

8. Cluster Analysis

Knowledge Discovery in Databases with Exercises

Dominik Probst, dominik.probst@fau.de

Computer Science 6 (Data Management), Friedrich-Alexander-Universität Erlangen-Nürnberg

Summer semester 2025

1. **Basic Concepts**
2. **Partitioning Methods**
3. **Hierarchical Methods**
4. **Density-based Methods**
5. **Grid-based Methods**
6. **Evaluation of Clustering**
7. **Summary**

Basic Concepts

- **Cluster:** A collection of data objects within a larger set that are:
 - Similar (or related) to one another within the same group and,
 - dissimilar (or unrelated) to the objects outside the group.
- **Cluster analysis (or clustering, data segmentation, . . .):**
 - Define similarities among data based on the characteristics found in the data (input from user!).
 - Group similar data objects into clusters.
- **Unsupervised learning:**
 - No predefined classes.
 - I.e., learning by observation (vs. learning by examples: supervised).
- **Typical applications:**
 - As a stand-alone tool to get insight into data distribution.
 - As a preprocessing step for other algorithms.

- **Biology:**
 - Taxonomy of living things: kingdom, phylum, class, order, family, genus, and species.
- **Land use:**
 - Identification of areas of similar land use in an earth-observation database.
- **Marketing:**
 - Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs.
- **City planning:**
 - Identifying groups of houses according to their house type, value, and geographical location.
- **Earthquake studies:**
 - Observed earthquake epicenters should be clustered along continent faults.
- **Climate:**
 - Understanding earth climate, find patterns of atmosphere and ocean.

- **A good clustering method will produce high-quality clusters.**
 - **High intra-class similarity:**
 - Cohesive within clusters.
 - **Low inter-class similarity:**
 - Distinctive between clusters.
- **The **quality** of a clustering method depends on:**
 - the **similarity measure** used by the method,
 - its implementation, and
 - its ability to discover some or all of the hidden patterns.

- **Dissimilarity/similarity metric:**

- Similarity is expressed in terms of a distance function, typically a metric: $d(x, y)$.
- The definitions of distance functions are usually rather different for interval-scaled, boolean, categorical, ordinal, ratio, and vector variables.
- **Weights** should be associated with different variables based on applications and data semantics.

- **Quality of clustering:**

- There is usually a separate **"quality" function** that measures the "goodness" of a cluster.
- It is hard to define "similar enough" or "good enough."
- The answer is typically highly case dependent.

- **Partitioning criteria:**
 - Single level vs. hierarchical partitioning.
 - Often, multi-level hierarchical partitioning is desirable.
- **Separation of clusters:**
 - Exclusive (e.g., one customer belongs to only one region) vs.
 - Non-exclusive (e.g., one document may belong to more than one class).
- **Similarity measure:**
 - Distance-based (e.g., Euclidean, road network, vector) vs.
 - Connectivity-based (e.g., density or contiguity).
- **Clustering space:**
 - Full space (often when low-dimensional) vs.
 - Subspaces (often in high-dimensional clustering).

- **Scalability:**
 - Clustering all the data instead of only the samples.
- **Ability to deal with different types of attributes:**
 - Numerical, binary, categorical, ordinal, linked, and mixture of these.
- **Constraint-based clustering:**
 - User may give inputs on constraints.
 - Use domain knowledge to determine input parameters.
- **Interpretability and usability.**
- **Others:**
 - Discovery of clusters with arbitrary shape.
 - Ability to deal with noisy data.
 - Incremental clustering and insensitivity to input order.
 - High dimensionality.

- **Partitioning approach:**

- Construct various partitions and then evaluate them by some criterion.
- E.g., minimizing the sum of square errors.
- Typical methods: k-means, k-medoids, CLARA, CLARANS.

- **Hierarchical approach:**

- Create a hierarchical decomposition of the set of data (or objects) using some criterion.
- Typical methods: AGNES, DIANA, BIRCH, CHAMELEON.

- **Density-based approach:**

- Based on connectivity and density functions.
- Typical methods: DBSCAN, OPTICS, DENCLUE.

- **Grid-based approach:**

- Based on a multiple-level granularity structure.
- Typical methods: STING, WaveCluster, CLIQUE.

- **Model-based approach:**

- 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 approach:**

- Based on the analysis of frequent patterns.
- Typical methods: p-Cluster.

