

Machine Learning Course - CS-433

K-Means Clustering

Nov 19, 2020

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

Last updated on: November 17, 2020



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 : 1-hot vector $\in \mathbb{R}^K$

$z_{nk} = \begin{cases} 1 & n \text{ is assigned to cluster } k \\ 0 & \text{o.w.} \end{cases}$

→ 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 μ_k with assignment = k

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

→ discrete constraint

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

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

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

Is this optimization problem easy?

NP-Hard

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

① For all n , compute z_n given μ .

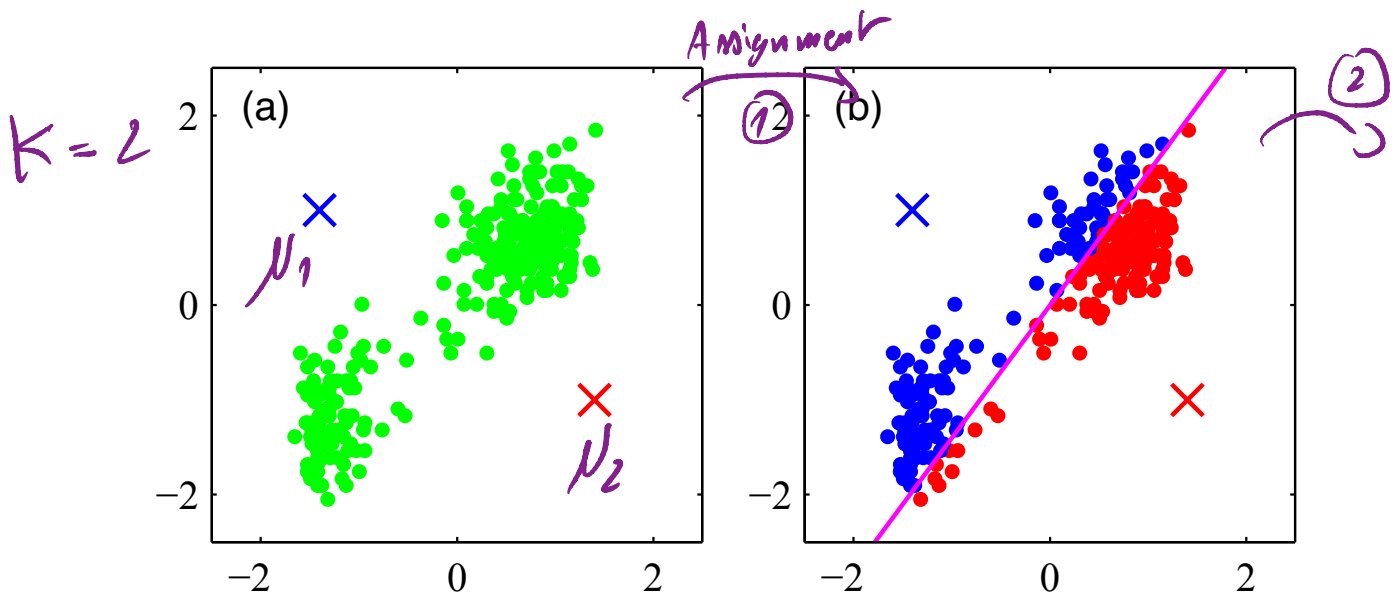
② For all k , compute μ_k given z .

Assignment step

center update



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



①:

