

Ludwig-Maximilians-Universität München Institut für Informatik Lehr- und Forschungseinheit für Datenbanksysteme



#### Lecture notes Knowledge Discovery in Databases

#### Summer Semester 2012

#### **Lecture 9: Clustering III**

Lecture: Dr. Eirini Ntoutsi Tutorials: Erich Schubert

http://www.dbs.ifi.lmu.de/cms/Knowledge\_Discovery\_in\_Databases\_I\_(KDD\_I)





- Previous KDD I lectures on LMU (Johannes Aßfalg, Christian Böhm, Karsten Borgwardt, Martin Ester, Eshref Januzaj, Karin Kailing, Peer Kröger, Jörg Sander, Matthias Schubert, Arthur Zimek)
- Tan P.-N., Steinbach M., Kumar V., *Introduction to Data Mining*, Addison-Wesley, 2006
- Jiawei Han, Micheline Kamber and Jian Pei, *Data Mining: Concepts and Techniques, 3rd ed.,* Morgan Kaufmann, 2011.
- Christoph Lippert | Data Mining in Bioinfortics | Clustering tutorial, http://agbs.kyb.tuebingen.mpg.de/wikis/bg/tutorial GMM.pdf





- Introduction
- A categorization of major clustering methods
- Density-based methods cont'
- Grid-based methods
- Model-based methods
- An overview of clustering
- Things you should know
- Homework/tutorial



## Major clustering methods I

- Partitioning approaches:
  - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
  - Typical methods: k-means, k-medoids, CLARANS
- Hierarchical approaches:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - Typical methods: Diana, Agnes, BIRCH, ROCK, CHAMELEON
- Density-based approaches:
  - Based on connectivity and density functions
  - Typical methods: DBSCAN, OPTICS, DenClue











## Major clustering methods II

- Grid-based approaches:
  - based on a multiple-level granularity structure
  - Typical methods: STING, WaveCluster, CLIQUE
- Model-based approaches:
  - A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
  - Typical methods: EM, SOM, COBWEB
- Frequent pattern-based approaches:
  - Based on the analysis of frequent patterns
  - Typical methods: pCluster
- User-guided or constraint-based approaches:
  - Clustering by considering user-specified or application-specific constraints
  - Typical methods: COD (obstacles), constrained clustering









- Introduction
- A categorization of major clustering methods
- Density-based methods cont'
- Grid-based methods
- Model-based methods
- An overview of clustering
- Things you should know
- Homework/tutorial



## **Density-based clustering**



- Clusters are regions of high density surrounded by regions of low density (noise)
- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
  - Discover clusters of arbitrary shape
  - Handle noise
  - One scan
  - Need density parameters as termination condition
- Several interesting studies:
  - DBSCAN: Ester, et al. (KDD'96)
  - OPTICS: Ankerst, et al (SIGMOD'99).
  - DENCLUE: Hinneburg & D. Keim (KDD'98)
  - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)





## DBSCAN (Ester et al, KDD'96) (previous lecture)



- Two parameters:
  - Eps (or  $\varepsilon$ ): Maximum radius of the neighbourhood
  - MinPts: Minimum number of points in an Eps-neighbourhood of that point
- Eps-neighborhood of a point p in D
  - N<sub>Eps</sub>(p): {q belongs to D | dist(p,q) <= Eps}</p>







A cluster is a maximal set of density-connected points



## **OPTICS (Ankerst et al, SIGMOD'99)**



- OPTICS: <u>Ordering Points To Identify the Clustering Structure</u>
  - Extension of DBSCAN
- It does not produce a clustering of a dataset explicitly, instead it produces a special ordering of the database w.r.t. its density-based clustering structure
- This cluster-ordering contains information that is equivalent to the densitybased clusterings corresponding to a broad range of parameter settings
- Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
- Can be represented graphically







- In many cases the intrinsic cluster structure cannot be characterized by *global* density parameters
- Different *local* densities may be needed to reveal clusters in different regions of the data space



Global densities would result in:

• A,B,C clusters

or

• C<sub>1</sub>, C<sub>2</sub>, C<sub>3</sub> clusters

- Use a hierarchical clustering?
  - Difficult to interpret the dendrogram for large datasets
- Use DBSCAN with different parameters?
  - Infinite number of possible parameters





 Observation: for a constant minPts value, density-based clusters w.r.t. to a lower value for ε are completely contained in density-based clusters with a higher value for ε



C is a cluster w.r.t.  $\varepsilon_1$ C<sub>1</sub>, C<sub>2</sub> are clusters w.r.t.  $\varepsilon_2 < \varepsilon_1$ C contains C<sub>1</sub>, C<sub>2</sub>

- Idea: extend DBSCAN so as several distance parameters are processed at the same time
- This way, density-based clusters w.r.t. different densities are constructed simultaneously



## **OPTICS** basic notions I



#### • Core distance of an object p

Let  $N_{\epsilon}(p)$  be the  $\epsilon$ -neighborhood of p in D and let MinPtsdistance(p) be the distance from p to its MinPts' neighbor.

