1) Introduction to clustering

2) Partitioning Methods
   - K-Means
   - K-Medoid
   - Choice of parameters: Initialization, Silhouette coefficient

3) <u>Expectation Maximization: a statistical approach</u>

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

Clustering

Statistical approach for finding maximum likelihood estimates of parameters in probabilistic models
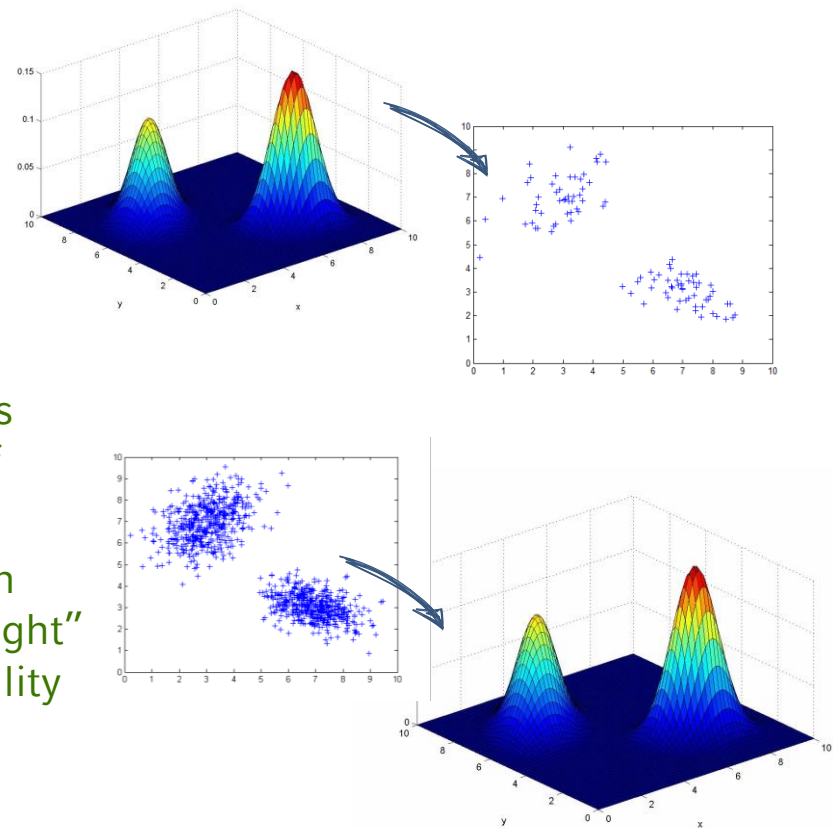
Here: using EM as clustering algorithm

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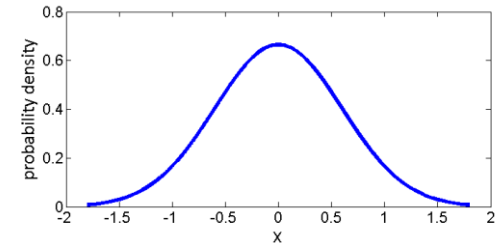
Note: EM is not restricted to Gaussian distributions, but they will serve as example in this lecture.

## Gaussian distribution:

– Univariate: single variable x ∈ ℝ:

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

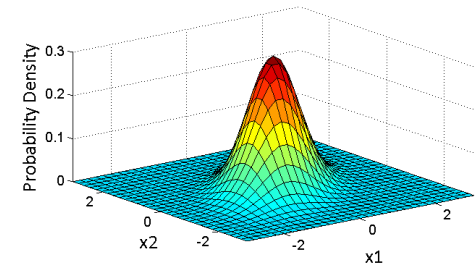$mean \in \mathbb{R}$    $variance \in \mathbb{R}$

