

## Clustering

$\mathcal{D}$ : given  $\mathcal{D}$ , find groups.

Partition of  $\mathcal{D} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$

$$\vec{x}_i \in \mathbb{R}^d$$

$n$  points in  $d$ -dim space

find  $k$ -parts

user-specified # of clusters

$C_1, C_2, \dots, C_k$

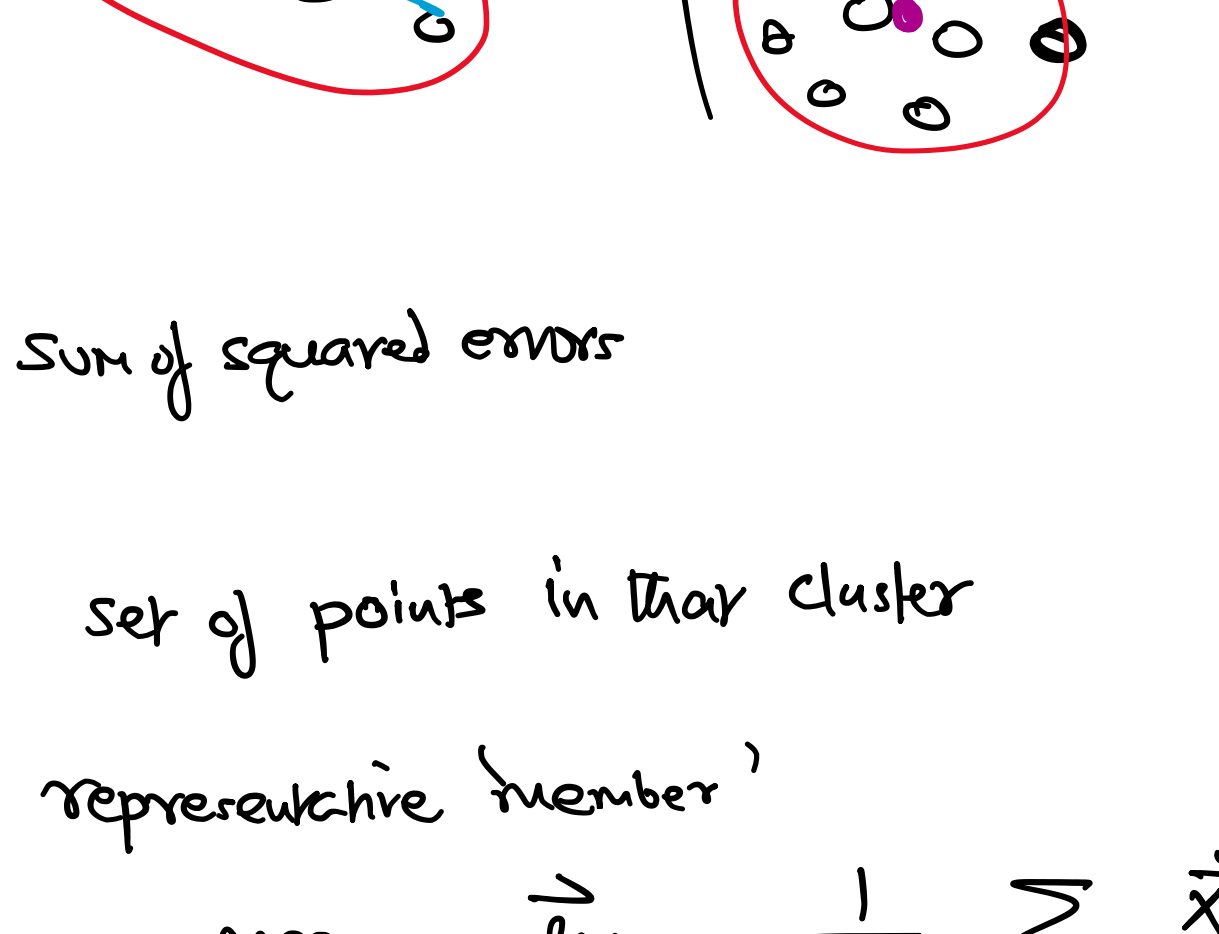
$C_i$  is a set of points

①  $C_i \subseteq \mathcal{D}$

②  $C_i \cap C_j = \emptyset$  ← disjoint clusters

③  $\bigcup_{i=1}^k C_i = \mathcal{D}$  ← every point belongs to some cluster

Partition



Similarity vs distance

1) minimize the intra cluster distances (max the intra cluster similarities)

