

CHAPTER 7:

# CLUSTERING

**Stella Grasshof**

# Overview

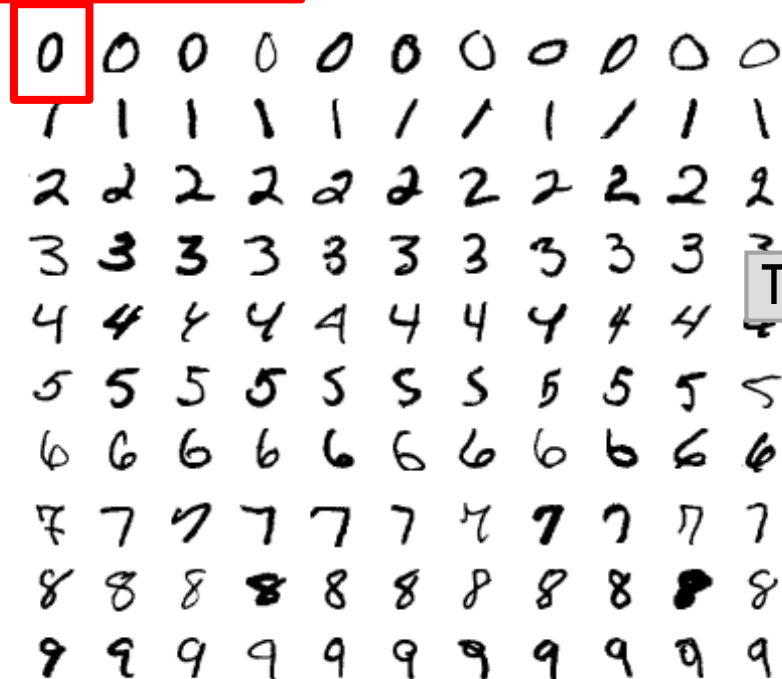
- 1) Intro
- 2) Hierarchical Clustering
- 3) K-means
- 4) Spectral Clustering
- 5) Expectation-Maximization Algorithm (EM)
- 6) Considerations
- 7) Mean-Shift
- 8) Outro

# Feature Extraction

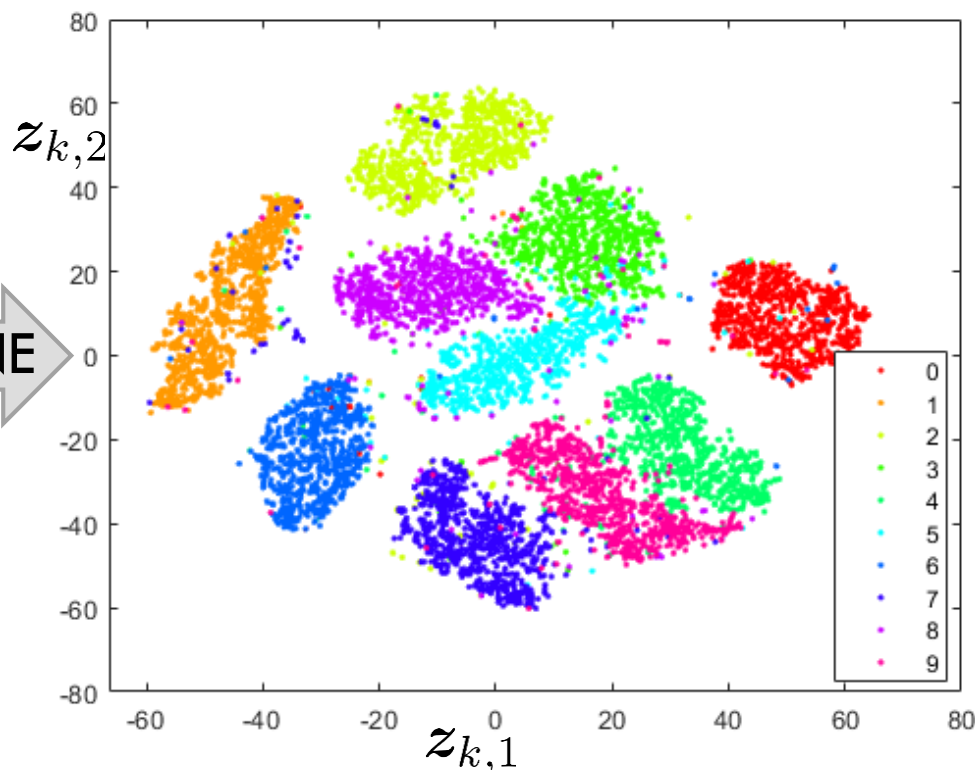
Consider high-dimensional data is given, we want:

- Compact representation of the data
- Extract **most relevant** information

$$x_k \in \mathbb{R}^{64 \cdot 64} \mapsto z_k \in \mathbb{R}^2$$



T-SNE



Sources:

[1] MNIST, wikipedia.org, Josef Steppan [CC BY-SA 4.0 (<https://creativecommons.org/licenses/by-sa/4.0/>)]

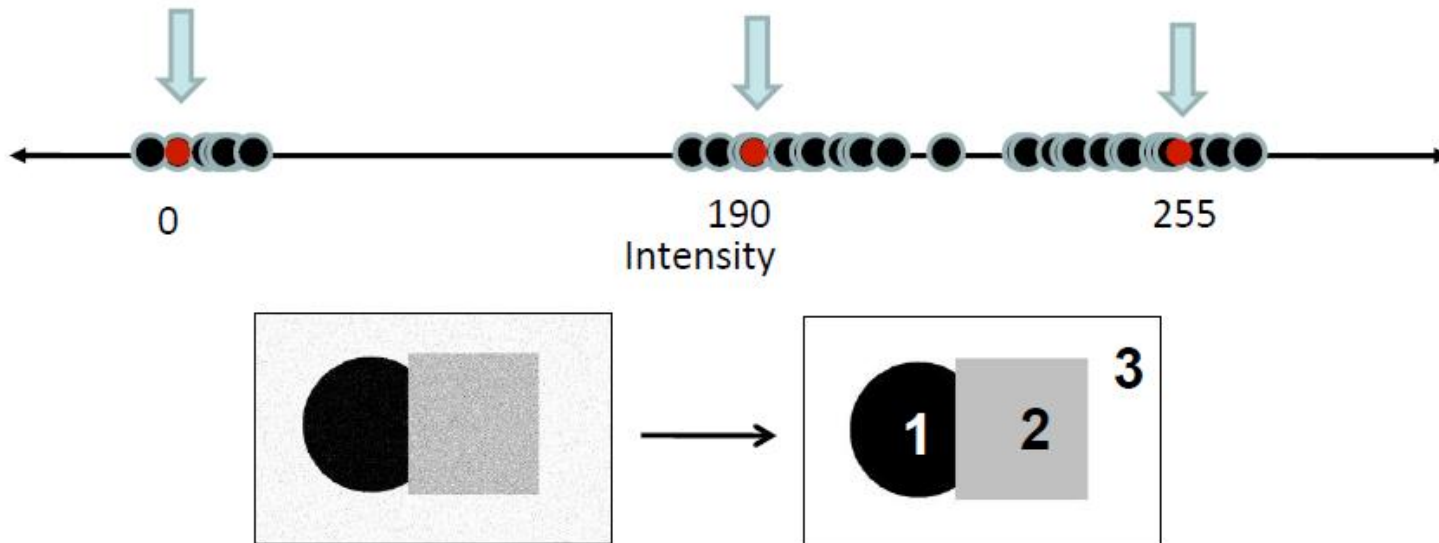
[2] MathWorks, <https://se.mathworks.com/help/stats/visualize-high-dimensional-data-using-t-sne.html>

