

Unsupervised Learning :

Requires minimum human supervision

does not require any knowledge of HUMAN-LABELLED DATA

2 different types of unsupervised learning:

- clustering,
- dimensionality reduction

Examples:

- recommendation system: to recommend customers:
online shopping purchases
movie recommendations

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_k - y_k)^2}$$

$d(\vec{x}, \vec{y})$

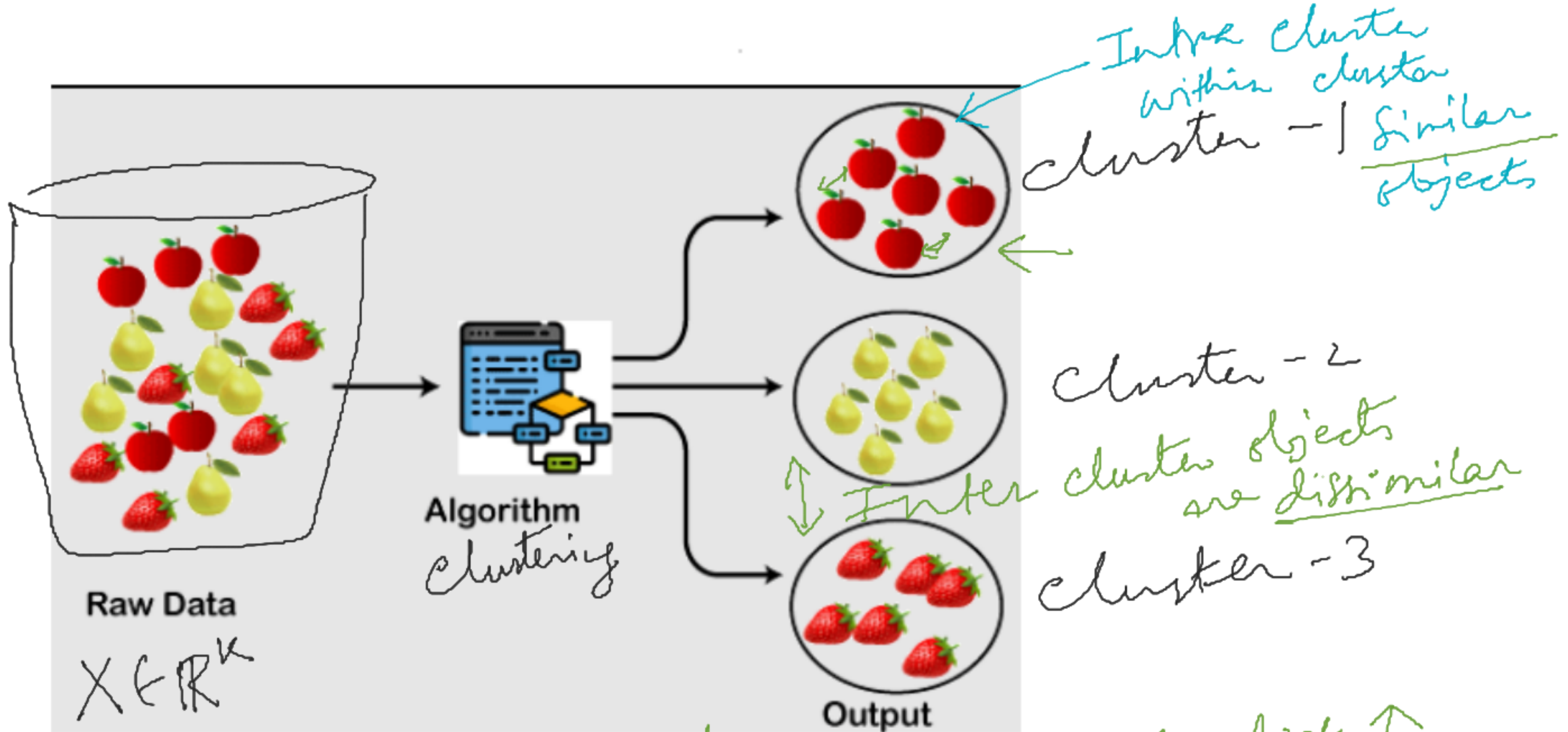
$$X \in \mathbb{R}^k$$

$$X$$

$$y \in \mathbb{R}^k$$

$$\begin{matrix} x_1 & x_2 \\ x_3 & \dots \end{matrix}$$

$$\begin{matrix} y_1 & y_2 \\ y_3 & \dots \end{matrix}$$



Intra cluster dist ↓

Inter cluster dist ↑

Clustering is an unsupervised learning technique

Objective of clustering:

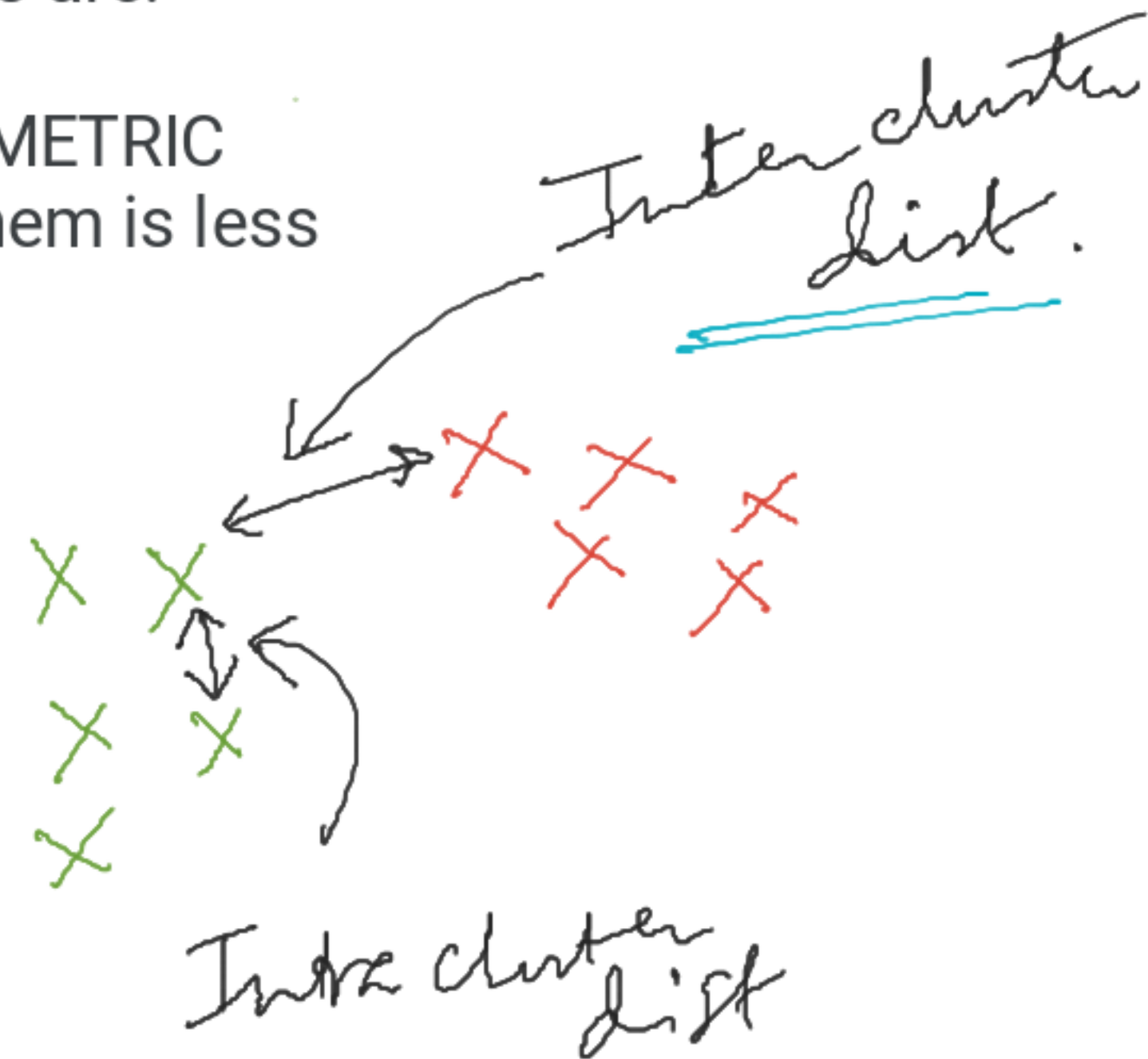
to discover overall distribution patterns
correlations among data distribution

Given the features of sample data to be clustered, objectives are:

- to put similar samples in the same cluster
- the similarity b/w 2 samples is measured by DISTANCE METRIC
- the similarity of 2 samples is more, if the distance b/w them is less

A clustering alog is considered to be good, if:

INTER-cluster distance b/w different clusters is more
INTRA-cluster distance of the same cluster is less



Inter-cluster dist. is the dis. b/w 2 sample data belonging to 2 diff. clusters

Different measures:

