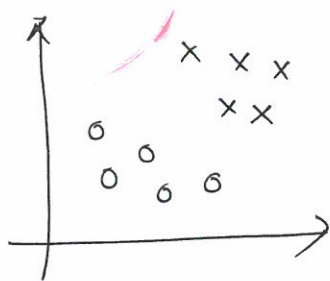


①

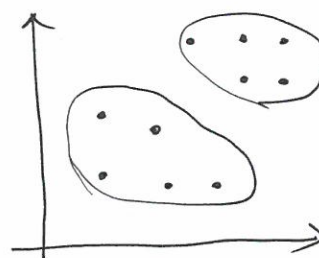
## Clustering - unsupervised learning Algo.

In Supervised learning



Training Data:  $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

In unsupervised learning



clusters  
↓  
clustering Algorithms.

Training Data:  $\{(x_1), (x_2), (x_3), \dots, (x_N)\}$   
↓  
No labels

Applications:

- To find structure/patterns in data

1. Market Segmentation
2. Social Network Analysis
3. Organize Computing clusters
4. Space data analysis.

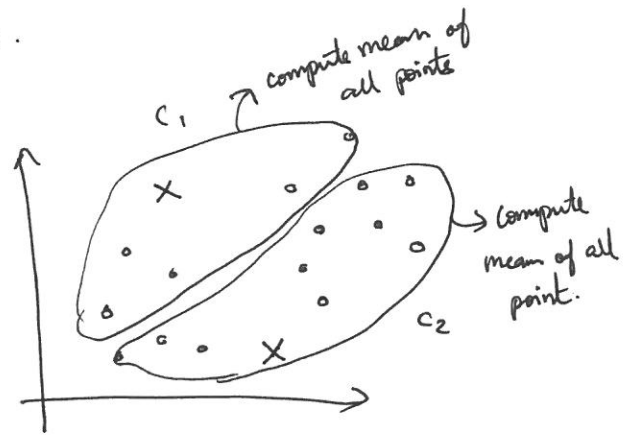
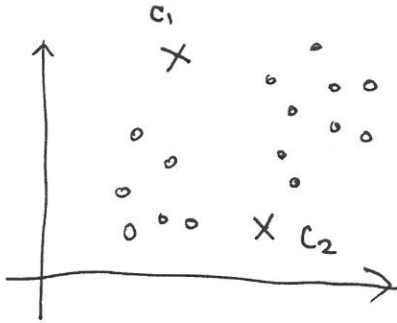
K-Means Clustering.



②

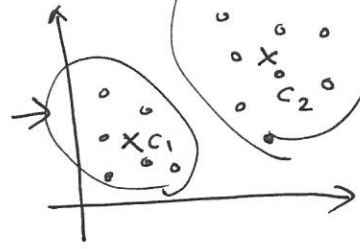
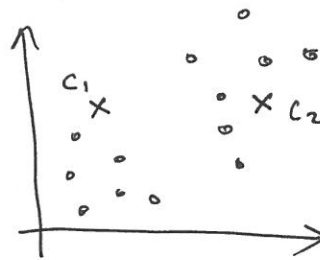
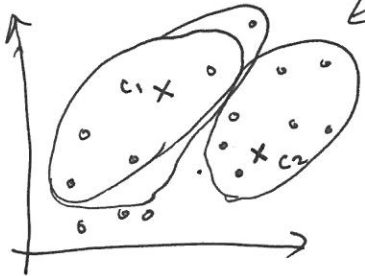
# K-Means Clustering Algo.

Cluster Assignment  
Step.



- randomly initialize cluster centroid assignment.

Move Centroid  
Step.



## K-Means Clustering:

Input:

- $K$  (number of clusters)
  - Training set  $(x_1, x_2, \dots, x_N) \Rightarrow N \times d = 1 \times$
- ↓
- $(x_{11}, x_{12}, \dots, x_{1d})$
- $x_i \in \mathbb{R}^d$

1. Randomly initialize  $K$ -cluster centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^d$

2. Repeat {

for  $i=1$  to  $N$

$c_i := \left\{ \begin{array}{l} \text{index (from 1 to } K) \text{ of cluster centroid} \\ \text{close to } x_i \end{array} \right\}$

$\min_k \|x_i - \mu_k\|^2$

↓

$c_i$

for  $k=1$  to  $K$

$\mu_k := \left\{ \begin{array}{l} \text{Average (mean) of points assigned to} \\ \text{cluster } k. \end{array} \right\}$

}

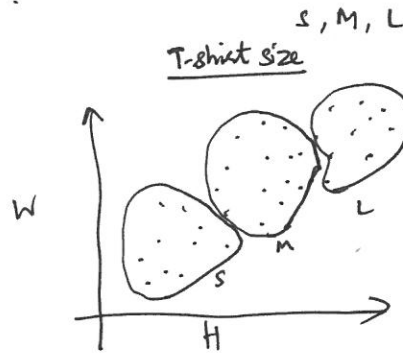
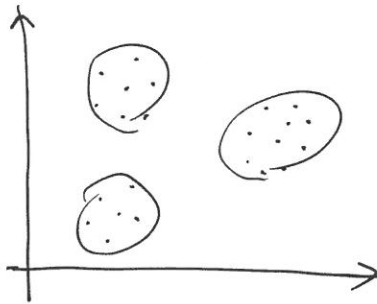
Cluster Assignment  
Step

Move centroid  
step

2

Note: if no points assigned to a cluster, then eliminate that cluster.