- **User-guided or constraint-based approach:**

- Clustering by considering user-specified or application-specific constraints.
- Typical methods: COD (obstacles), constrained clustering.

- **Link-based clustering:**

- Objects are often linked together in various ways.
- Massive links can be used to cluster objects: SimRank, LinkClus.

Partitioning Methods

- **Partitioning method:**

- Partition a database D of n objects $o_j, j \in \{1, \dots, n\}$ into a set of k -clusters $C_i, 1 \leq i \leq k$ such that the sum of squared distances to c_i is minimized (where c_i is the **centroid** or **medoid** of cluster C_i):

$$\min \sum_{i=1}^k \sum_{j=1}^n d(o_j, c_i)^2.$$

- **Given k , find a partition of k clusters that optimizes the chosen partitioning criterion.**

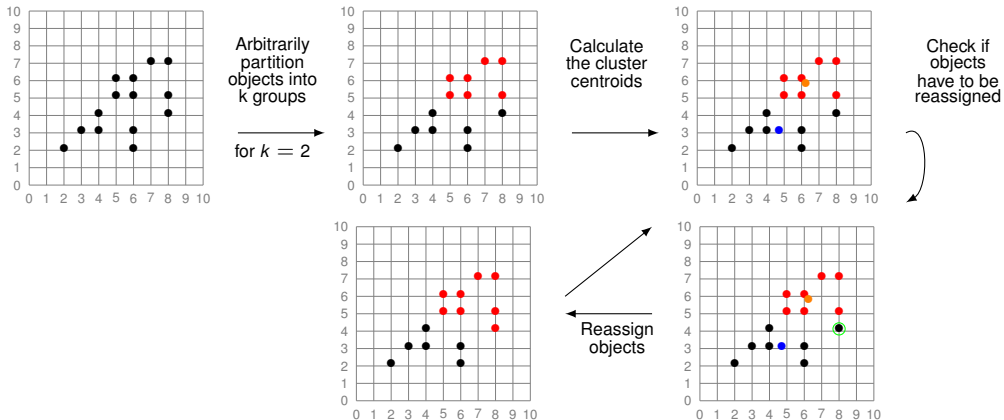
- Globally optimal: exhaustively enumerate all partitions.
- Heuristic methods: k-means and k-medoids algorithms.
- **k-means**¹:
 - Each cluster is represented by the center of the cluster.
- **k-medoids** or PAM (Partition around medoids)²:
 - Each cluster is represented by one of the objects in the cluster.

¹J. Macqueen, "Some methods for classification and analysis of multivariate observations," in *In 5-th Berkeley Symposium on Mathematical Statistics and Probability*, 1967, pp. 281–297.

²L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, 1990, ISBN: 978-0-47187876-6. DOI: 10.1002/9780470316801.

- **Given k , the k -means algorithm is implemented in four steps:**
 1. Partition the database into k non-empty subsets.
 - E.g. the first $\frac{n}{k}$ objects, then the next $\frac{n}{k}$ objects, . . .
 2. Compute the **centroids** of the clusters of the current partitioning.
 - The centroid is the center, i.e. mean point, of the cluster.
 - For each attribute (or dimension), calculate the average value.
 3. Assign each object to the cluster with the nearest centroid.
 - That is, for each object calculate distance to each of the k centroids and pick the one with the smallest distance.
 4. If any object has changed its cluster, go back to step 2. Otherwise stop.
- **Variant:**
 - Start with arbitrarily chosen k objects as initial centroids in step 1.
 - Continue with step 3.

An Example of k -means Clustering



- **Strength:**

- Efficient: $\mathcal{O}(tkn)$, where n is # objects, k is # of clusters, and t is the # of iterations.
Normally: $k, t \ll n$.
- Comparing: PAM: $\mathcal{O}(k(n-k)^2)$, CLARA: $\mathcal{O}(ks^2 + k(n-k))$.
- Comment: Often terminates at a local optimum.

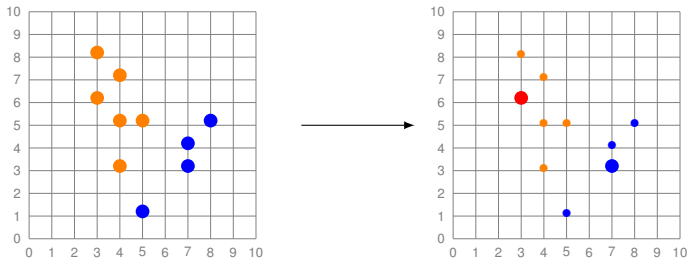
- **Weakness:**

- Applicable only to objects in a continuous n -dimensional space.
 - Using the k -modes method for categorical data.
 - In comparison, k -medoids can be applied to a wide range of data.
- Need to specify k , the number of clusters, in advance.
 - There are ways to automatically determine the best k^3 .
- Sensitive to noisy data and outliers.
- Not suitable to discover clusters with non-convex shapes.

