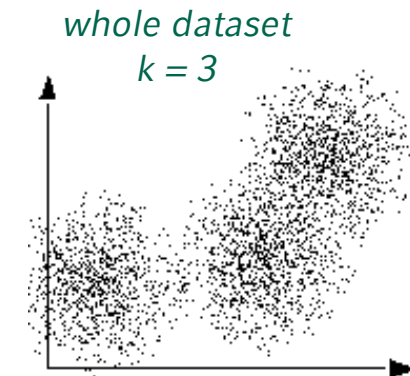


- 1) Introduction to clustering
- 2) Partitioning Methods
  - K-Means
  - Variants: K-Medoid, K-Mode, K-Median
  - Choice of parameters: Initialization, Silhouette coefficient
- 3) Probabilistic Model-Based Clusters: Expectation Maximization
- 4) Density-based Methods: DBSCAN
- 5) Hierarchical Methods
  - Agglomerative and Divisive Hierarchical Clustering
  - Density-based hierarchical clustering: OPTICS
- 6) Evaluation of Clustering Results
- 7) Further Clustering Topics
  - Scaling Up Clustering Algorithms

Just two examples:


[naïve]

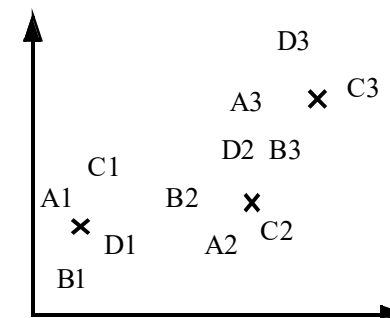
- Choose sample  $A$  of the dataset
- Cluster the sample and use centers as initialization



[Fayyad, Reina, and Bradley 1998]

- Choose  $m$  different (small) samples  $A, \dots, M$  of the dataset
- Cluster each sample to get  $m$  estimates for  $k$  representatives  
 $A = (A_1, A_2, \dots, A_k), B = (B_1, \dots, B_k), \dots, M = (M_1, \dots, M_k)$
- Then, cluster the set  $DS = A \cup B \cup \dots \cup M$   $m$  times. Each time use the centers of  $A, B, \dots, M$  as respective initial partitioning
- Use the centers of the best clustering as initialization for the partitioning clustering of the whole dataset

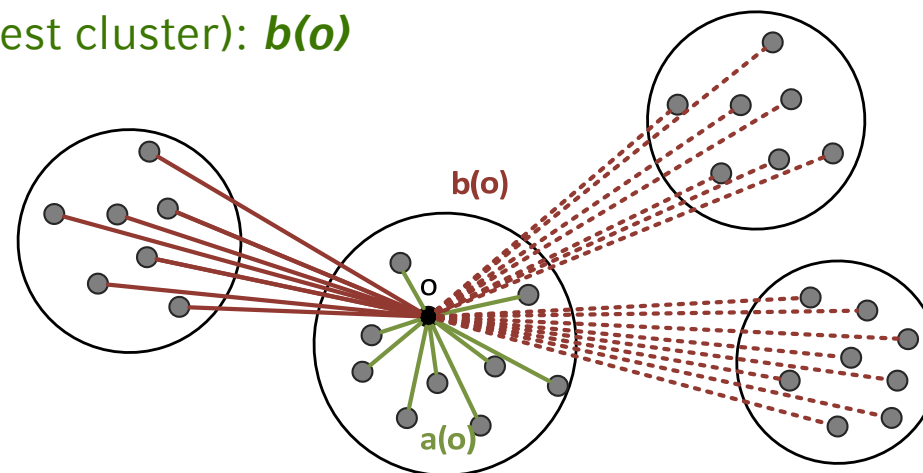
$m = 4$  samples  $A, B, C, D$   
 true cluster centers



Fayyad U., Reina C., Bradley P. S., „Initialization of Iterative Refinement Clustering Algorithms“, *In KDD 1998*), pp. 194—198.

