

annotated
version

Machine Learning Course - CS-433

K-Means Clustering

Nov 23, 2022

Martin Jaggi

Last updated on: November 20, 2022

credits to Mohammad Emtiyaz Khan & Rüdiger Urbanke

The logo of the École Polytechnique Fédérale de Lausanne (EPFL) is displayed in a bold, red, sans-serif font. The letters are stylized, with the 'E' and 'F' having a unique geometric design.

Clustering

Clusters are groups of points whose inter-point distances are small compared to the distances outside the cluster.

$$\mathbf{z}_n = \text{1-hot vector } \mathbb{R}^K$$

$$z_{nk} = \begin{cases} 1 & \text{if data point } n \text{ is assigned to group } k \\ 0 & \text{otherwise} \end{cases}$$

The goal is to find “prototype” points $\mu_1, \mu_2, \dots, \mu_K$ and cluster assignments $z_n \in \{1, 2, \dots, K\}$ for all $n = 1, 2, \dots, N$ data vectors $\mathbf{x}_n \in \mathbb{R}^D$.

K-means clustering

Assume K is known.

distance of \mathbf{x}_n to μ_k

$$\min_{\mathbf{z}, \mu} \mathcal{L}(\mathbf{z}, \mu) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$$

$$\text{s.t. } \mu_k \in \mathbb{R}^D, z_{nk} \in \{0, 1\}, \sum_{k=1}^K z_{nk} = 1,$$

$$\text{where } \mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nK}]^\top$$

$$\mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N]^\top$$

$$\mu = [\mu_1, \mu_2, \dots, \mu_K]^\top$$

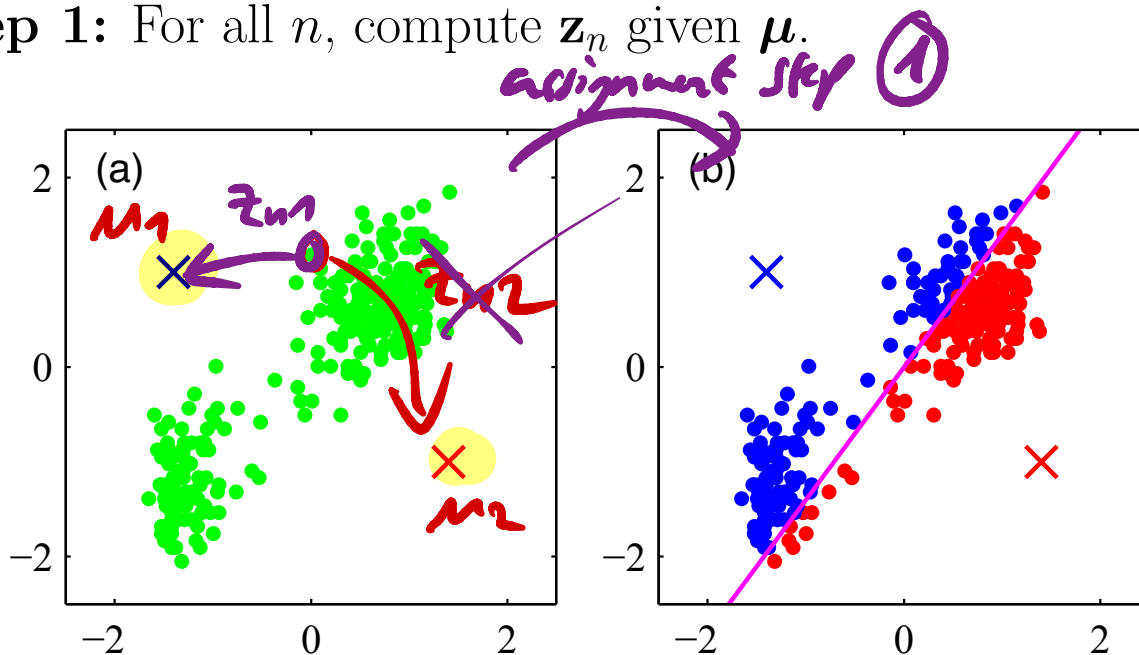
Is this optimization problem easy?

