



Machine Learning 1

Lecture 9.2 - Unsupervised Learning
K-Means Clustering

Erik Bekkers

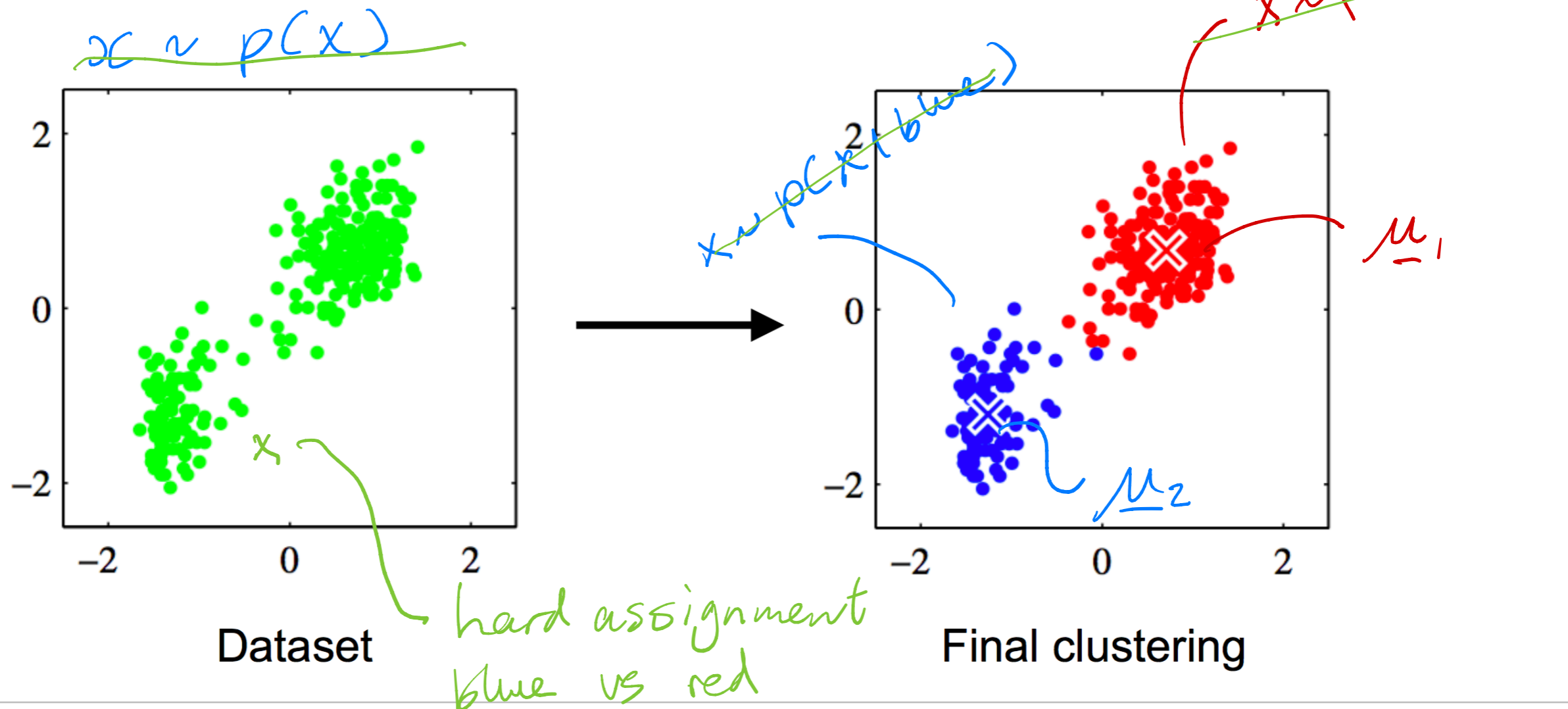
(Bishop 9.1)



Clustering with K-means

No probabilistic interpretation for now

- ▶ Data: a sample of points \mathcal{X} (without a target)
- ▶ Goal: every single data point is assigned to a cluster – a **discrete latent variable** $z = \{\text{blue}, \text{red}\}$



K-means clustering as minimization problem

▶ Data: $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

▶ Goal: partition into K clusters by minimizing

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

$z_{nk} \in \{0, 1\}$

A handwritten diagram in blue ink. It shows a vector $\underline{z}_n = \begin{pmatrix} z_{n1} \\ z_{n2} \\ \vdots \\ z_{nk} \end{pmatrix}$ followed by an equals sign and another vector $\begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \end{pmatrix}$. A blue arrow points from the z_{nk} term in the equation above to the z_{nk} term in the first vector.