2) maximize the inter-cluster distances (minimize the inter-cluster similarities)

SSE = sum of squared errors

$C_i$  = set of points in that cluster

→ representative member

$$\text{Mean } \vec{\mu}_i = \frac{1}{|C_i|} \sum_{\vec{x}_j \in C_i} \vec{x}_j$$

$$\text{Min SSE} = \sum_{i=1}^k \left( \sum_{j=1}^{|C_i|} \|\vec{x}_j - \vec{\mu}_i\|^2 \right)$$

Unknown parameters:  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k$

to find optimal solution, i.e.

$\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k$  is NP-hard!



Naive Algo:

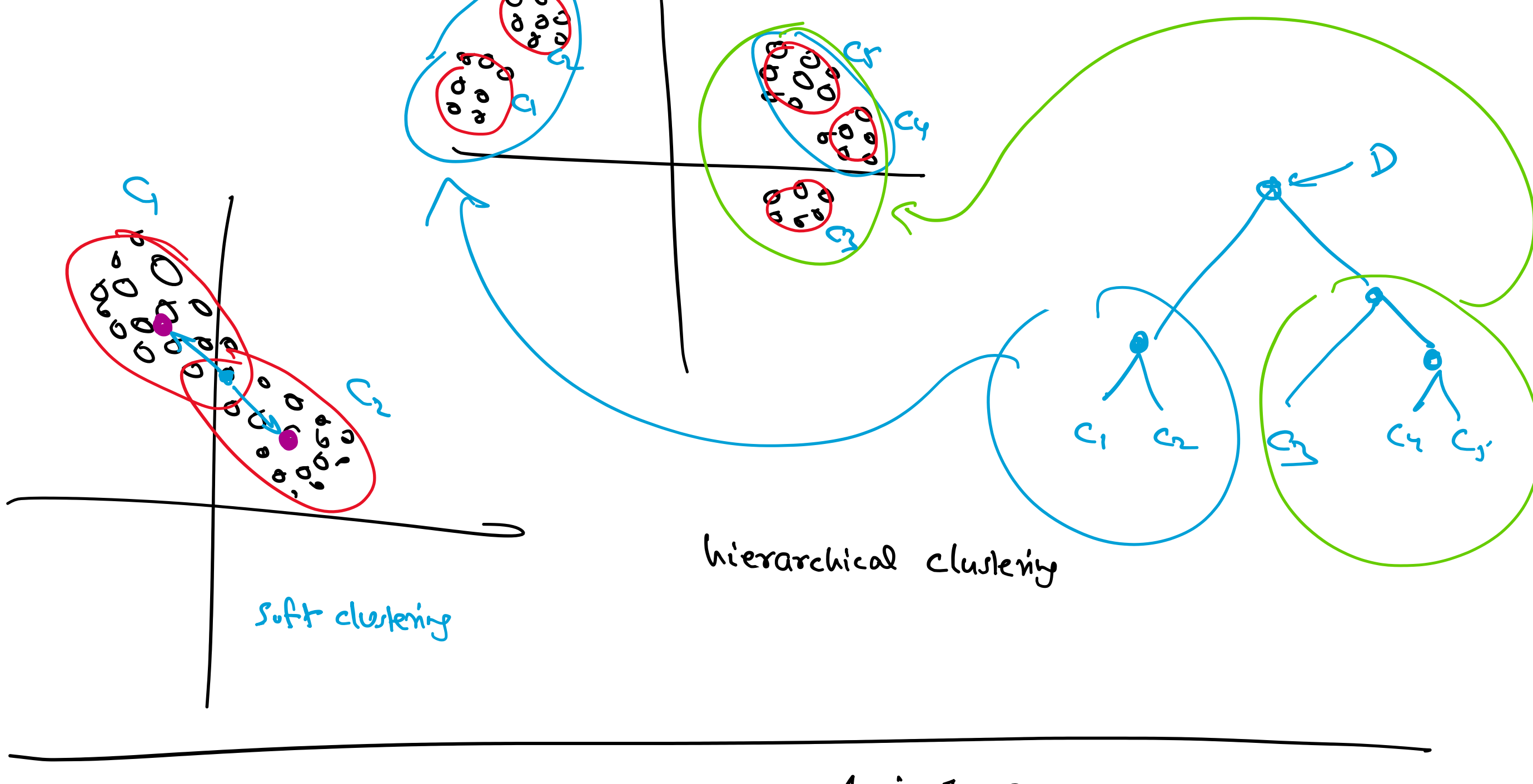
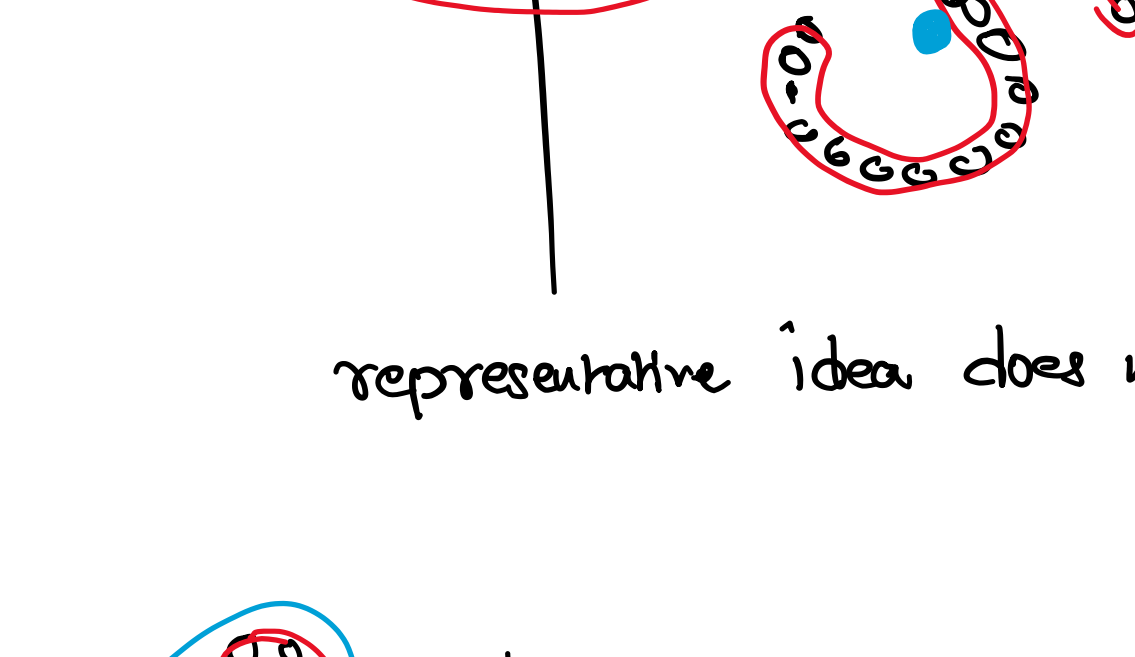
1) Enumerate all partitions of  $\mathcal{D}$  into  $k$  groups

2) evaluate the SSE objective

3) return the one with least SSE value

exponential # of these

Branch-and-bound



k-means : greedy algorithm (iterative)

