# Clustering

## Machine Learning

```
           Machine Learning
              /        \
      Supervised      Unsupervised
      ML algorithm     ML algorithm
       /     \            |
                          ├─ ① K-means Clustering
  Regression  Classification        (K-means++)
                          |
                          ├─ ② Hierarchical
                          |        Clustering
                          |
                          └─ ③ Dbscan
```

In case of Supervised

| Height (I/P) | weight | BMI (O/P) |
|---|---|---|
| 180 | 70 | 21 |
| 160 | 80 | 23 |
| 150 | 75 | 23 |
| 140 | 90 | 24 |
| 130 | 100 | 25 |
| 170 | 60 | 19 |

In case of Unsupervised

| Height | weight | BMI |
|---|---|---|
| 180 | 70 | normal ① |
| 160 | 45 | under weight ② |
| 150 | 70 | normal ① |
| 140 | 90 | normal ① |
| 130 | 100 | over weight ③ |
| 170 | 60 | normal. ① |

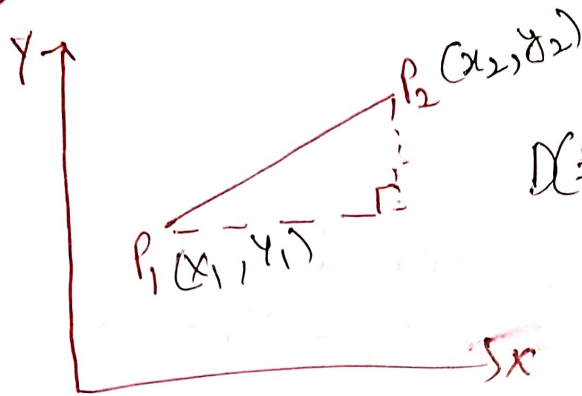→ we can • make clusterc by using BMI values.

-) ③ groups

① K-means

→ Unsupervised algorithm works by calculating
similarity score.

-) ① Eucledean distance.  ③ cosine similarity.
② Manhatten distance.

① Euclidean distance



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

② Manhatten distance

$$D(P_1 - P_2) = |x_1 - x_2| + |y_1 - y_2|$$

③ cosine similarity

$$\cos(\theta) = \frac{A \cdot B}{|A| |B|}$$

| Height | weight | | $\rightarrow$ Clustering. using K-means. |
|---|---|---|---|
| 185 | 72 | $C_1$ | |
| 170 | 56 | $C_2$ | |
| ③ 168 | 60 | | |
| ④ 179 | 68 | | |
| ⑤ 182 | 72 | | |
| ⑥ 188 | 77 | | |
| ⑦ 180 | 71 | | |
| ⑧ 183 | 84 | | |
| ⑨ 180 | 88 | | |
| ⑩ 180 | 67 | | |
| ⑪ 167 | 76 | | |

## Step ① Intialization of centroid

In case of K-means centroid intialization is random. Let

Let $K = 2$ ( $K$ = no. of centroid )

Lets assume 1st and 2nd row as $\textcircled{C_1}$ and $\textcircled{C_2}$

## Step ② Euclidean distance.

So to make cluster we have to find euclidean distance to other point from $\textcircled{C_1}$ and $\textcircled{C_2}$ (two cluster)
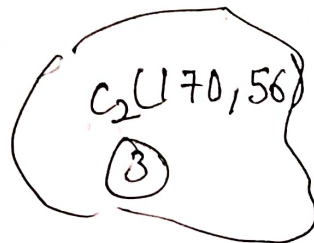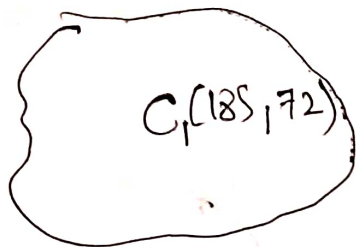
$$D(C_1, 3) = \sqrt{(185-168)^2 + (72-60)^2}$$
$$= \sqrt{(17)^2 + (12)^2} = 20.8$$

$$D(c_2, 3) = \sqrt{(168-170)^2 + (60-56)^2}$$

$$= \sqrt{2^2 + 4^2} = 4.4$$
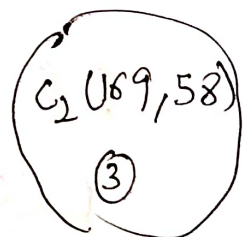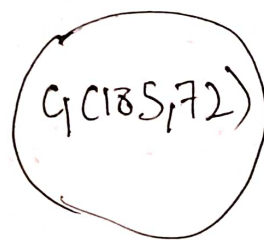
so $\quad D(c_2, 3) < D(c_1, 3)$

$C_1(185, 72)$

$C_2(170, 56)$

③

As ③ belongs to $C_2$ (cluster-2) $\quad C_2$ need to be updated.