NP-hard

Algorithm: Initialize $\mu_k \forall k$,
then iterate:

1. For all n , compute z_n given μ . assignment step z
2. For all k , compute μ_k given z . update representative μ

Step 1: For all n , compute z_n given μ .



assignment

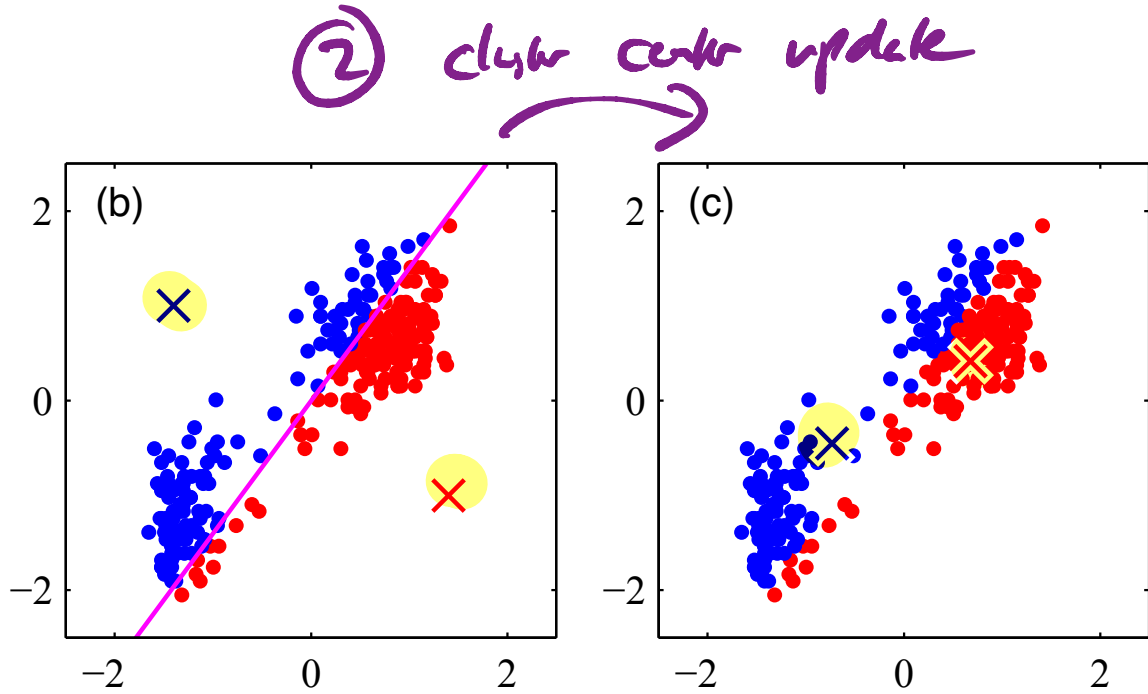
$$z_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_{j=1,2,\dots,K} \|\mathbf{x}_n - \mu_j\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

Step 2: For all k , compute μ_k given z .
Take derivative w.r.t. μ_k to get:

$$\mu_k = \frac{\sum_{n=1}^N z_{nk} \mathbf{x}_n}{\sum_{n=1}^N z_{nk}}$$

← points assigned to group k
← # points → . . .

Hence, the name 'K-means'.



Summary of K-means

Initialize $\mu_k \forall k$, then iterate:

1. For all n , compute z_n given μ .

$$\mathcal{O}(N \cdot K \cdot D)$$

$$z_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

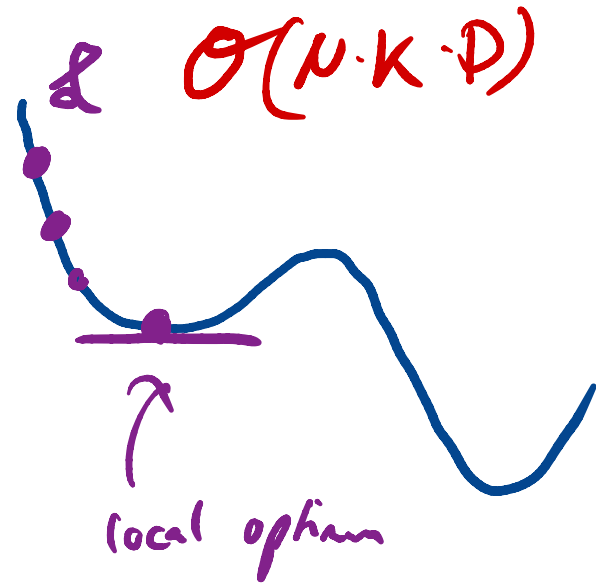
2. For all k , compute μ_k given z .