- Idea for a method:
  - Determine a clustering for each  $k = 2, \dots, K_{max} \leq n-1$
  - Choose the “best” clustering
- But how to measure the quality of a clustering?
  - A measure should not be monotonic over  $k$  because the measures for the compactness of a clustering SSE and TD are monotonously decreasing with increasing value of  $k$ .
- Silhouette-Coefficient [Kaufman & Rousseeuw 1990]
  - Measure for the quality of a  $k$ -means or a  $k$ -medoid clustering that is not monotonic over  $k$ .

- Basic idea:
  - How good is the clustering = how appropriate is the mapping of objects to clusters
  - Elements in cluster should be „similar“ to their representative  
→ measure the average distance of objects to their representative:  $a(o)$
  - Elements in different clusters should be „dissimilar“  
→ measure the average distance of objects to alternative clusters (i.e. second closest cluster):  $b(o)$



- $a(o)$ : average distance between object  $o$  and the objects in its cluster  $A$

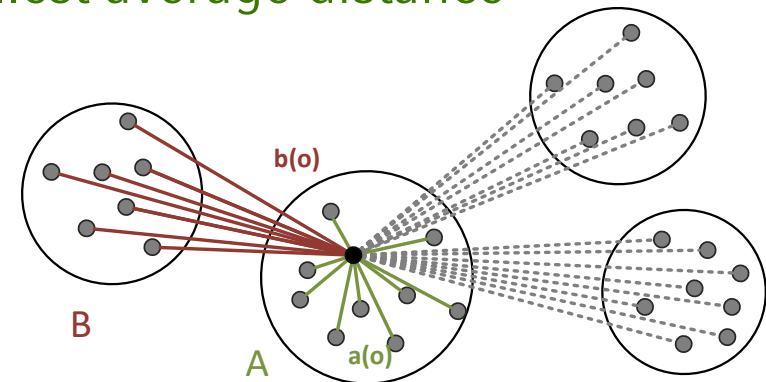
$$a(o) = \frac{1}{|C(o)|} \sum_{p \in C(o)} \text{dist}(o, p)$$

- $b(o)$ : for each other cluster  $C_i$  compute the average distance between  $o$  and the objects in  $C_i$ . Then take the smallest average distance

$$b(o) = \min_{C_i \neq C(o)} \left( \frac{1}{|C_i|} \sum_{p \in C_i} \text{dist}(o, p) \right)$$

- The silhouette of  $o$  is then defined as

$$s(o) = \begin{cases} 0 & \text{if } a(o) = 0, \text{ e.g. } |C_i| = 1 \\ \frac{b(o) - a(o)}{\max\{a(o), b(o)\}} & \text{else} \end{cases}$$



- The values of the silhouette coefficient range from  $-1$  to  $+1$

- The silhouette of a cluster  $C_i$  is defined as:

$$\text{silh}(C_i) = \frac{1}{|C_i|} \sum_{o \in C_i} s(o)$$

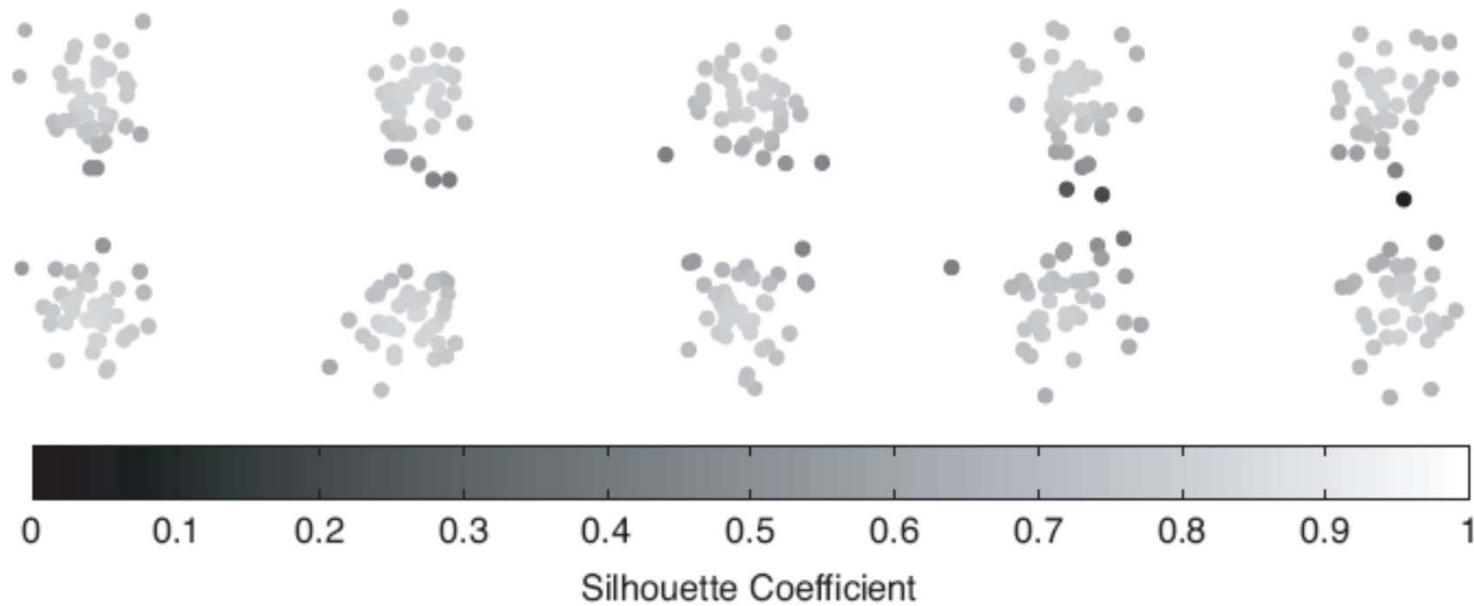
- The silhouette of a clustering  $\mathcal{C} = (C_1, \dots, C_k)$  is defined as:

$$\text{silh}(\mathcal{C}) = \frac{1}{|D|} \sum_{o \in D} s(o),$$

where  $D$  denotes the whole dataset.

- „Reading“ the silhouette coefficient:  
Let  $a(o) \neq 0$ .
  - $b(o) \gg a(o) \Rightarrow s(o) \approx 1$ : good assignment of  $o$  to its cluster  $A$
  - $b(o) \approx a(o) \Rightarrow s(o) \approx 0$ :  $o$  is in-between  $A$  and  $B$
  - $b(o) \ll a(o) \Rightarrow s(o) \approx -1$ : bad, on average  $o$  is closer to members of  $B$
  
- Silhouette Coefficient  $s_c$  of a clustering: average silhouette of all objects
  - $0.7 < s_c \leq 1.0$  strong structure,  $0.5 < s_c \leq 0.7$  medium structure
  - $0.25 < s_c \leq 0.5$  weak structure,  $s_c \leq 0.25$  no structure

### Silhouette Coefficient for points in ten clusters



in: Tan, Steinbach, Kumar: Introduction to Data Mining (Pearson, 2006)



- 1) Introduction to clustering
- 2) Partitioning Methods
  - K-Means
  - K-Medoid
  - Choice of parameters: Initialization, Silhouette coefficient
- 3) Expectation Maximization: a statistical approach
- 4) Density-based Methods: DBSCAN
- 5) Hierarchical Methods
  - Agglomerative and Divisive Hierarchical Clustering
  - Density-based hierarchical clustering: OPTICS
- 6) Evaluation of Clustering Results
- 7) Further Clustering Topics
  - Ensemble Clustering
  - Discussion: an alternative view on DBSCAN

Statistical approach for finding maximum likelihood estimates of parameters in probabilistic models