³T. Hastie et al., *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (Springer Series in Statistics). Springer New York, 2009, ISBN: 978-0-387-84857-0.

- **Most of the variants of the k -means differ in:**
 - Selection of the initial k subsets (or centroids).
 - Dissimilarity calculations.
 - Strategies to calculate cluster centroids.
- **Handling categorical data: k -modes:**
 - Replacing centroids with modes.
 - See Chapter 2: mode = value that occurs most frequently in the data.
 - Using new dissimilarity measures to deal with categorical objects.
 - Using a frequency-based method to update modes of clusters.
 - A mixture of categorical and numerical data: k -prototype method.

- The k -means algorithm is **sensitive to outliers!**
 - Since an object with an extremely large value may substantially distort the distribution of the data.
- **k -medoids:**
 - Instead of taking the mean value of the objects in a cluster as a reference point, medoids can be used, which is the most centrally located object in a cluster.

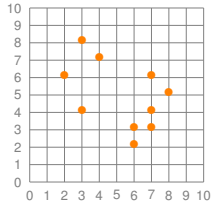


- **k -medoids clustering:**
 - Find representative objects (**medoids**) in clusters.
 - **PAM:**
 - Starts from an initial set of k medoids and iteratively replaces one of the medoids by one of the non-medoids, if it improves the total distance of the resulting clustering.
 - PAM works effectively for small data sets, but does not scale well for large data sets (due to the computational complexity).
- **Efficiency improvement on PAM:**
 - CLARA⁴: PAM on samples.
 - CLARANS⁵: Randomized re-sampling.

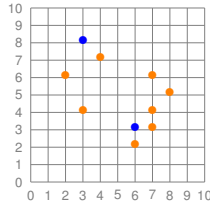
⁴L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, 1990, ISBN: 978-0-47187876-6. DOI: 10.1002/9780470316801.

⁵R. T. Ng and J. Han, "Efficient and effective clustering methods for spatial data mining," in *VLDB'94, Proceedings of 20th International Conference on Very Large Data Bases, September 12-15, 1994, Santiago de Chile, Chile*, J. B. Bocca et al., Eds., Morgan Kaufmann, 1994, pp. 144–155. [Online]. Available: <http://www.vldb.org/conf/1994/P144.PDF>.

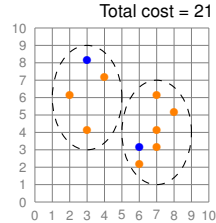
PAM: A Typical k -medoids Algorithm



Arbitrarily
choose k
object as
initial
medoids
for $k = 2$



Assign
each
remaining
object to
nearest medoid

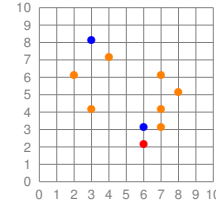


Randomly
select
non-medoid
object o_{random}

Swapping o
and o_{random}
if quality is
improved



Compute
total cost of
swapping

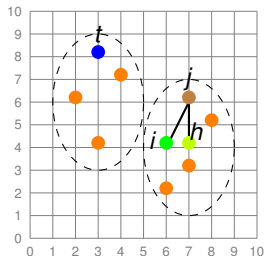


- **Use real objects to represent the clusters.**
- **Algorithm:**
 1. Arbitrarily choose k objects as the initial medoids.
 2. Repeat.
 3. Assign each remaining object to the cluster with the nearest medoid o_i .
 4. Randomly select a non-medoid object o_h .
 5. Compute the total cost TC_{ih} of swapping o_i with o_h .
 6. If $TC_{ih} < 0$ then swap o_i with o_h to form the new set of k medoids.
 7. Until no change.
- $TC_{ih} = \sum_j C_{jih}$
 - with C_{jih} as the cost for object o_j if o_i is swapped with o_h .
 - That is, distance to new medoid minus distance to old medoid.

- **Case 1:**

- o_j currently belongs to medoid o_i . If o_i is replaced with o_h as a medoid, and o_j is closest to o_h , then o_j is reassigned to o_h (same cluster, different distance).