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

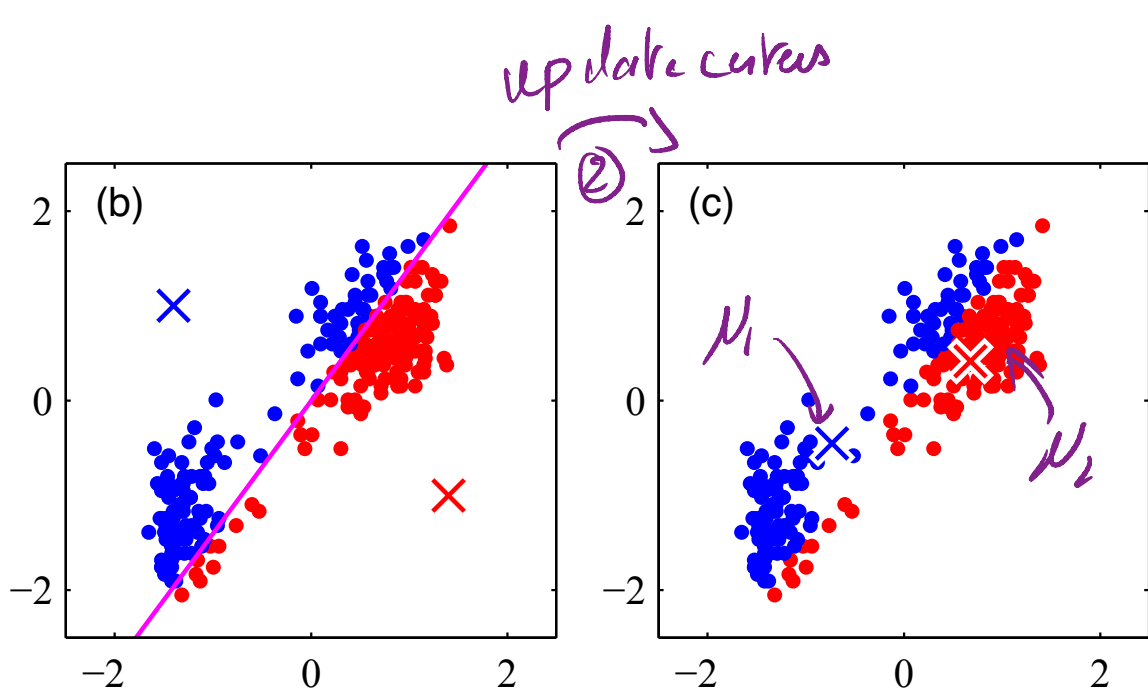
②

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

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

\leftarrow # points assigned to k

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 .

$$z_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|_2^2 \\ 0 & \text{otherwise} \end{cases} \quad \mathcal{O}(N \cdot K \cdot D)$$

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

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

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

$$\nabla Z(\mathbf{z}, \mu) \stackrel{!}{=} 0$$

Coordinate descent

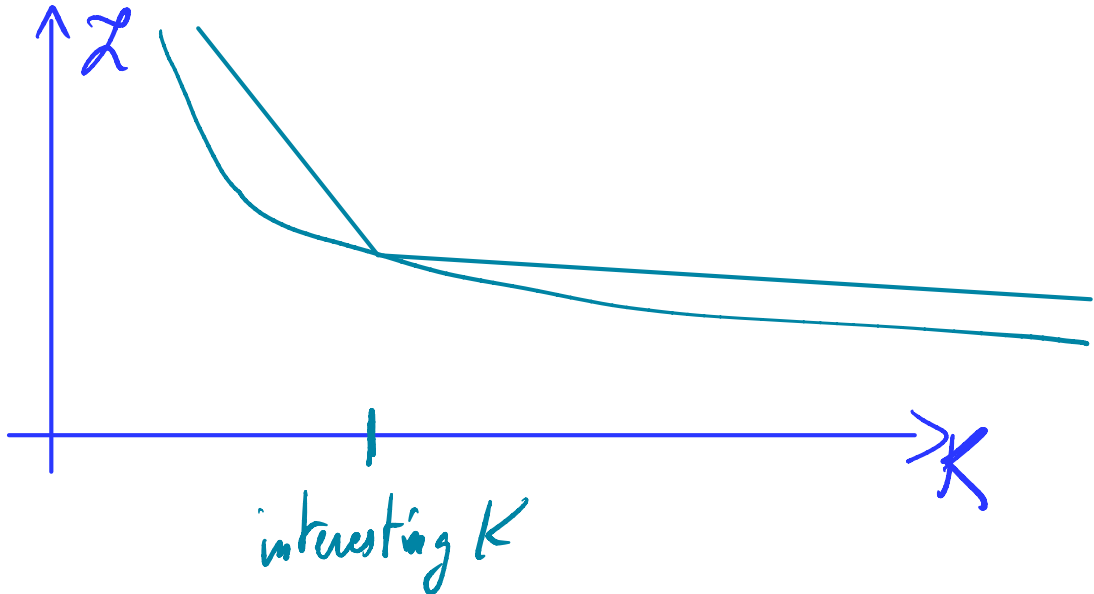
K-means is a coordinate descent algorithm, where, to find $\min_{\mathbf{z}, \mu} \mathcal{L}(\mathbf{z}, \mu)$, we start with some $\mu^{(0)}$ and repeat the following:

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

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

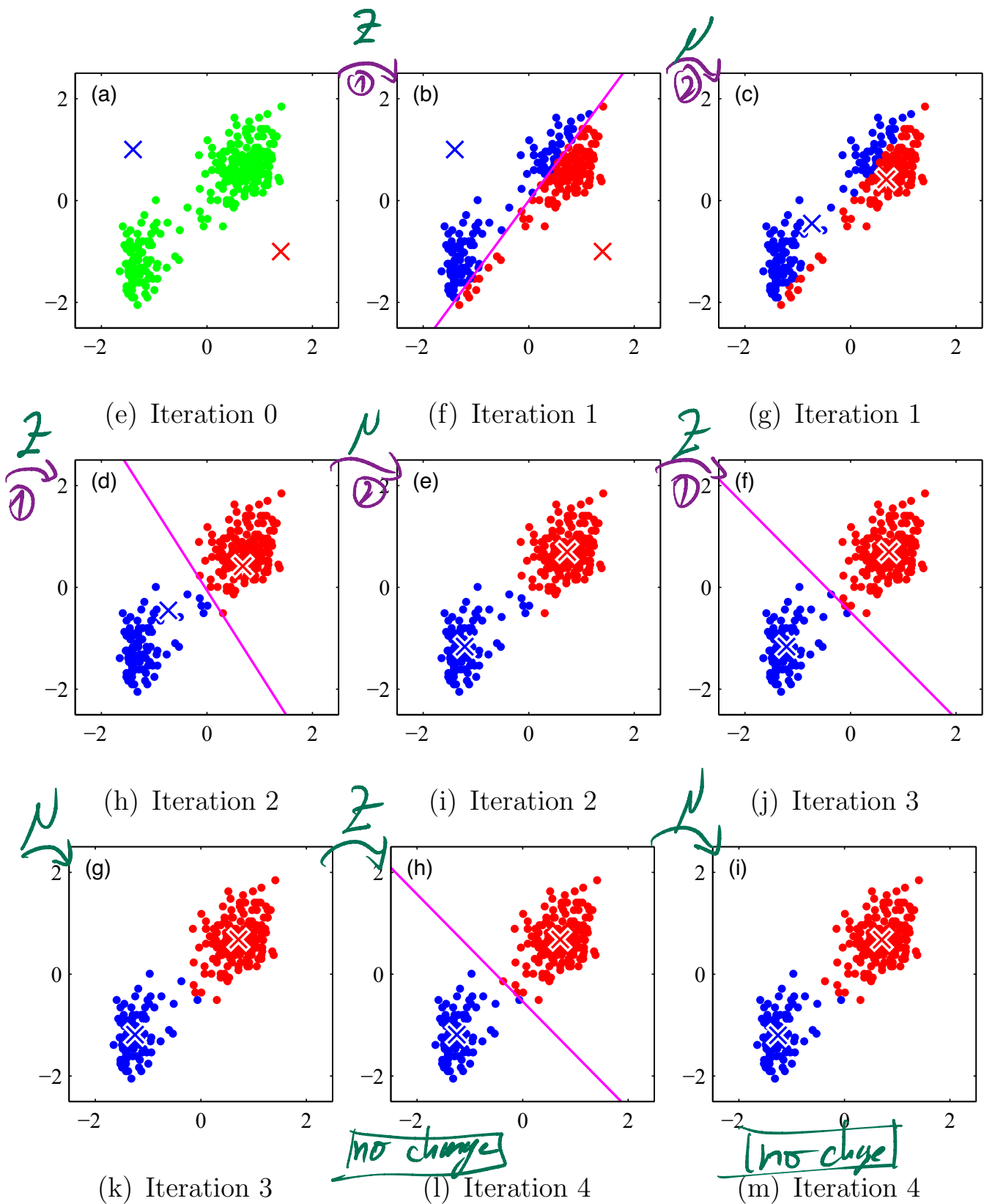
$\Rightarrow \mu$ -update of K-means

How to set K ?