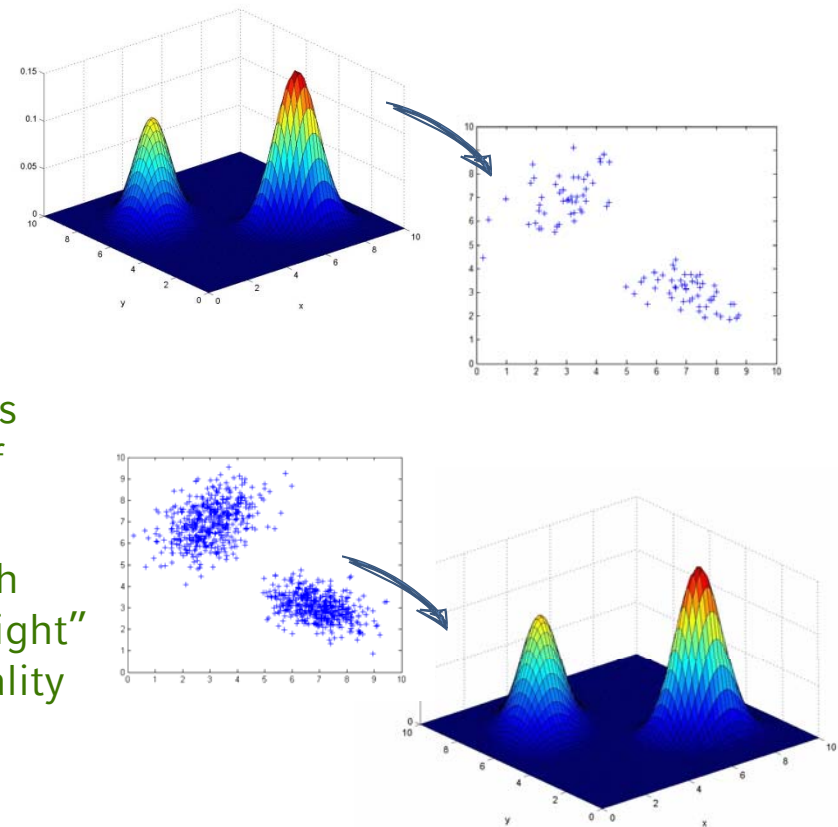
Here: using EM as clustering algorithm

Approach:

Observations are drawn from one of several components of a mixture distribution.

Main idea:

- Define clusters as probability distributions  
→ each object has a certain probability of belonging to each cluster
- Iteratively improve the parameters of each distribution (e.g. center, “width” and “height” of a Gaussian distribution) until some quality threshold is reached



Additional Literature: C. M. Bishop „Pattern Recognition and Machine Learning“, Springer, 2009

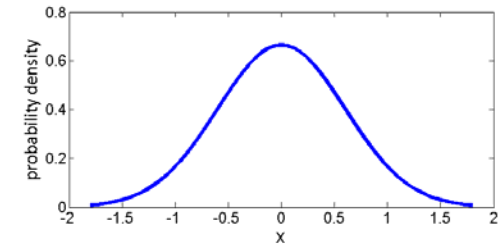
Note: EM is not restricted to Gaussian distributions, but they will serve as example in this lecture.

Gaussian distribution:

- Univariate: single variable  $x \in \mathbb{R}$ :

$$p(x|\mu, \sigma^2) = \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2\sigma^2} \cdot (x-\mu)^2}$$

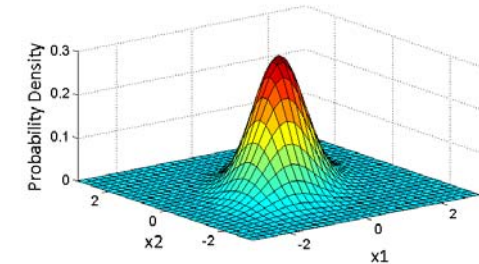
$\nearrow$  mean  $\in \mathbb{R}$       $\nwarrow$  variance  $\in \mathbb{R}$



- Multivariate:  $d$ -dimensional vector  $\mathbf{x} \in \mathbb{R}^d$ :

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \cdot e^{-\frac{1}{2} \cdot (\mathbf{x}-\boldsymbol{\mu})^T \cdot (\boldsymbol{\Sigma})^{-1} \cdot (\mathbf{x}-\boldsymbol{\mu})}$$

$\nearrow$  mean vector  $\in \mathbb{R}^d$       $\nwarrow$  covariance matrix  $\in \mathbb{R}^{d \times d}$

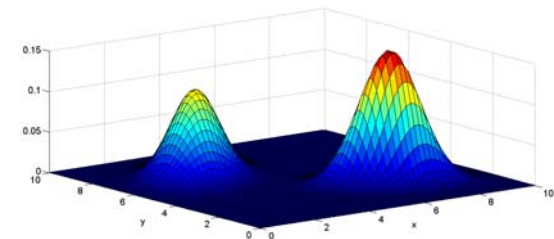


Gaussian mixture distribution with  $K$  components:

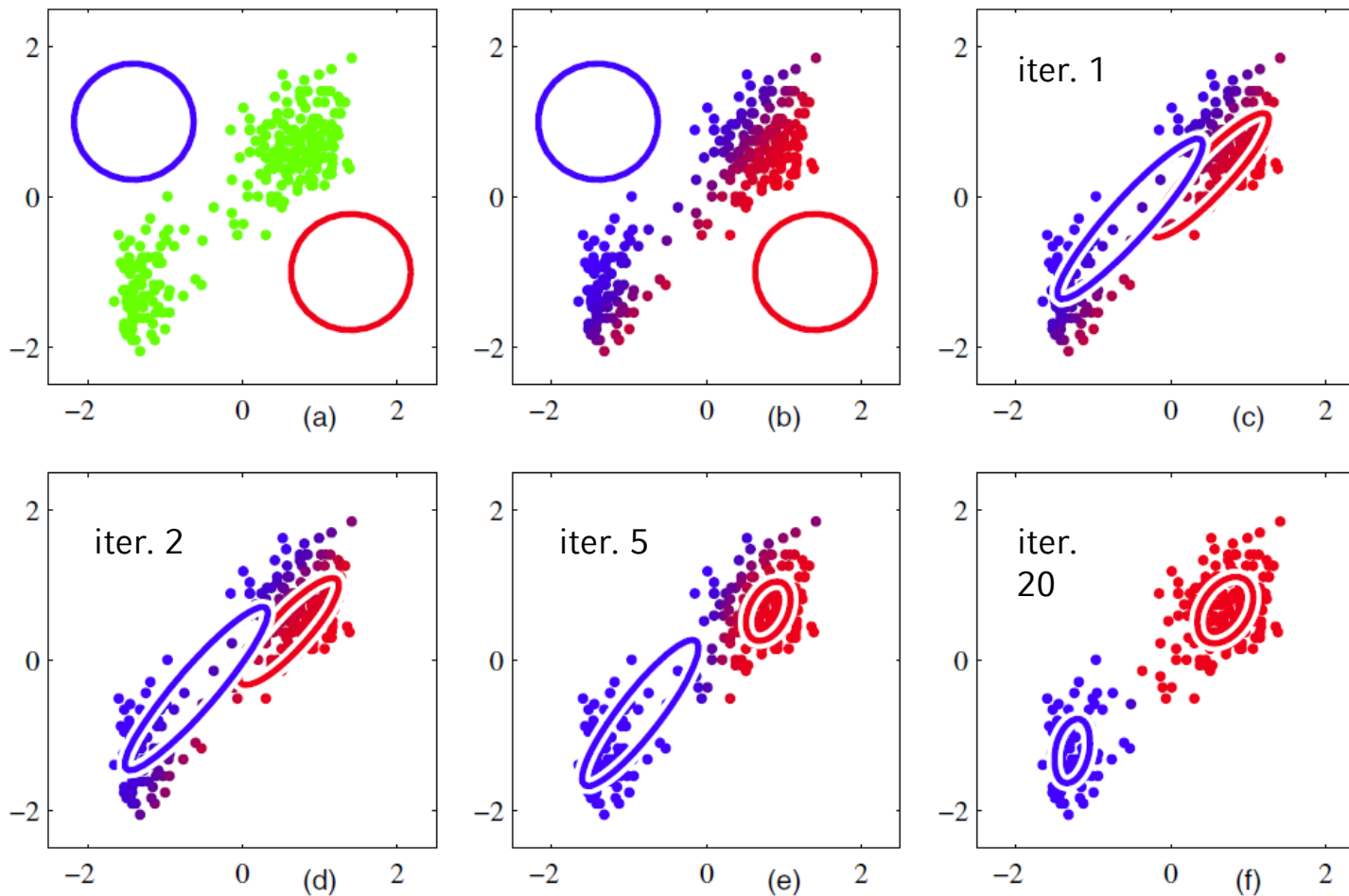
- $d$ -dimensional vector  $\mathbf{x} \in \mathbb{R}^d$ :

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$\uparrow$   
 mixing coefficients  $\in \mathbb{R} : \sum_k \pi_k = 1$  and  $0 \leq \pi_k \leq 1$



Example taken from: C. M. Bishop „Pattern Recognition and Machine Learning“, 2009



Note: EM is not restricted to Gaussian distributions, but they will serve as example in this lecture.

