⊛ Consider a set of unlabeled data points:

$$\mathcal{X} = \{\, \underline{x}_1, \underline{x}_2, \ldots, \underline{x}_N \,\} \quad \text{where} \quad \underline{x}_i \in \mathbb{R}^\ell, \; \forall j \in [N]$$
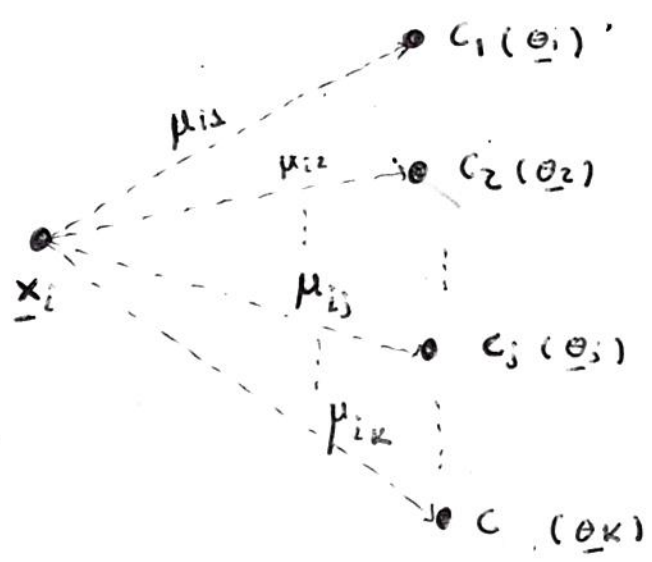
⊛ The clustering provides a partitioning of the original dataset into M clusters such that:

$$\mathcal{X} = \bigcup_{j=1}^{K} C_j \quad \text{with} \quad C_r \cap C_m = \phi, \; \forall r \neq m$$

⊛ **Problem Definition:** ● Assign each datapoint $\underline{x}_i$ to a unique cluster $C_j$ such that the data points which are assigned to the j-th cluster exhibit minimum distance towards the <u>cluster</u> <u>representative</u> given by $\underline{\theta}_j$.

● The total number of clusters is K and each cluster is represented by a vector $\underline{\theta}_j \in \mathbb{R}^\ell$.

Ⓐ Let $\mu_{ij} \in \{0,1\}$ denote the membership status of the i-th datapoint relative to the j-th cluster, $\{1 \leq i \leq N \text{ and } 1 \leq j \leq K\}$.



Each datapoint is assigned to a single cluster.

$$\sum_{j=1}^{K} \mu_{ij} = 1, \; \forall i \in [N]$$

⊛ Consider that all membership values are organized in a membership matrix $\underline{\underline{M}} = [\mu_{ij}]$ such that $1 \le i \le N$ and $1 \le j \le K$.

✪ In this setting, we may define the clustering problem as an optimization problem where the objective/cost function is given as:

$$
\begin{aligned}
&\min_{(\underline{\Theta}, \underline{\underline{M}})} \quad J(\underline{\Theta}; \underline{\underline{M}}) = \sum_{i=1}^{N} \sum_{j=1}^{K} \mu_{ij} \cdot \| x_i - \Theta_j \|^2 \\
&\text{s.t.} \quad \sum_{j=1}^{K} \mu_{ij} = 1, \quad \forall i \in [N] \\
&\text{where} \quad \underline{\Theta} = [\Theta_1, \Theta_2, \dots, \Theta_K]
\end{aligned}
$$

$\boxed{\text{OP. I}}$

S.O.S: The objective function $J(\underline{\Theta}; \underline{\underline{M}})$ is not differentiable since it's not continuous given that $\mu_{ij} \in \{0, 1\}$.

(combinatorial P)

⊛ This is a hard mixed-integer optimization problem since the number of possible clusterings of $N$ datapoints into $K$ clusters will be given by the so-called Stirling numbers of the second type:

$$
S(N, K) = \frac{1}{K!} \sum_{j=0}^{K} (-1)^{K-j} \binom{K}{j} j^N
$$

Ⓐ Therefore, we cannot use straightforward optimization techniques.

## Heuristic Optimization Scheme:

(α): Considering that $\Theta_j$ with $1 \leq j \leq K$ are fixed.

Since for each vector $\underline{x}_i$ only one $\mu_{ij}$ is $1$ and all the others are $0$, it is straightforward to see that $J(\underline{\Theta} ; \underline{\mu})$ is minimized by if we assign each $\underline{x}_i$ to each closest cluster-representative:

$$\mu_{ij} = \begin{cases} 1, & j = \arg\min_{r \in [K]} \| \underline{x}_i - \underline{\Theta}_r \|^2 ; \\ \\ 0, & \text{otherwise} \end{cases} \qquad [A]$$
$$\forall i \in [N]$$

(β): Considering that $\mu_{ij}$'s are fixed with $1 \leq i \leq N$ and $1 \leq j \leq K$.