### step-3

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

$$= (169, 58)$$

$C_1(185, 72)$

$C_2(169, 58)$

③

Again $\quad$ step ②

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

$$= \sqrt{6^2 + 4^2} = 7.07$$

$$D(c_2, 4) = \sqrt{(169 - 179)^2 + (58 - 68)^2}$$
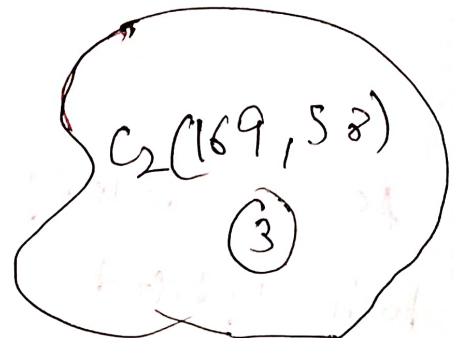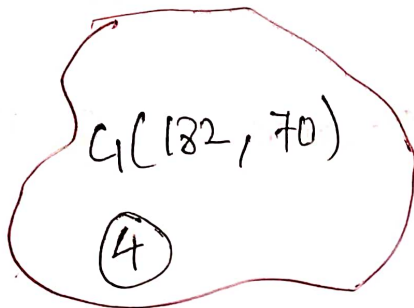
$$= \sqrt{10^2 + 10^2} = 14.14$$

As $\quad D(c_1, 4) < D(c_2, 4)$

∴ ④ belongs to $C_1$ (centroid -1)

## Step-3

As 4 belongs to $c_1$, $c_1$ centroid needs to update

$$\text{new } c_1 = \frac{185 + 179}{2}, \frac{72 + 68}{2}$$

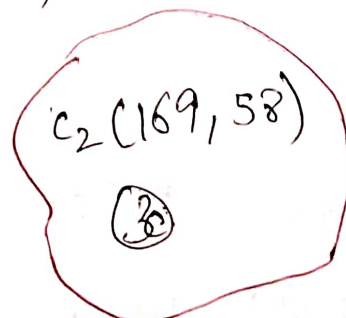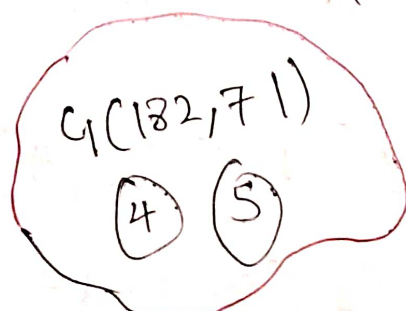$$= (182, 70)$$

$c_1(182, 70)$ ④

$c_2(169, 58)$ ③

## Again step ②

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

$$D(c_2, S) = \sqrt{(182-169)^2 + (72-58)^2}$$

$$= \sqrt{(13)^2 + (14)^2} = 19.1$$

As $D(c_1, S) < D(c_2, S)$

↓ step ③ Repeted.