# Why Clustering?

- Classification:  $x_n \mapsto C_k$
- Dimensionality Reduction:

$$x_n \in \mathbf{X} \in \mathbb{R}^{N \times D} \rightsquigarrow z_n \in \mathbf{Z} \in \mathbb{R}^{N \times E}, \quad E < D$$

- Clustering:
  - “estimate unknown class labels”
  - Analyze structure in data



# Semiparametric Density Estimation

Data probability density function (pdf):  $p(\boldsymbol{x})$

- **Parametric:** one model, (Ch. 4 and 5) e.g.

$$p(\boldsymbol{x}) \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

- **Semiparametric:** a mixture of densities

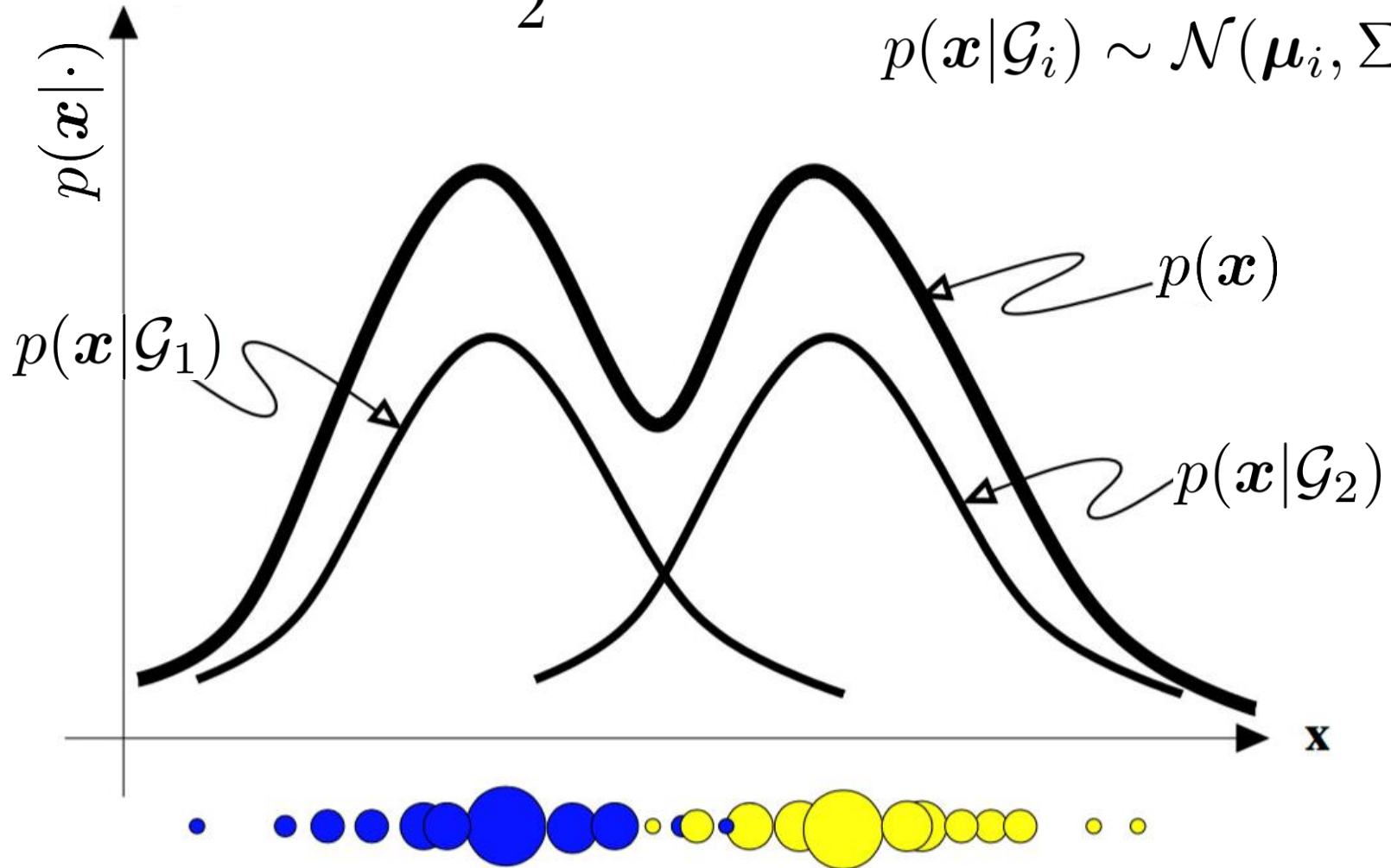
$$p(\boldsymbol{x}) = \sum_i P(\mathcal{G}_i) p(\boldsymbol{x}|\mathcal{G}_i)$$

- **Nonparametric:** No model (Ch. 8)

# Mixture of Gaussians

$$p(\mathbf{x}) = \frac{1}{2} (p(\mathbf{x}|\mathcal{G}_1) + p(\mathbf{x}|\mathcal{G}_2))$$

$$p(\mathbf{x}|\mathcal{G}_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$



# Classes vs. Clusters

**Supervised:** data with class labels  $\mathbf{x}_n \in C_k$

$$\mathbb{1} \{ \mathbf{x}_n \in C_k \} = \begin{cases} 1 & , \mathbf{x}_n \in C_k \\ 0 & , \text{else} \end{cases}, \quad N_k := \sum_{n=1}^N \mathbb{1} \{ \mathbf{x}_n \in C_k \},$$

$$p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}|C_k) P(C_k)$$

$$p(\mathbf{x}|C_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\hat{P}(C_k) = \frac{N_k}{N}$$

$$\boldsymbol{\theta} = \{P(C_k), \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$$

$$\hat{\boldsymbol{\mu}}_k = \mathbf{m}_k = \frac{1}{N_i} \sum_{n=1}^N \mathbb{1} \{ \mathbf{x}_n \in C_k \} \mathbf{x}_n$$

$$\hat{\boldsymbol{\Sigma}}_k = \mathbf{S}_k = \frac{1}{N_i} \sum_{n=1}^N \mathbb{1} \{ \mathbf{x}_n \in C_k \} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T$$

# Classes vs. Clusters

**Unsupervised:** data **without** class labels  $\mathbf{x}_n$

but assume there are groupings in the data  $\mathcal{G}_k \Rightarrow$  Labels?

$$\mathbb{1} \{ \mathbf{x}_n \in \mathcal{G}_k \} = \begin{cases} 1 & , \mathbf{x}_n \in \mathcal{G}_k \\ 0 & , \text{else} \end{cases} ,$$

$$p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x} | \mathcal{G}_k) P(\mathcal{G}_k)$$

$$p(\mathbf{x} | \mathcal{G}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\boldsymbol{\theta} = \{P(\mathcal{G}_k), \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$$