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

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

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

## Gaussian mixture distribution with $K$ components:

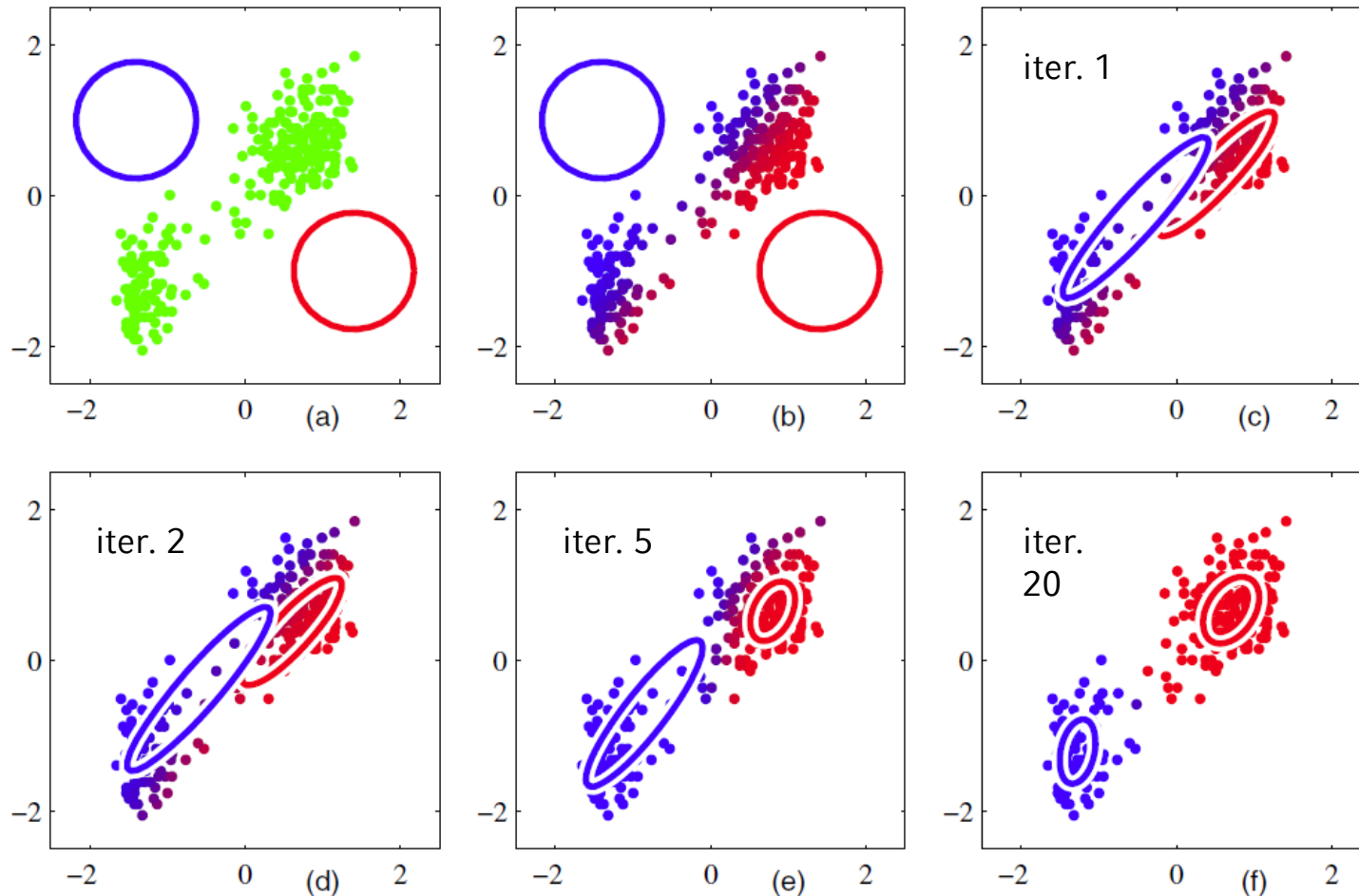– $d$-dimensional vector $\boldsymbol{x} \in \mathbb{R}^d$:

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

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

Example taken from: C. M. Bischop „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* $\mathcal{M} = \{C_1, \ldots, C_K\}$ is represented by a mixture distribution with parameters $\Theta = \{\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \ldots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K\}$ :

