

★ Unsupervised Learning

- ① K-mean
- ② DB Scan
- ③ Hierarchical

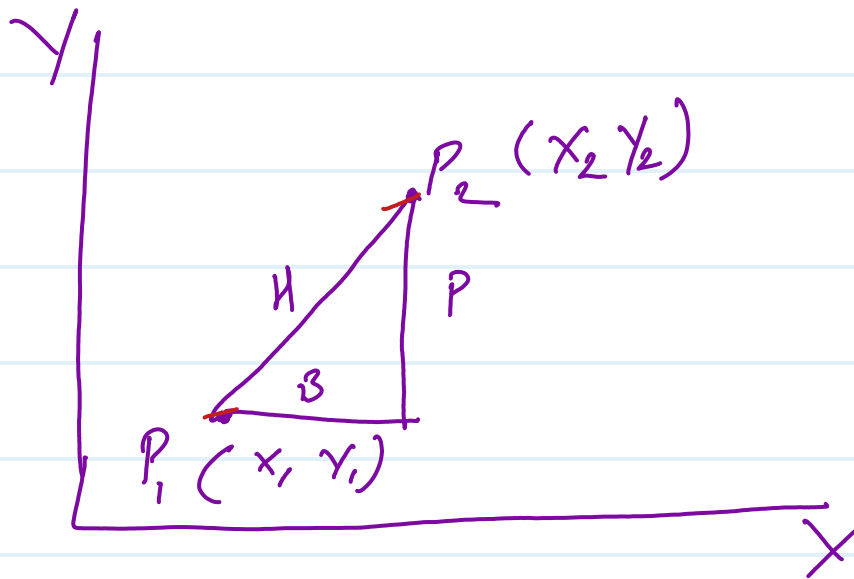
K-mean clustering

	Height	weight	cluster
①	185.	72.	
①	170.	56.	
	168	60	
	179	68	
	182	72	
	188	77	
	180	71	
	160	70	
	183	84	
	180	88	

* Centroid Based approach

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

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



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

*

- ① Centroid ✓
- ② Distance ✓
- ③ mean ✓

To find centroid we two different method

- ① Elbow
- ② WCSS

* Evaluation matrix

- ① Dunn Index
- ② silhote coeft.

For now we will take random centroid from the dataset

- ① 185 72 ✓
- ② 170 56 ✓

1st point

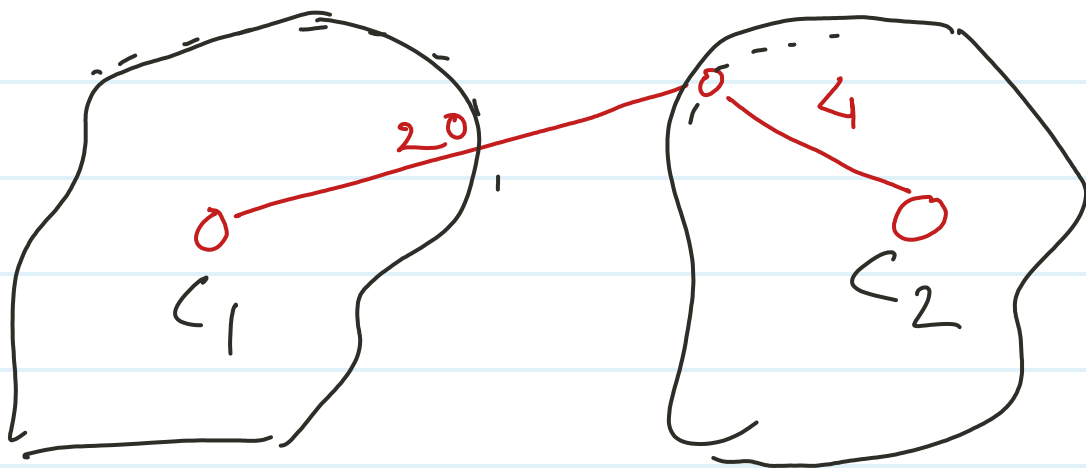
$$168 - 60$$

$$\begin{aligned} C_1 &= \sqrt{(168 - 185)^2 + (60 - 72)^2} \\ &= \sqrt{(-17)^2 + (-12)^2} \\ &= \sqrt{289 + 144} \\ &\Rightarrow \sqrt{433} \end{aligned}$$