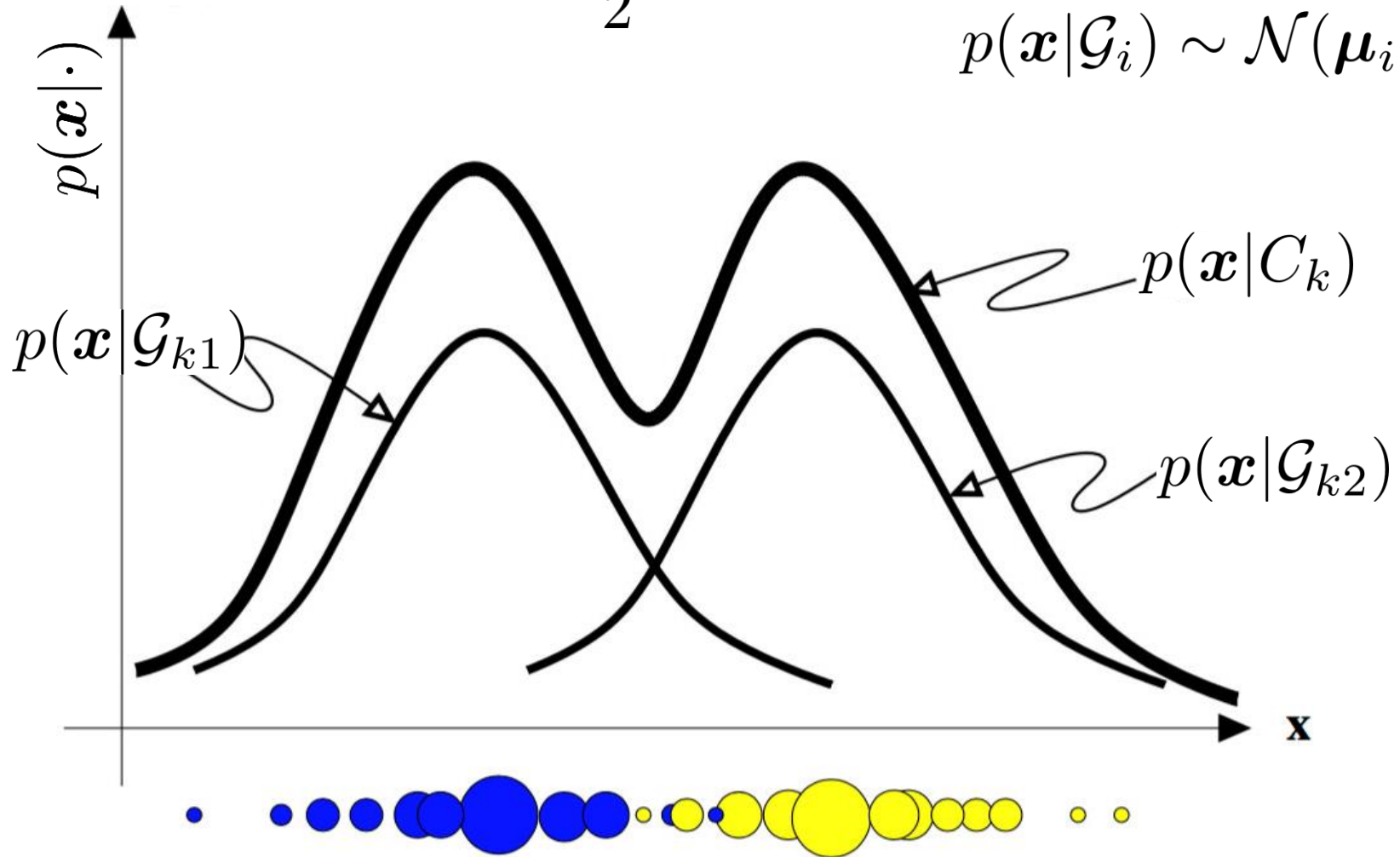
Classes: prior known category

Cluster: e.g. sub-groups per class



# Mixture of Mixtures

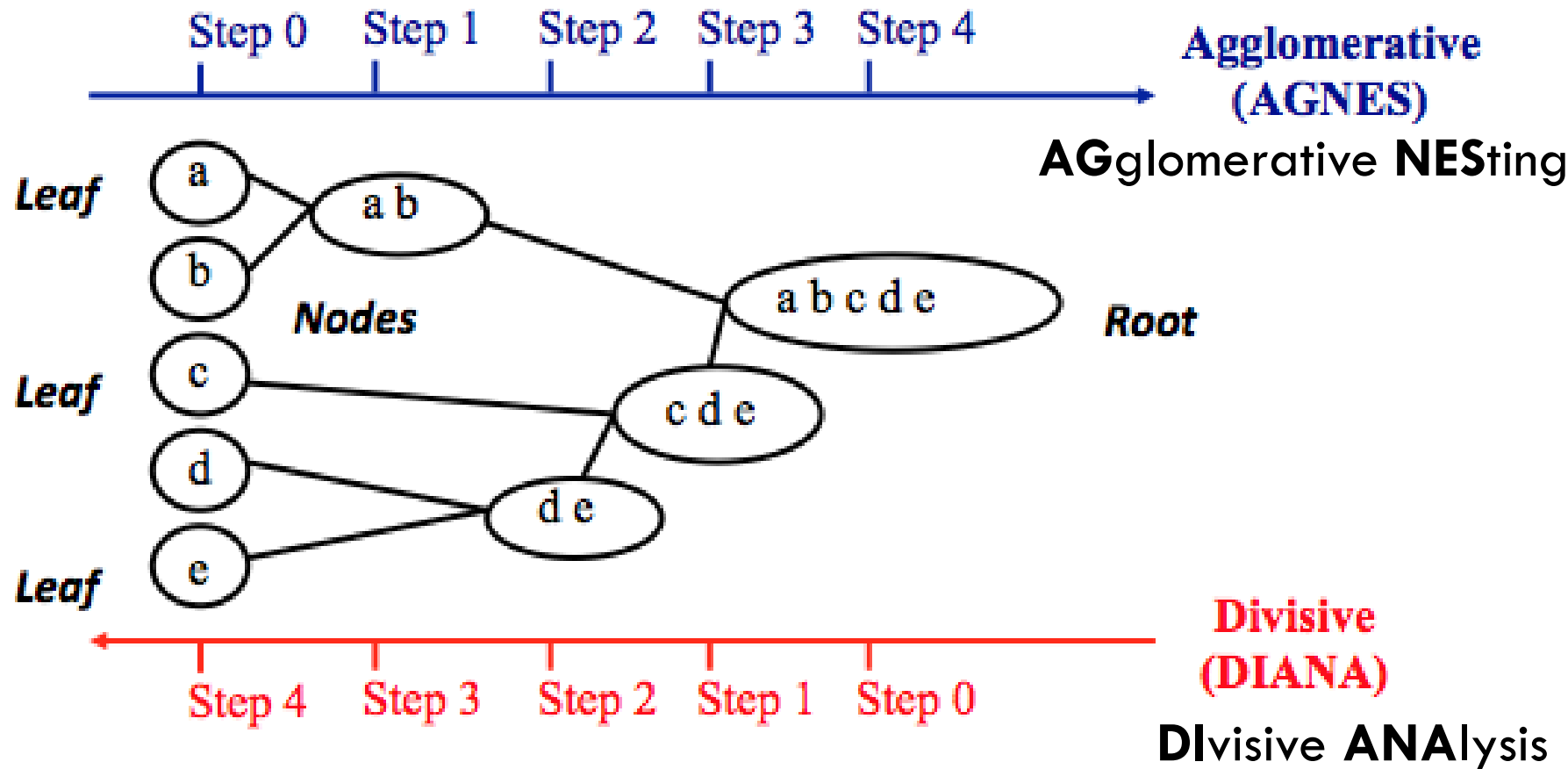
$$p(\mathbf{x}|C_k) = \frac{1}{2} (p(\mathbf{x}|\mathcal{G}_{k1}) + p(\mathbf{x}|\mathcal{G}_{k2}))$$
$$p(\mathbf{x}|\mathcal{G}_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \Sigma_i)$$



# Hierarchical Clustering

- Remember: Feature Selection
- **Agglomerative Clustering:**
  - Start with  $N$  cluster: each point is one cluster
  - Merge several cluster
- **Divisive Clustering:**
  - Start with 1 cluster: All points are 1 cluster
  - Divide cluster

# Hierarchical Clustering



# Hierarchical Clustering

How to decide merging and splitting?

