

# K-means Clustering

Mohammad Emtiyaz Khan  
EPFL

Nov 3, 2015



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

# 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  $r_n \in \{1, 2, \dots, K\}$  for all  $n = 1, 2, \dots, N$  data vectors.

## K-means clustering

Assume  $K$  is known.

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

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

$$\text{where } \mathbf{r}_n = [r_{n1}, r_{n2}, \dots, r_{nK}]^T$$

$$\mathbf{r} = [\mathbf{r}_1^T, \mathbf{r}_2^T, \dots, \mathbf{r}_N^T]^T$$

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

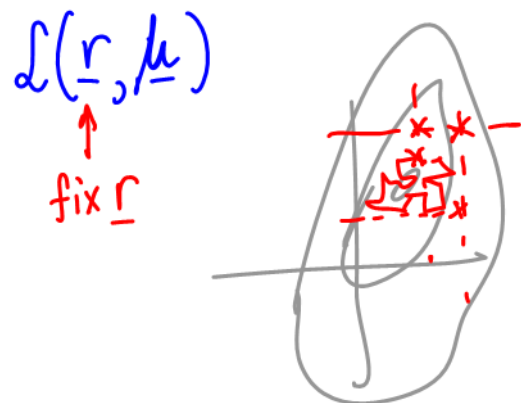
Is this optimization problem easy?

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

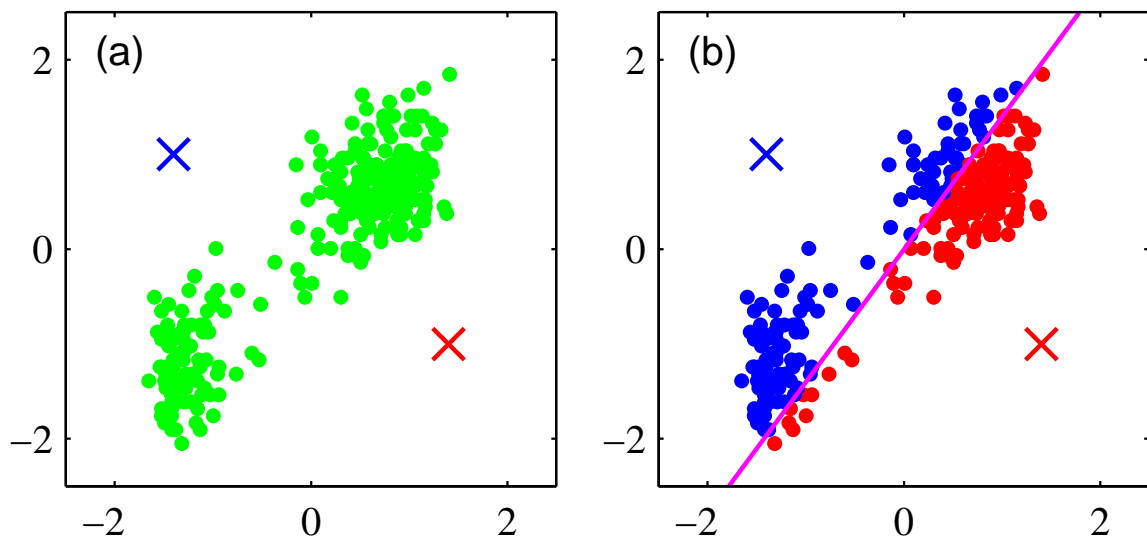
1. For all  $n$ , compute  $\mathbf{r}_n$  given  $\mu$ .
2. For all  $k$ , compute  $\mu_k$  given  $\mathbf{r}$ .

Handwritten notes and diagrams illustrating K-means clustering:

- Top: Three clusters of points labeled  $r_n=1$ ,  $r_n=2$ , and  $r_n=3$ . A point in the first cluster is circled in red and labeled  $\mu_1$ .
- Middle: Objective function:  $\sum_{n:r_n=1} \|\mathbf{x}_n - \mu_1\|_2^2 + \sum_{n:r_n=2} \|\mathbf{x}_n - \mu_2\|_2^2 + \sum_{n:r_n=3} \|\mathbf{x}_n - \mu_3\|_2^2$
- Below:  $\|\mathbf{x}_n - \mu_k\|_2^2$  is  $L_2$  norm  $= (\mathbf{x}_n - \mu_k)'(\mathbf{x}_n - \mu_k)$
- Bottom:  $\mu_1, \mu_2, \mu_3 \in \mathbb{R}^D, r_n \in \{1, 2, 3\}$