SLD

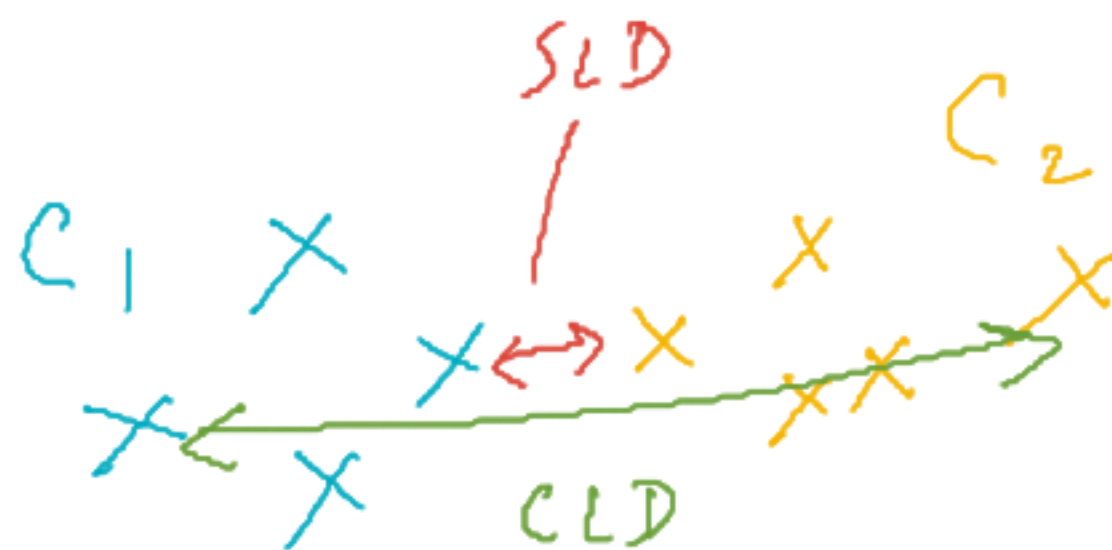
Single linkage distance: The minimum dist. b/w 2 sample data belonging to 2 diff. clusters

$$d(C_1, C_2) = \min_{\substack{x \in C_1, \\ y \in C_2}} d(x, y)$$

CLD

Complete linkage distance: The maximum dist. b/w 2 sample data belonging to 2 diff. clusters

$$d(C_1, C_2) = \max_{\substack{x \in C_1, \\ y \in C_2}} d(x, y)$$



Average linkage dist ... Avg. ...

$$\underline{d(c_1, c_2) = \frac{1}{|c_1| \times |c_2|} \sum_{\substack{x \in c_1, \\ y \in c_2}} d(x, y)}$$