- Cluster based on similarities/distances
- Distance between single points  $x_i, x_j$
- Distance between cluster/groups  $\mathcal{G}_k, \mathcal{G}_l$  containing several points

# Hierarchical Clustering

- Distance between points  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^D$   
 $\mathbf{x}_i = (x_{i1}, \dots, x_{iD})^T$

Minkowski ( $L_p$ )

$$d(\mathbf{x}_i, \mathbf{x}_j) = \left[ \sum_{d=1}^D |x_{id} - x_{jd}|^p \right]^{1/p}$$

- $p=2$ : Euclidean norm
- $p=1$ : City-block distance (Manhattan distance)
- ...

# Hierarchical Clustering

## Distance between two groups $\mathcal{G}_k, \mathcal{G}_l$

### □ Single-link:

$$d(\mathcal{G}_k, \mathcal{G}_l) = \min_{\mathbf{x}_i \in \mathcal{G}_k, \mathbf{x}_j \in \mathcal{G}_l} d(\mathbf{x}_i, \mathbf{x}_j)$$

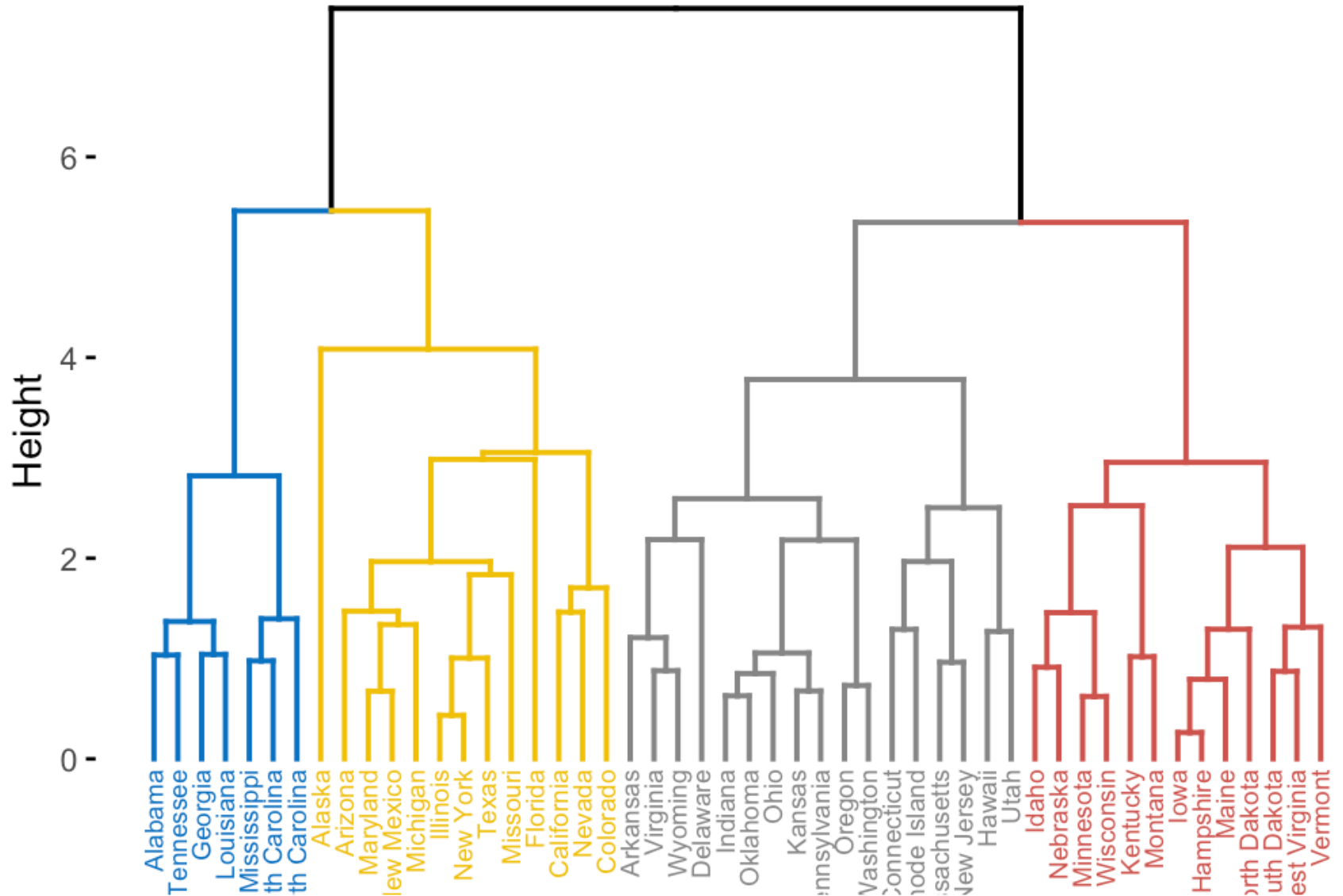
### □ Complete-link:

$$d(\mathcal{G}_k, \mathcal{G}_l) = \max_{\mathbf{x}_i \in \mathcal{G}_k, \mathbf{x}_j \in \mathcal{G}_l} d(\mathbf{x}_i, \mathbf{x}_j)$$

