

Annotated  
Version

Machine Learning Course - CS-433

# K-Means Clustering

Nov 25, 2021

changes by Martin Jaggi 2021, 2020, 2019, changes by Rüdiger Urbanke 2018, changes by  
Martin Jaggi 2016, 2017 ©Mohammad Emtiyaz Khan 2015

Last updated on: November 23, 2021

**EPFL**

# Clustering

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

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  $\mathbf{x}_n \in \mathbb{R}^D$ .

$z_n$   $n$ -hot vector  $\mathbb{R}^K$   
 $z_{nk} = \begin{cases} 1 & \text{if data } n \\ & \text{assigned to } k \end{cases}$   
otherwise 0  
assignment

## K-means clustering

Assume  $K$  is known.

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

distance of  $\mathbf{x}_n$  to  $\mu_k$

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

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

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

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

discrete var

Is this optimization problem easy?

NP-hard

# k-means algo

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

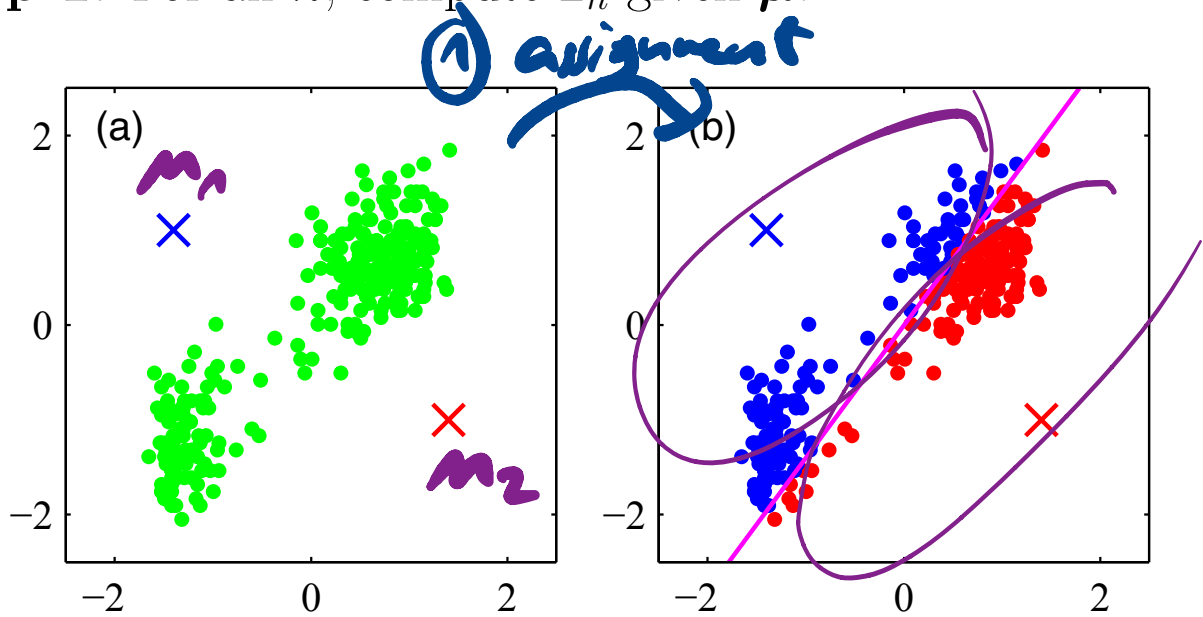
1. For all  $n$ , compute  $z_n$  given  $\mu$ .

assignment step  $z$

2. For all  $k$ , compute  $\mu_k$  given  $z$ .

update center  $\mu$

Step 1: For all  $n$ , compute  $z_n$  given  $\mu$ .