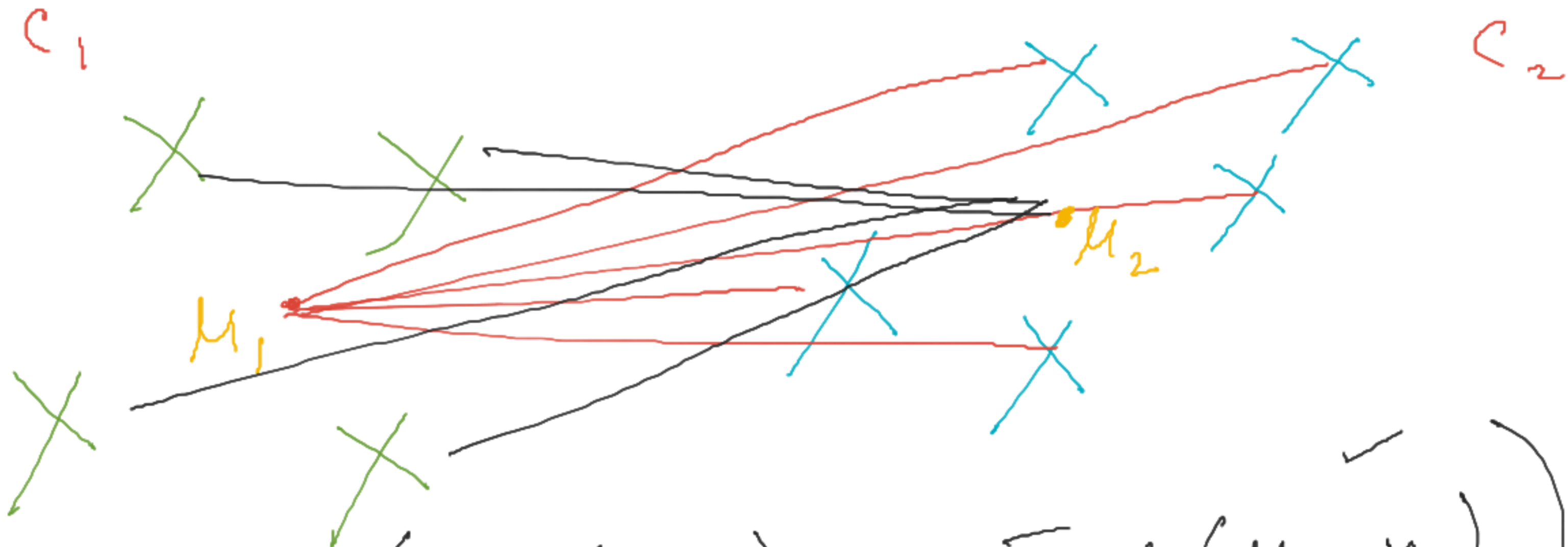


Centroid linkage distance: The dist. b/w 2 centroids belonging to 2 diff. clusters

$$\underline{d(c_1, c_2) = d(\mu_1, \mu_2)}$$

where,

$$\mu_1 = \frac{1}{|c_1|} \sum_{x \in c_1} x \quad \Bigg| \quad \mu_2 = \frac{1}{|c_2|} \sum_{y \in c_2} y$$

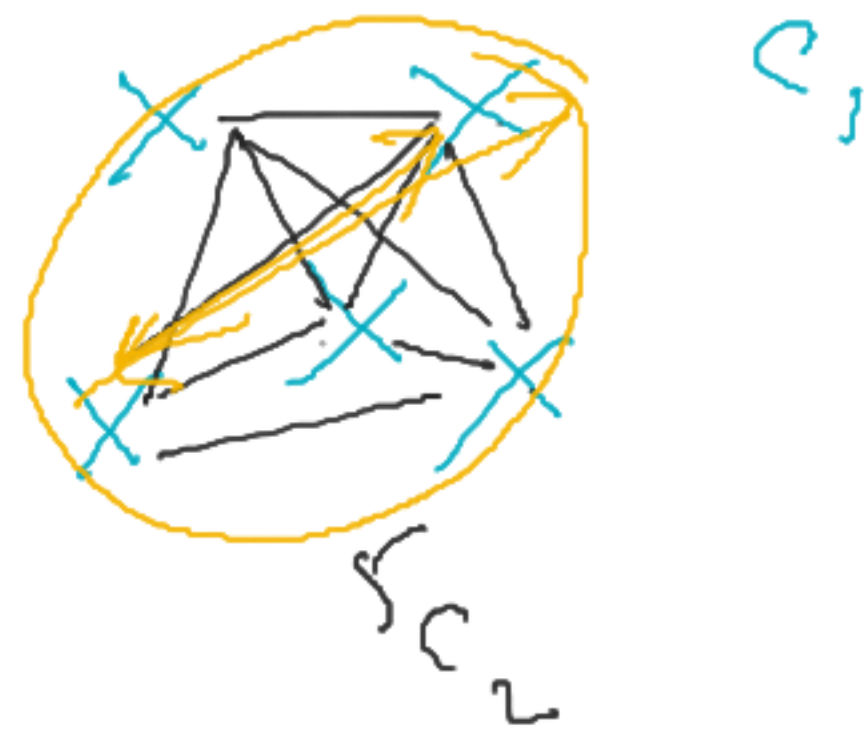


$$d(C_1, C_2) = \frac{1}{|C_1| + |C_2|} \left(\sum_{x \in C_1} d(x, \mu_2) + \sum_{y \in C_2} d(\mu_1, y) \right)$$

Average Centroid Linkage Distance

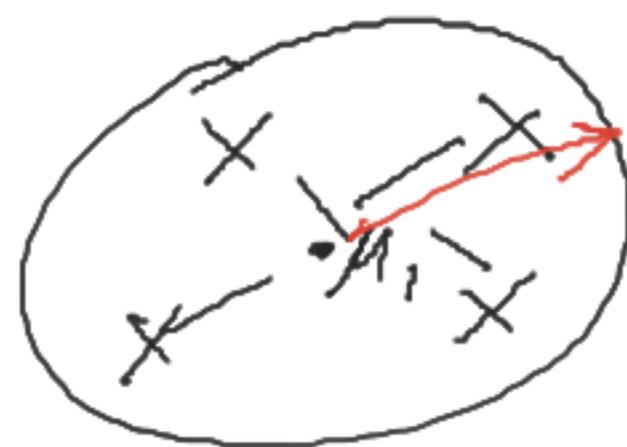
Intra-cluster dist: dist. b/w 2 data belonging to the SAME cluster