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

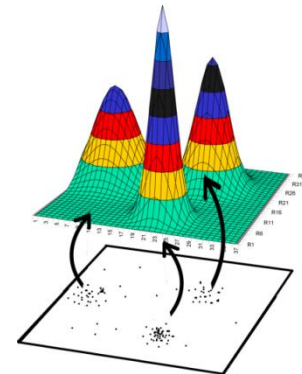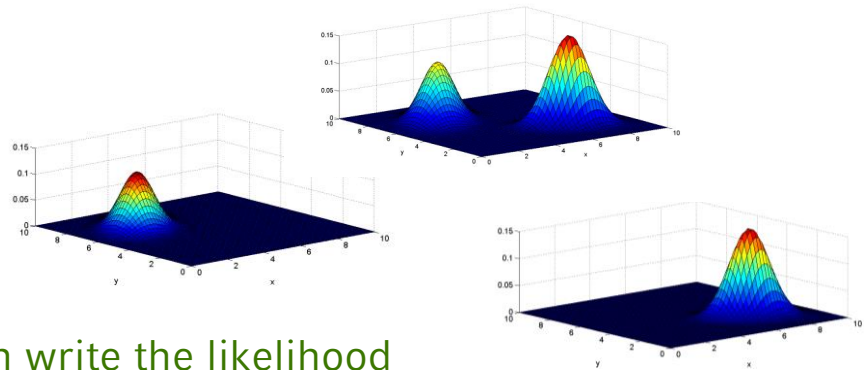Each *cluster* is represented by one component of the mixture distribution:

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

Given a dataset $\mathbf{X} = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\} \subseteq \mathbb{R}^d$, 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:

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

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:

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

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

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

- Define

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

$\gamma_j(x_n)$ is the probability that component $j$ generated the object $x_n$.

Maximization w.r.t. the means yields:

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

Maximization w.r.t. the covariance yields:

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

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

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

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

$$\boldsymbol{\mu}_j = \frac{\sum_{n=1}^{N} \gamma_j(\boldsymbol{x}_n)\, \boldsymbol{x}_n}{\sum_{n=1}^{N} \gamma_j(\boldsymbol{x}_n)} \quad \text{and} \quad \gamma_j(\boldsymbol{x}_n) = \frac{\pi_j \cdot \mathcal{N}\left(\boldsymbol{x}_n \big| \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right)}{\sum_{k=1}^{K} \pi_k \cdot \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

- 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 $\boldsymbol{\mu}_j$ and $\gamma_j(\boldsymbol{x}_n)$ 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:

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

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

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

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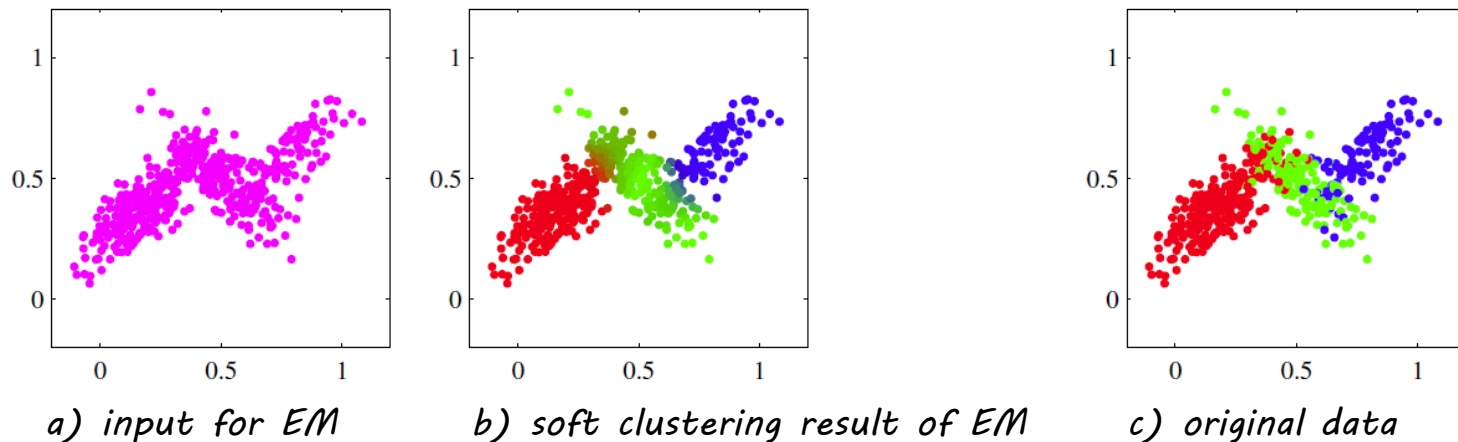
$$\pi_j^{new} = \frac{\sum_{n=1}^{N} \gamma_j^{new}(\boldsymbol{x}_n)}{\sum_{k=1}^{K} \sum_{n=1}^{N} \gamma_k^{new}(\boldsymbol{x}_n)}$$