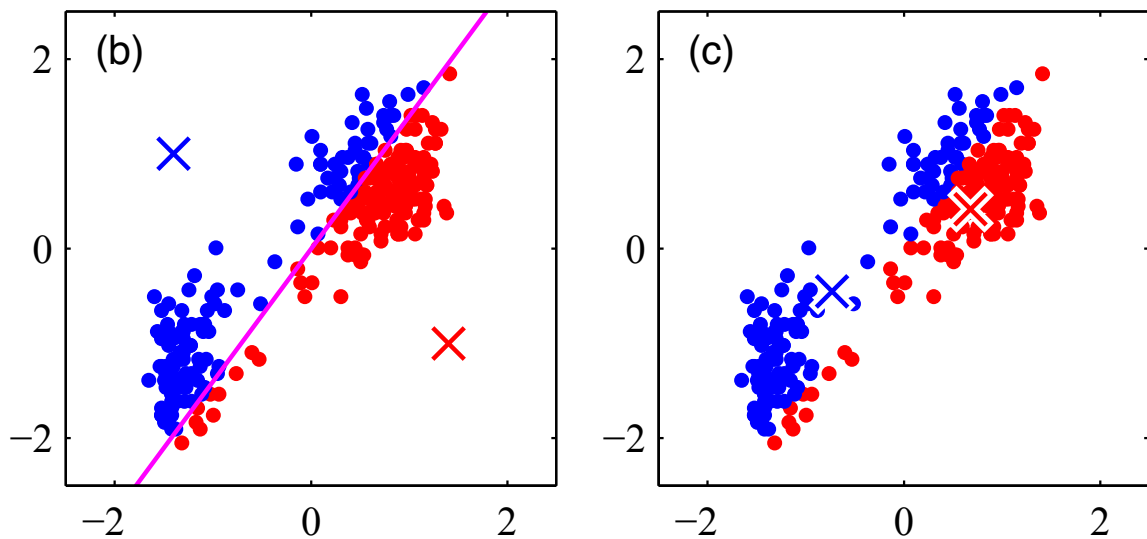
$$z_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_{j=1,2,\dots,K} \|\mathbf{x}_n - \mu_j\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

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

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

← # points assigned to cluster  $k$

Hence, the name 'K-means'.



## Summary of K-means

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

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

1. For all  $n$ , compute  $\mathbf{z}_n$  given  $\mu$ .

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

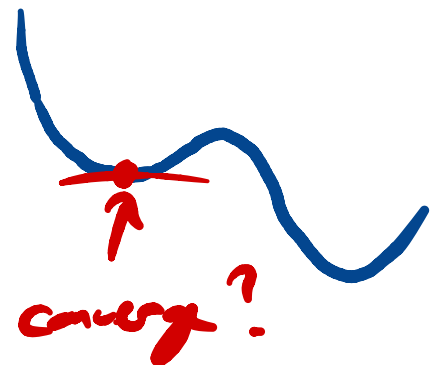
2. For all  $k$ , compute  $\mu_k$  given  $\mathbf{z}$ .