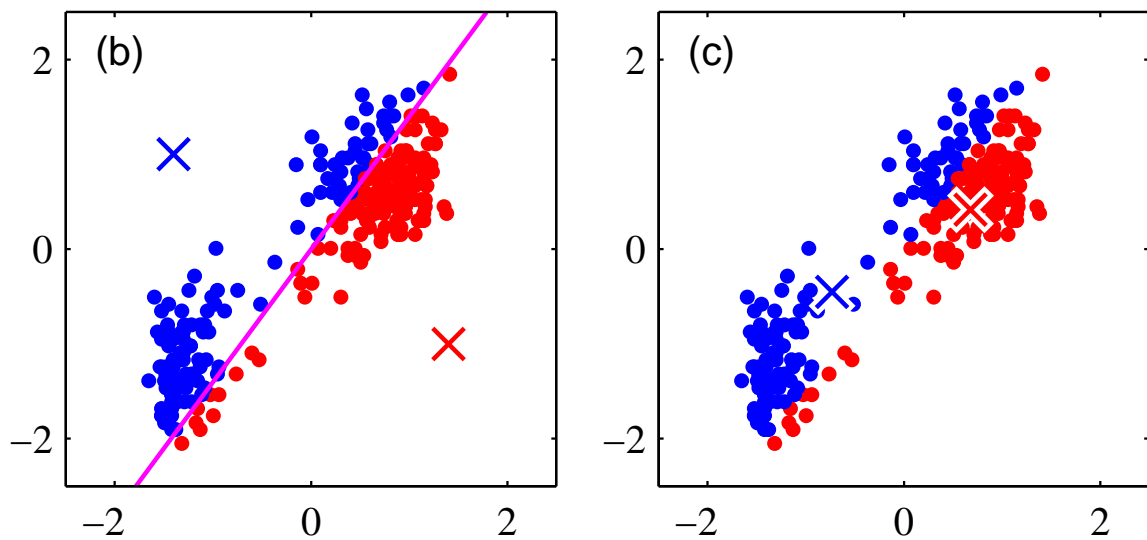


Step 1: For all  $n$ , compute  $\mathbf{r}_n$  given  $\boldsymbol{\mu}$ .



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

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



Take derivative w.r.t.  $\boldsymbol{\mu}_k$  to get:  $\mathcal{L}(\boldsymbol{\mu}_i) = \sum_{n=1}^N r_{ni} (\mathbf{x}_n - \boldsymbol{\mu}_i)^T (\mathbf{x}_n - \boldsymbol{\mu}_i) + \text{const}$

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

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_i} = - \sum_{n=1}^N r_{ni} (\mathbf{x}_n - \boldsymbol{\mu}_i) = 0$$

Hence, the name 'K-means'.

What is K-median? 2

$$\Rightarrow - \sum_n r_{ni} \mathbf{x}_n + \sum_n r_{ni} \boldsymbol{\mu}_i = 0$$

# Summary of K-means

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

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

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

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

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

Example

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

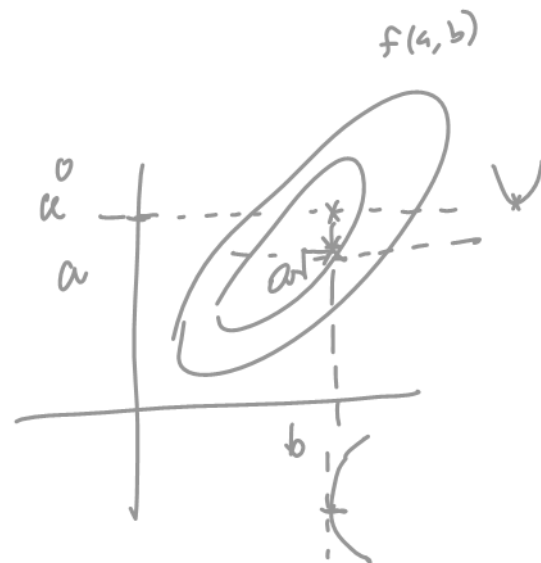
## Coordinate descent

K-means is a coordinate descent algorithm where, to find  $\min_{\mathbf{a}, \mathbf{b}} f(\mathbf{a}, \mathbf{b})$ , we start with some  $\mathbf{b}_0$  and repeat the following:

$$\mathbf{a}_{k+1} = \arg \min_b f(\mathbf{a}, \mathbf{b}_k)$$

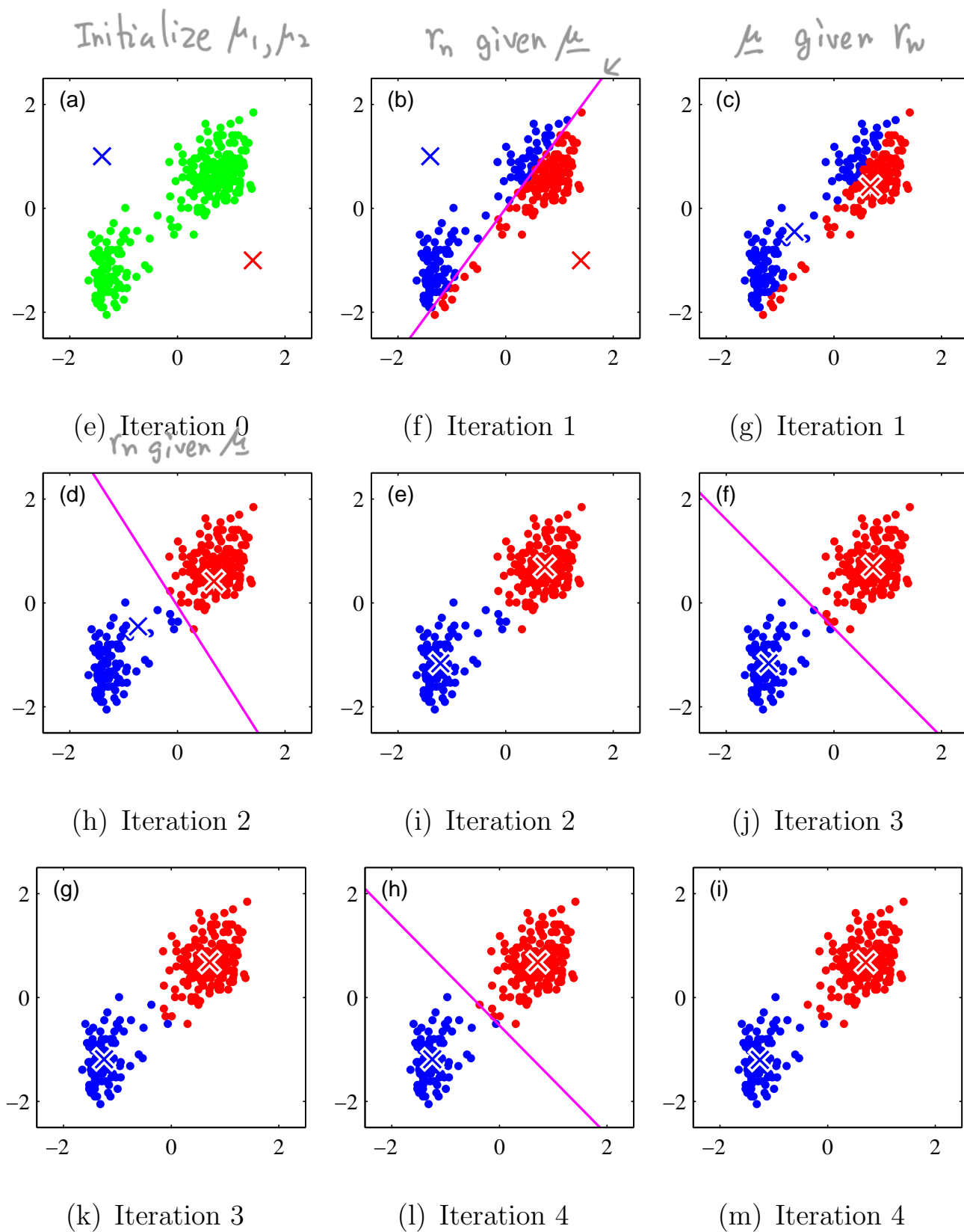
$$\mathbf{b}_{k+1} = \arg \min_a f(\mathbf{a}_{k+1}, \mathbf{b})$$

Convergence is assured when both subproblems have a unique minimum.