 $coredistance_{\varepsilon,MinPts}(p) = \begin{cases} UNDEFINED, & if |N_{\varepsilon}(p)| < MinPts\\ MinPtsdistance(p), & otherwise \end{cases}$ 

- Is simply the smallest distance  $\varepsilon'$  between p and an object q in its  $\varepsilon$ -neighborhood such that p would be a core point w.r.t.  $\varepsilon'$  if q was contained in N<sub> $\varepsilon$ </sub>(p). Otherwise is undefined.







• Reachability distance of an object p w.r.t. an object o Let  $N_{\epsilon}(o)$  be the  $\epsilon$ -neighborhood of o in D.

reachability distance<sub> $\varepsilon,MinPts</sub>(p,o) = \begin{cases} UNDEFINED & if |N_{\varepsilon}(o)| < MinPts \\ max{coredistance}(o), dist(o, p)}, & otherwise \end{cases}$ </sub>

- Is the smallest distance so as p is directly density-reachable from o, if o is core.
  It cant be smaller than coredistance(o) because otherwise o will not be core
- it depends on the core object o w.r.t. which it is calculated.











### **OPTICS pseudocode II**



ExpandClusterOrder(SetOfObjects, Object, ε, MinPts, OrderedFile); neighbors := SetOfObjects.neighbors(Object,  $\varepsilon$ ); Object.Processed := TRUE; Object.reachability\_distance := UNDEFINED; Object.setCoreDistance(neighbors, ε, MinPts); OrderedFile.write(Object); IF Object.core\_distance <> UNDEFINED THEN OrderSeeds.update(neighbors, Object); WHILE NOT OrderSeeds.empty() DO currentObject := OrderSeeds.next(): neighbors:=SetOfObjects.neighbors(currentObject, ε); currentObject.Processed := TRUE; currentObject.setCoreDistance(neighbors, ε, MinPts); OrderedFile.write(currentObject); IF currentObject.core\_distance<>UNDEFINED THEN OrderSeeds.update(neighbors, currentObject); END; // ExpandClusterOrder

Output: OPTICS outputs the points in a particular ordering. Each point is accompanied with its coredistance and its smallest reachability distance. The objects contained in OrderSeeds are sorted by their reachabilitydistance to the closest core object, Object, from which they have been directly density-reachable.

the object having the smallest reachability-distance in the seed-list is selected, currentObject

OrderSeeds::update(neighbors, CenterObject); c\_dist := CenterObject.core\_distance; FORALL Object FROM neighbors DO IF NOT Object.Processed THEN new\_r\_dist:=max(c\_dist,CenterObject.dist(Object)); IF Object.reachability\_distance=UNDEFINED THEN Object.reachability\_distance := new\_r\_dist; insert(Object, new\_r\_dist); ELSE // Object already in OrderSeeds IF new\_r\_dist<Object.reachability\_distance THEN Object.reachability\_distance := new\_r\_dist; decrease(Object, new\_r\_dist); END; // OrderSeeds::update





- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3



#### seed list:

![](_page_16_Picture_0.jpeg)

![](_page_16_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_16_Figure_4.jpeg)

#### seed list: (B,40) (I, 40)

![](_page_17_Picture_0.jpeg)

![](_page_17_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_17_Picture_4.jpeg)

![](_page_17_Figure_5.jpeg)

#### seed list: (I, 40) (C, 40)

![](_page_18_Picture_0.jpeg)

![](_page_18_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_18_Figure_4.jpeg)

seed list: (J, 20) (K, 20) (L, 31) (C, 40) (M, 40) (R, 43)

![](_page_19_Picture_0.jpeg)

![](_page_19_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_19_Figure_4.jpeg)

![](_page_19_Figure_5.jpeg)

seed list: (L, 19) (K, 20) (R, 21) (M, 30) (P, 31) (C, 40)