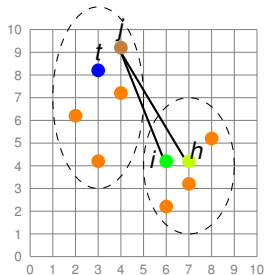
$$C_{jih} = d(o_j, o_h) - d(o_j, o_i).$$



- **Case 2:**

- o_j currently belongs to medoid o_t , $t \neq j$. If o_i is replaced with o_h as a medoid, and o_j is still closest to o_t , then the assignment does not change.

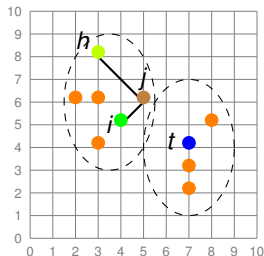
$$C_{jih} = d(o_j, o_h) - d(o_j, o_t) \geq 0.$$



- **Case 3:**

- o_j currently belongs to medoid o_i . If o_i is replaced with o_h as a medoid, and o_j is closest to medoid o_t of one of the other clusters, then o_j is reassigned to o_t (new cluster, different distance).

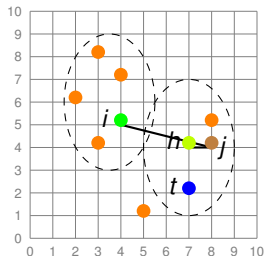
$$C_{jih} = d(o_j, o_t) - d(o_j, o_h) < 0.$$



- **Case 4:**

- o_j currently belongs to medoid o_t , $t \neq j$. If o_i is replaced with o_h as a medoid (from a different cluster!), and o_j is closest to o_h , then o_j is reassigned to o_h (new cluster, different distance).

$$C_{jih} = d(o_j, o_h) - d(o_j, o_t) < 0.$$



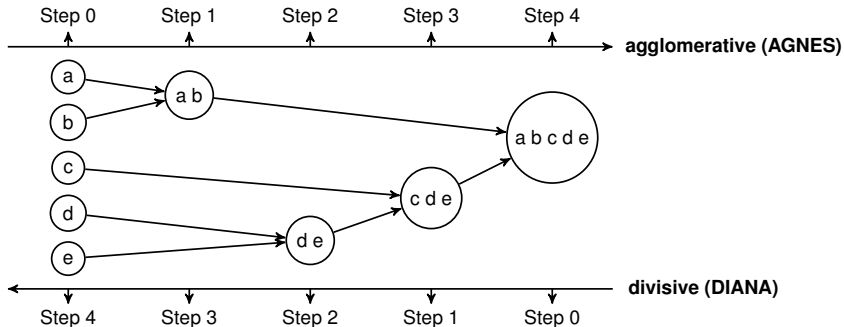
- **Draws multiple samples of the data set, applies PAM on each sample, and gives the best clustering as the output.**
- **Strength:**
 - Deals with larger data sets than PAM.
- **Weakness:**
 - Efficiency depends on the sample size.
 - A good clustering based on samples will not necessarily represent a good clustering of the whole data set if the sample is biased.

- **A Clustering Algorithm based on Randomized Search**
- **Samples:**
 - Drawn dynamically with some randomness in each step of the search.
- **Clustering process:**
 - Can be presented as searching a graph where each node is a potential solution, that is, a set of k medoids.
- **If local optimum found,**
 - start with new randomly selected node in search for a new local optimum.
- **More efficient and scalable than both PAM and CLARA.**
- **Focusing techniques and spatial access structures may further improve its performance.**⁶

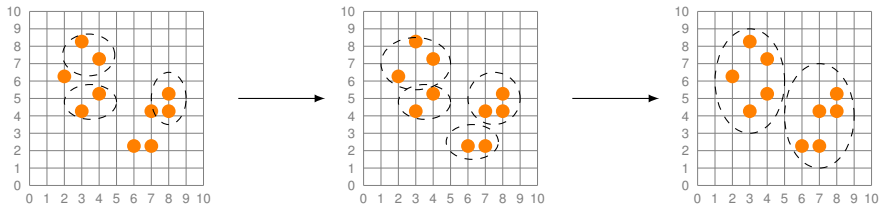
⁶M. Ester et al., "Knowledge discovery in large spatial databases: Focusing techniques for efficient class identification," in *Advances in Spatial Databases, 4th International Symposium, SSD'95, Portland, Maine, USA, August 6-9, 1995. Proceedings*, M. J. Egenhofer and J. R. Herring, Eds., ser. Lecture Notes in Computer Science, vol. 951, Springer, 1995, pp. 67–82.