Complete Diameter Dist ✓
$$\Delta(c_1) = \max_{\substack{(x,y) \in c_1 \\ x \neq y}} d(x,y)$$



Centroid diameter dist

$$\Delta(c_1) = 2 \times \text{avg}_{x \in c_1} d(x, M_1)$$



Avg. diameter dist
$$\Delta(c_1) = \text{avg}_{\substack{(x,y) \in c_1 \\ x \neq y}} d(x,y)$$

A good clustering algo:

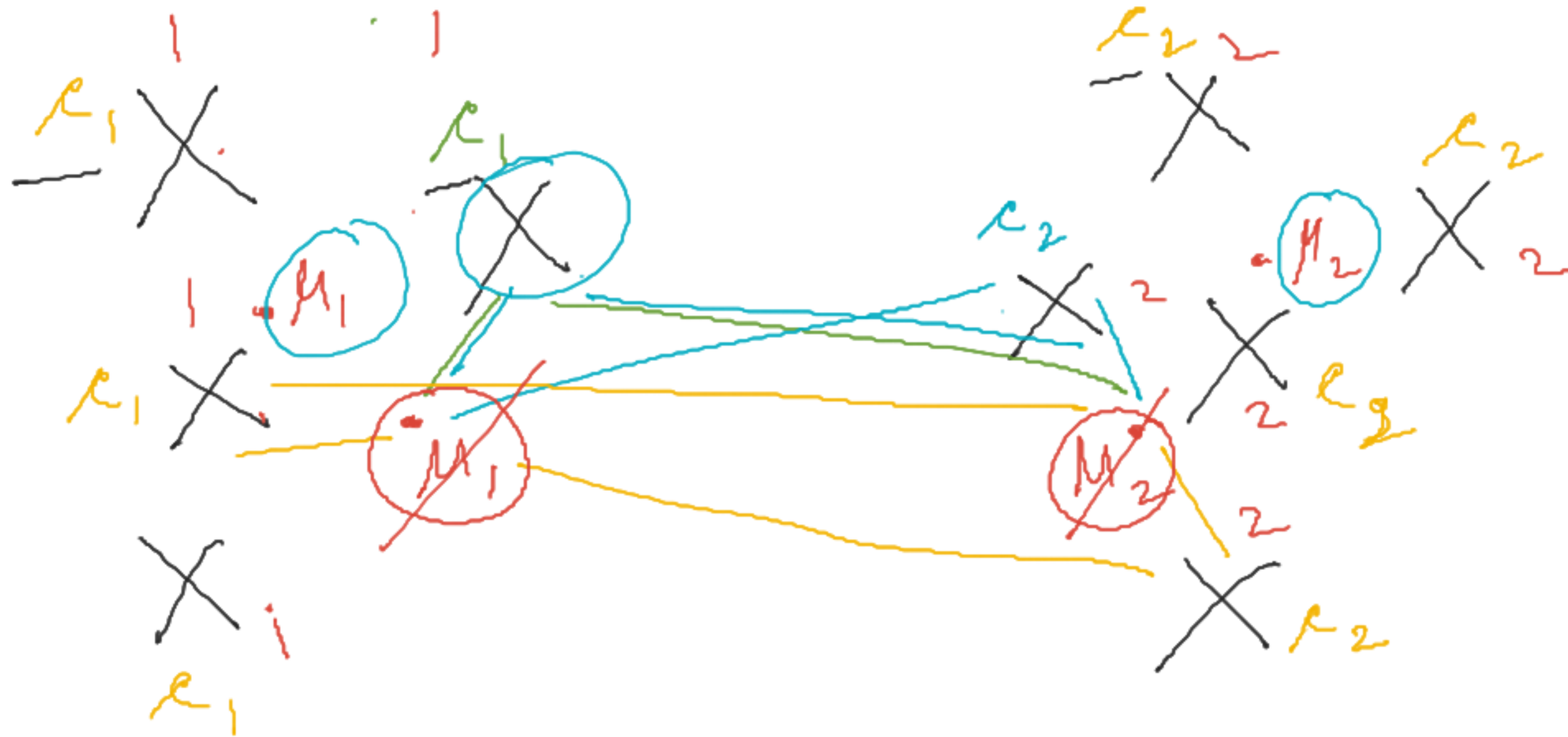
- A) more complete diameter dist
- B) less complete linkage dist
- C) less avg. diameter dist
- D) more centroid diameter dist

k-means

i/p: no. of cluster
(k)