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



a) input for EM        b) soft clustering result of EM        c) original data

Modification to obtain a *partitioning* variant

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

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

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 principals of model selection (see next slide)



k-Means Clustering



EM Clustering

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 $K$

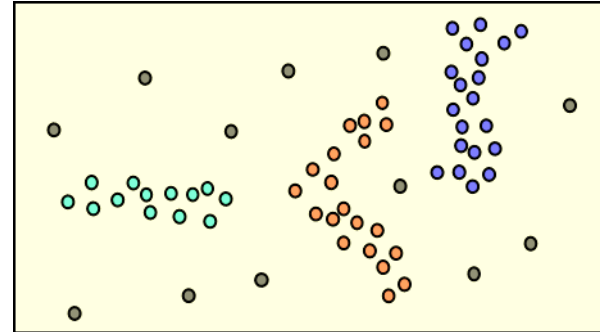Solution: deterministic or stochastic *model selection* methods[MP'00]
which try to balance the goodness of fit with simplicity.

- Deterministic: $qual(\Theta_K) = \log p(\mathbf{X}|\Theta_K) + \mathcal{P}(K)$
  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.

# Contents

Clustering

# Density-Based Clustering

- Basic Idea:

  - Clusters are dense regions in the data space, separated by regions of lower object density



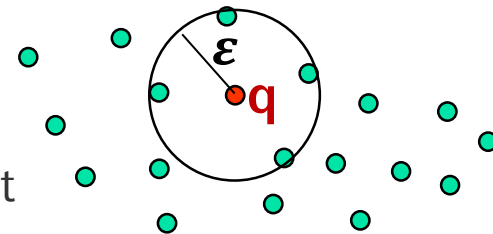- Why Density-Based Clustering?



Results of a *k*-medoid algorithm for *k*=4

- Different density-based approaches exist (see Textbook & Papers)
  Here we discuss the ideas underlying the DBSCAN algorithm
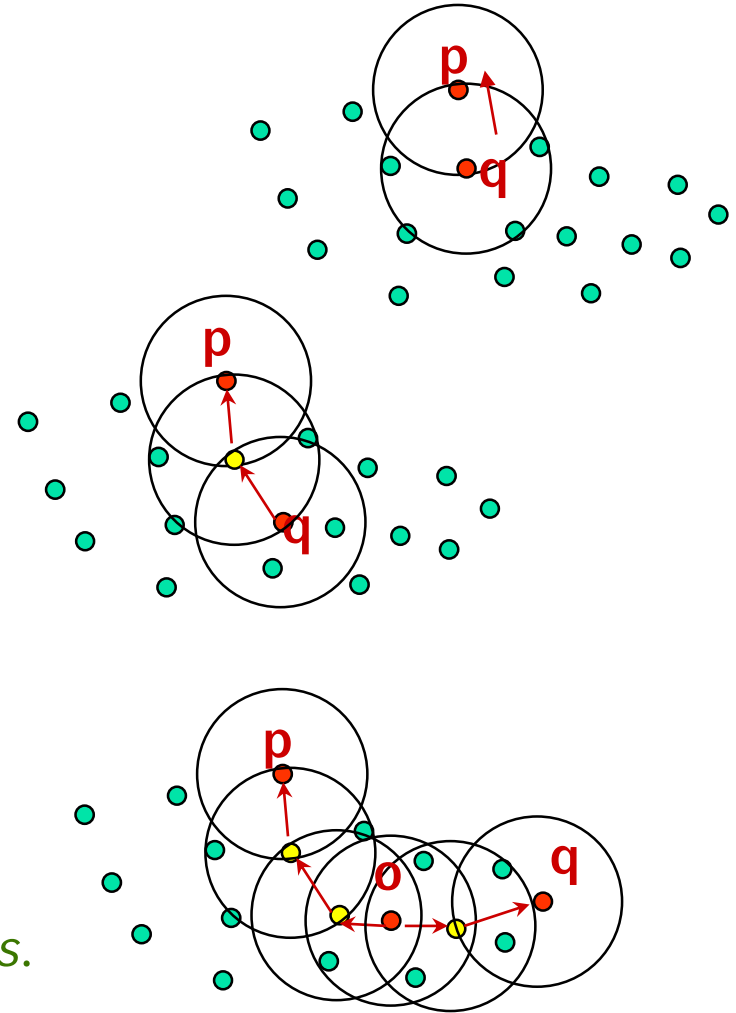
- Intuition for the formalization of the basic idea