$$\mu_k = \frac{\sum_{n=1}^N z_{nk} \mathbf{x}_n}{\sum_{n=1}^N z_{nk}}$$

$$\mathcal{O}(N \cdot K \cdot D)$$

Convergence to a local optimum is assured since each step decreases the cost (see Bishop, Exercise 9.1).

$$\mathcal{L}(\mu, z)$$



Coordinate descent

K-means is a coordinate descent algorithm, where, to find $\min_{\mathbf{z}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{z}, \boldsymbol{\mu})$, we start with some $\boldsymbol{\mu}^{(0)}$ and repeat the following:

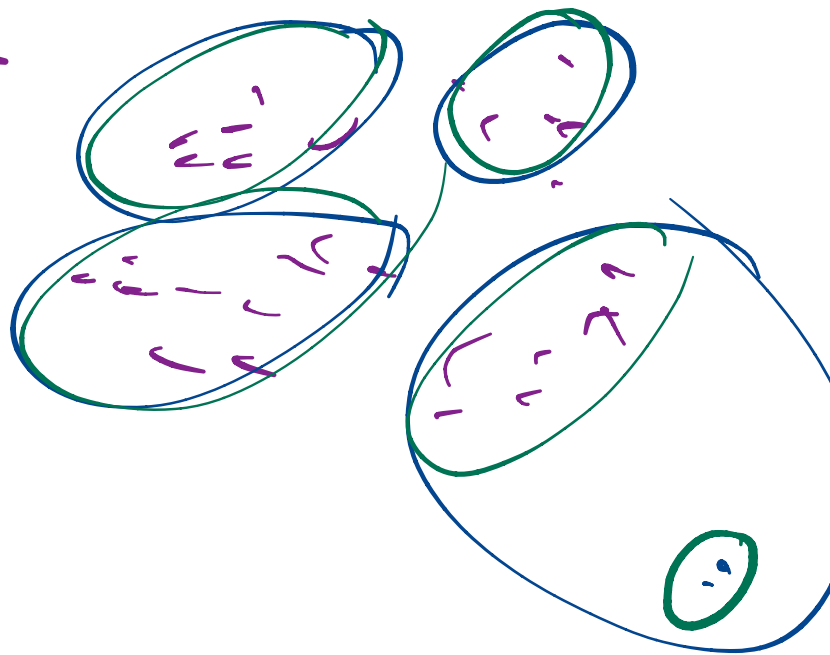
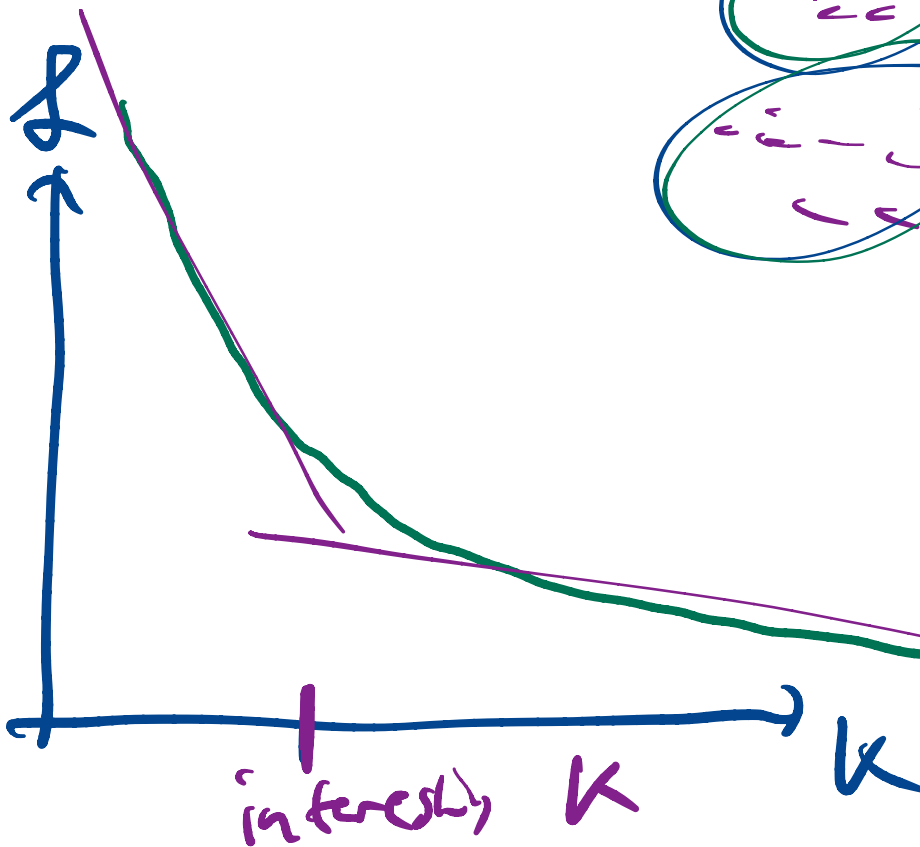
$$\mathbf{z}^{(t+1)} := \arg \min_{\mathbf{z}} \mathcal{L}(\mathbf{z}, \boldsymbol{\mu}^{(t)})$$