$k=2$

$k=2$



K-means : i/p: ① Training data: $X = \{x_1, x_2, \dots, x_n\}$
where $\forall x_i \in \mathbb{R}^m$
② K no. of clusters $\{C_1, C_2, \dots, C_K\}$

Algo : Randomly initialize 'K' cluster centroids
 $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^m$

Repeat (until terminate) \leftarrow no update in centroid

{ for all $x_i \in X$:

$\rightarrow \underline{c_i} = \underset{j \in \{1, 2, \dots, K\}}{\operatorname{argmin}} \|x_i - \mu_j\|^2$ // cluster assignment step

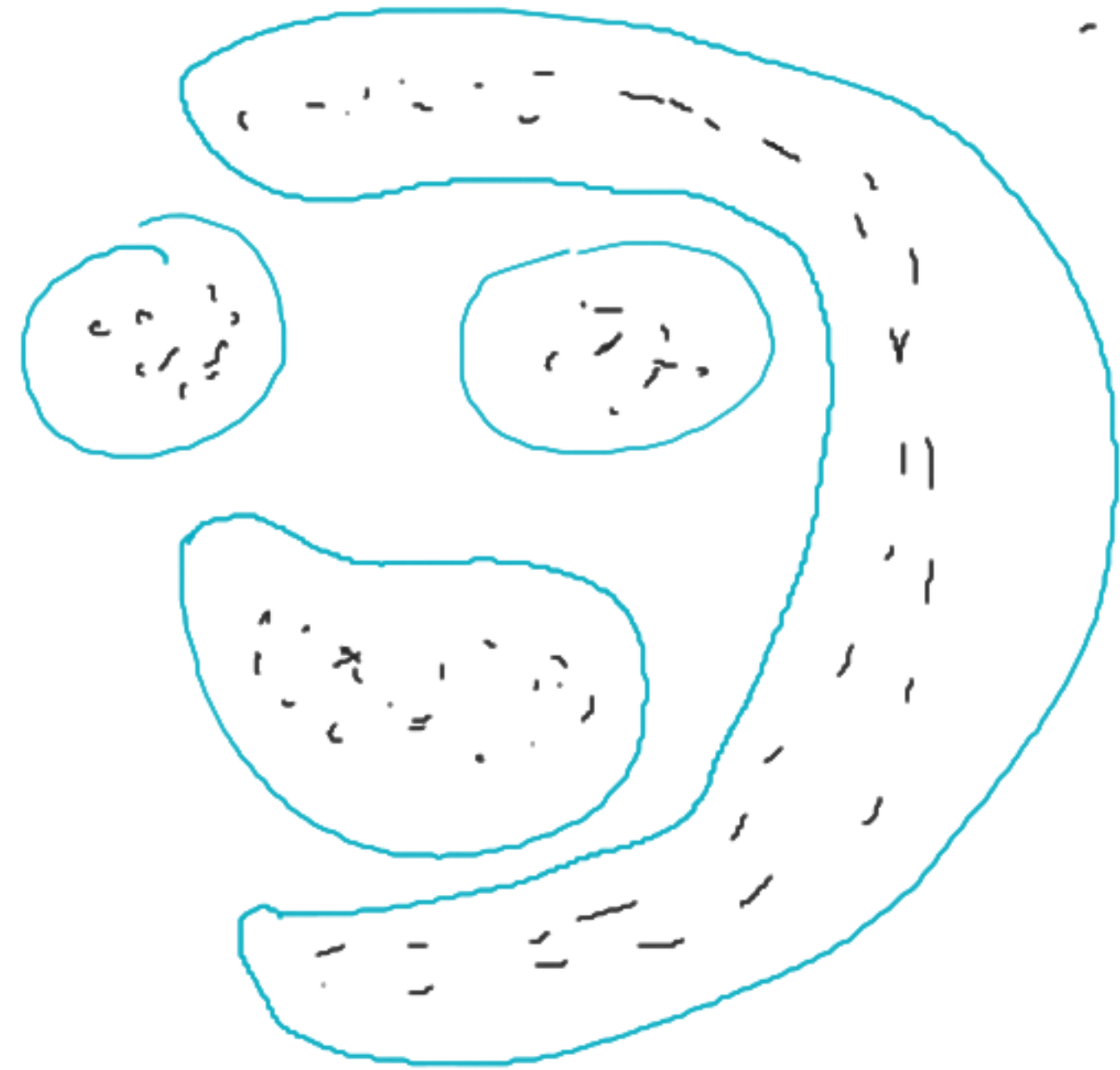
for all $\mu_i \in \{\mu_1, \mu_2, \dots, \mu_K\}$