A clustering  $\mathcal{M} = \{C_1, \dots, C_K\}$  is represented by a mixture distribution with parameters  $\Theta =$

$\{\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K\}$  :

$$p(\mathbf{x}|\Theta) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Each *cluster* is represented by one component of the mixture distribution:

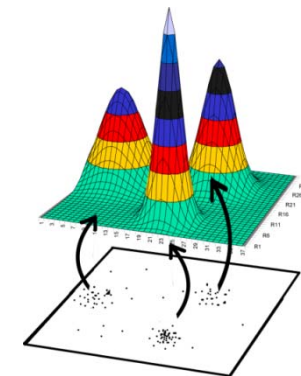
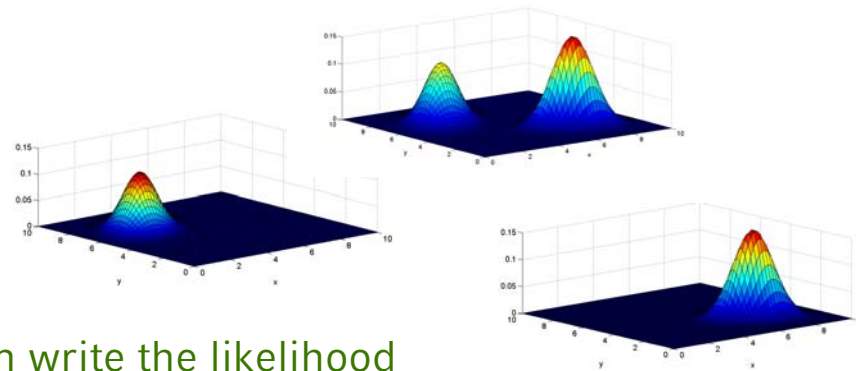
$$p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Given a dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subseteq \mathbb{R}^d$ , we can write the likelihood that all data points  $\mathbf{x}_n \in \mathbf{X}$  are generated (independently) by the mixture model with parameters  $\Theta$  as:

$$\log p(\mathbf{X}|\Theta) = \log \prod_{n=1}^N p(\mathbf{x}_n|\Theta)$$

Goal: Find the parameters  $\Theta_{ML}$  with **maximal (log-)likelihood estimation (MLE)**

$$\Theta_{ML} = \arg \max_{\Theta} \{\log p(\mathbf{X}|\Theta)\}$$



- Goal: Find the parameters  $\Theta_{ML}$  with the **maximal (log-)likelihood estimation!**

$$\Theta_{ML} = \arg \max_{\Theta} \{\log p(\mathbf{X}|\Theta)\}$$

$$\log p(\mathbf{X}|\Theta) = \log \prod_{n=1}^N \sum_{k=1}^K \pi_k \cdot p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \cdot p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Maximization with respect to the means:

$$\frac{\partial}{\partial \boldsymbol{\mu}_j} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_j) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

$$\frac{\partial \log p(\mathbf{X}|\Theta)}{\partial \boldsymbol{\mu}_j} = \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n|\Theta)}{\partial \boldsymbol{\mu}_j} = \sum_{n=1}^N \frac{\frac{\partial p(\mathbf{x}_n|\Theta)}{\partial \boldsymbol{\mu}_j}}{p(\mathbf{x}_n|\Theta)} = \sum_{n=1}^N \frac{\frac{\partial \pi_j \cdot p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_j}}{\sum_{k=1}^K p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} = \sum_{n=1}^N \frac{\pi_j \cdot \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_j) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

$$\frac{\partial \log p(\mathbf{X}|\Theta)}{\partial \boldsymbol{\mu}_j} = \boldsymbol{\Sigma}_j^{-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_j) \frac{\pi_j \cdot \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \stackrel{\text{def}}{=} \mathbf{0}$$

- Define

$$\gamma_j(\mathbf{x}_n) := \pi_j \cdot \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j).$$

$\gamma_j(\mathbf{x}_n)$  is the probability that component  $j$  generated the object  $\mathbf{x}_n$ .

