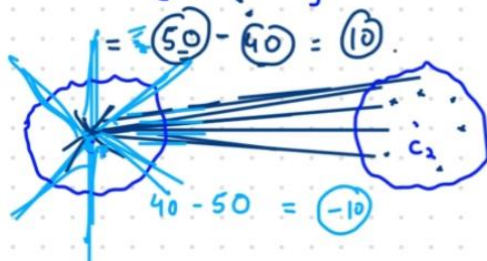
how to validate cluster  $K=5$

Distance

- ① dunn index.
- ② silhouette score

① dunn index =  $\frac{\max \text{dist}(x_i, x_j)}{\max \text{dist}(y_i, y_j)}$

② silhouette score =  $\frac{|b_i - q_i|}{\max(b_i, q_i)}$  Same cluster



Linreg

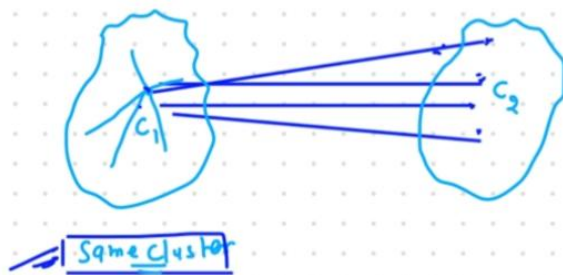
Regression

$\rightarrow R^2 / \text{adj-}R^2$   $R^2 = [0, 1]$

classification worst [0] best [1]

$\rightarrow \text{Roc AUC} / \text{Confusion mat}$

worst best  
 $-1 \quad +1$



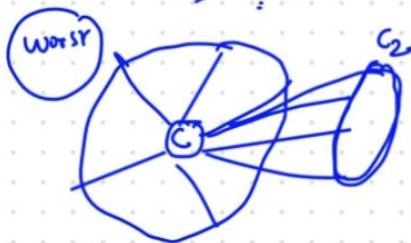
K22

$a_i$

$b_i$

$$\equiv a_i > b_i$$

+ve  
-1



① unsupervised.

② k-means.  $\rightarrow$   $K=?$

③ how to choose optimal  $K=?$

① random centroid.

② Distance, min.

③ update the centroid

ELBOW WCSS

④  $K25 \rightarrow$  how to validate.

Silhouette.  $\rightarrow [-1, +1]$   
 $\uparrow \quad \uparrow$   
 worst best

② hierarchical.

③ DBSCAN

④ Partial.

① you have to revise

② k means v/s kmeans++

③ Dunn index.

}

# How to make a best model or optimize sol<sup>n</sup>.

Custom learning or Custom model = Supervised + Unsupervised.  
(Semi Supervised)