![](_page_20_Picture_0.jpeg)

![](_page_20_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_20_Figure_4.jpeg)

seed list: (M, 18) (K, 18) (R, 20) (P, 21) (N, 35) (C, 40)

![](_page_21_Picture_0.jpeg)

![](_page_21_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_21_Figure_4.jpeg)

seed list: (K, 18) (N, 19) (R, 20) (P, 21) (C, 40)

![](_page_22_Picture_0.jpeg)

![](_page_22_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_22_Figure_4.jpeg)

![](_page_22_Figure_5.jpeg)

#### seed list: (N, 19) (R, 20) (P, 21) (C, 40)

![](_page_23_Picture_0.jpeg)

![](_page_23_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_23_Figure_4.jpeg)

seed list: (R, 20) (P, 21) (C, 40)

![](_page_24_Picture_0.jpeg)

![](_page_24_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_24_Figure_4.jpeg)

#### seed list: (P, 21) (C, 40)

![](_page_25_Picture_0.jpeg)

![](_page_25_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_25_Picture_4.jpeg)

seed list: (C, 40)

![](_page_26_Picture_0.jpeg)

![](_page_26_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_26_Figure_4.jpeg)

#### seed list: (D, 22) (F, 22) (E, 30) (G, 35)

![](_page_27_Picture_0.jpeg)

![](_page_27_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_27_Figure_4.jpeg)

#### seed list: (F, 22) (E, 22) (G, 32)

![](_page_28_Picture_0.jpeg)

![](_page_28_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_28_Figure_4.jpeg)

#### seed list: (G, 17) (E, 22)

![](_page_29_Picture_0.jpeg)

![](_page_29_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_29_Figure_4.jpeg)

#### seed list: (E, 15) (H, 43)

![](_page_30_Picture_0.jpeg)

![](_page_30_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_30_Figure_4.jpeg)

![](_page_31_Picture_0.jpeg)

![](_page_31_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_31_Figure_4.jpeg)

#### seed list: -

![](_page_32_Picture_0.jpeg)

![](_page_32_Picture_1.jpeg)

- Example Database (2-dimensional, 16 points)
- $\varepsilon$ = 44, MinPts = 3

![](_page_32_Figure_4.jpeg)

![](_page_33_Picture_0.jpeg)

#### **Reachability plot**

![](_page_33_Picture_2.jpeg)

- A 2D plot of objects ordering (x-axis) and reachability distance (y-axis)
- Clusters correspond to valleys in the plot since their cluster members have a low reachability distance to their nearest neighbor.

![](_page_33_Figure_5.jpeg)

![](_page_33_Figure_6.jpeg)

![](_page_34_Picture_0.jpeg)

#### **Reachability plot: another example**

![](_page_34_Picture_2.jpeg)

![](_page_34_Figure_3.jpeg)

![](_page_35_Picture_0.jpeg)

## Reachability plot: an example with hierarchical clusters

![](_page_35_Picture_2.jpeg)

![](_page_35_Figure_3.jpeg)

![](_page_35_Figure_4.jpeg)

**Cluster ordering of the objects** 

![](_page_36_Picture_0.jpeg)

## Reachability plot: how to obtain a clustering?

![](_page_36_Picture_2.jpeg)

• Draw a horizontal line to obtain a clustering.

![](_page_36_Figure_4.jpeg)

Depending on the data distribution, the lower the line is the more clusters would emerge

![](_page_36_Figure_6.jpeg)

![](_page_37_Picture_0.jpeg)

#### **Parameters effect**

![](_page_37_Picture_2.jpeg)

![](_page_37_Picture_3.jpeg)

MinPts = 10,  $\varepsilon$  = 10

![](_page_37_Figure_5.jpeg)

optimal parameters

![](_page_37_Figure_7.jpeg)

![](_page_37_Picture_8.jpeg)

smaller  $\boldsymbol{\epsilon}$ 

MinPts = 2,  $\varepsilon$  = 10

![](_page_37_Picture_11.jpeg)

smaller MinPts

- $\bullet$  Cluster ordering is robust to the parameter values  $\epsilon$  and MinPts
  - Good results when parameter values are "large enough"
  - the smaller  $\boldsymbol{\epsilon},$  the more objects may have undefined reachability distance
  - the smaller MinPts, the more jagged the plot looks

• Also, cluster ordering is independent from the dimension of the dataset

#### Knowledge Discovery in Databases I: Clustering III

![](_page_38_Picture_0.jpeg)