Hierarchical Methods

- Does not require the number of clusters k as an input, but needs a termination condition.

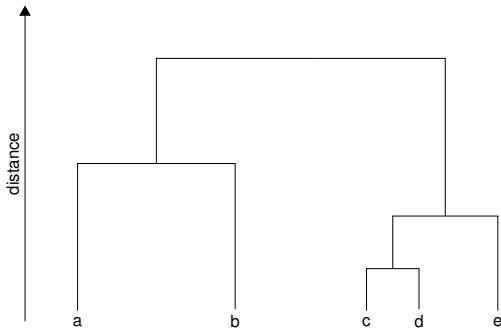


- Introduced by Kaufman et al.⁷
- Use the single-link method. (see below)
- Merge nodes that have the least dissimilarity.
- Go on in a non-descending fashion.
- Eventually all nodes belong to the same cluster.

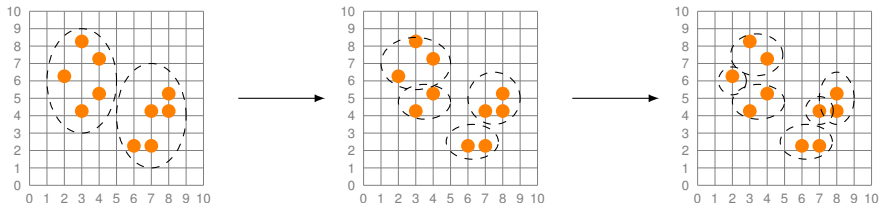


⁷L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, 1990, ISBN: 978-0-47187876-6. DOI: 10.1002/9780470316801.

- Decompose data objects into a several levels of nested partitioning (tree of clusters), called a **dendrogram**.
- A clustering of the data objects is obtained by **cutting** the dendrogram at the desired level, then each connected component forms a cluster.



- Introduced by Kaufman et al.⁸
- Inverse order of AGNES.
- Eventually each node forms a cluster of its own.



⁸L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, 1990, ISBN: 978-0-47187876-6. DOI: 10.1002/9780470316801.

- **Minimum distance:**

- Smallest distance between an object in one cluster and an object in the other, i.e.,
 $\text{dist}_{\min}(C_i, C_j) = \min_{o_{ip} \in C_i, o_{jq} \in C_j} d(o_{ip}, o_{jq}).$

- **Maximum distance:**

- Largest distance between an object in one cluster and an object in the other, i.e.,
 $\text{dist}_{\max}(C_i, C_j) = \max_{o_{ip} \in C_i, o_{jq} \in C_j} d(o_{ip}, o_{jq}).$

- **Average distance:**

- Average distance between an object in one cluster and an object in the other, i.e.,
 $\text{dist}_{\text{avg}}(C_i, C_j) = \frac{1}{n_i \cdot n_j} \sum_{o_{ip} \in C_i, o_{jq} \in C_j} d(o_{ip}, o_{jq}).$

- **Mean distance:**

- Distance between the centroids of two clusters, i.e., $\text{dist}_{\text{mean}}(C_i, C_j) = d(c_i, c_j).$

- **Nearest-neighbor clustering algorithm:**
 - Uses **minimum distance** to measure distance between clusters.
- **Single-linkage algorithm:**
 - Terminates if distance between nearest clusters exceeds user-defined threshold.
- **Minimal spanning-tree algorithm:**
 - View objects (data points) as nodes of a graph.
 - Edges form a path between nodes in a cluster.
 - Merging of two clusters corresponds to adding an edge between the nearest pair of nodes.
 - Because edges linking clusters always go between distinct clusters, resulting graph will be a tree.
 - Thus, agglomerative hierarchical clustering that uses minimum distance produces minimal spanning tree.

- **Farthest-neighbor clustering algorithm:**
 - Uses **maximum distance** to measure distance between clusters.
- **Complete-linkage algorithm:**
 - Terminates if maximum distance between nearest clusters exceeds user-defined threshold.
 - Good if true clusters are rather compact and approx. equal in size.