$$\text{SSE} = \sum_{i=1}^k \sum_{j=1}^{|C_i|} \|\vec{x}_j - \vec{\mu}_i\|^2$$

Unknown parameters:  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k$

k-means, one per cluster

0) Random initialization

Pick any  $k$  points  $\in \mathcal{D}$  and designate as the  $k$  means

e.g.,  $\vec{\mu}_1 = \vec{x}_{10}$

$\vec{\mu}_2 = \vec{x}_{99}$

$\vec{\mu}_3 = \vec{x}_5$

$n=100$

$k=3$

Iterative improvement

1) given  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k$

assign each point to the closest mean

partitioning step

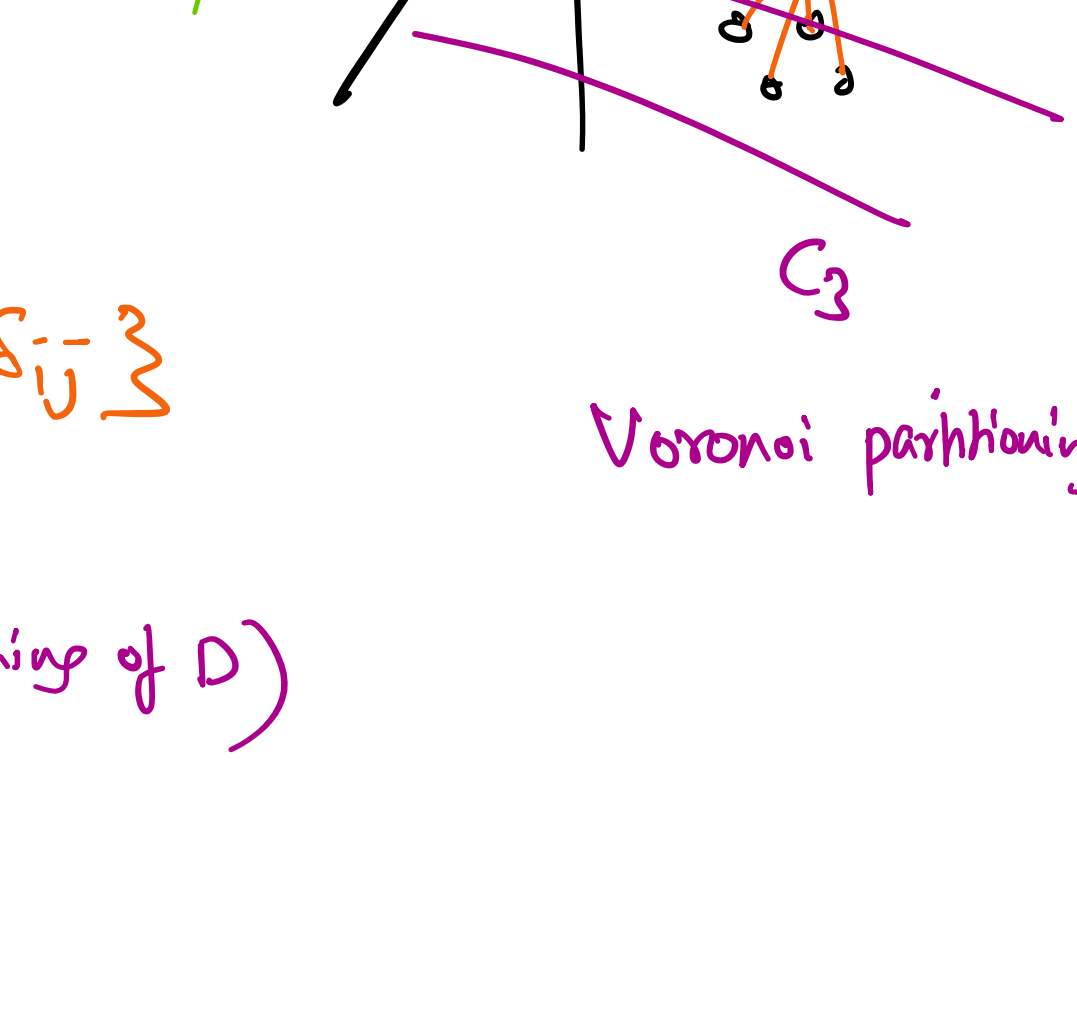
$\forall \vec{x}_j \in \mathcal{D}$

$\forall \vec{\mu}_i$

find  $\delta_{ij} = \|\vec{x}_j - \vec{\mu}_i\|^2$

assign  $\vec{x}_j$  to  $C_{i^*}$

where  $i^* = \arg \min_{i=1}^k \delta_{ij}$



result in  $C_1, C_2, \dots, C_k$  (partitioning of  $\mathcal{D}$ )