![](_page_38_Picture_1.jpeg)

- Introduction
- A categorization of major clustering methods
- Density-based methods cont'
- Grid-based methods
- Model-based methods
- An overview of clustering
- Things you should know
- Homework/tutorial

![](_page_39_Picture_0.jpeg)

## **Grid-based methods**

![](_page_39_Picture_2.jpeg)

- A grid structure is used to capture the density of the dataset.
  - A cluster is a set of connected dense cells
  - STING (a STatistical INformation Grid approach) by Wang, Yang and Muntz (VLDB'97)
  - WaveCluster by Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
  - CLIQUE: Agrawal, et al. (SIGMOD'98)
    - for high-dimensional data

- Appealing features
  - No assumption on the number of clusters
  - Discovering clusters of arbitrary shapes
  - Ability to handle outliers

![](_page_39_Picture_13.jpeg)

![](_page_40_Picture_0.jpeg)

![](_page_40_Picture_1.jpeg)

- The spatial area is divided into rectangular cells
- There are several levels of cells corresponding to different levels of resolution

![](_page_40_Figure_4.jpeg)

![](_page_41_Picture_0.jpeg)

![](_page_41_Picture_1.jpeg)

- Each cell at a high level is partitioned into a number of smaller cells in the next lower level
- Statistical info of each cell is pre-computed and stored beforehand for query answering
- Parameters of higher level cells can be easily calculated from parameters of lower level cells
  - Count, mean, standard deviation, min, max
  - Type of distribution—normal, uniform, etc

![](_page_42_Picture_0.jpeg)

![](_page_42_Picture_1.jpeg)

- A top-down approach
- Start from a pre-selected layer—typically with a small number of cells
- Remove the irrelevant cells from further consideration
- When finish examining the current layer, proceed to the next lower level
- Repeat this process until the bottom layer is reached

![](_page_43_Picture_0.jpeg)

![](_page_43_Picture_1.jpeg)

- Advantages:
  - Query-independent, easy to parallelize, incremental update
  - O(K), where K is the number of grid cells at the lowest level
- Disadvantages:
  - All the cluster boundaries are either horizontal or vertical, and no diagonal boundary is detected

![](_page_44_Picture_0.jpeg)

![](_page_44_Picture_1.jpeg)

- Introduction
- A categorization of major clustering methods
- Density-based methods cont'
- Grid-based methods
- Model-based methods
- An overview of clustering
- Things you should know
- Homework/tutorial

![](_page_45_Picture_0.jpeg)

## Model-based clustering

![](_page_45_Picture_2.jpeg)

- Assumption: data have been generated by a statistical process
- Goal: find the statistical model that best fits the data
- Each cluster can be seen as one distribution
  - e.g., Gaussian distribution
- Objects are assumed to be independent samples from their cluster distribution
- A particular kind of statistical model: Gaussian mixture models
- Procedure: decide on the model and find the parameters of that model from the data

![](_page_46_Picture_0.jpeg)

## Gaussian Mixture Models

![](_page_46_Picture_2.jpeg)

- Objects are points  $x = (x_1, ..., x_d)$  in a Euclidean vector space
- Data are independent and identically distributed samples from a mixture of k distributions
- Each cluster is a multivariate Gaussian distribution
- Each cluster is represented by
  - Mean (centroid)  $\mu_c$
  - d x d covariance matrix  $\boldsymbol{\Sigma}_c$  for the points in cluster c
- Probability density function of a Gaussian distribution

$$P(x \mid c) = \frac{1}{\sqrt{(2\pi)^{d} \mid \Sigma_{c} \mid}} e^{-\frac{1}{2} \cdot (x - \mu_{c})^{T} \cdot \Sigma_{c}^{-1} \cdot (x - \mu_{c})}$$

![](_page_47_Picture_0.jpeg)

## **Multivariate normal distribution**

![](_page_47_Picture_2.jpeg)

![](_page_47_Figure_3.jpeg)

![](_page_48_Picture_0.jpeg)

![](_page_48_Picture_2.jpeg)

• Probability of a cluster c

$$\mathbf{P}(c_l) = \frac{1}{N} \sum_{i=1}^{N} \mathbf{P}(c_l | x_i)$$

• Probability of observing an object x<sub>i</sub>

$$\mathbf{P}(x_i) = \sum_{l=1}^k \mathbf{P}(c_l) \mathbf{P}(x_i | c_l)$$

where P(x<sub>i</sub>|c<sub>i</sub>) is given by the probability density function of the Gaussian distribution