In this setting, we need to evaluate the F.O.Cs with respect to the parameters $\underline{\Theta}_r$, $\boxed{\forall r \in [K]}$ as:

$$\frac{\partial J}{\partial \underline{\Theta}_r} = \underline{0} \in \mathbb{R}^\ell \quad \text{or} \quad \frac{\partial J}{\partial \underline{\Theta}_r} = \frac{\partial}{\partial \underline{\Theta}_r} \left\{ \sum_{i=1}^{N} \sum_{j=1}^{K} \mu_{ij} \| \underline{x}_i - \underline{\Theta}_j \|^2 \right\} = \underline{0}$$

Ⓐ Since we are differentiating w.r.t $\underline{\Theta}_r$, all the other $\underline{\Theta}_m$'s with $m \neq r$ as we do not contribute to the differentiation process. They are constant as far as the partial differentiation process is concerned.

Ⓐ Thus, we may write that:

$$\frac{\partial J}{\partial \underline{\Theta}_r} = \sum_{i=1}^{N} \sum_{j=1}^{K} \frac{\partial}{\partial \underline{\Theta}_r} \left[ \mu_{ij} \| \underline{x}_i - \underline{\Theta}_j \|^2 \right] = \sum_{i=1}^{N} \frac{\partial}{\partial \underline{\Theta}_r} \left[ \mu_{ir} \| \underline{x}_i - \underline{\Theta}_r \|^2 \right] = \underline{0}$$
$$\forall r \in [K] \Rightarrow$$

$\circledast$ $\displaystyle\sum_{i=1}^{N} \mu_{ir} \frac{\partial}{\partial \underline{\theta}_r} \| \underline{x}_i - \underline{\theta}_r \|^2 = \underline{0}$ $\Rightarrow$ $\displaystyle\sum_{i=1}^{N} \mu_{ir} (-2\underline{x}_i + 2\underline{\theta}_r) = \underline{0}$ $\Rightarrow$

$\circledA$

$\| \underline{x}_i - \underline{\theta}_r \|^2 = (\underline{x}_i - \underline{\theta}_r)^T (\underline{x}_i - \underline{\theta}_r) = (\underline{x}_i^T - \underline{\theta}_r^T)(\underline{x}_i - \underline{\theta}_r) =$

$\qquad = \underline{x}_i^T \underline{x}_i - \underline{x}_i^T \underline{\theta}_r - \underline{\theta}_r^T \underline{x}_i + \underline{\theta}_r^T \underline{\theta}_r =$

$\qquad = \underline{x}_i^T \underline{x}_i - 2\underline{x}_i^T \underline{\theta}_r + \underline{\theta}_r^T \underline{\theta}_r .$

$\circledA$ $\dfrac{\partial}{\partial \underline{\theta}_r} \left[ \underline{x}_i^T \underline{x}_i - 2\underline{x}_i^T \underline{\theta}_r + \underline{\theta}_r^T \underline{\theta}_r \right] = -2\underline{x}_i + 2\underline{\theta}_r$

$\circledA$ $\dfrac{\partial \underline{x}^T \underline{x}}{\partial \underline{x}} = \dfrac{\partial \underline{x}^T I \underline{x}}{\partial \underline{x}} = (I + I^T) \underline{x} = 2 I \underline{x} = 2\underline{x}$

$\circledast$ $-2 \displaystyle\sum_{i=1}^{N} \mu_{ir} \underline{x}_i + 2 \sum_{i=1}^{N} \mu_{ir} \underline{\theta}_r = \underline{0}$ $\Rightarrow$

$\underline{\theta}_r \circ \displaystyle\sum_{i=1}^{N} \mu_{ir} = \sum_{i=1}^{N} \mu_{ir} \underline{x}_i$ $\Rightarrow$ $\boxed{\underline{\theta}_r = \dfrac{\displaystyle\sum_{i=1}^{N} \mu_{ir} \underline{x}_i}{\displaystyle\sum_{i=1}^{N} \mu_{ir}}}$ [B]

$\circledast$ Eqs. (A) and (B) can pave the way in order to formulate the Generalized Hard Clustering Scheme.

Step 1 : Choose $\underline{\theta}_j(0)$ as initial estimates for $\underline{\theta}_j$, $\forall j \in [K]$.

Step 2 :    $t = 0$

Step 3: Repeat :

Step 3.1: $\forall i \in [N]$: Update $\underline{\mu}_i$ according to:

$$\mu_{ij}(t) = \begin{cases} 1, & j = \arg\min_{r \in [K]} \|\underline{x}_i - \underline{\theta}_r(t)\|^2; \\ \\ 0, & \text{otherwise.} \end{cases}$$

(r)

Step 3.2: $t = t+1$.

Step 3.3: $\forall j \in [K]$: Update $\underline{\theta}_j$ according to:

$$\underline{\theta}_j(t) = \frac{\sum_{i=1}^{N} \mu_{ij}(t-1)\underline{x}_i}{\sum_{i=1}^{N} \mu_{ij}(t-1)}$$

(r)

Until $\|\underline{\theta}(t) - \underline{\theta}(t-1)\| < \varepsilon$

Ⓐ Consider $\underline{\underline{\theta}} = \{\underline{\theta}_1, \underline{\theta}_2, \ldots, \underline{\theta}_K\}$ and $\underline{\underline{\mu}} = \begin{bmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \\ \vdots \\ \underline{\mu}_N \end{bmatrix}$

where $\underline{\theta}_j$ the $j$-th cluster representative and $\underline{\mu}_i$ the $i$-th datapoint membership vector.