Unsupervised ML

to discover
groups in the
data

Datasets

→ target column ——→ missing
(unlabelled data)

clustering

Semi-supervised Technique

↳ 1000 Records

real time application

① Customer Segmentation
(Marketing, EdTech,
fintech)

400 Records

↳ labeled data

600 Records
↳ unlabeled
data

② Image Compression

③ Anomaly Detection

④ Document/topic grouping
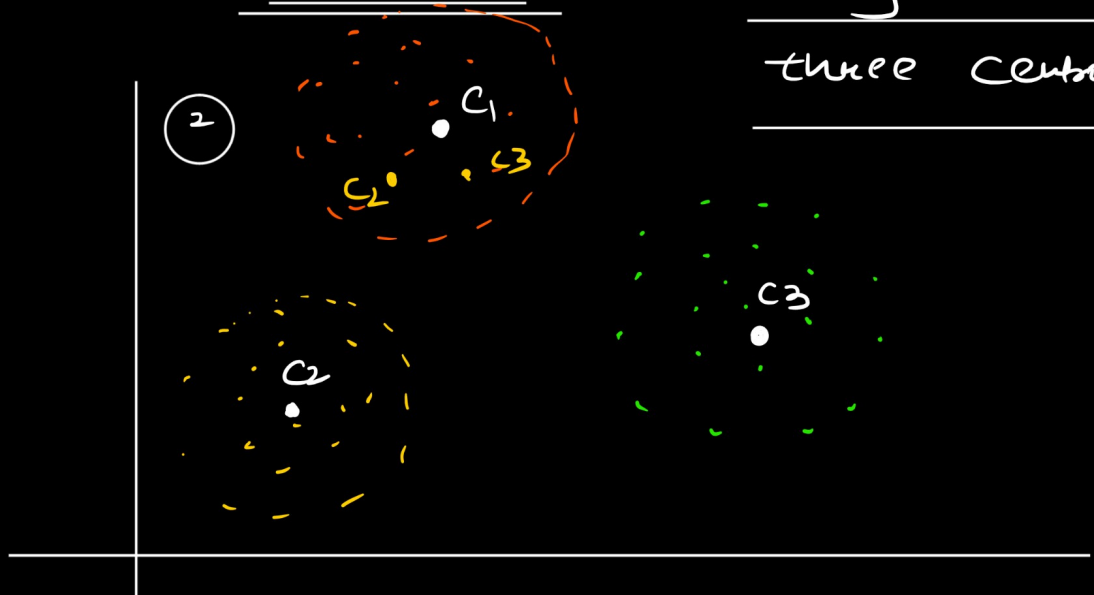
Clustering

←← Labels

Supervised
ML

Distance measures :-

① Euclidean distance
② Manhattan distance

Numeric Points

③ Hamming distance

④ Cosine similarity (for textual data)

K-Means Clustering ⟶ customer segmentation

( # Clusters

① $k = 3$ ⟶ Randomly initialize three centroids
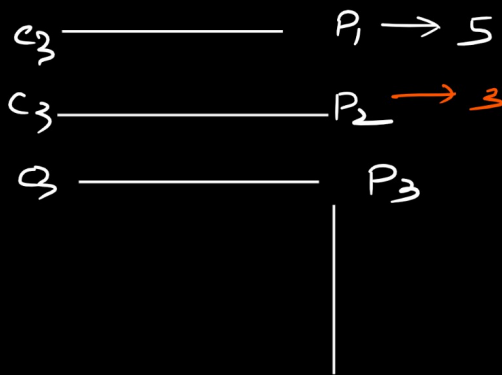
②



$C_1$
$C_2$
$C_3$

$C_2$

$C_3$

③ Re-evaluation of centroids $(C_1, C_2, C_3)$

$$c_1 = \frac{P_1 + P_3 + P_5}{3}$$

$$c_2 = \frac{P_2 + P_4 + P_6}{3}$$

$$c_3 = \frac{P_7 + P_8 + P_9 + P_{10}}{4}$$

$c_1$ —————— $P_1 \to 2$         $C_2$ —————— $P_1 \to 1$

$C_1$ —————— $P_2 \to 1$         $C_2$ —————— $P_2 \to 2$

$C_1$ —————— $P_3$              $C_2$ —————— $P_3$

$C_1$ ————                      ————

$C_1$ —————— $P_{10}$           $C_2$ —————— $P_{10}$

$c_3$ —————— $P_1 \to 5$

$C_3$ —————— $P_2 \to 3$

$C_3$ ———— $P_3$

$C_3$ —————— $P_{10}$

$$\left( \frac{x_1 + x_3 + x_4 + x_7}{4} , \right.$$

$\xleftarrow{\quad}$ $C_1$ $\xrightarrow{\quad}$ $\underline{P_1 , P_3 , P_4 , P_7}$   $(x_1, y_1)$ $(x_3, y_3)$ $(x_4, y_4)$ $(x_7, y_7)$

$$\left. \frac{y_1 + y_3 + y_4 + y_7}{4} \right)$$

$C_2$ $\xrightarrow{\quad}$ $\underline{P_2 , P_5 , P_6 , P_8}$

$C_3$ $\xrightarrow{\quad}$ $\underline{P_9 , P_{10}, P_{11} , P_{12}}$

## Steps of KMeans

1. Choose number of clusters k ✔ k.
2. Initialize centroids randomly.
3. Assign points to nearest centroid.
4. Update centroid as mean of assigned points.
5. Repeat until convergence.