$\rightarrow \mu_i =$ update the value from newly generated clusters // centroid update

}

Random initialization of centroid

Choosing "k"



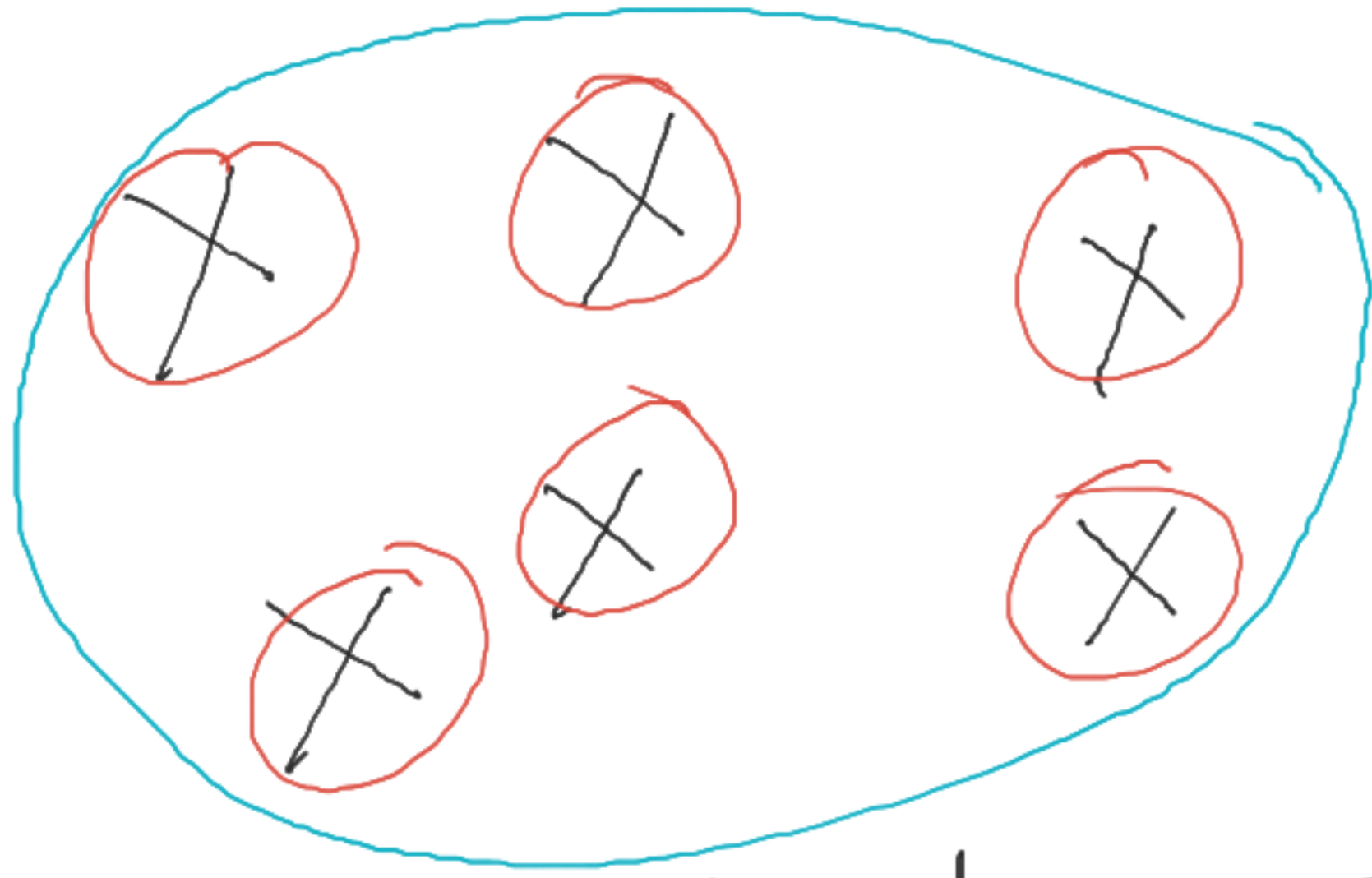
Globular data X

outlier X

→ (DBSCAN)

$$n = 100$$

$$k \leq \sqrt{n} = 10$$



K-mode

min # cluster = 1

max # cluster = n = 6

k-means
s/p: $\forall x_i \in X, (x_i, c_i); c_i \in \{C_1, C_2, \dots, C_k\}$
is the cluster index

