**Machine Learning Course - CS-433**

# K-Means Clustering

Nov 5, 2019

**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 $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \ldots, \boldsymbol{\mu}_K$ and cluster assignments $z_n \in \{1, 2, \ldots, K\}$ for all $n = 1, 2, \ldots, N$ data vectors $\mathbf{x}_n \in \mathbb{R}^D$.

$$x_1, \ldots, x_N$$

specify
$K = \# groups$

# K-means clustering

Assume $K$ is known.

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

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

distance to own representative

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

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

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

$\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \ldots, \boldsymbol{\mu}_K]^\top$
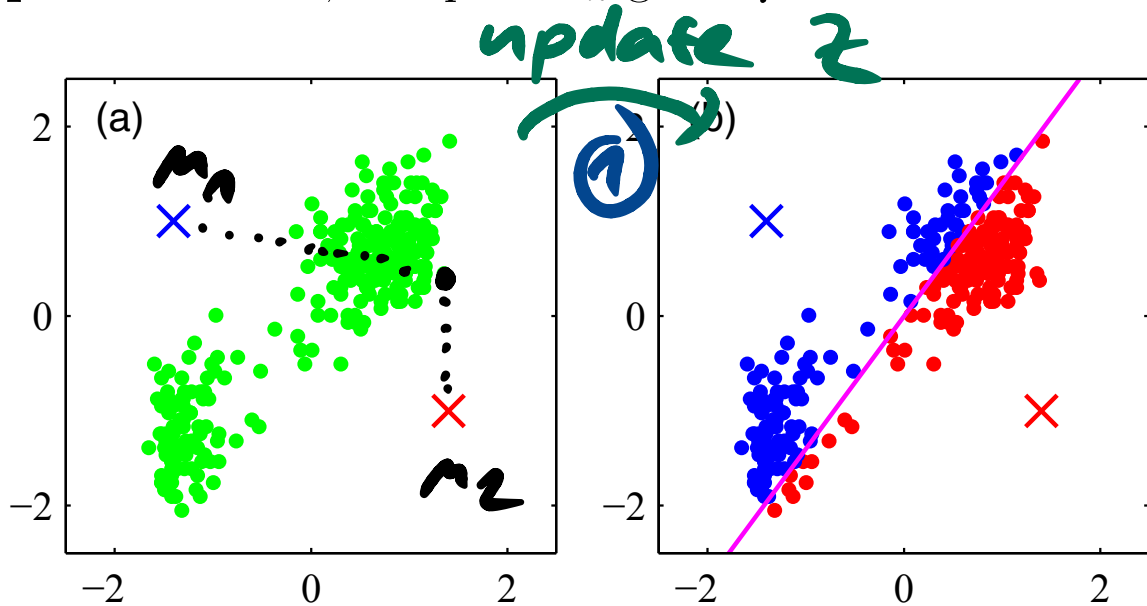
Is this optimization problem easy?

non-convex

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

2. For all $k$, compute $\boldsymbol{\mu}_k$ given $\mathbf{z}$.

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

update z



(a)  (b)