## Advantages

↳ simple yet

1. efficient for large datasets.

2. Works well with spherical clusters

## Disadvantages

→ Elbow Method

1. Requires # clusters (k) beforehand.

2. Sensitive to outliers

3. Randomly initialize centroids (within 1 cluster)
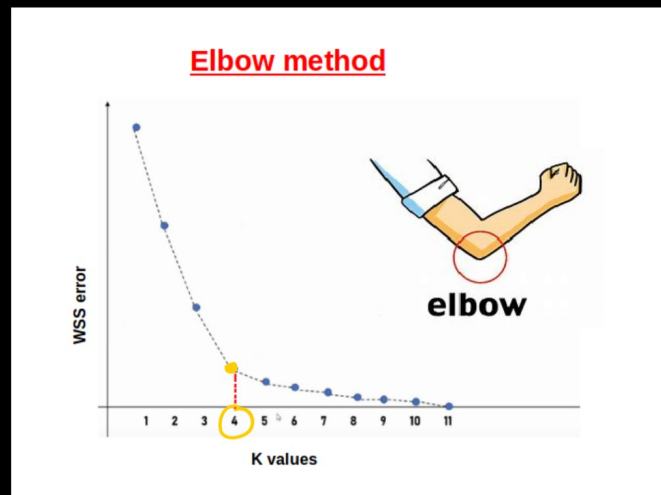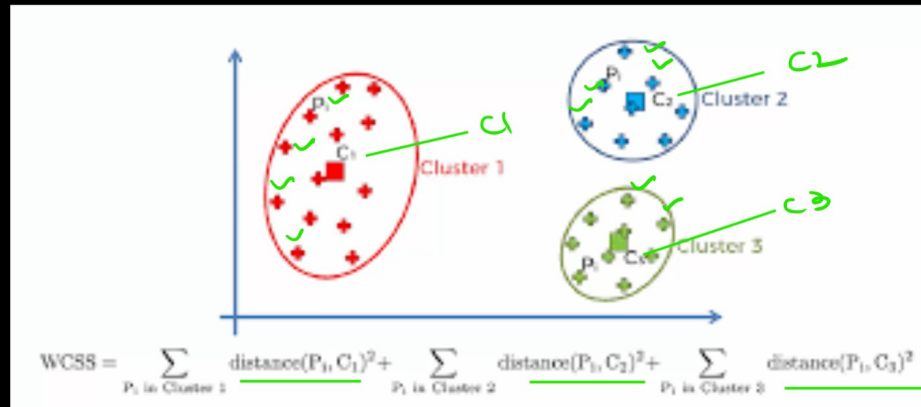
Task ← K means ++

# Elbow Method

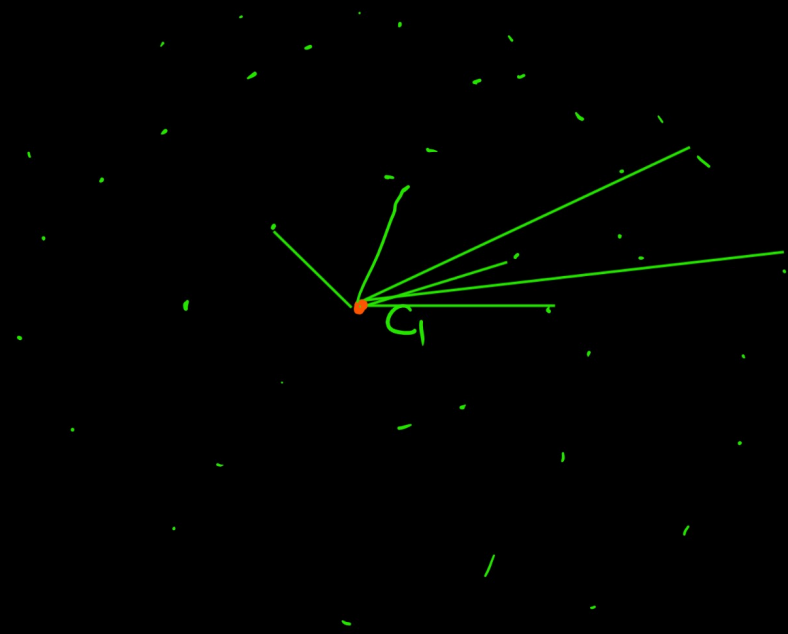WCSS $\longrightarrow$ Within Cluster Summation of Squares

or

Inertia



$$WCSS = \sum_{P_i \text{ in Cluster 1}} distance(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster 2}} distance(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster 3}} distance(P_i, C_3)^2$$



Elbow method

$K = 1$ ——————————— $WCSS = max \, (Peak)$

$K = n$ ——————————— $WCSS = 0$

$$n = 4$$

$c_1 \quad c_2 \quad c_3 \quad c_4$