▶ Find cluster **assignments** (latent var) $z_{nk} \in \{0, 1\}$

and cluster means $\boldsymbol{\mu}_k \in \mathbb{R}^D$

▶ **1-hot-encoding**: $z_{nk} = 1$ if and only if point n is assigned to cluster k

Minimize J (EM algorithm)

▶ Initialize with a random $\mu_k \in \mathbb{R}^D$

▶ Repeat until convergence:

▶ Find the **assignment** (fixed means) – **E-step**

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

▶ **Find the means** (fixed assignments) – **M-step**

$$\mu_k = \frac{\sum_n z_{nk} \mathbf{x}_n}{\sum_n z_{nk}}$$

"Expectation"

"Maximization"

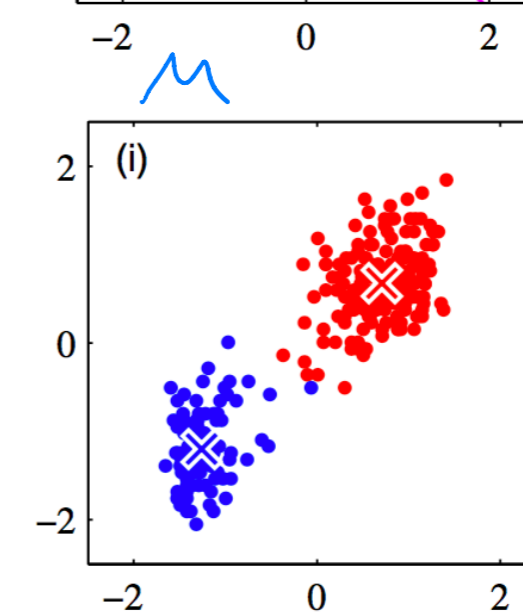
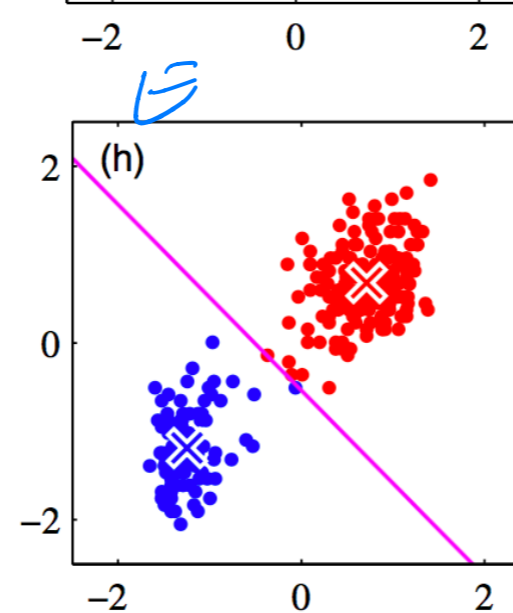
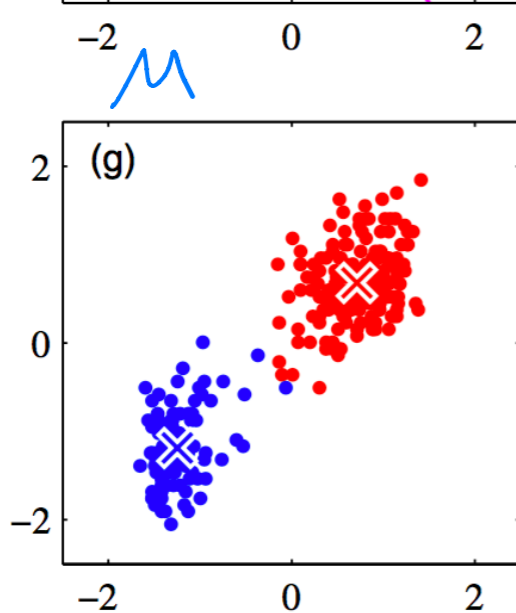
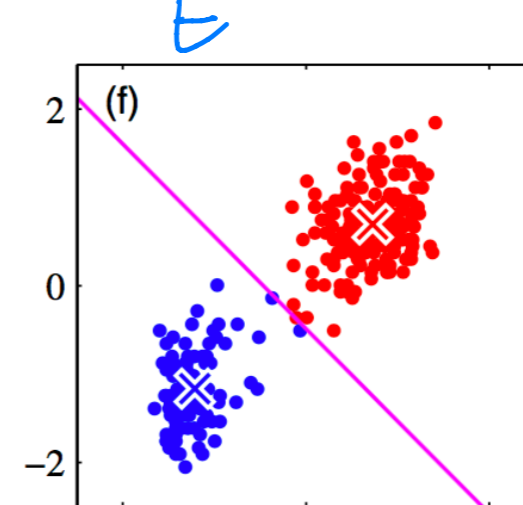
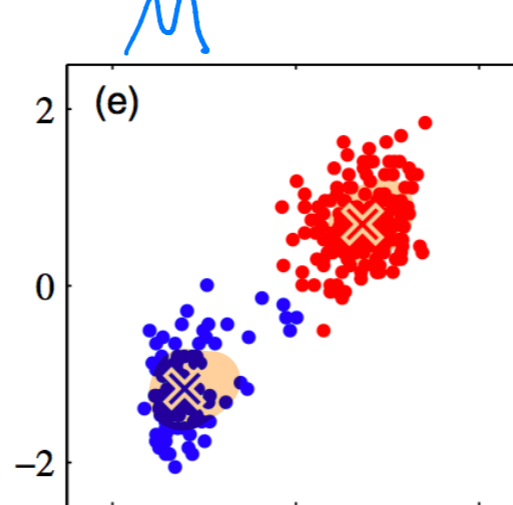
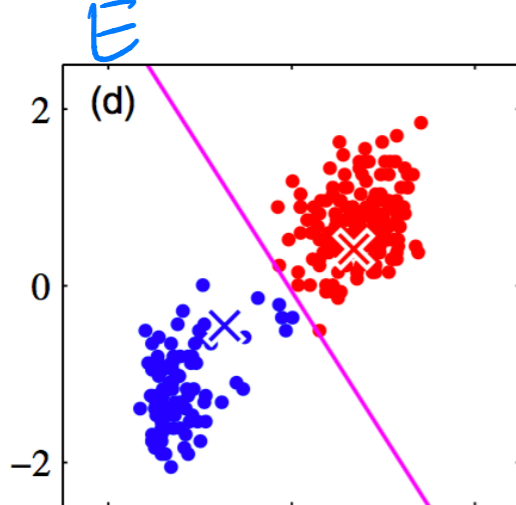
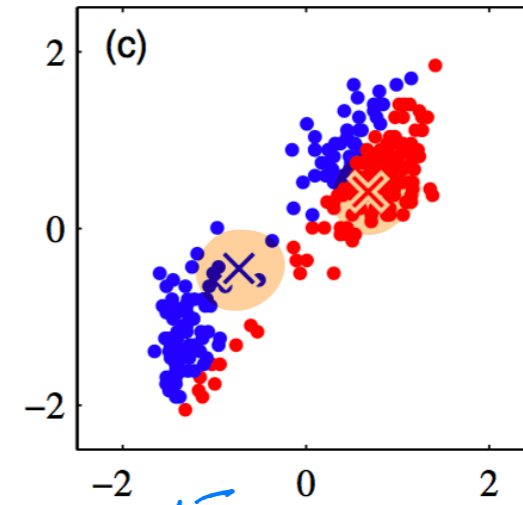
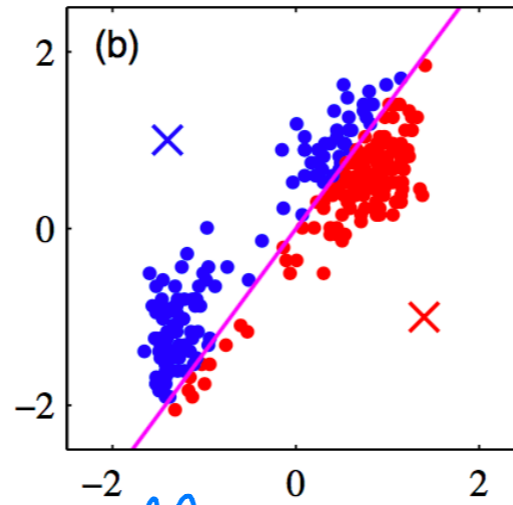
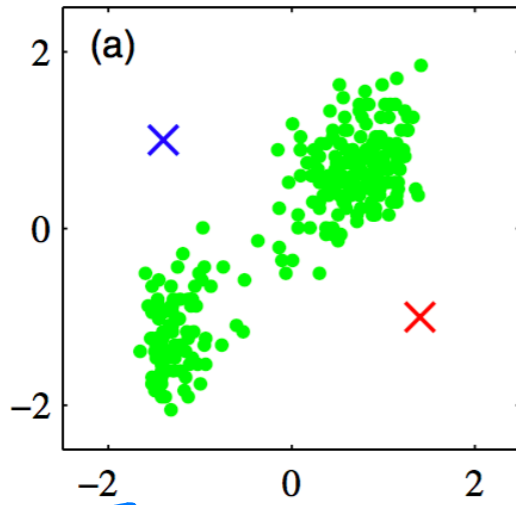
Example

