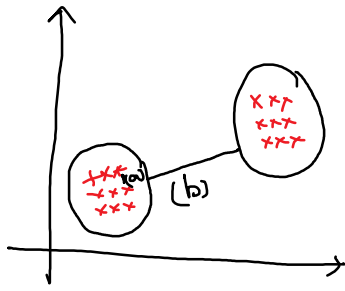# Unsupervised learning

$D \Rightarrow \{x_i, y_i\} \rightarrow$ supervised learning

$D = \{x_i\} \rightarrow$ unsupervised learning

Clustering
↓
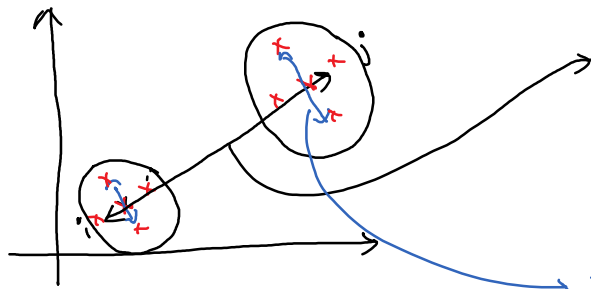K Means

Applications:  e-commerce:  group customers → location, income, gender
etc.

⇒ Review Analysis :     Amazon Reviews

+ve                    — ve

→ <u>Image Segmentation:</u>



<u>Metrics</u>

$\Rightarrow$ Intercluster distance (b)

$\Rightarrow$ Intracluster distance

Characteristics:   a) $\rightarrow$ intracluster distance should be small

(good cluster)   b) $\rightarrow$ intercluster distance should be large

Dunn's Index $\Rightarrow$   $\dfrac{\max d(i,j) \longrightarrow \text{intercluster distance}}{\max d'(k) \longrightarrow \text{intracluster distance}}$



$\max d(i,j) \Rightarrow$ distance b/w farthest points in different cluster.

$\max d'(k) \Rightarrow$ distance b/w farthest points in same cluster

b) Silhouette's Score $\Rightarrow$ $\dfrac{b-a}{\max(b,a)} \longrightarrow$ $b \Rightarrow$ avg intercluster distance

$[-1, +1]$   $a \Rightarrow$ avg intracluster distance

Case 1:   $a \Rightarrow \min \Rightarrow 0$   $b \Rightarrow b$ (best scenario)

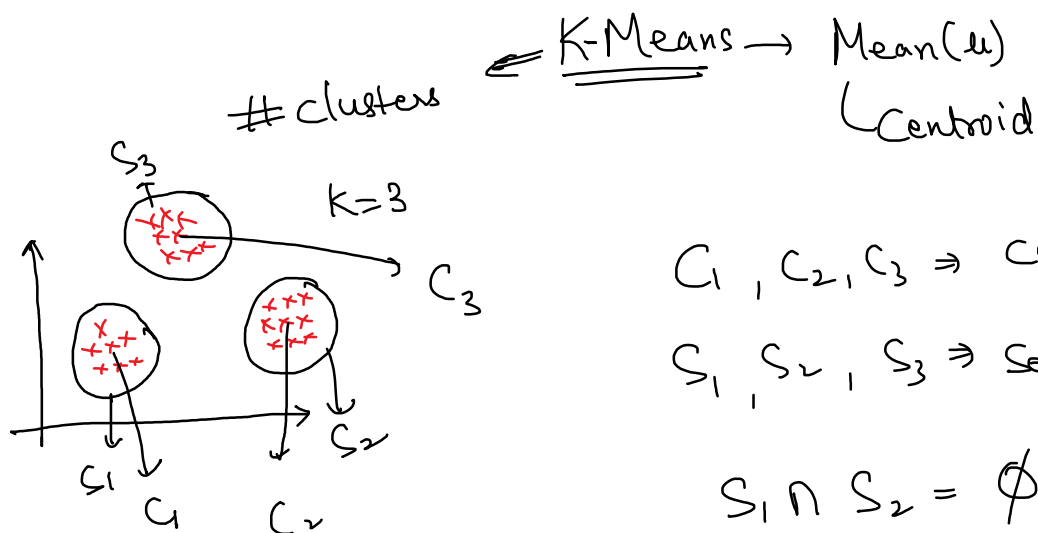$$SS = \frac{b-0}{\max(b,0)} = \frac{b}{b} = 1$$

Case 2:   $b < a$, $b = 0$, $a = a$      completely wrong

$$SS = \frac{0-a}{\max(0,a)} = \frac{-a}{a} = -1$$

Case 3:   $a = b$

$$SS = \frac{a-a}{\max(a,a)} = \frac{0}{a} = 0$$

#clusters  $\Leftarrow$ K-Means $\rightarrow$ Mean($\mu$)
$\qquad\qquad\qquad\qquad\qquad$ $\hookrightarrow$ Centroid



$K = 3$

$C_3$

$C_1, C_2, C_3 \Rightarrow$ Centroids

$S_1, S_2, S_3 \Rightarrow$ Sets

$S_1 \cap S_2 = \emptyset \quad S_3 \cap S_1 = \emptyset$

$S_2 \cap S_3 = \emptyset$

$$C = \frac{1}{n} \sum_{i=1}^{n} x_i$$
$$x_i \in S_i$$

MOF $\Rightarrow$ $C^* = \underset{C_1, C_2, C_3 \cdots C_k}{\arg\min} \sum_{i=1}^{k} \sum_{x \in S_i} \| x_i - C_i \|^2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\hookrightarrow$ intracluster distance

$\qquad\qquad\qquad\qquad\qquad\qquad$ $x_i \in S_i$

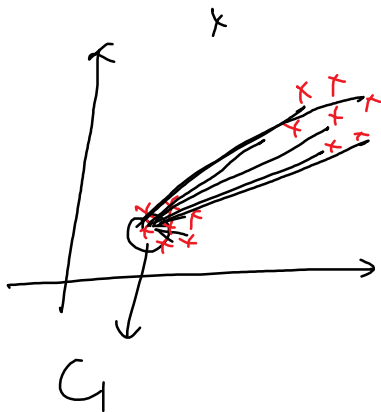$\qquad\qquad\qquad\qquad\qquad\qquad$ $S_i \cap S_j = \emptyset$ $\downarrow$

# LLoyd's Algorithm:

① Randomly choose k datapoints as centroids

② Assignment: for each point, select the nearest centroid with the help of distance & add that point to the corresponding cluster.

③ Updation: Recalculate centroids,

$$C_i = \frac{1}{S_i} \sum_{i=1}^{n} x_i$$

$$x_i \in S_i$$

④ Repeat step ② & ③ till convergence

# KMeans ++



| datapoints | distance |
|---|---|
| $x_1$ | $d_1$ |
| $x_2$ | $d_2$ |
| $\mid$ | $\mid$ |
| $\mid$ | $\mid$ |
| $x_n$ | $d_n$ |

distance $\propto$ prob. of being picked as
centroid

KMeans $\Rightarrow$ affected by outliers