Examples

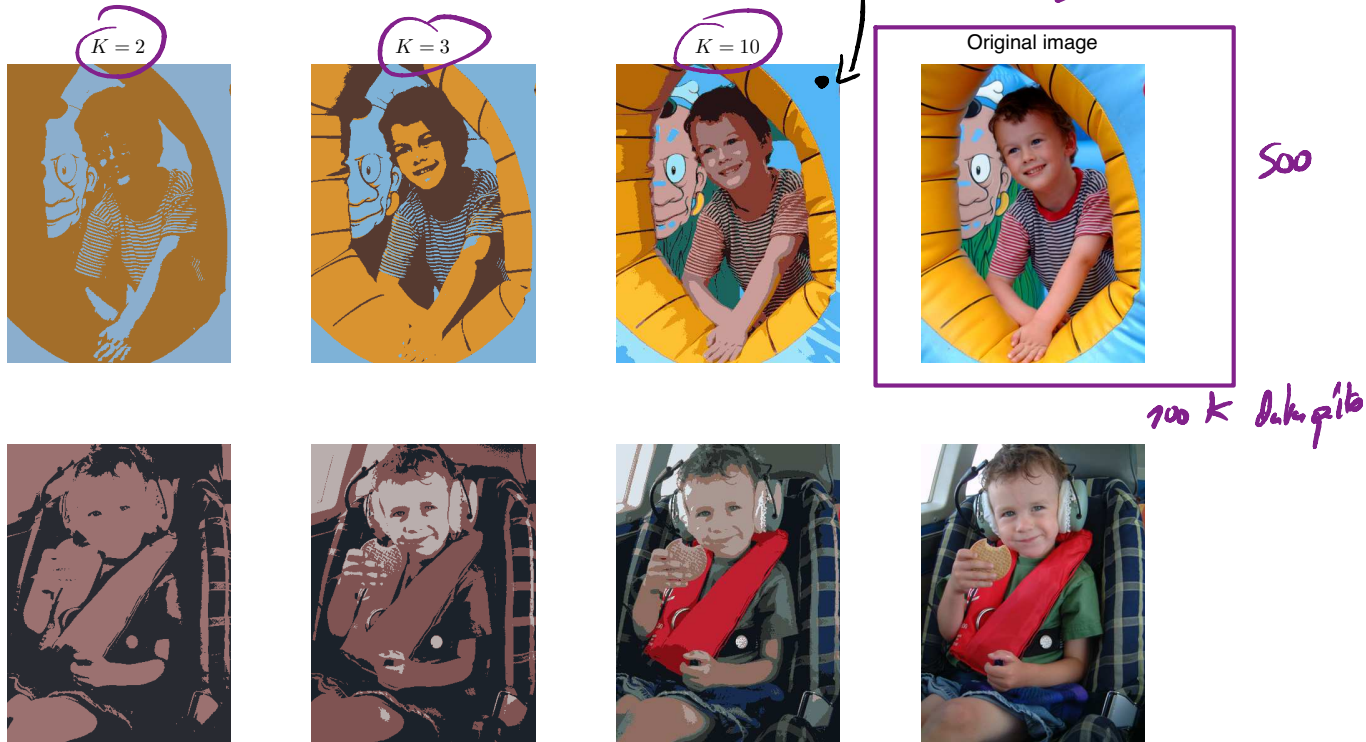
K-means for the “old-faithful” dataset (Bishop’s Figure 9.1)



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

$$X_n: \mathbb{R}^3 \rightarrow \mu_k \in \mathbb{R}^3$$

200



Probabilistic model for K-means

likelihood of X given param. μ, Z

$$\begin{aligned}
 P(X|\mu, Z) &= \prod_{n=1}^N \mathcal{N}(x_n | \mu_k, I) \\
 &= \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(x_n | \mu_k, I)^{z_{nk}} \\
 &= \prod_{n=1}^N \prod_{k=1}^K c \cdot e^{-\frac{1}{2} \|x_n - \mu_k\|^2 \cdot z_{nk}} \\
 -\log(P(X|\mu, Z)) &= \underbrace{\sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} \|x_n - \mu_k\|^2 z_{nk}}_{Z(\mu, Z)} + C'
 \end{aligned}$$

$P(X|\mu, Z)$ \uparrow all dataset

μ_k \uparrow assigned to k for x_n

z_{nk} \leftarrow spherical

K-means as a Matrix Factorization

Recall the objective

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

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

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

$$\min_{\mathbf{M}, \mathbf{Z}} f(\mathbf{M}, \mathbf{Z}^\top)$$

$$\begin{aligned}\mathbf{M} &= (\boldsymbol{\mu}_1 \dots \boldsymbol{\mu}_K)_{D \times K} \\ \mathbf{Z} &= \begin{pmatrix} z_{11} \\ \vdots \\ z_{N1} \end{pmatrix}_{N \times K}\end{aligned}$$

Issues with K-means

1. Computation can be heavy for large N , D and K .

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

2. Clusters are forced to be spherical (e.g. cannot be elliptical).

3. Each example can belong to only one cluster ("hard" cluster assignments).

→ Gaussian Mixture Models solve these.