(assignment)

E

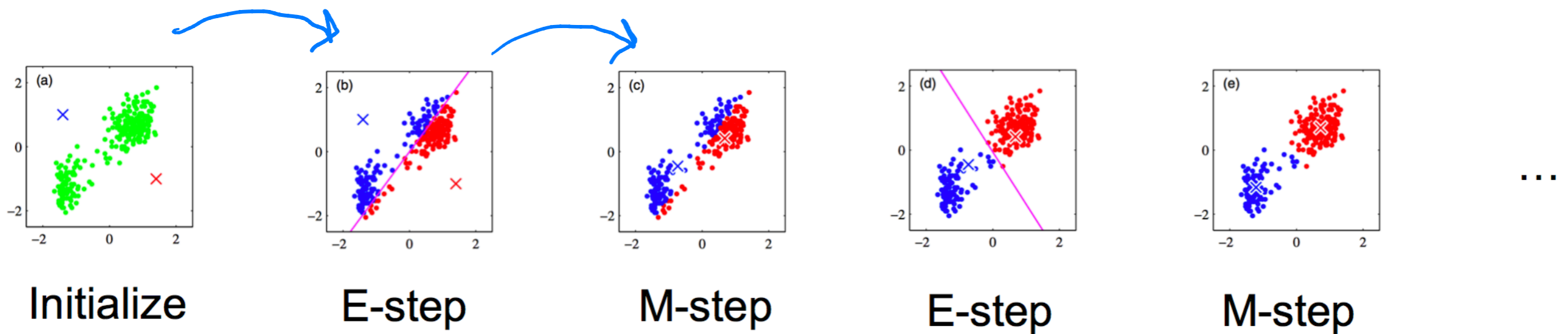
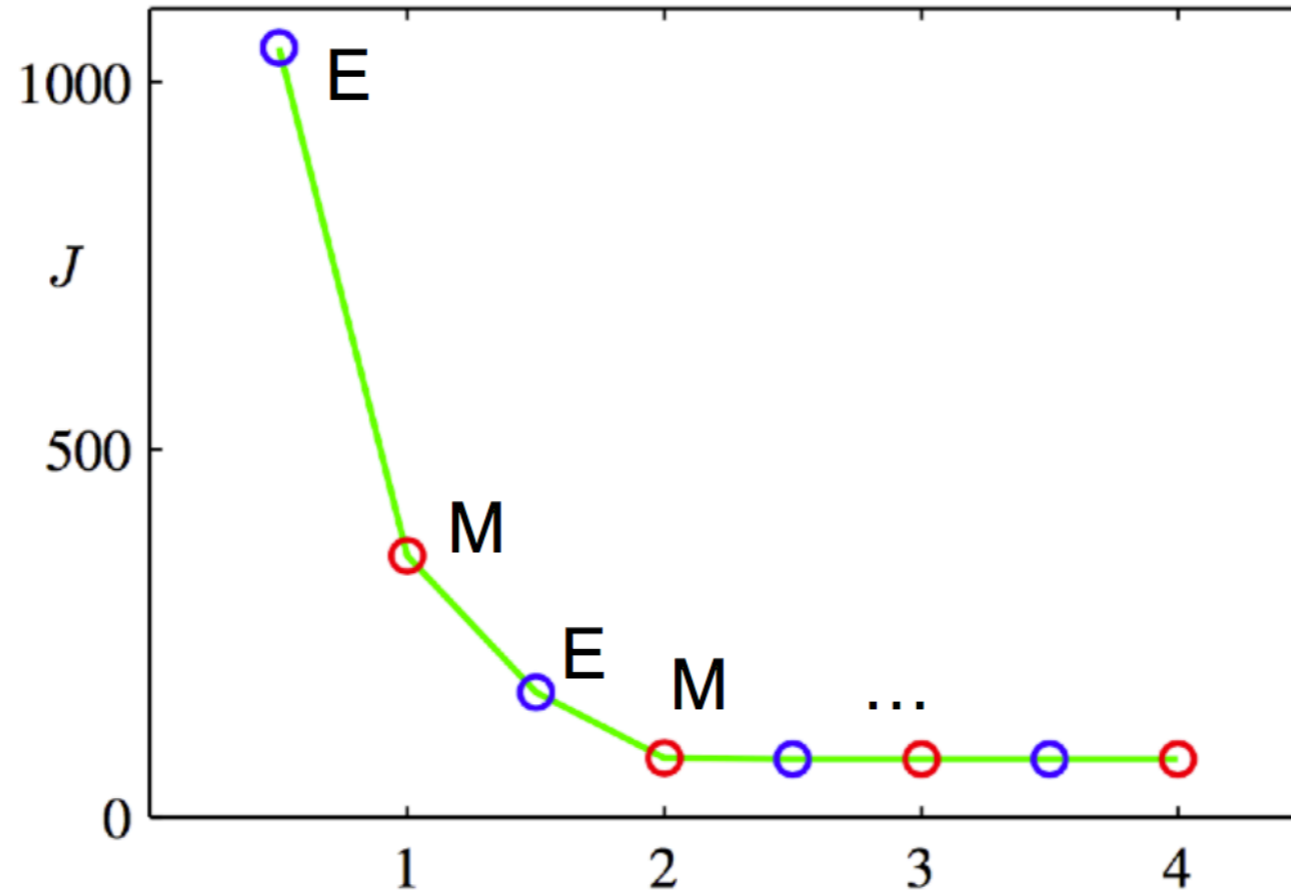
(update means)