2) Given  $C_1, C_2, \dots, C_k$

update the means

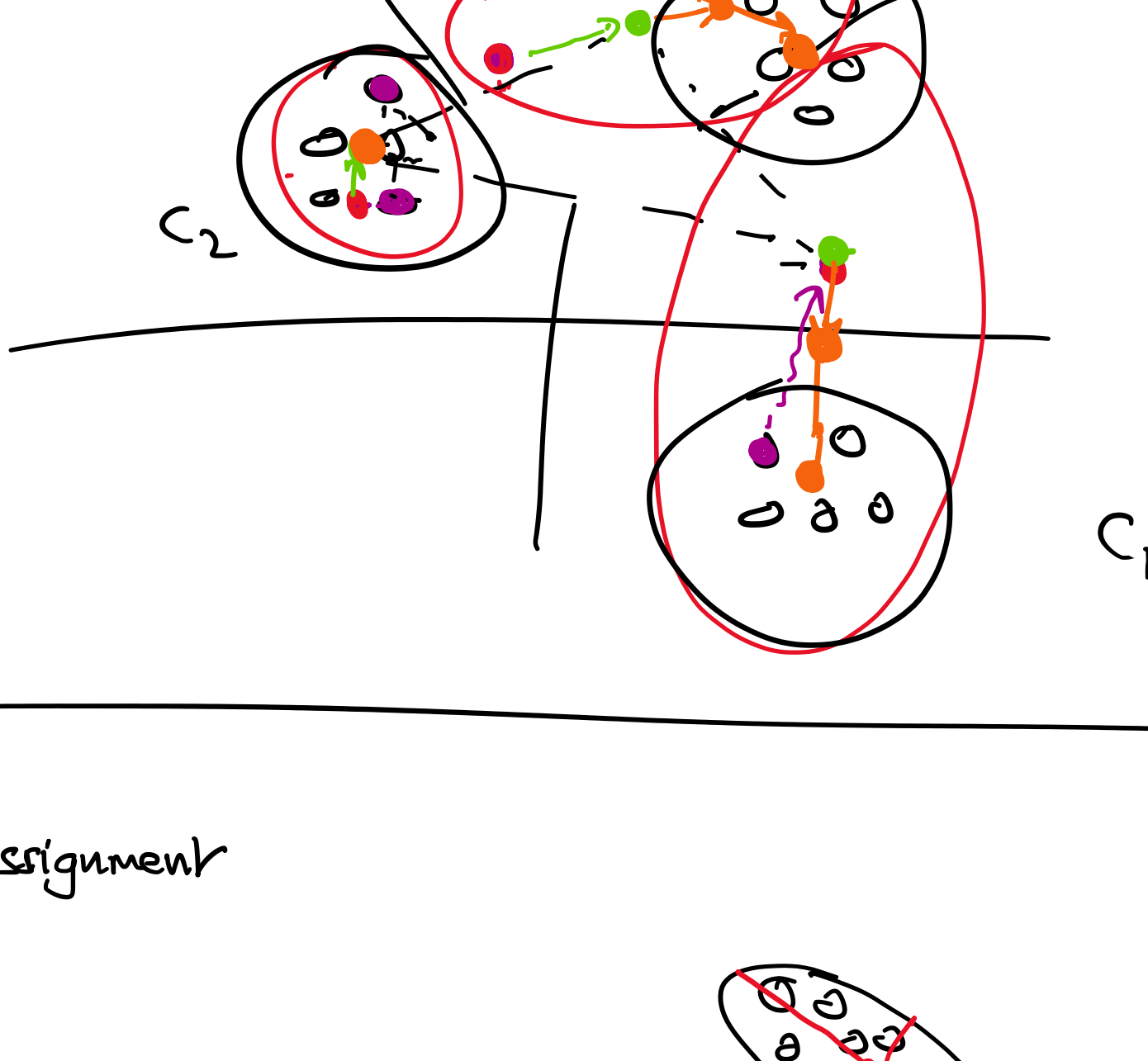
$$\vec{\mu}_i = \frac{1}{|C_i|} \sum_{\vec{x}_j \in C_i} \vec{x}_j$$

result in new estimates for the  $k$  means

3) repeat 1) & 2) until convergence

$$\mathcal{E} = \sum_{i=1}^k \|\vec{\mu}_i^t - \vec{\mu}_i^{t-1}\|^2 \quad t \text{ is step}$$