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

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

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

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



# Coordinate descent

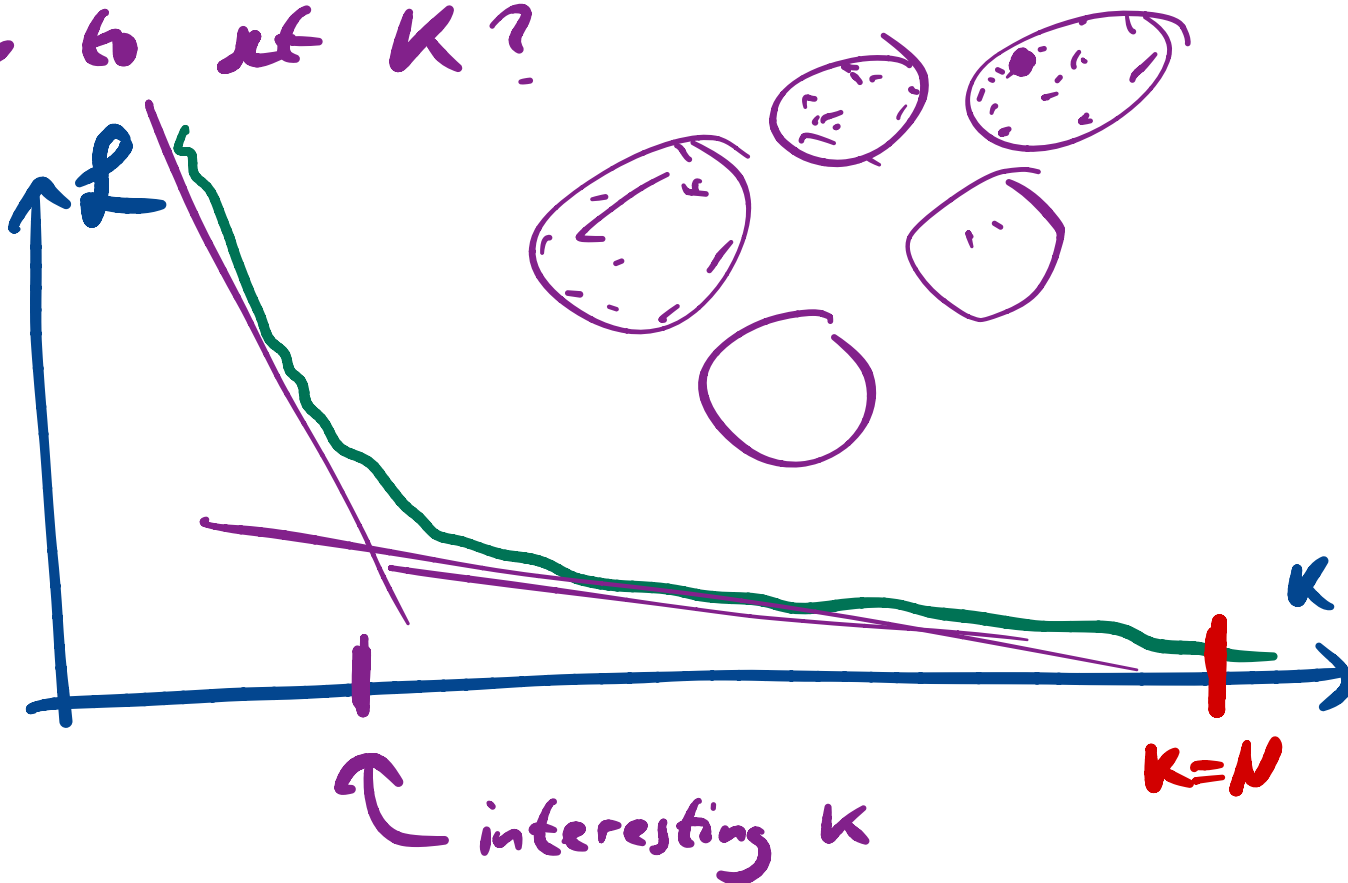
K-means is a coordinate descent algorithm, where, to find  $\min_{\mathbf{z}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{z}, \boldsymbol{\mu})$ , we start with some  $\boldsymbol{\mu}^{(0)}$  and repeat the following:

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

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

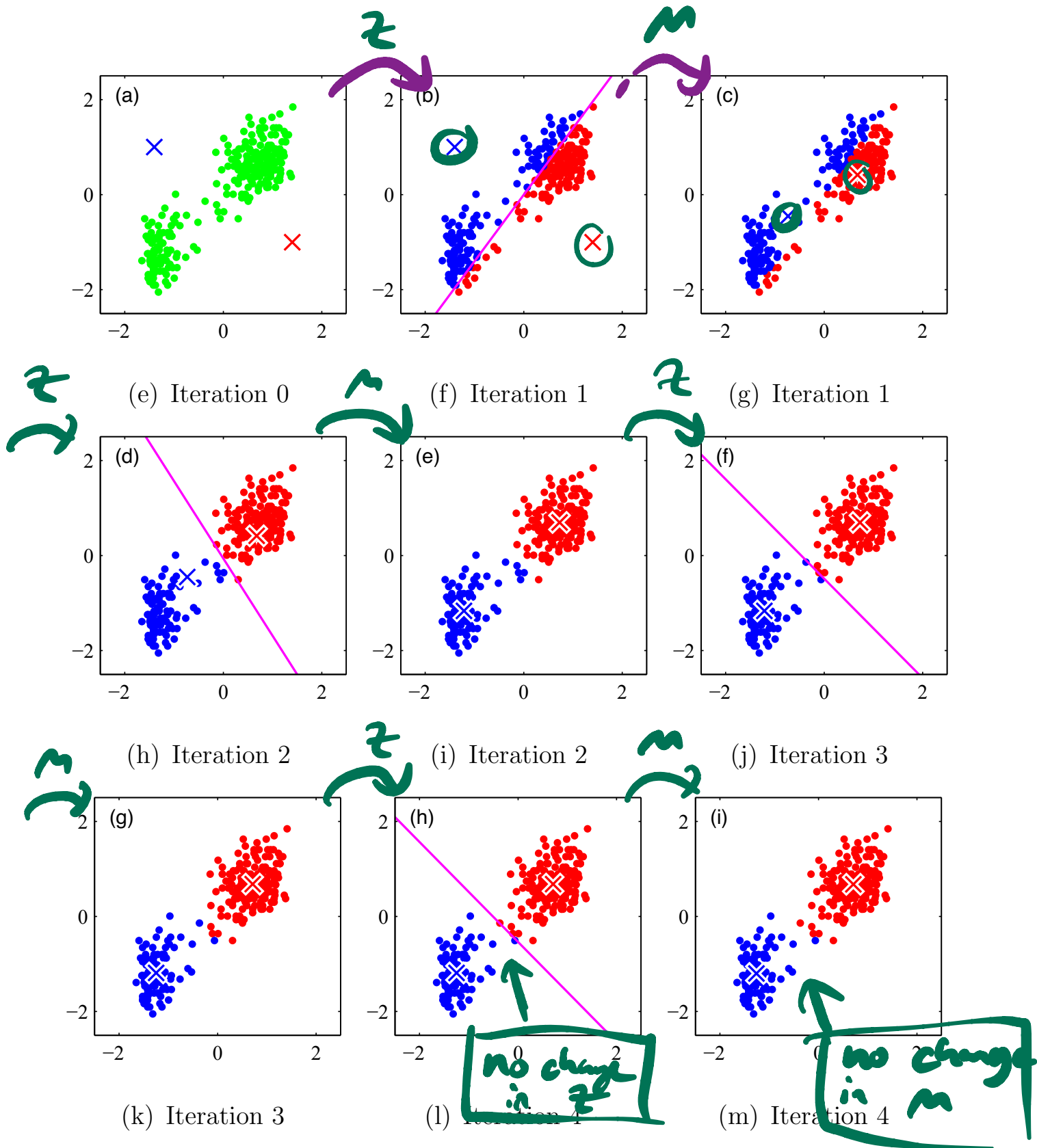
perfect  $\Leftrightarrow \mu$ -update  
min.

How to pick  $K$ ?



# Examples

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



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

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



$K=2$



$K=3$



$K=10$

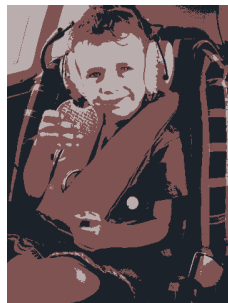


Original image

500

200

100k datapoints



vector quantization

Probabilistic model for K-means

Likelihood of  $x$  given  $\mu, z$

$$p(x_n | \mu, z) = \prod_{n=1}^N \mathcal{N}(x_n | \mu_k, I)$$

$$p(X | \mu, z)$$

$$= \prod_n \prod_{k=1}^K \mathcal{N}(x_n | \mu_k, I)^{z_{nk}}$$

$$= \prod_n \prod_{k=1}^K c e^{-\frac{1}{2} \|x_n - \mu_k\|^2 \cdot z_{nk}}$$

$$-\log p(X, \mu, z) = \underbrace{\sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} \|x_n - \mu_k\|^2 \cdot z_{nk} + c}_{\mathcal{L}(\mu, z)}$$

# K-means as a Matrix Factorization

Recall the objective

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

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

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

product of  
two matrices

$$\mathbf{M} = \begin{pmatrix} | & & | \\ \boldsymbol{\mu}_1 & \dots & \boldsymbol{\mu}_K \\ | & & | \end{pmatrix}_{D \times K}$$

$$\mathbf{Z} = \begin{pmatrix} - & z_{11} & - \\ & & \\ - & z_{N1} & - \\ & & \\ - & z_{Nk} & - \end{pmatrix}_{N \times K}$$

## 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).