M



Convergence

$$J = \sum_n \sum_k z_{nk} \| \underline{x}_n - \underline{\mu}_k \|^2$$



But *global* convergence?

- ◆ J is non-convex for μ_k and z_{nk} together and k-means converges to a **local minimum**
- ◆ Can we do better? Random restarts with different initial cluster means and then select the clusters with minimal J.

Derivation of the M-step

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

◆ It is a convex function in $\boldsymbol{\mu}_k$ (fixed z_{nk})

◆ Find the minimum by setting gradient = 0

(Handwritten note: $\frac{\partial}{\partial \boldsymbol{\mu}_k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 = 2(\mathbf{x}_n - \boldsymbol{\mu}_k)$)

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\sum_{n=1}^N \sum_{l=1}^K z_{nl} \|\mathbf{x}_n - \boldsymbol{\mu}_l\|^2 \right) = \sum_{n=1}^N z_{nk} \frac{\partial}{\partial \boldsymbol{\mu}_k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

≠ 0 only when l=k

$$= -2 \sum_{n=1}^N z_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T = 0$$

◆ Hence

$$\sum_{n=1}^N z_{nk} \mathbf{x}_n - \boldsymbol{\mu}_k \sum_{n=1}^N z_{nk} = 0 \implies \boldsymbol{\mu}_k = \frac{\sum_{n=1}^N z_{nk} \mathbf{x}_n}{\sum_{n=1}^N z_{nk}}$$

Application: image compression

Data points: $x_n \in (K, G, B)$ pixel values
 K clusters: color representation

$K = 2$

$K = 3$

$K = 10$

Original image



$\begin{matrix} \sim \sim \\ \sim \sim \\ \sim \sim \end{matrix} \left. \begin{matrix} 1 \\ 2 \\ 1 \end{matrix} \right\} \begin{matrix} 1 \\ 2 \\ 1 \end{matrix} \right)$
 X, Y integers $\in \{1, 2\}$