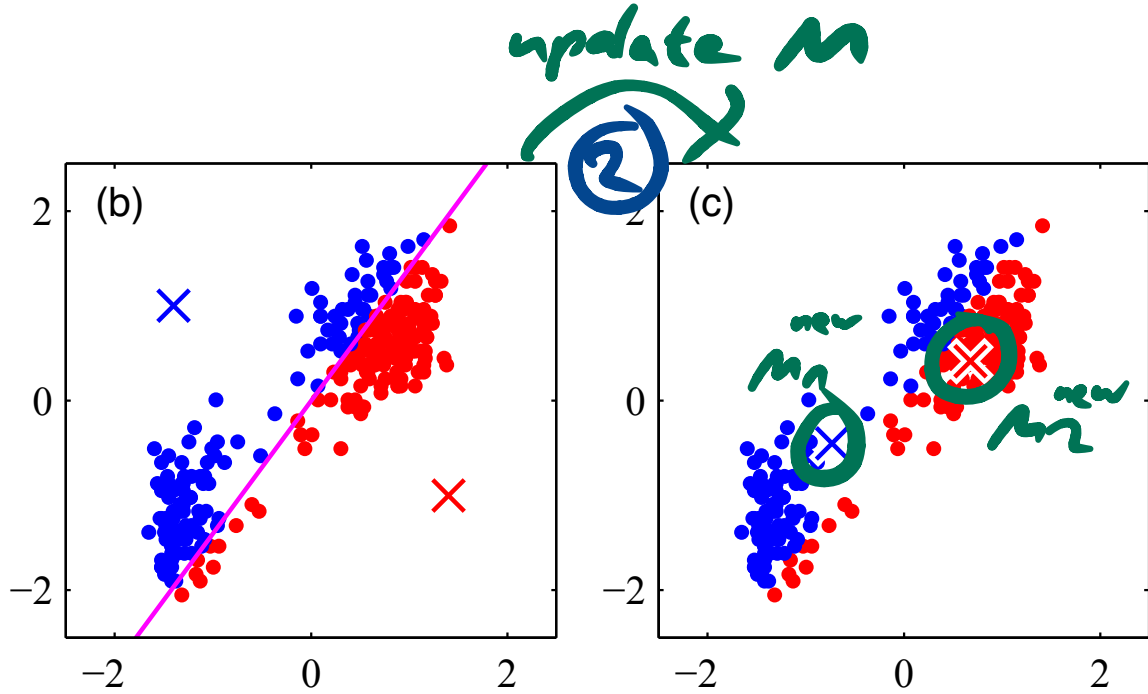① assign to coset $\boldsymbol{\mu}_k$

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

grap k

② **Step 2:** For all $k$, compute $\boldsymbol{\mu}_k$ given $\mathbf{z}$.
Take derivative w.r.t. $\boldsymbol{\mu}_k$ to get:

$$\boldsymbol{\mu}_k := \frac{\sum_{n=1}^{N} z_{nk}\mathbf{x}_n}{\sum_{n=1}^{N} z_{nk}}$$

Hence, the name 'K-means'.

update $\mu$

(b) (c)

$m_3^{new}$  $m_2^{new}$

## Summary of K-means

Initialize $\boldsymbol{\mu}_k \, \forall k$, then iterate:

iterations

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

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

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

cost $\mathcal{O}(D)$

2. For all $k$, compute $\boldsymbol{\mu}_k$ given $\mathbf{z}$.

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

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

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

$\mathcal{L} \geq 0$

$\nabla \mathcal{L}(z, \mu) \stackrel{!}{=} 0$

## Coordinate descent

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

*„Block-coordinate descent"*

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

② $\mu^{(t+1)} := \arg\min_{\mu} \mathcal{L}(z^{(t+1)}, \mu)$ ← $\nabla_{\mu} \mathcal{L} \doteq 0$

How to set K ?

# Examples

K-means for the "old-faithful" dataset (Bishop's Figure 9.1)



(e) Iteration 0      (f) Iteration 1      (g) Iteration 1

(h) Iteration 2      (i) Iteration 2      (j) Iteration 3

(k) Iteration 3      (l) Iteration 4      (m) Iteration 4

$M_n = \begin{bmatrix} R \\ G \\ B \end{bmatrix}$

$200$

$K = 2$     $K = 3$     $K = 10$     Original image

$500$

$100k$



## Probabilistic model for K-means

Likelihood of $X$ given parameters $M, z$

$$p(x_n | m, z) = \prod_{n=1} \mathcal{N}(x_n | M_k, I)$$

$\hookleftarrow$ assignment for $x_n$

$$p(X | m, z) = \prod_{n=1}^{N} \prod_{k=1}^{K} \mathcal{N}(x_n | M_k, I)^{z_{nk}}$$

$$= \prod_{n=1}^{N} \prod_{k=1}^{N} c \cdot e^{-\frac{1}{2} \| x_n - M_k \|^2 \cdot z_{nk}}$$

$$-\log p(X | m, z) = \sum_{n} \sum_{k} \frac{1}{2} \| x_n - M_k \|^2 z_{nk} + c'$$

$$= \mathcal{L}(M, z)$$

# K-means as a Matrix Factorization

Recall the objective

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

$$= \|\mathbf{X}^\top - \mathbf{M}\mathbf{Z}^\top\|_{\text{Frob}}^2$$

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

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

$f(M \cdot Z^\top)$

entry-wise

$$M = \begin{pmatrix} | & & | \\ m_n & \dots & m_k \\ | & & | \end{pmatrix}$$
$D \cdot K$

$$Z = \begin{pmatrix} - z_1 - \\ \\ - z_N - \end{pmatrix}$$
$N \cdot K$

## 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).