$$C_1 = 20.6$$

$$\begin{aligned} C_2 &= \sqrt{(168 - 170)^2 + (60 - 56)^2} \\ &\Rightarrow \sqrt{4 + 16} \\ &\Rightarrow \sqrt{20} \end{aligned}$$

$$C_2 = 4.4$$



For point 1 4.4 distance is from C_2 and 20.6 is from C_1 , so the least distance is 4.4, so that the 1st point will be group with C_2 .

* Update value of C_2

$$\text{new } C_2 = \frac{168 + 170}{2}, \frac{56 + 60}{2}$$

$$\text{new } C_2 = (169, 58)$$

centroid

formula of WCSS

$$WCSS = \sum_{k=1}^k \sum_{i=1}^{n_k} \left[\text{Distance}(x_i, \mu_k) \right]^2$$

k = number of clusters

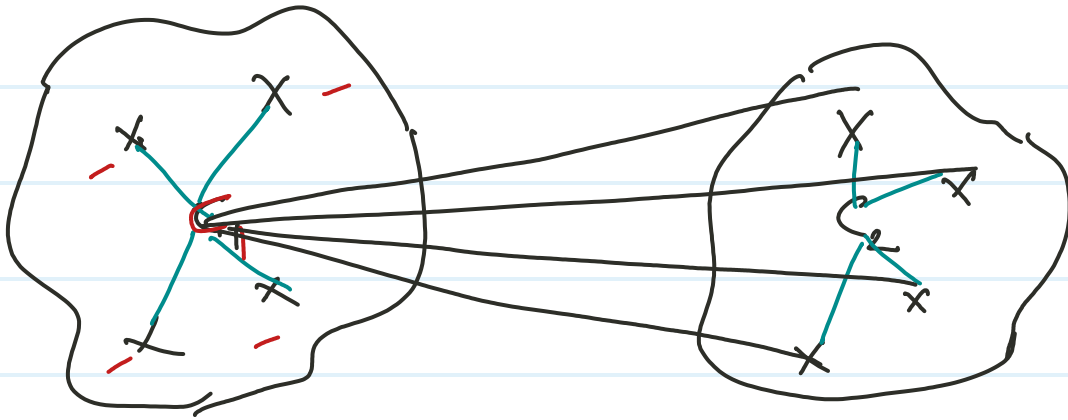
i = The datapoint in cluster k

μ_k = centroid of cluster k

n_k = The number of point in cluster k

$\text{Distance}(x_i, \mu_k)$ = The Euclidian distance b/w data point x_i and assign cluster μ_k

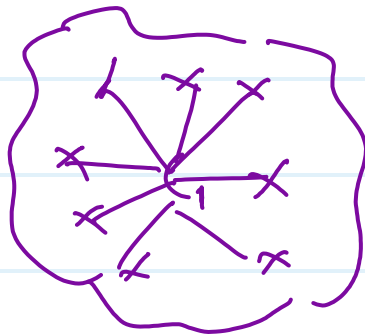
* WCSS (within cluster sum of square)



Intra cluster.

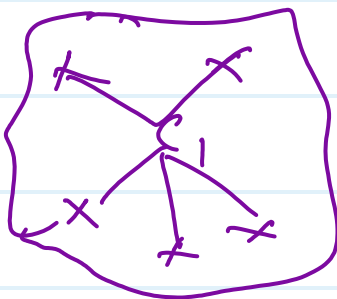
Inter cluster.