Maximization w.r.t. the means yields:

$$\boldsymbol{\mu}_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)} \quad (\text{weighted mean})$$

Maximization w.r.t. the covariance yields:

$$\boldsymbol{\Sigma}_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_j) (\mathbf{x}_n - \boldsymbol{\mu}_j)^T}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)}$$

Maximization w.r.t. the mixing coefficients yields:

$$\pi_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)}{\sum_{k=1}^K \sum_{n=1}^N \gamma_k(\mathbf{x}_n)}$$

Problem with finding the optimal parameters  $\Theta_{ML}$ :

$$\mu_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)} \quad \text{and} \quad \gamma_j(\mathbf{x}_n) = \frac{\pi_j \cdot \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}$$

- Non-linear mutual dependencies.
- Optimizing the Gaussian of cluster  $j$  depends on all other Gaussians.
- There is no closed-form solution!
- Approximation through iterative optimization procedures
- Break the mutual dependencies by optimizing  $\mu_j$  and  $\gamma_j(\mathbf{x}_n)$  independently



## EM-approach: iterative optimization

1. Initialize means  $\boldsymbol{\mu}_j$ , covariances  $\boldsymbol{\Sigma}_j$ , and mixing coefficients  $\pi_j$  and evaluate the initial log likelihood.
2. **E step:** Evaluate the responsibilities using the current parameter values:

$$\gamma_j^{new}(\mathbf{x}_n) = \frac{\pi_j \cdot \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

3. **M step:** Re-estimate the parameters using the current responsibilities:

$$\boldsymbol{\mu}_j^{new} = \frac{\sum_{n=1}^N \gamma_j^{new}(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j^{new}(\mathbf{x}_n)}$$

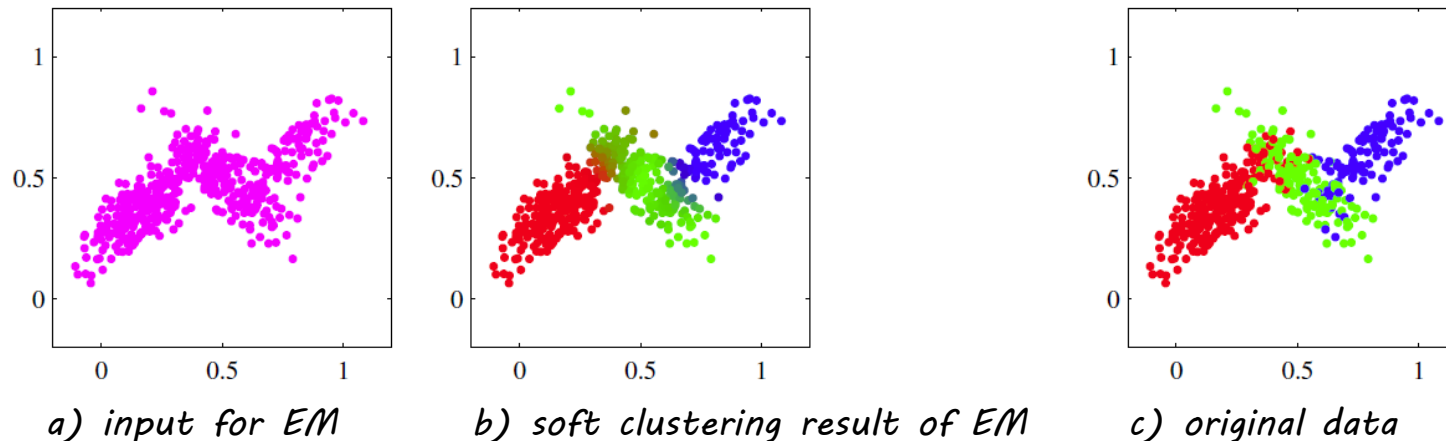
$$\boldsymbol{\Sigma}_j^{new} = \frac{\sum_{n=1}^N \gamma_j^{new}(\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_j^{new})(\mathbf{x}_n - \boldsymbol{\mu}_j^{new})^T}{\sum_{n=1}^N \gamma_j^{new}(\mathbf{x}_n)}$$

$$\pi_j^{new} = \frac{\sum_{n=1}^N \gamma_j^{new}(\mathbf{x}_n)}{\sum_{k=1}^K \sum_{n=1}^N \gamma_k^{new}(\mathbf{x}_n)}$$