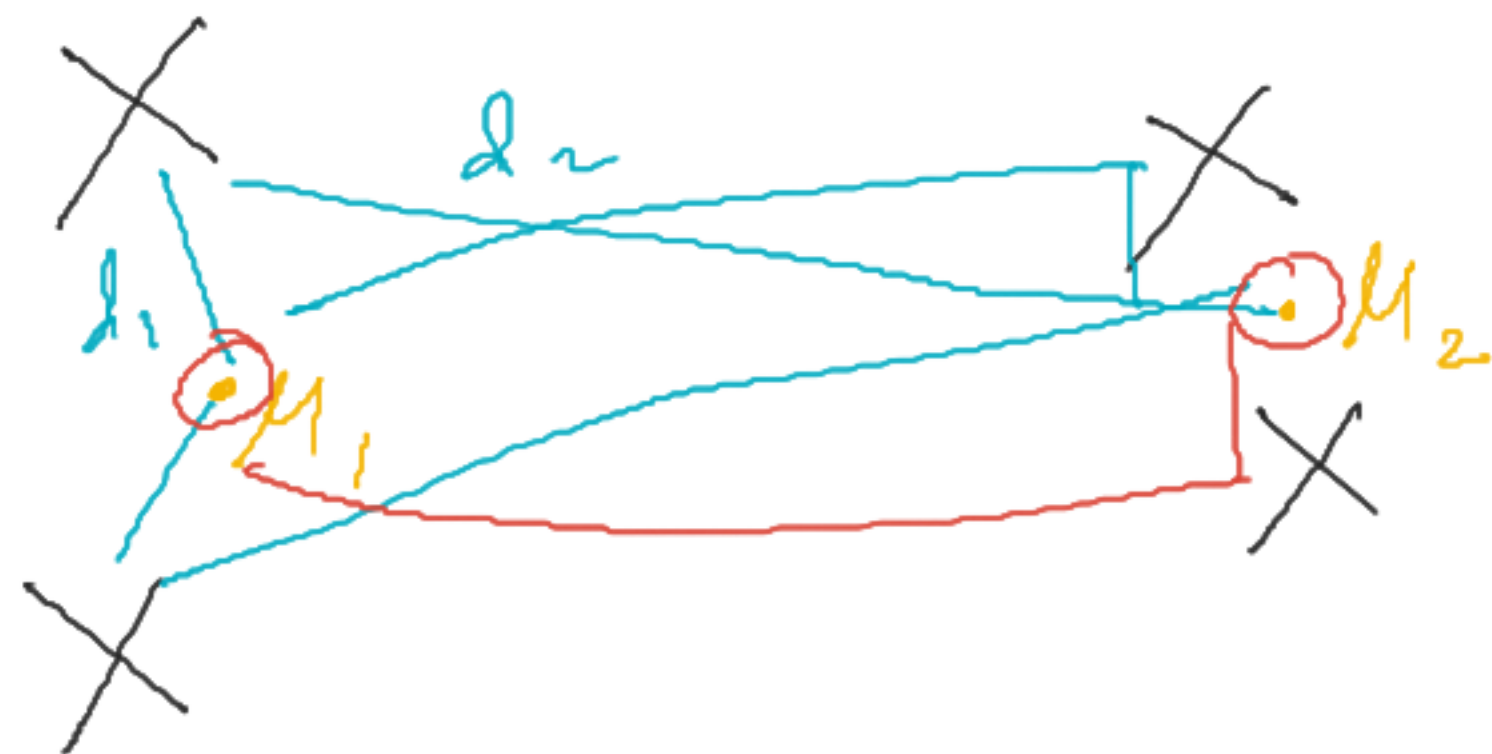
objective $\mu_i \in \mathbb{R}^m$ is the centroid of $C_i \in \{C_1, C_2, \dots, C_k\}$

optimization objective

minimize $J(x_1, x_2, \dots, x_n, \mu_1, \mu_2, \dots, \mu_k)$
 $n + k$ parameters

$$J(x_1, \dots, x_n, \mu_1, \dots, \mu_k) = \frac{1}{n} \sum_{i=1}^n \|x_i - \mu_{c_i}\|^2$$

$\mu_{c_i} \in \{\mu_1, \dots, \mu_k\}$ is the centroid of the cluster to which x_i belongs



$$J = \frac{1}{8} \left(d_1 + d_2 + d_3 + d_4 + d_5 + d_6 + d_7 + d_8 \right)$$

Suppose x_1 belongs to Cluster 3,
which are true

✓ A) $\kappa_1 = 3$

B) $\kappa_3 = 1$

C) $\kappa_1 = 3$ ✗

D) $\kappa_3 = 1$ ✗

$\kappa_1 = 3$

 x_1