

Machine Learning Course - CS-433

K-Means Clustering

Nov 8, 2016

©Mohammad Emtiyaz Khan 2015

minor changes by Martin Jaggi 2016



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Clustering

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

choose K

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

assigning n -th point to cluster k
 $z_n = (0, \dots, 1, \dots, 0) \in \mathbb{R}^K$

K-means clustering

Assume K is known.

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

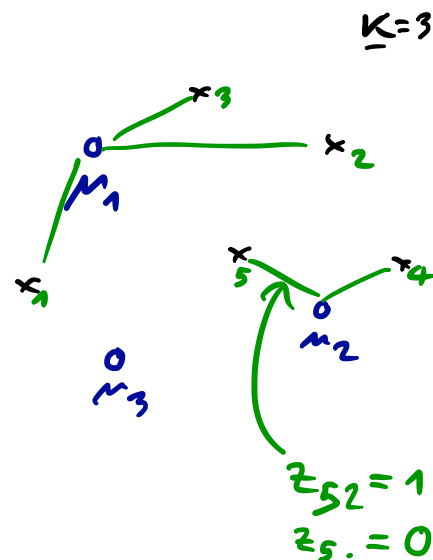
$$\text{s.t. } \boldsymbol{\mu}_k \in \mathbb{R}^D, z_{nk} \in \{0, 1\}, \sum_{k=1}^K z_{nk} = 1,$$

\mathbf{z} : assignment matrix

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

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

$$\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K]^\top$$



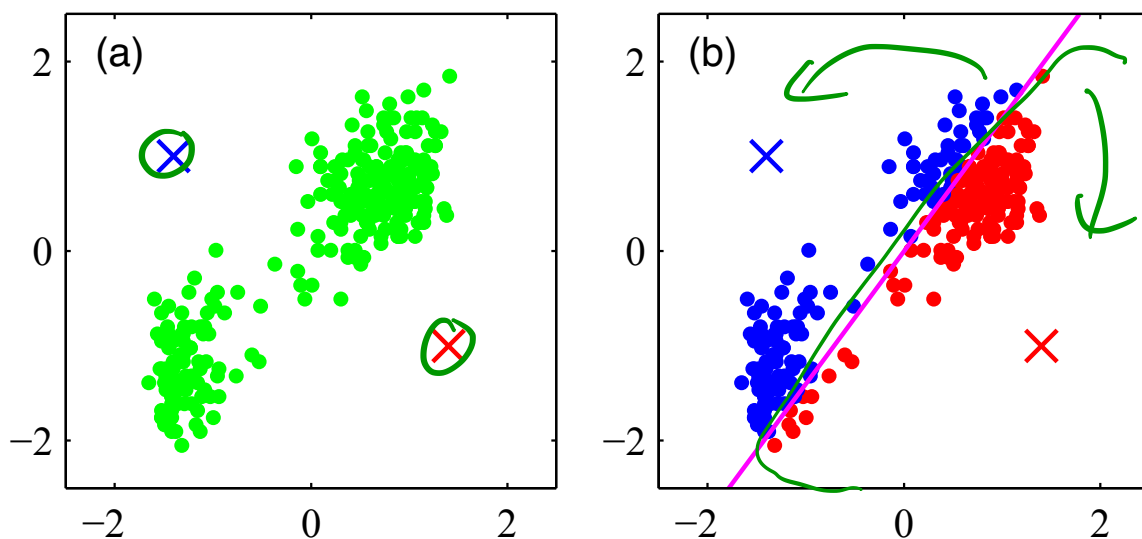
Is this optimization problem easy?

