

annotated
Version

Machine Learning Course - CS-433

K-Means Clustering

Nov 9, 2017

©Mohammad Emtiyaz Khan 2015

minor changes by Martin Jaggi 2016

minor changes by Martin Jaggi 2017

Last updated: November 7, 2017



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Clustering

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

The goal is to ^{\mathbb{R}^D} 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 $\mathbf{x}_n \in \mathbb{R}^D$. ^{means}

Specify # groups K

K-means clustering

Assume K is known.

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

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

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

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

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

Is this optimization problem easy?

non-convex

Binary Assignment

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

$$\mathbf{z}_n = (0, \dots, 1, \dots, 0)$$

actual k

$K=3$

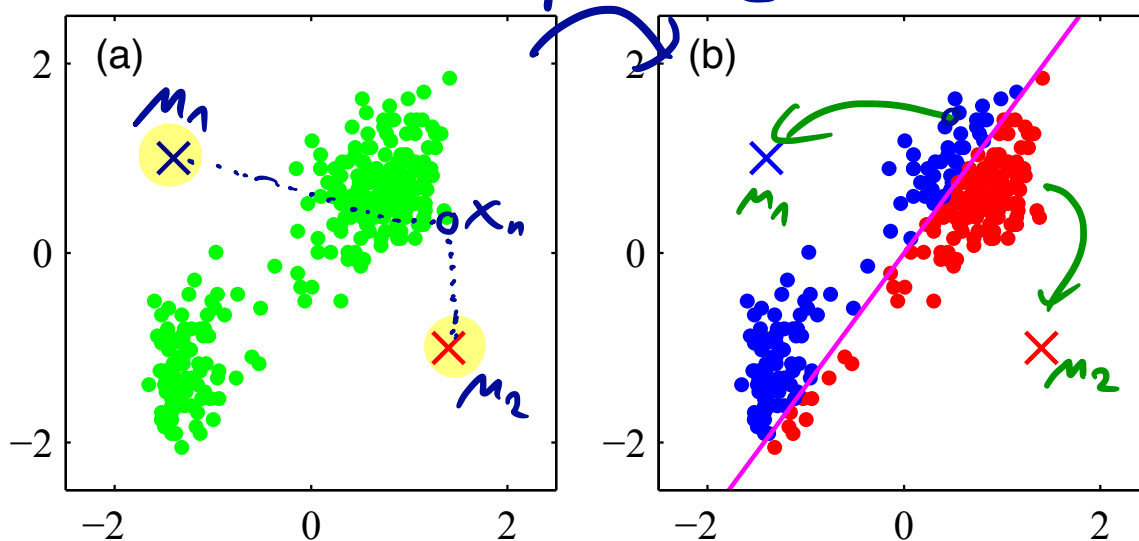


Algorithm: Initialize $\mu_k \forall k$,
then iterate:

- ① For all n , compute z_n given μ .
- ② For all k , compute μ_k given z .

$k=2$

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



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

for every n

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

Fix k

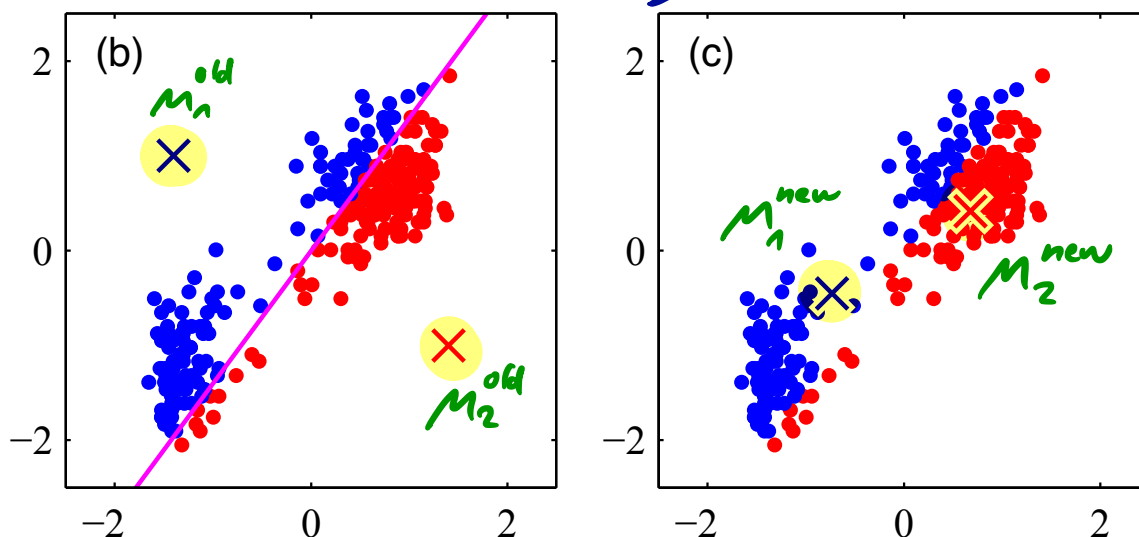
$$\mu_k := \frac{\sum_{n=1}^N z_{nk} \mathbf{x}_n}{\sum_{n=1}^N z_{nk}} = \text{Average Mean of all points of group } k$$