4. Evaluate the new log likelihood  $\log p(\mathbf{X} | \boldsymbol{\Theta}^{new})$  and check for convergence of parameters or log likelihood ( $|\log p(\mathbf{X} | \boldsymbol{\Theta}^{new}) - \log p(\mathbf{X} | \boldsymbol{\Theta})| \leq \epsilon$ ).  
If the convergence criterion is not satisfied, set  $\boldsymbol{\Theta} = \boldsymbol{\Theta}^{new}$  and go to step 2.

EM obtains a *soft* clustering (each object belongs to each cluster with a certain probability) reflecting the uncertainty of the most appropriate assignment.

Example taken from: C. M. Bishop „Pattern Recognition and Machine Learning“, 2009



Modification to obtain a *partitioning* variant

- Assign each object to the cluster to which it belongs with the highest probability

$$\text{Cluster}(\text{object}_n) = \underset{k \in \{1, \dots, K\}}{\text{argmax}} \{ \gamma(z_{nk}) \}$$

Superior to k-Means for clusters of varying size  
or clusters having differing variances  
→ more accurate data representation

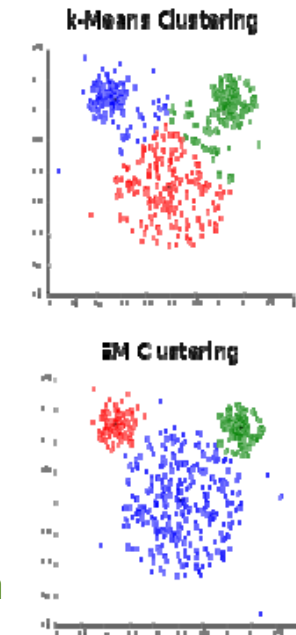
Convergence to (possibly local) maximum

Computational effort for  $N$  objects,  $K$  derived clusters, and  $t$  iterations:

- $O(t \cdot N \cdot K)$
- #iterations is quite high in many cases

Both - result and runtime - strongly depend on

- the initial assignment
  - do multiple random starts and choose the final estimate with highest likelihood
  - Initialize with clustering algorithms (e.g., K-Means usually converges much faster)
  - Local maxima and initialization issues have been addressed in various extensions of EM
- a proper choice of parameter  $K$  (= desired number of clusters)
  - Apply principles of model selection (see next slide)



Classical trade-off problem for selecting the proper number of components  $K$

- If  $K$  is too high, the mixture may overfit the data
- If  $K$  is too low, the mixture may not be flexible enough to approximate the data

Idea: determine candidate models  $\Theta_K$  for a range of values of  $K$  (from  $K_{min}$  to  $K_{max}$ ) and select the model  $\Theta_{K^*} = \max\{\text{qual}(\Theta_K) | K \in \{K_{min}, \dots, K_{max}\}\}$

- Silhouette Coefficient (as for  $k$ -Means) only works for partitioning approaches.
- The MLE (Maximum Likelihood Estimation) criterion is nondecreasing in  $K$

Solution: deterministic or stochastic *model selection* methods<sup>[MP'00]</sup>

which try to balance the goodness of fit with simplicity.

- Deterministic:  $\text{qual}(\Theta_K) = \log p(\mathbf{X} | \Theta_K) + \mathcal{P}(K)$   
where  $\mathcal{P}(K)$  is an increasing function penalizing higher values of  $K$
- Stochastic: based on Markov Chain Monte Carlo (MCMC)

[MP'00] G. McLachlan and D. Peel. *Finite Mixture Models*. Wiley, New York, 2000.