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

1. For all n , compute z_n given μ . update assignments
2. For all k , compute μ_k given z . update means/centers

Example
 $K=2$

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



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

$\forall n$

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

fix k

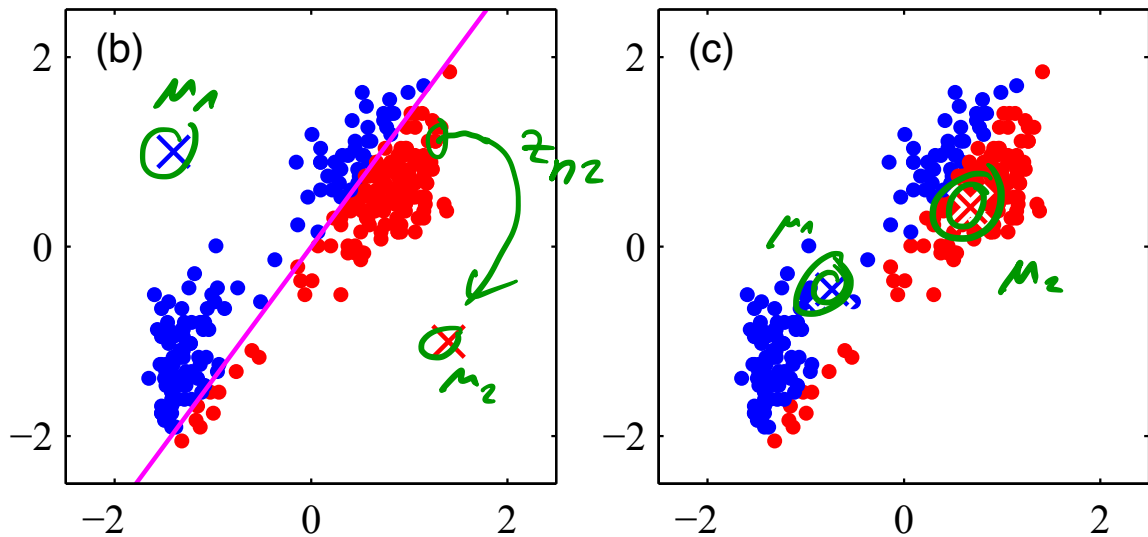
$$\mu_k = \frac{\sum_{n=1}^N z_{nk} \mathbf{x}_n}{\sum_{n=1}^N z_{nk}}$$

update means

Hence, the name 'K-means'.

$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

$$\nabla_{\mu} \mathcal{L} = \dots = 0$$



Summary of K-means

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

- ① For all n , compute \mathbf{z}_n given μ . *update assignments \mathbf{z}*

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

- ② For all k , compute μ_k given \mathbf{z} . *update means μ*

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

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

$$\mathcal{L} \geq 0$$

Coordinate descent

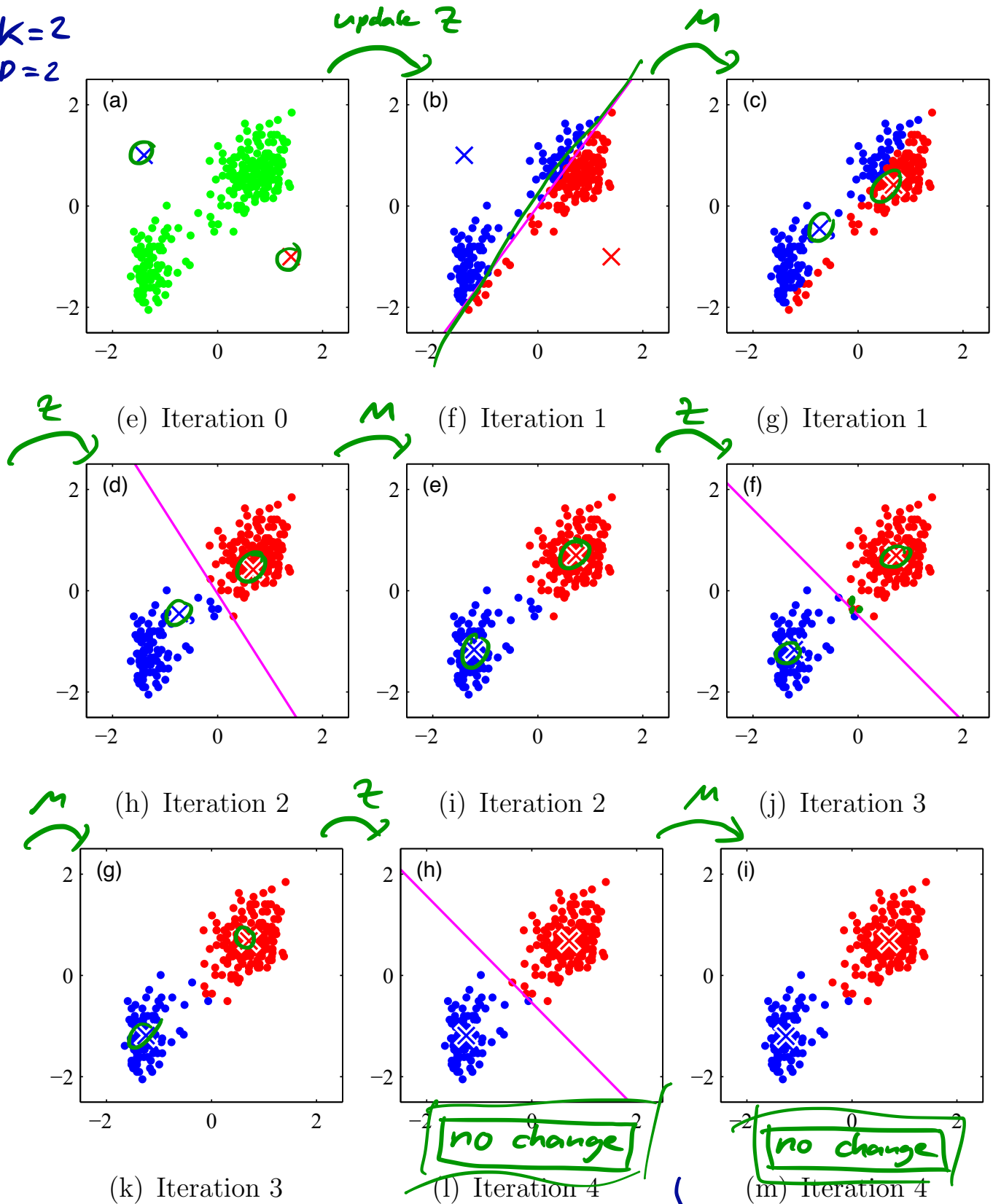
K-means is a coordinate descent algorithm, where, to find $\min_{\mathbf{z}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{z}, \boldsymbol{\mu})$, we start with some $\boldsymbol{\mu}^{(0)}$ and repeat the following:

$$\begin{aligned} \mathbf{z}^{(t+1)} &:= \arg \min_{\mathbf{z}} \mathcal{L}(\mathbf{z}, \boldsymbol{\mu}^{(t)}) && \Leftrightarrow \text{step 1: update } \mathbf{z} \text{ } \mathbf{z} \in \{0,1\}^n \\ \boldsymbol{\mu}^{(t+1)} &:= \arg \min_{\boldsymbol{\mu}} \mathcal{L}(\mathbf{z}^{(t+1)}, \boldsymbol{\mu}) && \Leftrightarrow \text{step 2: update } \boldsymbol{\mu} \end{aligned}$$

Examples

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

$K=2$
 $D=2$



How to choose K ?



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



Probabilistic model for K-means

likelihood of X given parameters μ, z

$$\begin{aligned}
 p(x_n | \mu, z) &= \prod_n \mathcal{N}(x_n | \mu_k, \mathbf{I}) \\
 p(\mathbf{x} | \mu, z) &= \prod_n \prod_k \mathcal{N}(x_n | \mu_k, \mathbf{I})^{z_{nk}} \\
 &= \prod_n \prod_k \left(\exp^{-\frac{1}{2} \|x_n - \mu_k\|_2^2} \right)^{z_{nk}} \\
 -\log p(\mathbf{x} | \mu, z) &= \sum_n \sum_k \frac{1}{2} \|x_n - \mu_k\|_2^2 z_{nk} \\
 &= \mathcal{L}(\mu, z) + \text{const}
 \end{aligned}$$

Diagram illustrating the probabilistic model for K-means. A box labeled X has a downward arrow pointing to it. A green circle contains the text $k \text{ s.t. } z_{nk}=1$. A green circle contains the text \sum . A green arrow points from the \sum circle to the expression $(x - \mu)^T \mathbf{1} (x - \mu)$.

$$\|A\|_{\text{Frob}}^2 = \sum_{i,j} (A_{ij})^2$$

K-means as a Matrix Factorization

Recall the objective

Recall the objective

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

Handwritten notes and definitions:

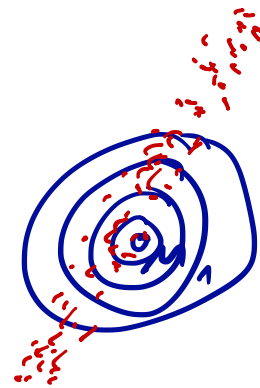
- $= \|\mathbf{x}_n - \mathbf{M} \mathbf{z}_{n:}^\top\|_2^2$ (with a blue arrow pointing from the inner product in the objective to this expression)
- $\mathbf{z}_{n:} = (0, \dots, 1, \dots, 0)$ (with a blue arrow pointing from the z_{nk} term in the objective to this expression)
- $\mathbf{z} = \begin{pmatrix} \mathbf{z}_{1:} \\ \vdots \\ \mathbf{z}_{N:} \end{pmatrix}_N$ (with a blue arrow pointing from the \mathbf{z} in the Frobenius norm to this expression)
- $\mathbf{M} = \mathbf{M} = \begin{pmatrix} | & & | \\ \mu_1 & \dots & \mu_K \\ | & & | \end{pmatrix}_D^K$ (with a blue arrow pointing from the $\mathbf{M} \mathbf{z}$ term in the Frobenius norm to this expression)
- $\Leftrightarrow \|\mathbf{X} - \mathbf{M} \mathbf{Z}^\top\|_{\text{Frob}}^2$ (with a blue circle around the equivalence symbol and a blue arrow pointing from the objective to this expression)

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

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

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).



ToDo

1. Understand the iterative algorithm for K-means. Why is the problem difficult to optimize and how does the iterative algorithm make it simpler?
2. What is the computational complexity of K-means?
3. Derive the probabilistic model associated with the cost function.