$$\boldsymbol{\mu}^{(t+1)} := \arg \min_{\boldsymbol{\mu}} \mathcal{L}(\mathbf{z}^{(t+1)}, \boldsymbol{\mu})$$

← ①

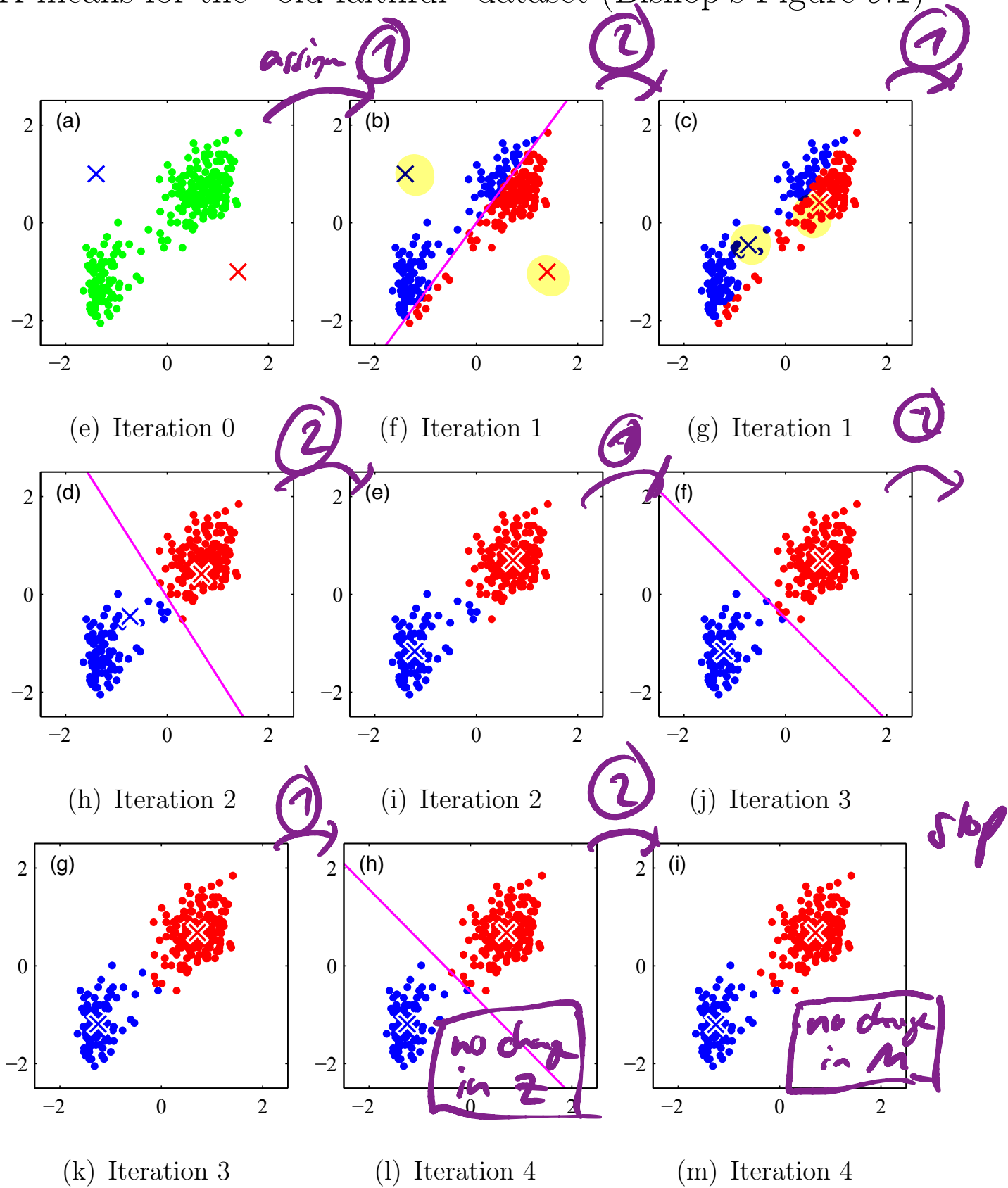
← equivalent to ②

How to set K



Examples

K-means for the “old-faithful” dataset (Bishop’s Figure 9.1)



$$\mu_k \in \mathbb{R}^3$$

Data compression for images (this is also known as vector quantization).

$$x_n = \begin{bmatrix} r \\ g \\ b \end{bmatrix} \in \mathbb{R}^3$$

$K=2$

$K=2$



$K=3$



$K=10$

$K=10$



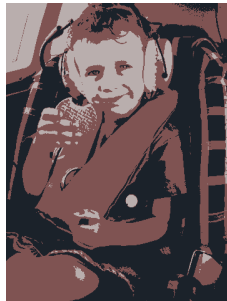
Original image



500

100K
datapoints

200



„vector quantization“

Probabilistic model for K-means

$$\begin{aligned}
 \prod_{n=1}^N p(x_n | \mu, z) &= \log \prod_{n=1}^N \mathcal{N}(x_n | \mu_k, I) \\
 \log p(x | \mu, z) &= \log \prod_n \prod_{k=1}^K \mathcal{N}(x_n | \mu_k, I)^{z_{nk}} \\
 &= \log \prod_n \prod_k c \cdot e^{-\frac{1}{2} \|x_n - \mu_k\|^2 \cdot z_{nk}} \\