![](_page_49_Picture_0.jpeg)

![](_page_49_Picture_2.jpeg)

 If the objects are generated in an independent manner, the probability of the whole set of objects X, |X|=N, is just the product of the probabilities of each x<sub>i</sub> in X:

$$\mathcal{L} = \prod_{i=1}^{N} P(x_i)$$
$$= \prod_{i=1}^{N} \sum_{l=1}^{k} P(c_l) P(x_i | c_l)$$

• Using statistical methods, we can estimate the parameters of these distributions from the data, and thus describe the clusters.

![](_page_50_Picture_0.jpeg)

![](_page_50_Picture_2.jpeg)

- How can we partition the data?
  - Choose the most likely cluster assignment of each object

$$\operatorname{argmax}_{l} P(c_{l}|x_{i}) = \operatorname{argmax}_{l} P(c_{l}) P(x_{i}|c_{l})$$

- How to estimate the efficient statistics of each cluster?
  - Use Expectation Maximization (EM) algorithm
  - Original algorithm by [Dempster, Laird and Rubin, 1977]
  - A general method for method for finding the maximum-likelihood estimate of a data distribution, when the data is partially missing or hidden.
    - In our case, data are fully observed
    - The cluster assignments of an object x<sub>i</sub> though can be seen as hidden variables

![](_page_51_Picture_0.jpeg)

![](_page_51_Picture_1.jpeg)

- Initialize cluster assignments
- Two alternating steps:
  - E-step

re-estimate the expected-values of the hidden data (cluster assignments) under the current estimate of the model

- M-step

re-estimate the model parameters such that the likelihood according to the current estimate of the complete data is maximized

Until convergence

$$\frac{\mathcal{L}_{new}}{\mathcal{L}_{old}} < 1 + \epsilon$$

![](_page_52_Picture_0.jpeg)

![](_page_52_Picture_1.jpeg)

• E-step: re-estimate the expected-values of the hidden data (cluster assignments) under the current estimate of the model

$$\mathbf{P}^{new}(c_l|x_i) = \mathbf{P}(c_l)\mathbf{P}(x_i|c_l)$$

![](_page_53_Picture_0.jpeg)

## **EM algorithm III**

![](_page_53_Picture_2.jpeg)

- M-step: re-estimate the model parameters such that the likelihood according to the current estimate of the complete data is maximized
  - Cluster densities:

 $\mathbf{P}^{new}(c_l) = \frac{1}{N} \sum_{i=1}^{N} \mathbf{P}^{new}(c_l | x_i)$  $\mu_l^{new} = \frac{\sum_{i=1}^N x_i \mathbf{P}^{new}(c_l | x_i)}{\sum_{i=1}^N \mathbf{P}^{new}(c_l | x_i)}$ **Cluster means:** 

Cluster covariances:

$$\Sigma_{l}^{new} = \frac{\sum_{i=1}^{N} (x_{i} - \mu_{l}^{new})(x_{i} - \mu_{l}^{new})' \mathbf{P}^{new}(c_{l}|x_{i})}{\sum_{i=1}^{N} \mathbf{P}^{new}(c_{l}|x_{i})}$$

![](_page_54_Picture_0.jpeg)

## EM and k-Means

![](_page_54_Picture_2.jpeg)

- EM is similar to the k-Means algorithm (previous lecture)
- k-Means for Euclidean data is a special case of EM for spherical Gaussian distributions with equal covariance matrices, but different means
- E-step (EM) → assign each object to a cluster step (k-Means)
  In EM each object is assigned to a cluster with a probability
- M-step (EM)  $\rightarrow$  compute cluster centroids step (k-Means)
  - In EM, the computation of the mean also considers the fact that each object belong to a distribution with a certain probability

![](_page_55_Picture_0.jpeg)

![](_page_55_Picture_2.jpeg)

- EM can be slow
- Not practical for models with a large number of components
- Problematic when clusters contain only a few points or if the points are nearly co-linear
- The choice of the exact model to use
- Difficulties with noise and outliers
- + More general than k-Means and fuzzy c-Means because they can use distributions of various types
- + Thus, it can find clusters of different sizes and elliptical shapes
- + It is easy to characterize the produced clusters

![](_page_56_Picture_0.jpeg)

![](_page_56_Picture_1.jpeg)