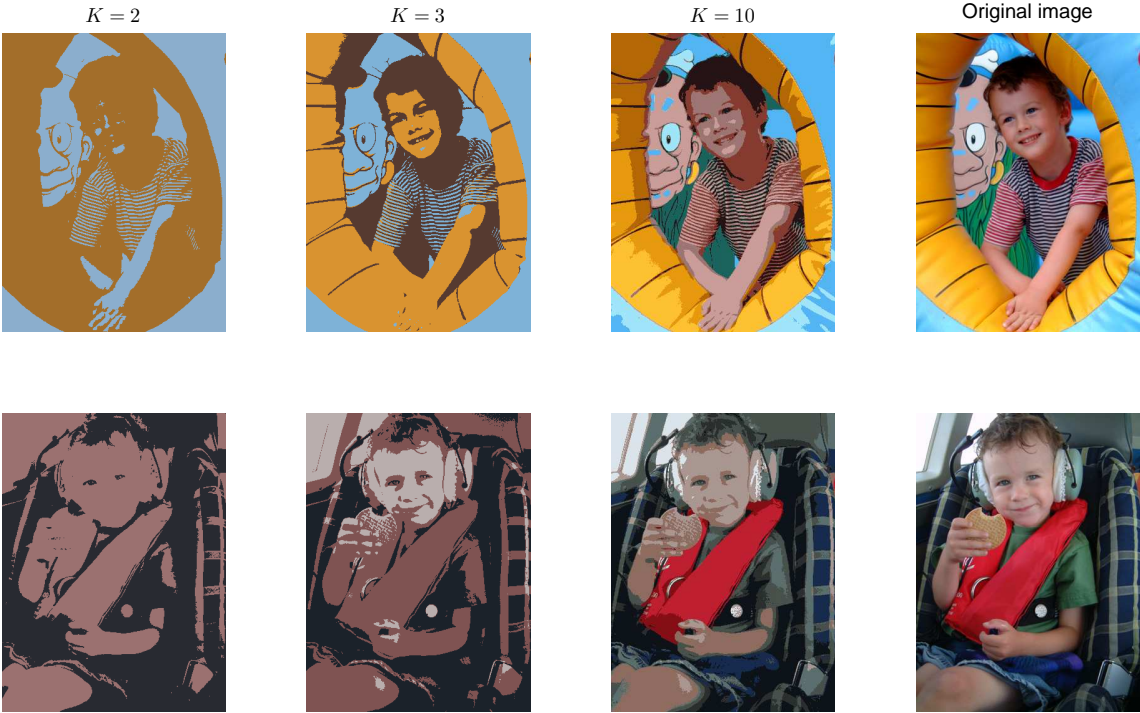


# Examples

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



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



$$-\frac{1}{2} \log p(\underline{r}, \underline{\mu}) = \mathcal{L}(\underline{r}, \underline{\mu}) + \text{const}$$

## Probabilistic model for K-means

$$\mathcal{L}(\underline{r}, \underline{\mu}) = \sum_{k=1}^K \sum_{n=1}^N r_{nk} \underbrace{\left( \frac{x_n}{\mu_k} \right)^2}$$

$$p(\underline{r}, \underline{\mu}) \propto \exp^{-\sum_k \sum_n r_{nk} \left( \frac{x_n}{\mu_k} \right)^2}$$

$$p(\underline{r}, \underline{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \left[ \mathcal{N} \left( \frac{x_n}{\mu_k}, \mathbf{I} \right)^{r_{nk}} \right] = \prod_n \prod_k \exp^{-\frac{1}{2} r_{nk} \left( \frac{x_n}{\mu_k} \right)^2}$$

$$= \prod_n \prod_k \left[ \exp^{-\frac{1}{2} \left( \frac{x_n}{\mu_k} \right)^2} \right]^{r_{nk}}$$

## Issues with K-means

1. Computation can be heavy for large  $N$  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).

Q: What is the computational complexity?  
(Do it as homework)

## To do

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. Derive the probabilistic model associated with the cost function.
3. What is the computational complexity of K-means?