if  $\mathcal{E} < 10^{-5}$  stop



Probabilistic / soft assignment

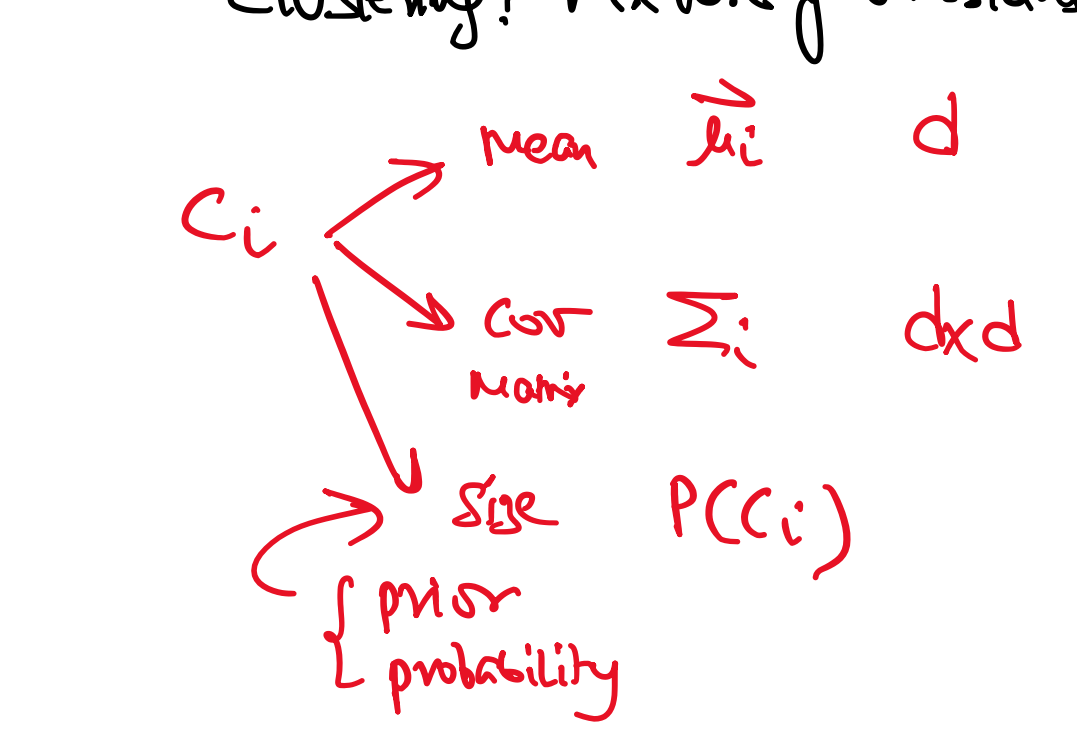
$\{C_1, C_2, \dots, C_k\}$

no longer a partition

$$\begin{cases} P(C_i | \vec{x}_j) \equiv w_{ij} \\ \text{probability that point } j (\vec{x}_j) \\ \text{belongs to cluster } i (C_i) \end{cases}$$

$\forall C_i, i=1 \dots k$

$$\sum_{i=1}^k w_{ij} = 1 \quad \forall \vec{x}_j$$



mean  $\vec{\mu}_i$   $d$

cov  $\Sigma_i$   $d \times d$

prior  $P(C_i)$

posterior probability

$$\Theta = \begin{cases} \vec{\mu}_1, \Sigma_1, P(C_1) \\ \vec{\mu}_2, \Sigma_2, P(C_2) \\ \vdots \\ \vec{\mu}_k, \Sigma_k, P(C_k) \end{cases} \quad \text{Unknown}$$

Expectation - Maximization Algorithm (EM)

0) Random initialization

a) Pick any  $k$  (random) points as  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k$

b) assign each point  $\vec{x}_j \in \mathcal{D}$  to closest mean

to obtain  $\{C_1, C_2, \dots, C_k\}$

c) for each  $C_i$ , compute  $\Sigma_i$

$$P(C_i) = \frac{|C_i|}{n}$$

1) Expectation step: (soft assignment)