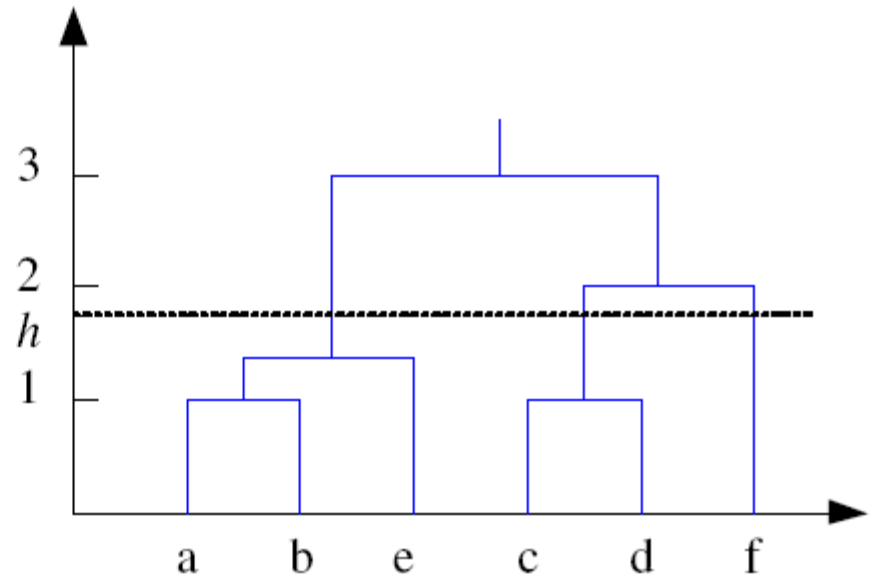
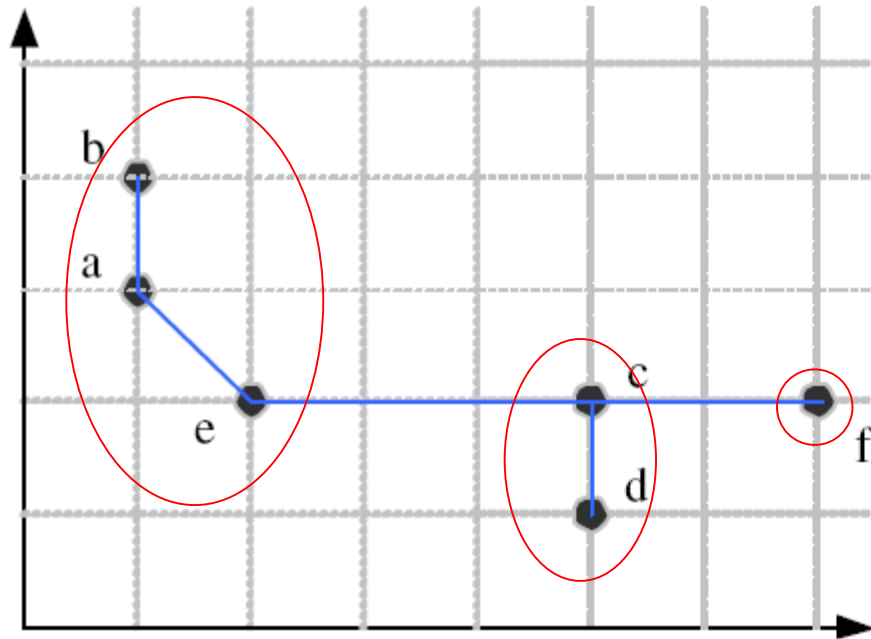
### □ Average-link, centroid

$$d(\mathcal{G}_k, \mathcal{G}_l) = \text{ave}_{\mathbf{x}_i \in \mathcal{G}_k, \mathbf{x}_j \in \mathcal{G}_l} d(\mathbf{x}_i, \mathbf{x}_j)$$

# Cluster Dendrogram



# Example: Single-Link Clustering



*Dendrogram*

How to decide number of clusters?



# Hierarchical Clustering

- Pro:

- Flexibility: free choice of distance measure
- Easy to implement and widespread
- Provides hierarchy of points and clusters

- Contra:

- Runtime
- No model
- Imbalanced: outlier will end up in separate small cluster

# K-Means Clustering

## Input:

- $N$  data points  $\mathbf{x}_n \in \mathbb{R}^D, n = 1, \dots, N$
- parameter  $K$  : number of classes

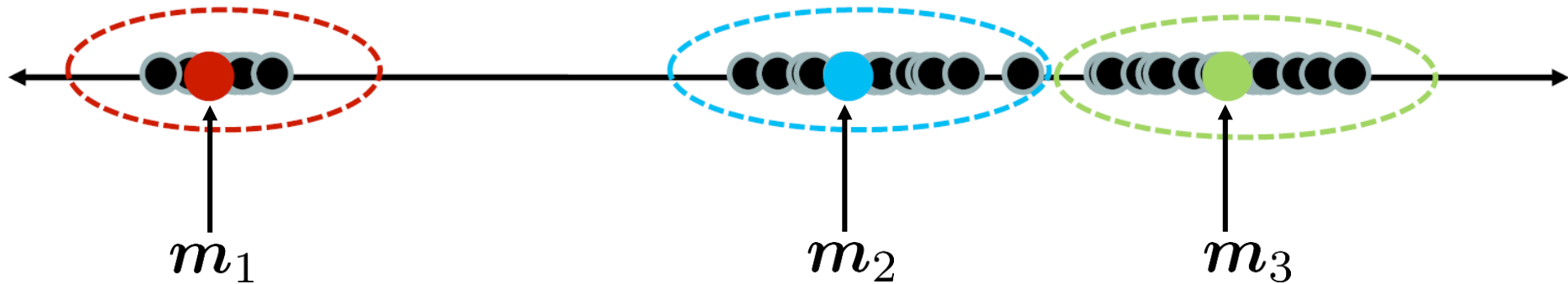
## Output:

- labels for each point  $\mathbf{x}_n \mapsto m_k$
- Codebook of reference vectors  $m_k, k = 1, \dots, K$   
Template Matching, Vector Quantization (VQ), ...

=>  $N$  points are matched on  $K$  references  
Compression from  $N$  to  $K$

# K-Means Clustering: $K=3$

- Assuming known  $m_k$ ,  $k = 1, \dots, K$   
cluster centers / templates / references
- Assign closest center to each data point:  $x_n \mapsto m_k$



- Problem:
  - No class labels given
  - $m_k$  are unknown

# K-Means Clustering

1) **Initialize** **K** cluster centers  $\mathbf{m}_k$ ,  $k = 1, \dots, K$

2) Assign each point to closest cluster center

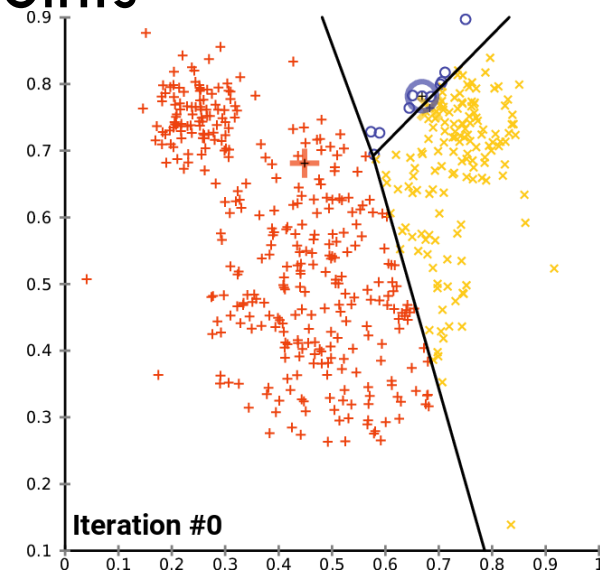
$$1(\mathbf{x}_n \mapsto \mathbf{m}_k) = \begin{cases} 1 & , \|\mathbf{x}_n - \mathbf{m}_k\|_2 = \min_i \|\mathbf{x}_n - \mathbf{m}_i\|_2 \\ 0 & , \text{else} \end{cases}$$

3) Update cluster center as mean of points

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n=1}^N \mathbf{x}_n 1(\mathbf{x}_n \mapsto \mathbf{m}_k)$$

$$N_k = \sum_{n=1}^N 1(\mathbf{x}_n \mapsto \mathbf{m}_k)$$

4) Repeat 2)-3) until **convergence**



# K-Means Clustering

## □ Initialization

- Randomly selected K data points
- Divide data space in intervals: choose mean of these
- ...