$$\text{new } c_1 = \left(\frac{182+182}{2}, \frac{70+72}{2}\right)$$

$$= (182, 71)$$

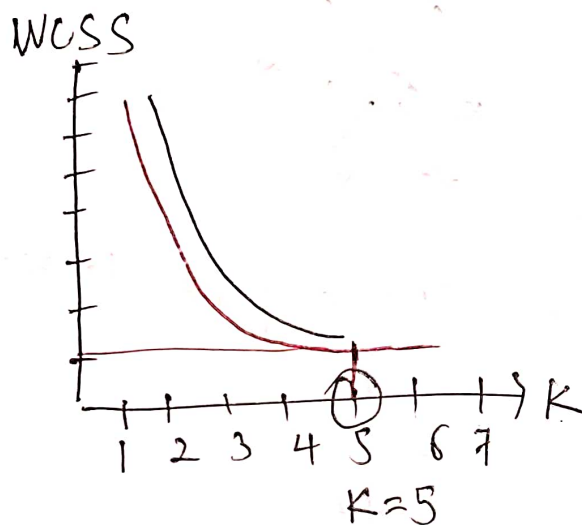$c_1(182, 71)$ ④ ⑤

$c_2(169, 58)$ ③

## Steps
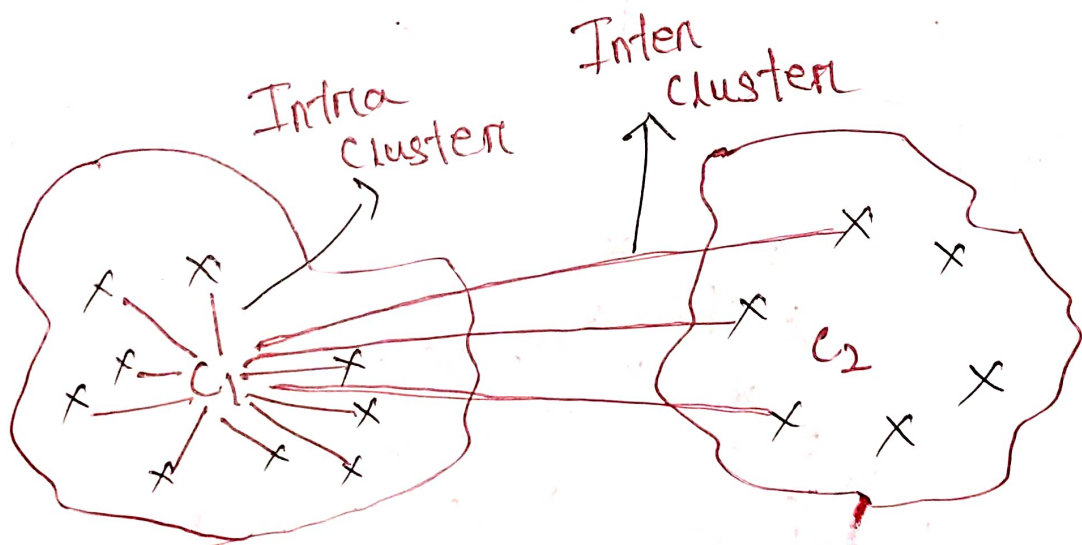
① centroid

② Distance ( compare minimum distance)

③ include point in clusters & update centroid.

Q. How to decide k value ?.

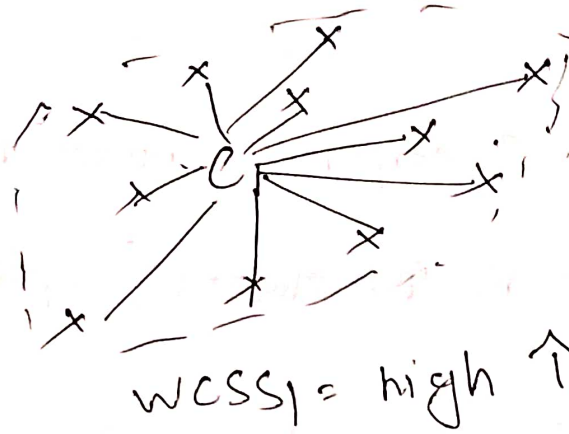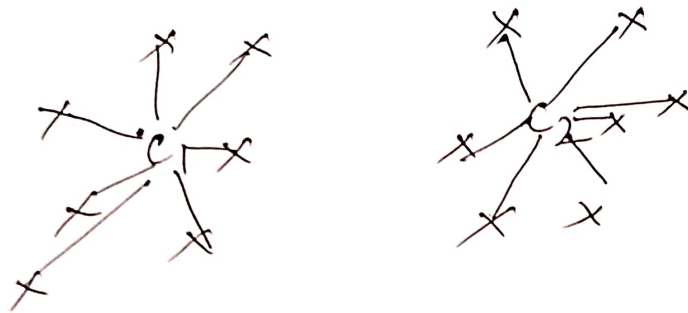Ans. ELBOW method.



$K = 5$
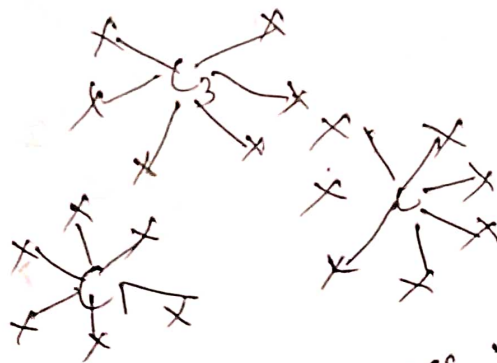
WCSS → within-cluster sum of square.

Let $k = 1$



$$WCSS_1 = high \uparrow$$

$k = 2$



$$WCSS_2 = Low \; (As \; compare \; to \; k=1)$$

$$WCSS = \sum_{i=1}^{n} d(c_i, x_i)^2$$  $WCSS_1 > WCSS_2$

$k = 3$



$$WCSS_1 > WCSS_2 > WCSS_3$$

wess is minimum at specific value of K
assume ( K = 5)

★ How to validate cluster value (K = 5]

① Dunn index.

② Silhouette score

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

② Silhouette score $= \dfrac{b_i - a_i}{\max (b_i - a_i)} \longrightarrow$ [−1 to 1]

In silhouette score

$a_i =$ intra cluster distance

$b_i =$ inter cluster distance.

If $a_i > b_i$

silhouett score −ve (wrost value)

→ In case of K-means ++ , random initialization
of K-means problem is solved.

→ Kmeans is a centroid based algorithm.