Given  $\vec{\mu}_i, \Sigma_i, P(C_i) \quad \forall i=1 \dots k$

compute  $w_{ij} = P(C_i | \vec{x}_j)$

$\forall \vec{x}_j \in \mathcal{D}, \forall i=1 \dots k$

Computing for each point  $\vec{x}_j$ , the prob that it belongs to each of the  $k$  clusters

2) Maximization step (update params)

Given  $w_{ij} = P(C_i | \vec{x}_j) \quad \forall \vec{x}_j \in \mathcal{D}, \forall C_i, i=1 \dots k$

compute updates to  $\Theta$

$$\vec{\mu}_i = \frac{\sum_{j=1}^n w_{ij} \vec{x}_j}{\sum_{j=1}^n w_{ij}}$$

weighted sum/mean

add up each point proportional to its weight/prob of belonging to  $C_i$

$$\Sigma_i = \frac{\sum_{j=1}^n w_{ij} (\vec{x}_j - \vec{\mu}_i) (\vec{x}_j - \vec{\mu}_i)^T}{\sum_{j=1}^n w_{ij}}$$

weighted outer products for cov

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{n} \quad \leftarrow \text{total 'weight' for } C_i \text{ over all the points}$$

3) stop when  $\mathcal{E} < 10^{-5}$

$$\mathcal{E} = \sum_{i=1}^k \|\vec{\mu}_i^t - \vec{\mu}_i^{t-1}\|^2$$

how to compute

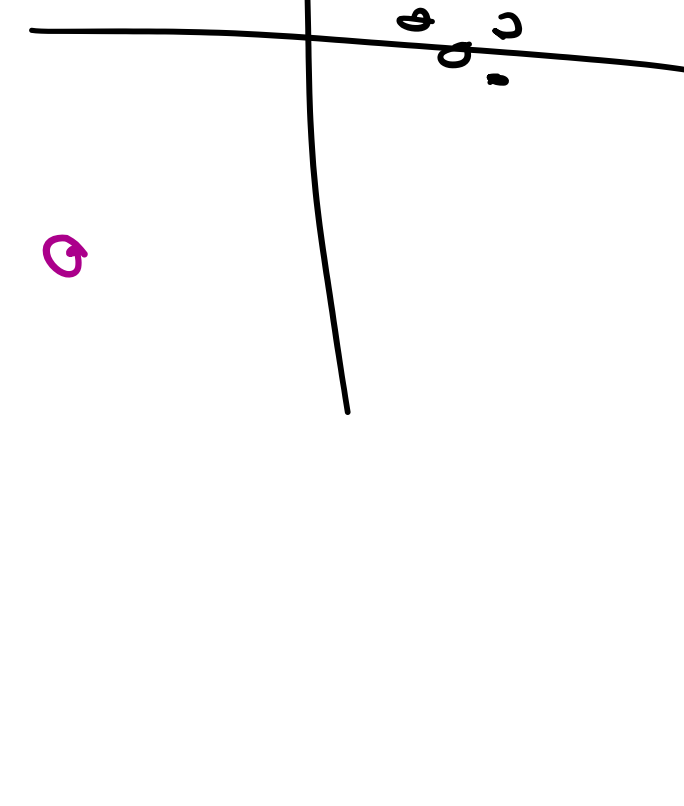
$w_{ij} = P(C_i | \vec{x}_j)$  Prob that  $\vec{x}_j$  belongs to  $C_i$

$$P(C_i | \vec{x}_j) = \frac{P(\vec{x}_j | C_i) \cdot P(C_i)}{P(\vec{x}_j)}$$

Posterior probability

$$= \frac{P(\vec{x}_j | C_i) \cdot P(C_i)}{\sum_{i=1}^k P(\vec{x}_j | C_i) \cdot P(C_i)}$$

$P(\vec{x}_j)$



$$P(C_i | \vec{x}_j) = \frac{N(\vec{x}_j | \vec{\mu}_i, \Sigma_i) \cdot P(C_i)}{P(\vec{x}_j)}$$

$w_{ij} \rightarrow$  hard clustering

Assign each point to the max prob cluster