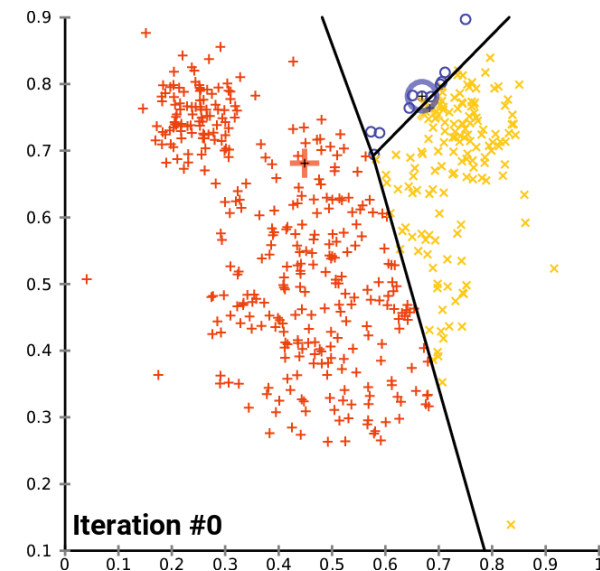
## □ Convergence

- If cluster centers do not change anymore
- ...

## □ How to choose K?

- Depends on application
- Re-run and inspect the results
- look at data !

[https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)



# K-Means: Image Segmentation

$K = 2$



$K = 3$



$K = 10$



Original image



# K-Means Clustering

- Different variants, e.g.
  - assign cluster, such that variance is minimized
  - Fuzzy K-means: assign probability instead of one cluster
- **Pro**
  - easy to implement
- **Contra**
  - Need to choose K
  - Sensitive to initial cluster centers
  - Local minima
  - Spherical structure of clusters (euclidean norm)

# Spectral Clustering

- First process data, then clustering
- Given pairwise similarities (or “ minus“ distance )

$$w_{ij} = w_{ji} = \text{similarity}(\mathbf{x}_i, \mathbf{x}_j)$$

- Steps:
  - I. Use these values for Laplacian Eigenmaps (Ch. 6) to map  $\mathbf{x}_n$  to a new space  $\mathbf{z}_n$
  - II. Use  $k$ -means on  $\mathbf{z}_n$  for clustering



# Expectation-Maximization (EM)

Probability density function (pdf) vs. Likelihood

$$x \sim \mathcal{N}(\mu, \Sigma) \quad p(x) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

- Pdf: evaluate on data **given parameters**

$$p(x|\theta), \quad \theta = (\mu, \Sigma)$$

- Likelihood: evaluate on parameters **given data**

$$l(\theta) = p(\theta|x), \quad \theta = (\mu, \Sigma)$$

- Given data:

- maybe missing values, no labels

- Unknown parameters of pdf

=> Complete likelihood unknown

# Expectation-Maximization (EM)

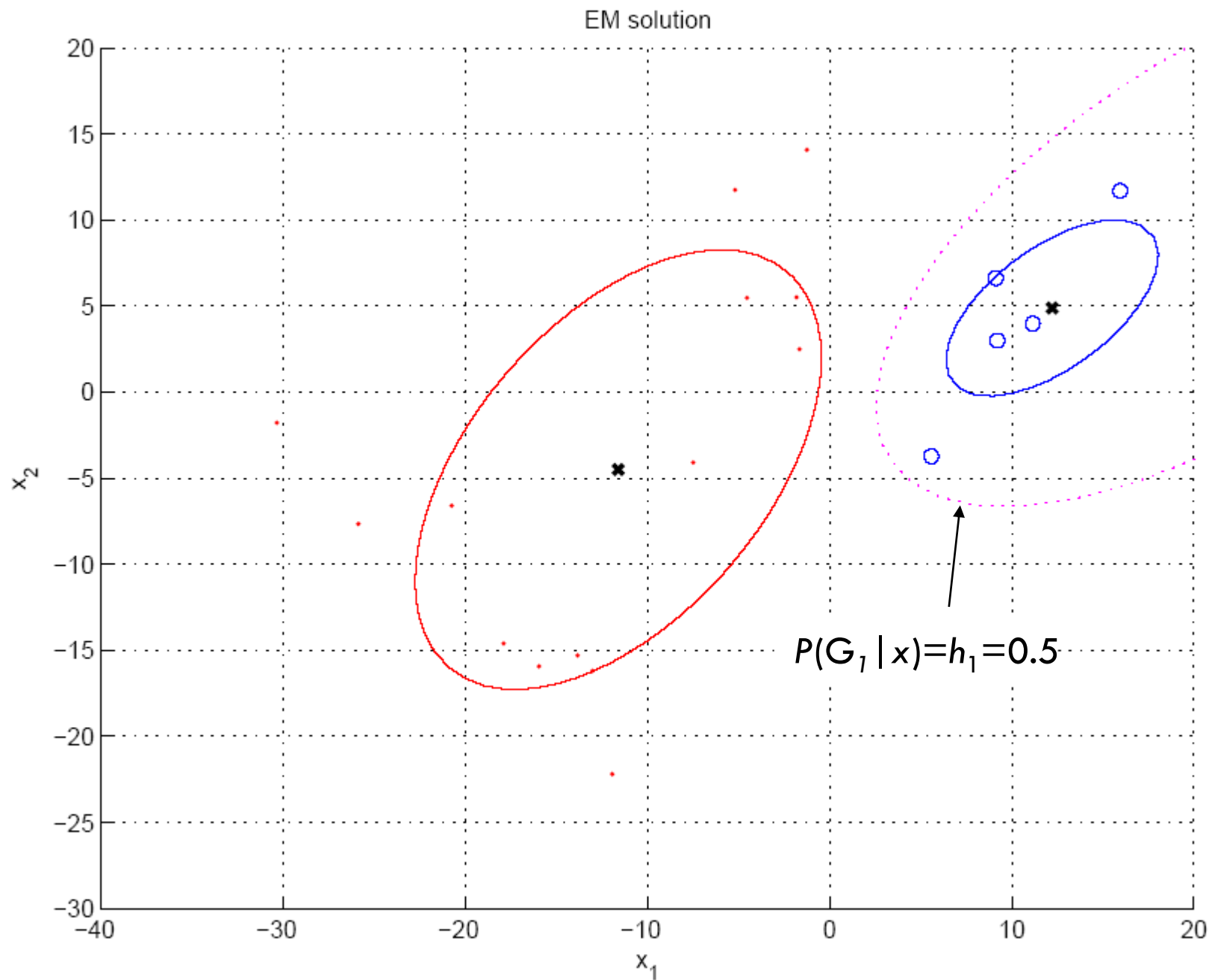
- Log likelihood with iid samples  $\mathbf{x}_n \in \mathbf{X}$

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}|\mathbf{X}) &= \log \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\theta}) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n|\mathcal{G}_k)P(\mathcal{G}_k)\end{aligned}$$

- Assume hidden variables  $\mathcal{Z}$   
when known, make optimization much simpler
- Complete likelihood  $\mathcal{L}_c(\boldsymbol{\theta}|\mathbf{X}, \mathcal{Z})$
- Incomplete likelihood  $\mathcal{L}(\boldsymbol{\theta}|\mathbf{X})$