Hence, the name 'K-means'.

$$\min_{\mu} \mathcal{L}(\mu, z) : \nabla_{\mu} \mathcal{L}(\mu, z) \stackrel{!}{=} 0$$

step 2

update μ



Summary of K-means

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

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

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

$O(N \cdot K \cdot D)$

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

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

$\forall k \quad O(N \cdot K \cdot D)$

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

$\mathcal{L} \geq 0$

repeat

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:

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

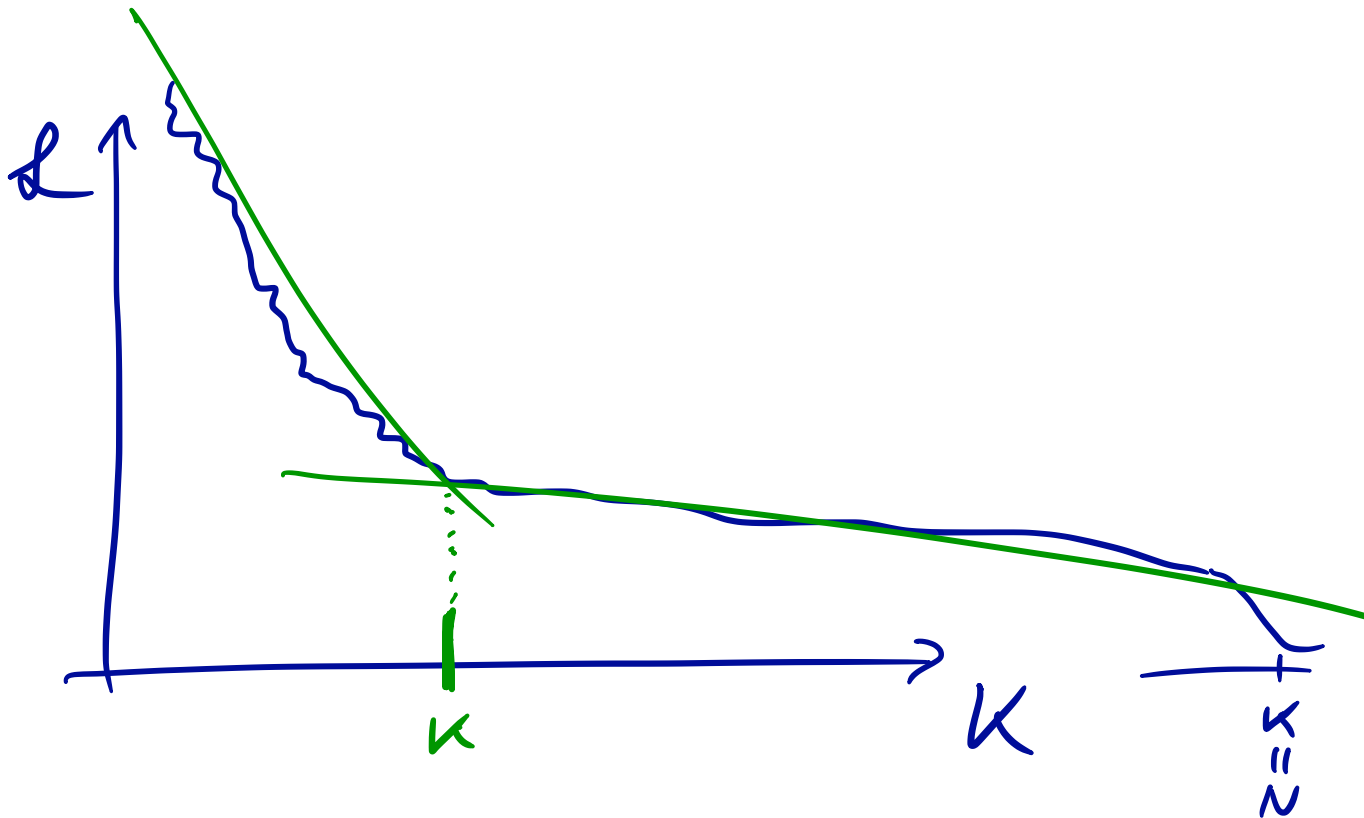
$$\mu^{(t+1)} := \arg \min_{\mu} \mathcal{L}(\mathbf{z}^{(t+1)}, \mu)$$