$k = 1$

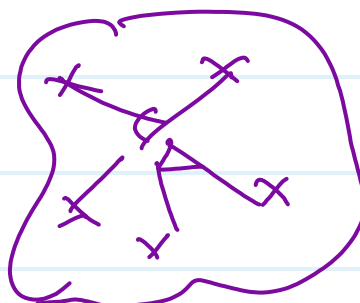


$$WCSS_1 = 5$$

$k = 2$



$WCSS_1$



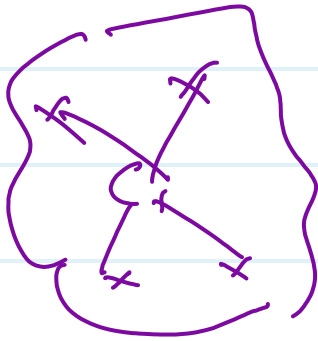
$WCSS_2$

$$WCSS_1 > WCSS_2$$

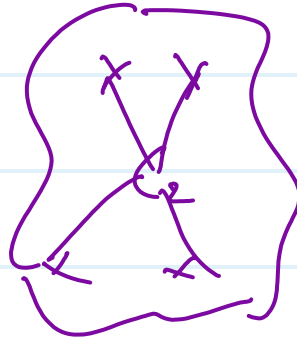
$$4.4 \quad 2.5$$

$$WCSS = \underline{\underline{2.5}}$$

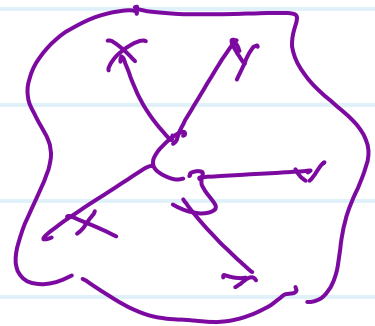
$$k = 3$$



W_1



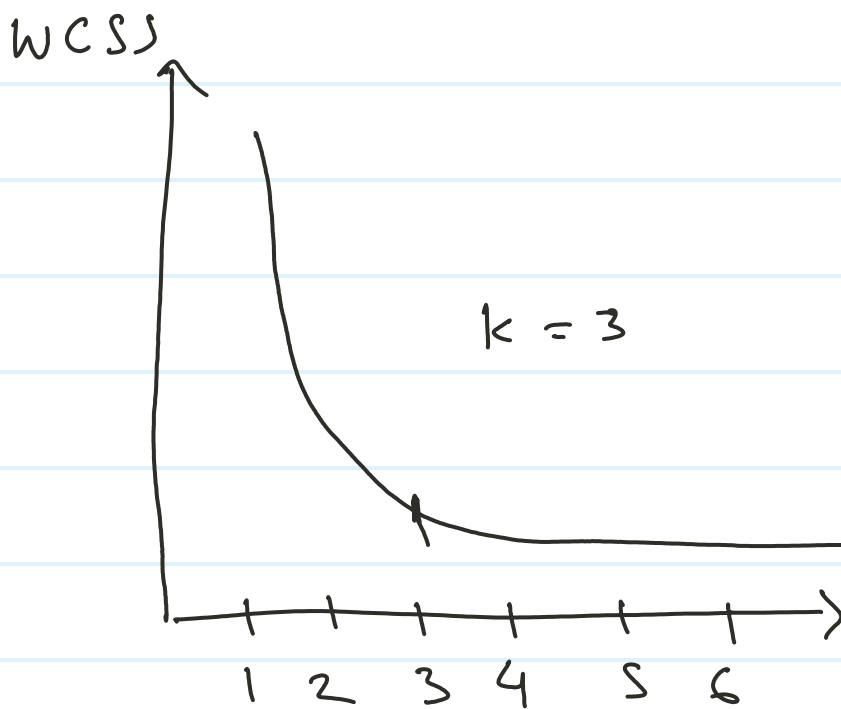
W_2



W_3

$$WCSS_1 > WCSS_2 > WCSS_3$$

$$WCSS_3 = 2$$



★ Dunn Indexing. -

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

★ silhouette coeff.

$$\Rightarrow \frac{b_i - a_i}{\max(b_i - a_i)}$$

a_i = intra cluster

b_i = inter cluster



intra

inter dust

Range of silhouette score is

-1 to +1

-1 = worst

+1 = best



0.4

-0.7

* Agglomerative clustering (Hierarchical clustering)

A, B, C, D, E, F, G,