K-means for Non-separated clusters:



K-Means Optimization Objective:

K-means :  $K$  = # clusters  
 $k = \{1, 2, \dots, K\}$  = index  
 $C_i$  = index of clusters  $(1, 2, \dots, K)$  to which  $x_i$  is currently assigned.  
 $\mu_k$  = cluster centroid  $k$ . ( $\mu_k \in \mathbb{R}^d$ )  
 $\mu_{C_i}$  = cluster centroid of cluster to which  $x_i$  has been assigned.  
 $x_i \rightarrow 5 \quad C_i = 5 \Rightarrow \mu_{C_i} = \mu_5$

Optimization Objective:

$$E(C_1, \dots, C_N, \mu_1, \dots, \mu_K) = \frac{1}{N} \sum_{i=1}^N \|x_i - \mu_{C_i}\|^2$$

$$\min_{\substack{C_1, \dots, C_N, \\ \mu_1, \dots, \mu_K}} E(C_1, \dots, C_N, \mu_1, \dots, \mu_K)$$

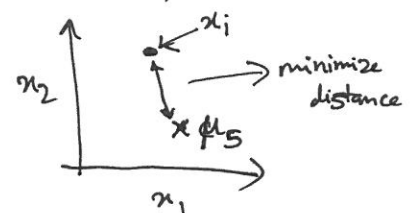
↓  
distortion cost fn.

Cluster Assignment Step:

$\min E(\dots)$  with  $(C_1, C_2, \dots, C_N)$   
and fixing  $(\mu_1, \mu_2, \dots, \mu_K)$

Moving Centroid Step:

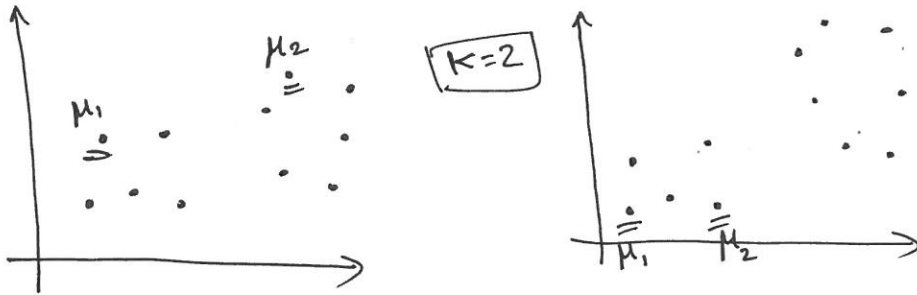
$\min E(\dots)$  with  $(\mu_1, \mu_2, \dots, \mu_K)$   
and fixing  $(C_1, C_2, \dots, C_N)$



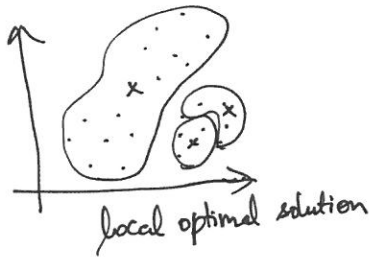
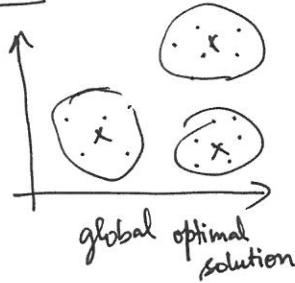
④

## Random Initialization:

- should have  $K < N$
- Randomly pick  $K$  training samples
- Set  $\mu_1, \dots, \mu_K$  equal to these  $K$  examples.



Local Optima  $\therefore \rightarrow$  K-mean can stuck at local optima.



$\downarrow$   
To avoid this, initialize and run K-means multiple times  
 $\downarrow$   
100-times.

$\swarrow$   
{ Randomly initialize and Run K-means } 100 times

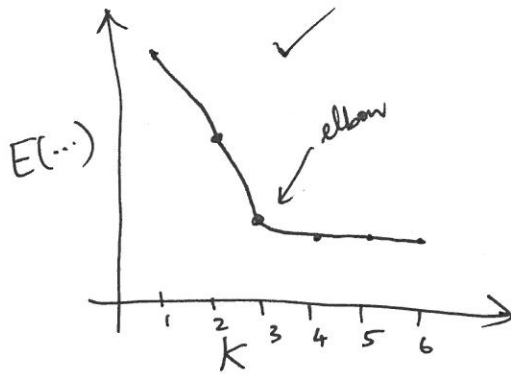
$\downarrow$   
Pick clustering that give lowest cost  $E(\dots)$

$\downarrow$   
 $K=2 \dots 10$  problem size.

⑤

Choose value of  $K$ ? → # clusters

- Choose  $K$  manually → depending on application → T-shirt sizes (S, M, L).
- Use Elbow method (XS, " XL)



(choose  $K=3$  as it is the elbow in the curve)

