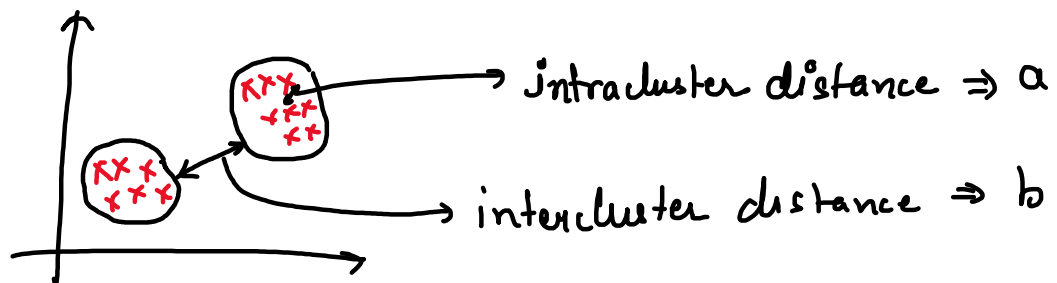Clustering → grouping
↓
$X_i$     └→ Unsupervised learning

Applications · ① e commerce ⇒ Customer segmentation

② Review Analysis

③ Image Segmentations

# Metrics



→ intracluster distance ⇒ a

→ intercluster distance ⇒ b

# Characteristics of a good cluster·

1) intracluster distance must be small
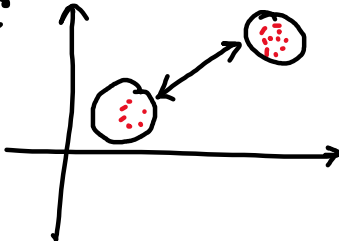
2) intercluster "      "    "   large

# Silhouette's Score:

b ⇒ avg intercluster

Silhouettes Score:

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

$b \Rightarrow$ avg intercluster distance
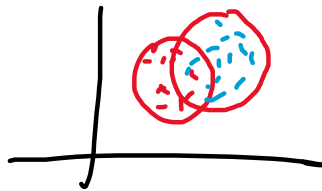$a \Rightarrow$ avg intracluster distance.

Case 1:



$a \Rightarrow 0 \Rightarrow \min \qquad b \Rightarrow b$

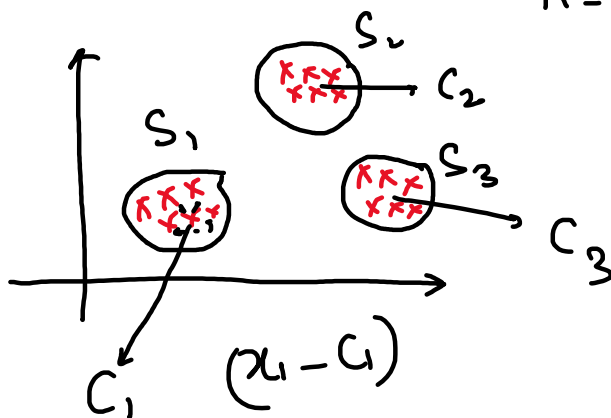$$SS = \frac{b - 0}{\max(b, 0)} = \frac{b}{b} = 1 \quad.$$

Case 2:



$a = b$

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

K Means $\longrightarrow$ Mean (Centroid)

\# clusters

$K = 3$



$C_1, C_2, C_3 \Rightarrow$ centroids

$S_1, S_2, S_3 \Rightarrow$ sets

$S_1 \cap S_2 = \emptyset$

$S_2 \cap S_3 = \emptyset$

$(x_1 - C_1)$

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

$$S_3 \cap S_1 = \emptyset$$

MOF: $C^* = \text{argmin} \sum_{i=1}^{\overline{S}} \sum_{x_i \in S_1} \|X - C_i\|^2$