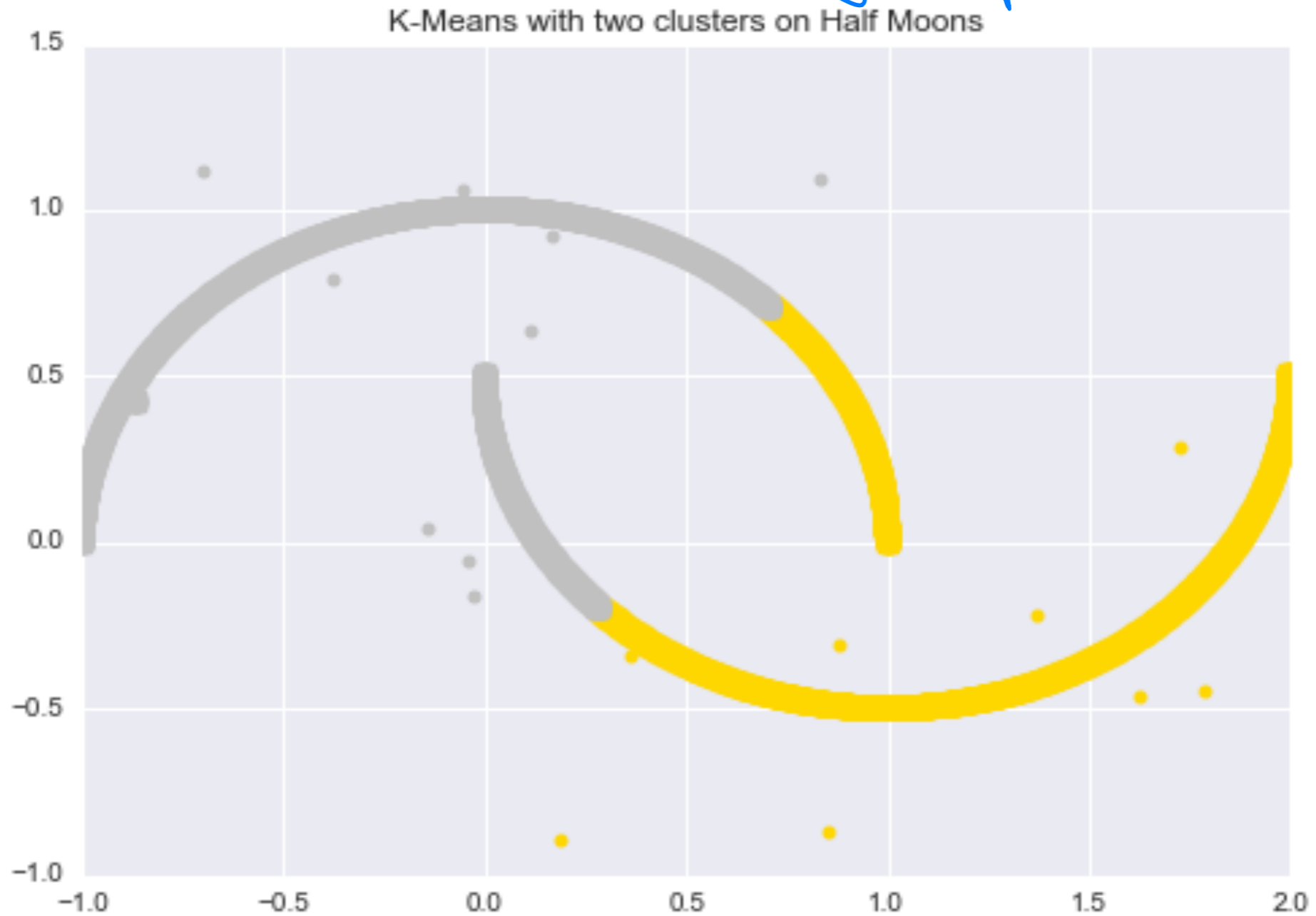
$\left(\begin{matrix} r \\ g \\ b \end{matrix} \right) \left(\begin{matrix} \hat{r} \\ \hat{g} \\ \hat{b} \end{matrix} \right)$
 X, Y, Z values $\in [0, 255]$