update ~~assignments~~

update means

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

How to set K

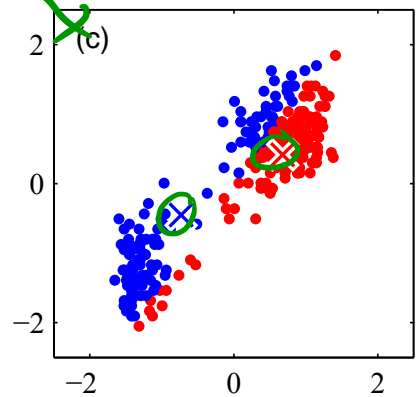
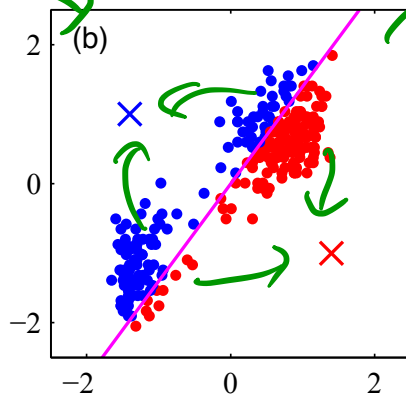
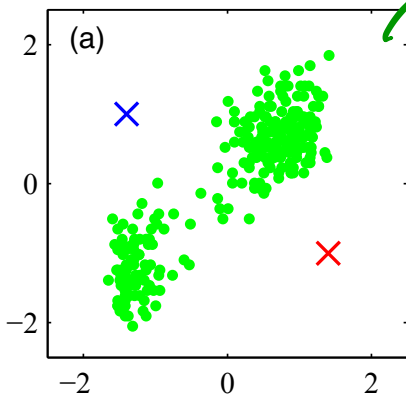


Examples

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

update z

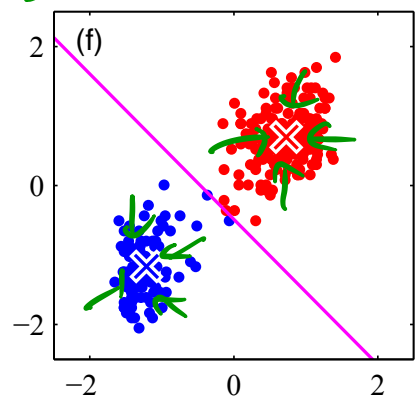
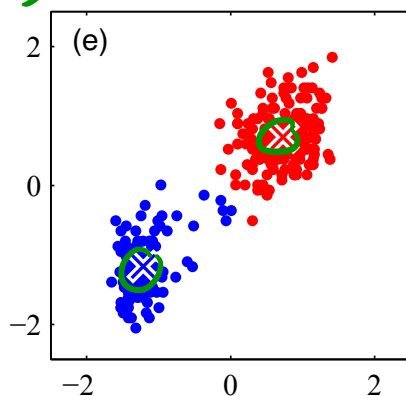
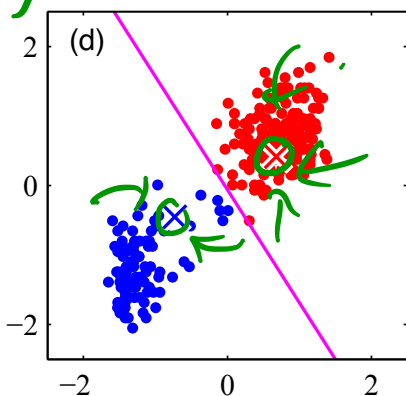
update μ