  – For any point in a cluster, the local point density around that point has to exceed some threshold

  – The set of points from one cluster is spatially connected

- Local point density at a point $q$ defined by two parameters

  – $\varepsilon$–radius for the neighborhood of point $q$:
    $N_\varepsilon(q) := \{p \in D | dist(p, q) \leq \varepsilon\}$     *! contains q itself !*

  – **MinPts** – minimum number of points in the given neighbourhood $N_\varepsilon(q)$

- $q$ is called a **core object** (or core point)
  w.r.t. $\varepsilon$, *MinPts* if $| N_\varepsilon(q) | \geq$ *MinPts*

$MinPts = 5$ → **q** is a core object

- $p$ **directly density-reachable** from $q$ w.r.t. $\varepsilon$, *MinPts* if
  1) $p \in N_\varepsilon(q)$  and
  2) $q$ is a core object w.r.t. $\varepsilon$, *MinPts*

- **density-reachable**: transitive closure of *directly* density-reachable

- $p$ is **density-connected** to a point $q$ w.r.t. $\varepsilon$, *MinPts* if there is a point $o$ such that both, $p$ and $q$ are density-reachable from $o$ w.r.t. $\varepsilon$, *MinPts*.

- **Density-Based Cluster**: non-empty subset $S$ of database $D$ satisfying:

  1) *Maximality*: if $p$ is in $S$ and $q$ is density-reachable from $p$ then $q$ is in $S$

  2) *Connectivity*: each object in $S$ is density-connected to all other objects in $S$

- **Density-Based Clustering** of a database $D$ : $\{S_1, ..., S_n; N\}$ where

  - $S_1, ..., S_n$ : all density-based clusters in the database $D$

  - $N = D \setminus \{S_1 \cup ... \cup S_n\}$ is called the **noise** (objects not in any cluster)
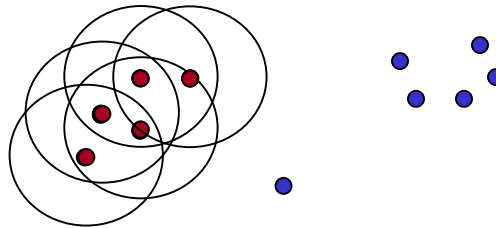


$$\varepsilon = 1.0$$
$$MinPts = 5$$

- Density Based Spatial Clustering of Applications with Noise
- Basic Theorem:
  - Each object in a density-based cluster C is density-reachable from any of its core-objects
  - Nothing else is density-reachable from core objects.

---

**for** each $o \in D$ **do**
    **if** $o$ is not yet classified **then**
        **if** $o$ is a core-object **then**
            collect all objects density-reachable from $o$
            and assign them to a new cluster.
        **else**
            assign $o$ to NOISE

---

  - density-reachable objects are collected by performing successive $\varepsilon$-neighborhood queries.

Ester M., Kriegel H.-P., Sander J., Xu X.: „A Density-Based Algorithm for Discovering Clusters in Large  Spatial Databases with Noise", *In KDD 1996*, pp. 226—231.
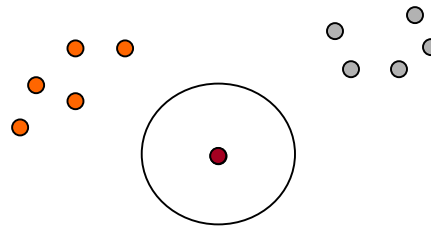
- Parameter
    - $\varepsilon = 2.0$
    - $MinPts = 3$



```
for each o ∈ D do
    if o is not yet classified then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```
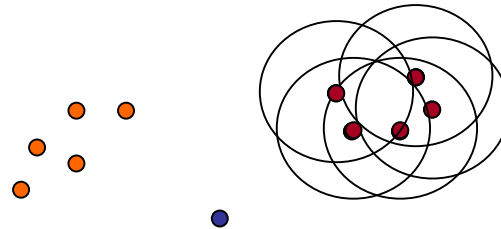
- Parameter
  - $\varepsilon = 2.0$
  - $MinPts = 3$



```
for each o ∈ D do
    if o is not yet classified then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```

■ Parameter

– $\varepsilon = 2.0$

– $MinPts = 3$



```
for each o ∈ D do
    if o is not yet classified then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```