 &= - \sum_n \sum_k \frac{1}{2} \|x_n - \mu_k\|^2 \cdot z_{nk} + c \\
 &= \mathcal{Q}(\mu, z)
 \end{aligned}$$

Annotations:
 $\mathcal{N}(x_n | \mu_k, I)$ is a *group of x_n* .
 z_{nk} is highlighted in green.
 \log is written in purple.

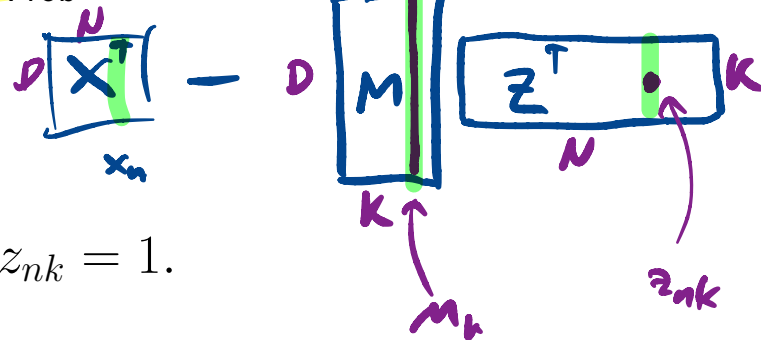
K-means as a Matrix Factorization

Recall the objective

$$\begin{aligned} \min_{\mathbf{z}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{z}, \boldsymbol{\mu}) &= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2 \\ &= \|\mathbf{X}^\top - \mathbf{M}\mathbf{Z}^\top\|_{\text{Frob}}^2 \end{aligned}$$

$$\text{s.t. } \boldsymbol{\mu}_k \in \mathbb{R}^D,$$

$$z_{nk} \in \{0, 1\}, \sum_{k=1}^K z_{nk} = 1.$$



Issues with K-means

1. Computation can be heavy for large N , D and K .
2. Clusters are forced to be spherical (e.g. cannot be elliptical).
3. Each example can belong to only one cluster ("hard" cluster assignments).