(e) Iteration 0

(f) Iteration 1

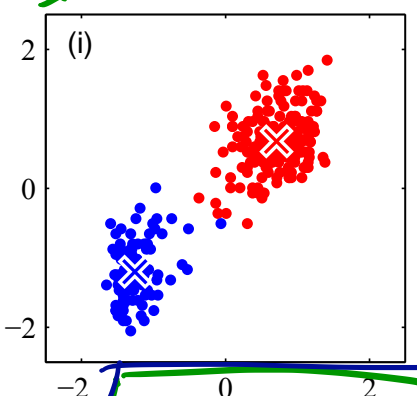
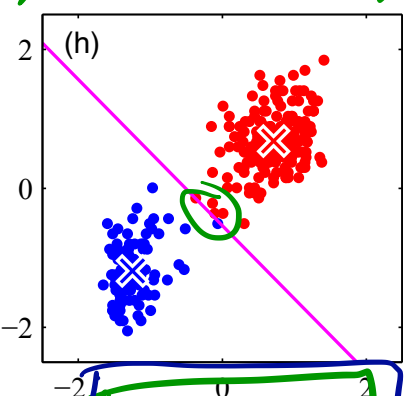
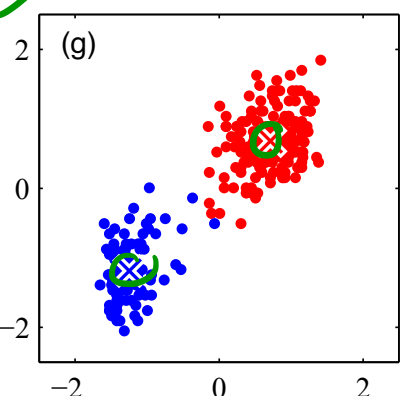
(g) Iteration 1



(h) Iteration 2

(i) Iteration 2

(j) Iteration 3



(k) Iteration 3

(l) Iteration 4

(m) Iteration 4

no change

no change

stop

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



Probabilistic model for K-means

Likelihood of X given parameters μ, z

$$\begin{aligned}
 p(x_n | \mu, z) &= \prod_{n=1}^N N(x_n | \mu_k, I) \\
 p(X | \mu, z) &= \prod_{n=1}^N \prod_{k=1}^K N(x_n | \mu_k, I)^{z_{nk}} \\
 &= \prod_{n=1}^N \prod_{k=1}^K \left(e^{-\frac{1}{2} \|x_n - \mu_k\|^2} \right)^{z_{nk}} = \prod_{n=1}^N \prod_{k=1}^K e^{-\frac{1}{2} (x_n - \mu_k)^T (x_n - \mu_k) z_{nk}}
 \end{aligned}$$

↑ k for which $z_{nk}=1$

$$\begin{aligned}
 -\log p(x_n | \mu, z) &= \sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} \|x_n - \mu_k\|^2 z_{nk} + c' \\
 -\log p(X | \mu, z) &= \sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} \|x_n - \mu_k\|^2 z_{nk} + c' \\
 &= \mathcal{L}(\mu, z) + \text{const}
 \end{aligned}$$

K-means as a Matrix Factorization

Recall the objective

$$\|x_n - M z_{n:}^T\|^2$$

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

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

Data:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_N \end{pmatrix}_{N \times D}$$

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

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

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

$$= \sum_j \|A_{:,j}\|^2 = \sum_i \|A_{i,:}\|^2$$

$$\mathbf{M} = \begin{pmatrix} \mu_{11} & \dots & \mu_{1K} \\ \vdots & & \vdots \\ \mu_{N1} & \dots & \mu_{NK} \end{pmatrix}_{D \times K}$$

$$\mathbf{Z} = \begin{pmatrix} z_{11} & \dots & z_{1K} \\ \vdots & & \vdots \\ z_{N1} & \dots & z_{NK} \end{pmatrix}_{N \times K}$$

$$z_{n:} = (0, \dots, 1, \dots, 0) \in \mathbb{R}^K$$

Issues with K-means

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

problematic for very large 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).

Exercises

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.