# Expectation-Maximization (EM)

- **Initialize** parameters, set:  $\theta^0$
- **Expectation:** given parameters:  $\theta^l$ 
  - ▣ Estimate hidden variables  $\mathcal{Z}$
  - ▣ Get expected value of likelihood
$$Q(\theta|\theta^l) = E[\mathcal{L}_c(\theta|X, \mathcal{Z})|\theta^l]$$
- **Maximization:**
  - ▣ Estimate parameters which maximize likelihood
$$\theta^{l+1} = \operatorname{argmax}_{\theta} Q(\theta|\theta^l)$$
- Repeat E-step and M-step



# Mixtures of Latent Variable Models

## Regularize clusters

1. Assume shared/diagonal covariance matrices
2. Use PCA/FA to decrease dimensionality:

Mixtures of PCA/FA

$$\text{FA: } \boldsymbol{x} = \boldsymbol{V}\boldsymbol{z} + \boldsymbol{\epsilon} \quad \text{Cov}(\boldsymbol{z}) = \boldsymbol{I}, \text{ Cov}(\boldsymbol{\epsilon}) = \boldsymbol{\Psi}$$

$$p(\boldsymbol{x}|\mathcal{G}_k) = \mathcal{N}(\boldsymbol{m}_k, \boldsymbol{V}_k \boldsymbol{V}_k^{\text{T}} + \boldsymbol{\Psi}_k)$$

# Choosing $K$

- Defined by the application, e.g., image quantization
- Plot data (after e.g. PCA) and check for clusters
- Incremental (leader-cluster) algorithm:  
Add one at a time until large change
- Manually check for meaning

# Mean Shift: no need to choose K

**Mean Shift** instead of K-means

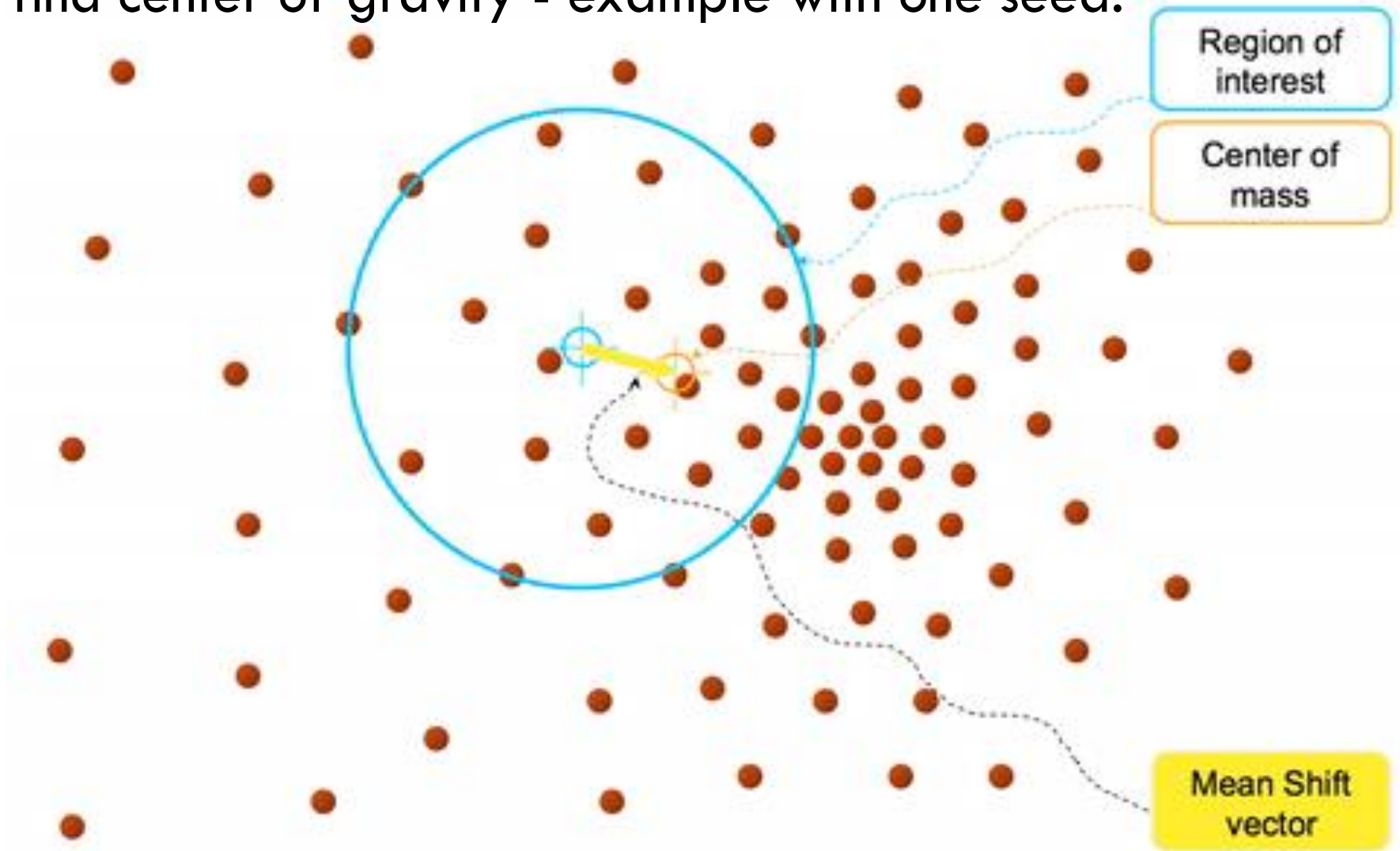
**Parameter:** window function and width

=> locates maxima

- 1) Initialize random seed  $m_{(0)}$  and window  $W$
- 2) Calculate new center of gravity (mean)  
 $m_{(l)}$  as (weighted) mean of all points in window  
 $x_n \in W(m_{(l-1)})$  centered at  $m_{(l-1)}$
- 3) Shift search window to new mean  $m_{(l)}$
- 4) Repeat 2-3