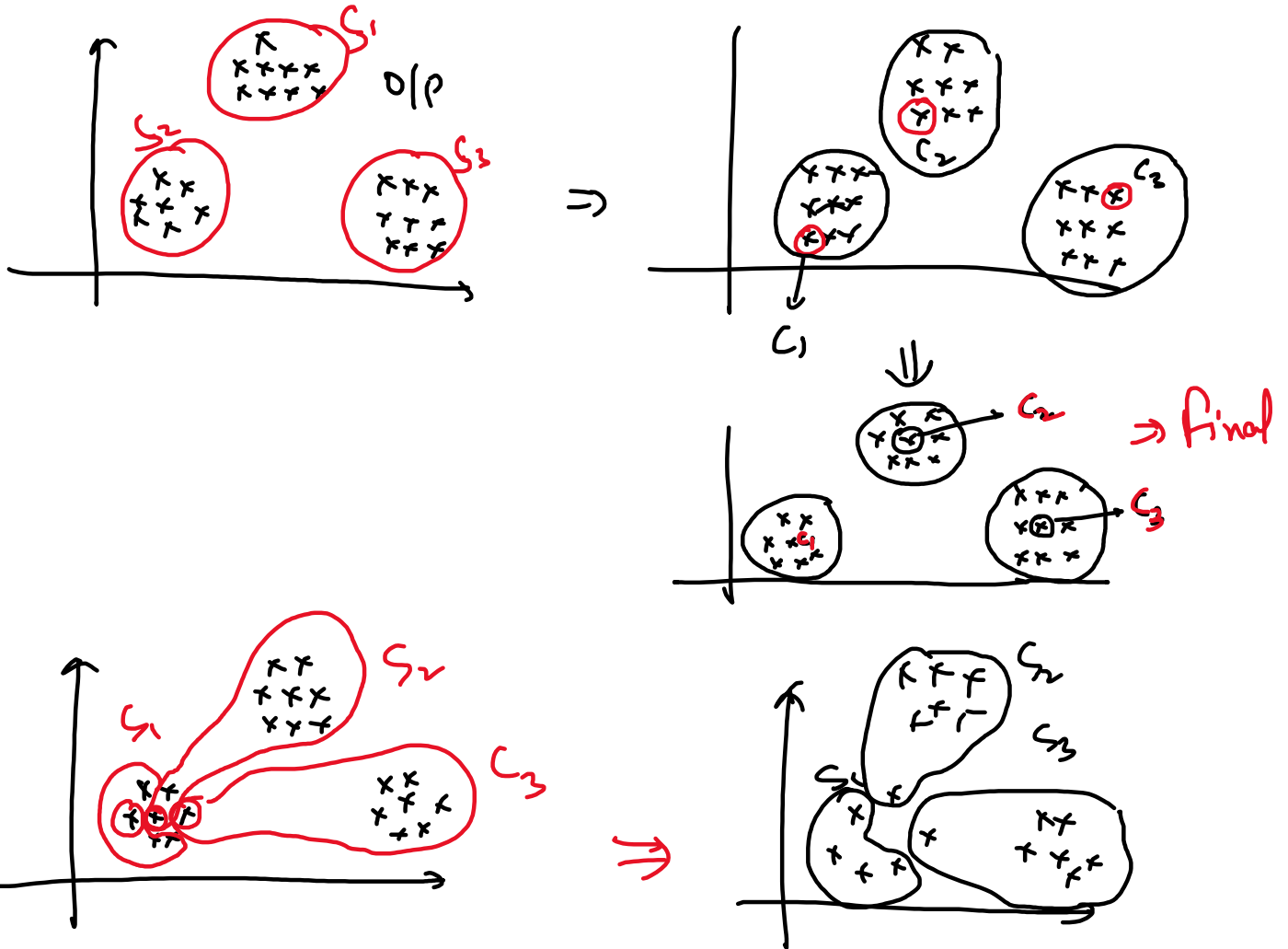$\hookrightarrow$ intracluster distance

np hard problem

# Lloyd's Algorithm

① Randomly choose pts as centroids.

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

③ Update. Recalculate Centroid
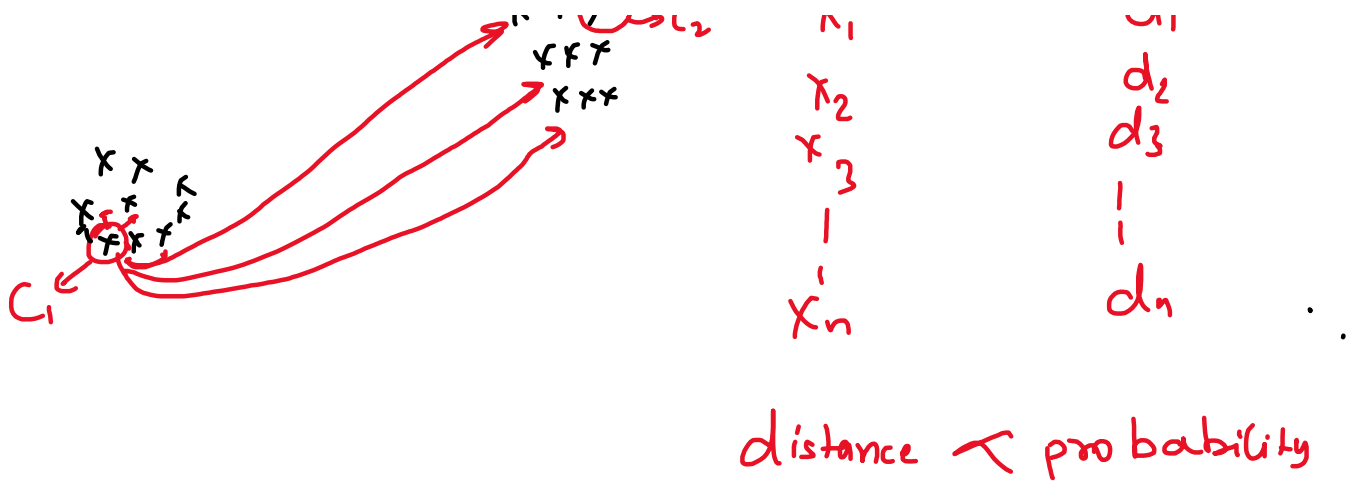
$$C_1 = \frac{1}{S_1} \sum_{i=1}^{n} x_i$$
$$x_1 \in S_1$$

**(4)** These **②** & **③** steps will be repeated till centroid stop changing



O/p

⇒

⇒ final

# KMeans++ :

**①** Choose only one random centroid

| datapoint | distance from $c_1$ |
|---|---|
| $x_1$ | $d_1$ |
| $x_2$ | $d_2$ |

$C_2$

$x_1$          $d_1$

$x_2$          $d_2$

$x_3$          $d_3$

$|$            $|$

$|$            $|$

$x_n$          $d_n$

distance $\curlywedge$ probability

KMeans is sensitive to outliers.