- Introduction
- A categorization of major clustering methods
- Density-based methods cont'
- Grid-based methods
- Model-based methods
- An overview of clustering
- Things you should know
- Homework/tutorial

![](_page_57_Picture_0.jpeg)

![](_page_57_Picture_1.jpeg)

- A cluster is a set of data objects that are similar to one another within the same cluster and dissimilar to the objects in other clusters
- Cluster analysis: Find similarities between data according to the characteristics found in the data and group similar data objects into clusters
- Key points in clustering
  - Similarity/ distance function
  - Learning algorithm
- An unsupervised learning task
  - No clues on the number of clusters, nor in the characteristics of these clusters
- Important DM task: as a stand-alone tool or as a preprocessing step
- A large amount of algorithms
  - Partitioning methods
  - Hierarchical methods
  - Density-based methods
  - Model-based methods

![](_page_58_Picture_0.jpeg)

## **Clustering methods I**

![](_page_58_Picture_2.jpeg)

#### Partitioning methods

- Construct a partition of a database D of n objects into a set of k clusters
  - Each object belongs to exactly one cluster (hard clustering)
  - The number of clusters k is given in advance
- The partition should optimize the chosen partitioning criterion
  - e.g., minimize the intra-cluster variance, i.e., the sum of the squared distances from each data point to its cluster center.
- k-Means: choose a set of k points  $\{c_1, c_2, ..., c_k\}$  in the d-dimensional space to form clusters  $\{C_1, C_2, ..., C_k\}$  such that the following quantity is minimized

$$Cost(C) = \sum_{i=1}^{k} \sum_{x \in C_i} (x - c_i)^2$$

![](_page_58_Picture_11.jpeg)

- k-Means (centroid) k-Medoids (medoid)
- Other methods that scale to large datasets, e.g. CLARA, CLARANS

![](_page_59_Picture_0.jpeg)

## **Clustering methods II**

# LMU

#### **Hierarchical methods**

- Create a hierarchical decomposition of the dataset. Not a single clustering but a set of nested clusters organized as a hierarchical tree (dendrogram)
- 2 ways: Agglomerative (bottom up) Divisive (top down)
- How to merge (split) and when to stop?
- Inter-cluster similarity:
  - single link,
  - complete link,
  - group average,
  - centroid,
  - Ward's method,

![](_page_60_Picture_0.jpeg)

## **Clustering methods III**

![](_page_60_Picture_2.jpeg)

**Density-based clustering** 

- Clusters are regions of high density surrounded by regions of low density (noise)
- Density is measured locally in Eps-neighborhood
- DBSCAN
  - minPts, Eps parameters
  - Core points, border points, noise points
  - Direct reachability, reachability, connectivity, cluster
- OPTICS
  - Cluster ordering
  - Core distance, reachability distance
  - Reachability plot

![](_page_60_Picture_14.jpeg)

![](_page_61_Picture_0.jpeg)

## **Clustering methods IV**

![](_page_61_Picture_2.jpeg)

#### Grid-based clustering

- A grid structure is used to capture the dataset distribution
- Work in the grid, after mapping the points to the grid

#### Model-based clustering

- Data have been generated by a statistical process, the goal is to find the statistical model that best fits the data
- EM algorithm

![](_page_62_Picture_0.jpeg)

![](_page_62_Picture_1.jpeg)

- Introduction
- A categorization of major clustering methods
- Density-based methods cont'
- Grid-based methods
- Model-based methods
- An overview of clustering
- Things you should know
- Homework/tutorial

![](_page_63_Picture_0.jpeg)

## Things you should know

![](_page_63_Picture_2.jpeg)

- Density-based methods cont'
  - OPTICS
    - o core-distance, reachability distance
    - o Reachability plot
- Grid-based methods
  - STING
- Model-based methods
  - EM

![](_page_64_Picture_0.jpeg)

## Homework/ Tutorial

![](_page_64_Figure_2.jpeg)

#### **Tutorial:** Tutorial this Thursday on clustering

#### <u>Homework</u>:

- Try OPTICS in Weka, Elki
  - Can you interpret the reachability plot?
  - What if you change some parameter, ε, MinPts?

#### Suggested reading:

- Tan P.-N., Steinbach M., Kumar V., Introduction to Data Mining, Addison-Wesley, 2006 (Chapter 9).
- Han J., Kamber M., Pei J. Data Mining: Concepts and Techniques 3rd ed., Morgan Kaufmann, 2011 (Chapter 10)