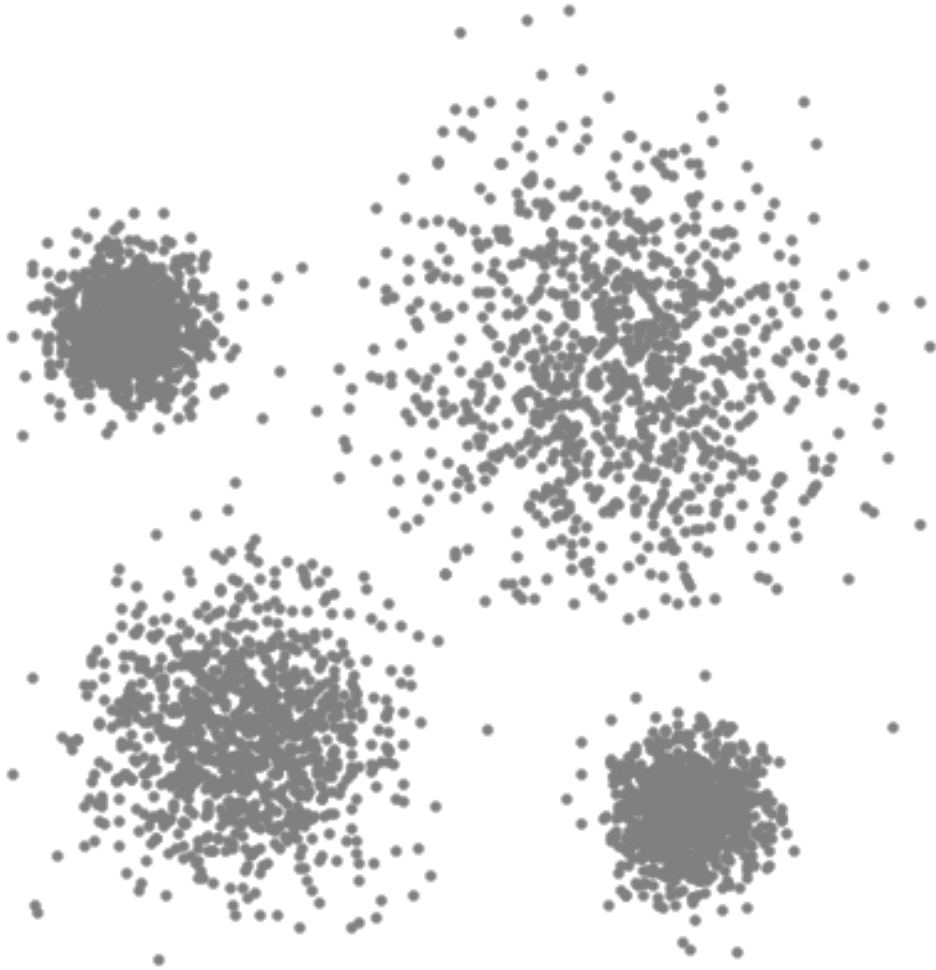
# Mean Shift

find center of gravity - example with one seed:





# Mean Shift



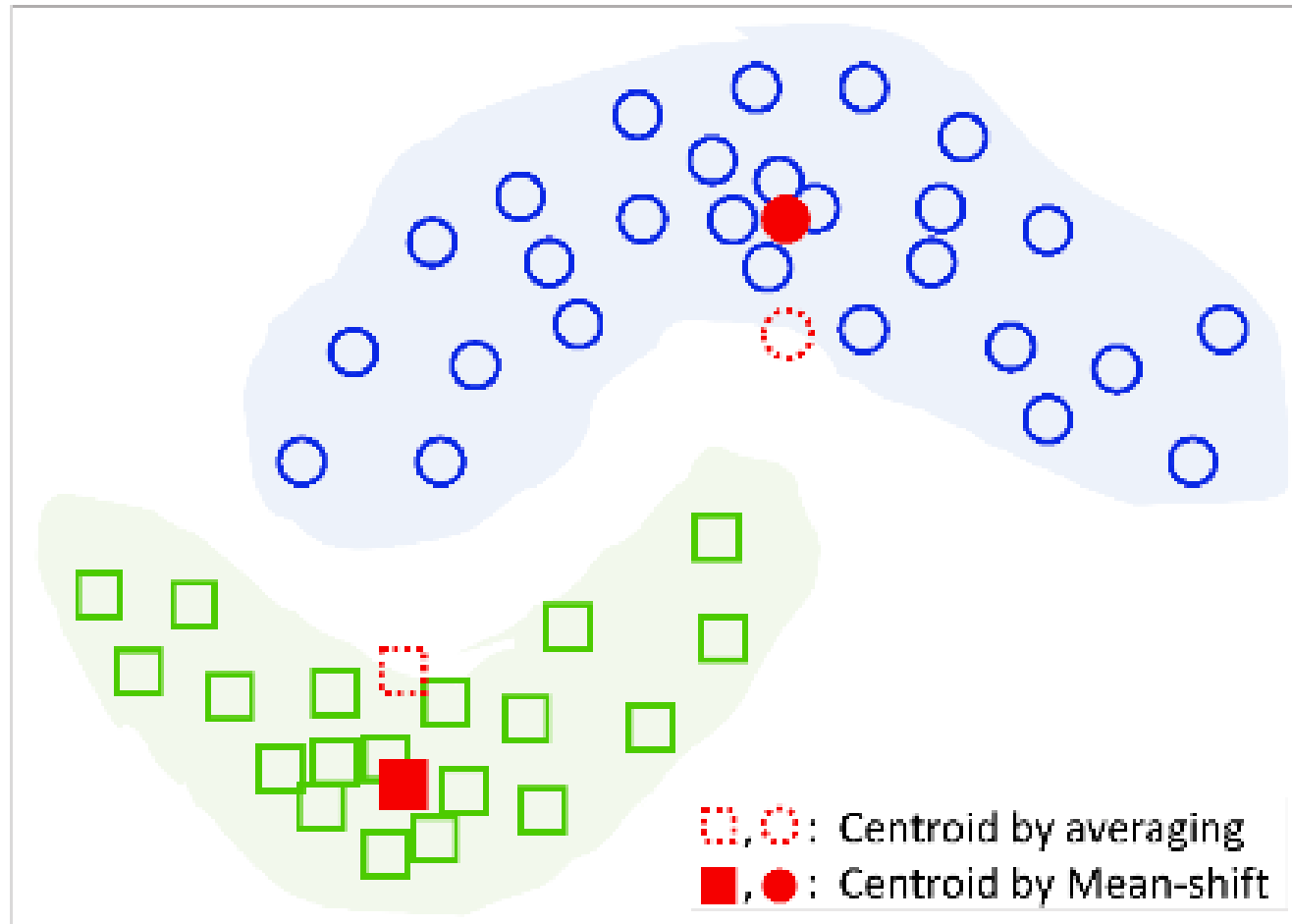
Initialize:

- random starting points
- Each data point

Mean Shift:

- for start each point:  
where does it end up?
- Different points share  
the same  
center of gravity =  
cluster center

# Mean Shift



# Mean Shift:

## **Pro:**

- ☐ No need to choose number of clusters  $K$   
=> finds variable number of clusters
- ☐ Only one parameter window size:  
has physical meaning
- ☐ Model-free
- ☐ Robust to outliers

## **Contra:**

- ☐ Result depends on window size (not trivial)
- ☐ Computationally expensive
- ☐ Attention in high dims!

# After Clustering

- **Dimensionality reduction**  
remove redundant information, e.g.
  - correlations between features and group features
- **Clustering methods** find similarities between instances and group instances  
=> Analysis
- knowledge extraction through
  - number of clusters
  - prior probabilities
  - cluster parameters, i.e., center, range of features.

Example: CRM, customer segmentation

# Clustering as Preprocessing

$$b_k^n := \mathbb{1} \{ \mathbf{x}_n \in \mathcal{G}_k \} = \begin{cases} 1 & , \mathbf{x}_n \in \mathcal{G}_k \\ 0 & , \text{else} \end{cases} \text{ vs. } h_i^k := P(\mathcal{G}_k | \mathbf{x}_i)$$

- Estimated group labels: **hard** vs. **soft**  
dimensions of a new K-dim. space  
=> learn discriminant or regressor

- **Local** representation

- only one  $b_i$  is 1, all others are 0
- only few  $h_i$  are nonzero

**vs**

**Distributed** representation, e.g.:  
after PCA all  $z_i$  are nonzero

# Summary

- Dimensionality reduction as preprocessing
- Clustering to:
  - ▣ analyze data (subgroups, etc.)
  - ▣ Estimate classes
- Classification
  - ▣ Given class labels

# APPENDIX