↑

→

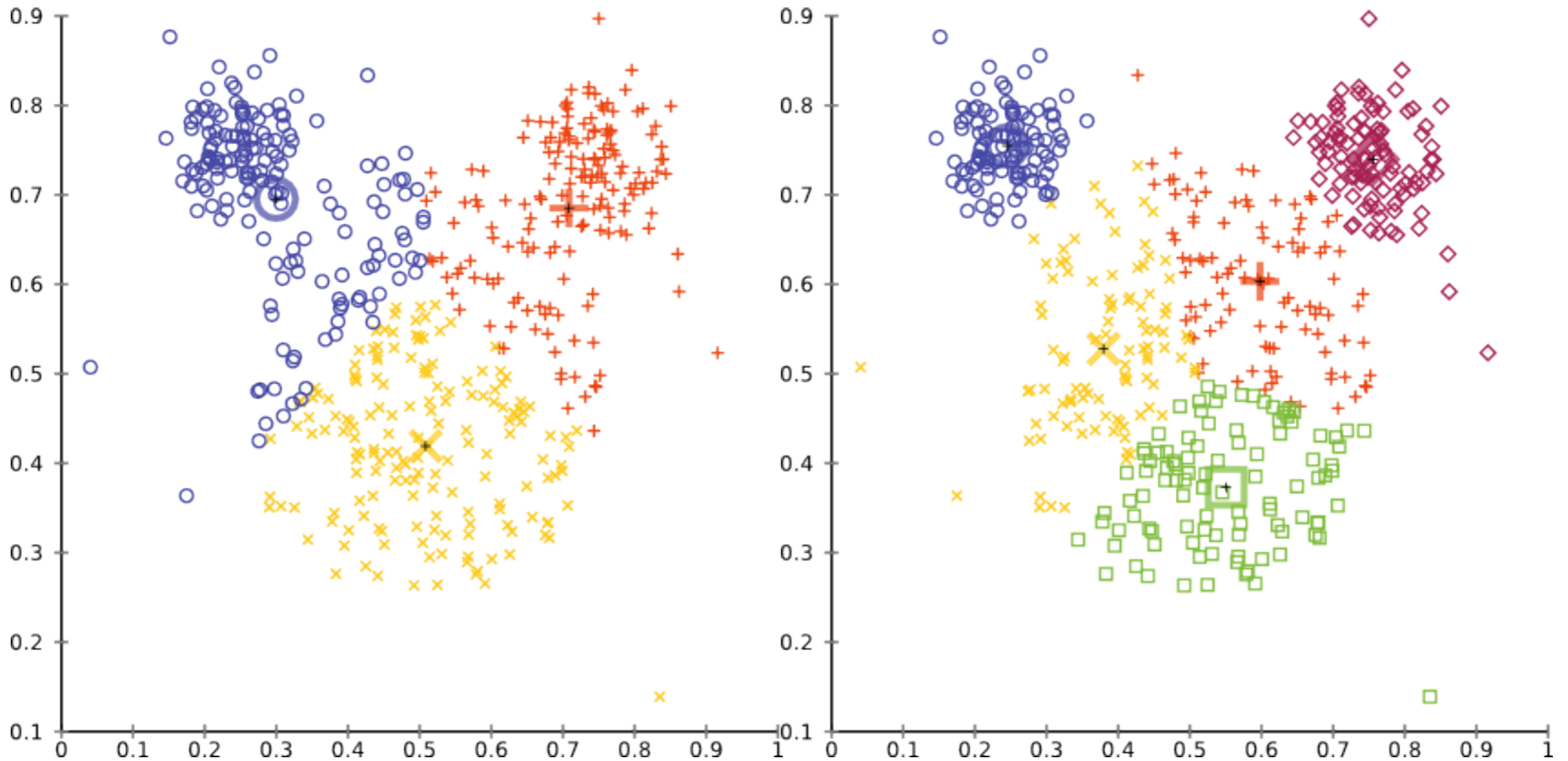
Failures of K-means

⇒ only spherical clusters



Failures of K-means (the mouse data!)

⇒ each cluster equal size



K = 3

K = 5

Improvements

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \| \mathbf{x}_n - \boldsymbol{\mu}_k \|^2$$

- ▶ Stochastic gradient for big data

- ▶ For each datapoint, update nearest cluster mean:

$$\begin{aligned} \boldsymbol{\mu}_k &= \boldsymbol{\mu}_k - \eta \left(\frac{\partial J}{\partial \boldsymbol{\mu}_k} \right)^T \\ &= \boldsymbol{\mu}_k + 2\eta(\mathbf{x}_n - \boldsymbol{\mu}_k) \end{aligned}$$

- ▶ Other distances between points (K-medoids)

- ▶ Euclidean not always appropriate (discrete data), sensitive to outliers

$$\tilde{J} = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$$

Pros & Cons

▶ Good

- Simple to implement
- Fast

▶ Bad

- Local minima
- Model only “spherical” clusters
- (*) • Sensitive to the features scale
- Number of clusters K to be chosen in advance
- Cluster assignments are “hard”, not probabilistic => next topic, Gaussian Mixture Model

via whitening operator