- **Major weakness of agglomerative clustering methods:**
 - Can never undo what was done previously.
 - Do not scale well: Time complexity of at least $\mathcal{O}(n^2)$, where n is the number of objects.
- **Integration of hierarchical and distance-based clustering:**
 - BIRCH: Uses CF-tree and incrementally adjusts the quality of sub-clusters.
 - CHAMELEON: Hierarchical clustering using dynamic modeling.

Incrementally construct a CF (Clustering Feature) tree:

- A hierarchical data structure for multiphase clustering.
- Phase 1: Scan DB to build an initial in-memory CF-tree.
 - A multi-level compression of the data that tries to preserve the inherent clustering structure of the data.
- Phase 2: Use an arbitrary clustering algorithm to cluster the leaf nodes of the CF-tree.

Scales linearly:

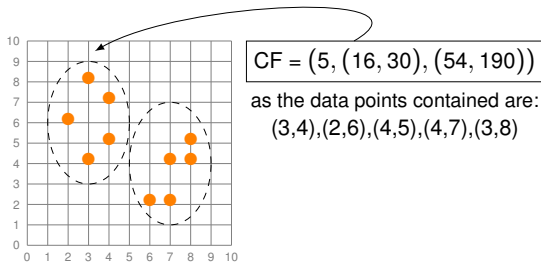
- Finds a good clustering with a single scan and improves the quality with a few additional scans.

Weakness:

- Handles only numerical data, and sensitive to the order of the data records.

Clustering Feature CF = (n, LS, SS):

- 3D vector summarizing statistics about clusters:
 - n : number of data points.
 - LS: linear sum of N points $\sum_{i=1}^n x_i$.
 - SS: square sum of N points $\sum_{i=1}^n x_i^2$.



- **Allows to derive many useful statistics of a cluster. E.g.:**

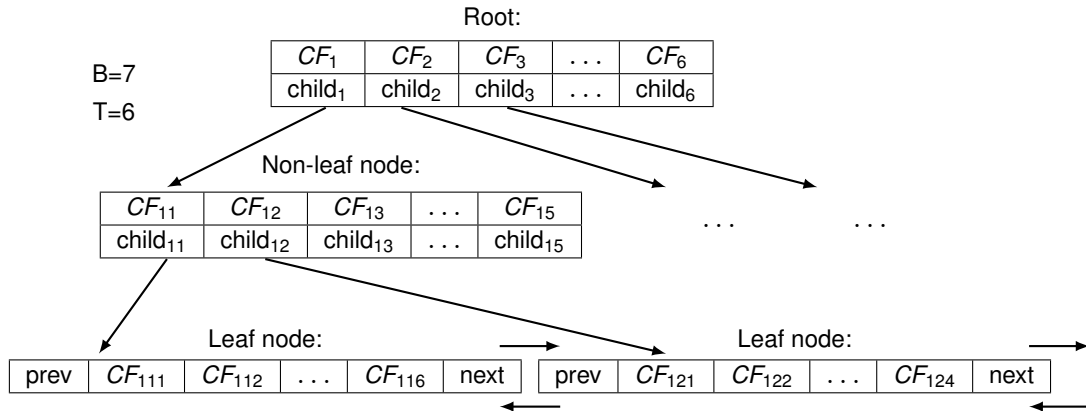
- centroid: $c = \frac{\sum_{i=1}^n o_i}{n} = \frac{LS}{n}$,
- radius: $R = \sqrt{\frac{\sum_{i=1}^n (x_i - c)^2}{n}} = \sqrt{\frac{SS}{n} - \left(\frac{LS}{n}\right)^2}$,
- diameter: $D = \sqrt{\frac{\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2}{n(n-1)}} = \sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}$

- **Additive:**

- For two disjoint clusters C_1 and C_2 with clustering features $CF_1 = (n_1, LS_1, SS_1)$ and $CF_2 = (n_2, LS_2, SS_2)$, the clustering feature of the cluster that is formed by merging C_1 and C_2 is simply: $CF_1 + CF_2 = (n_1 + n_2, LS_1 + LS_2, SS_1 + SS_2)$.

- **Height-balanced tree.**
- **Stores the clustering features for a hierarchical clustering.**
 - Non-leaf nodes store sums of the CFs of their children.
- **Two parameters:**
 - **Branching factor B :** max # of children.
 - **Threshold T :** maximum diameter of sub-clusters stored at leaf nodes.
 - Diameter D : average pairwise distance within a cluster, reflects the tightness of the cluster around the centroid:

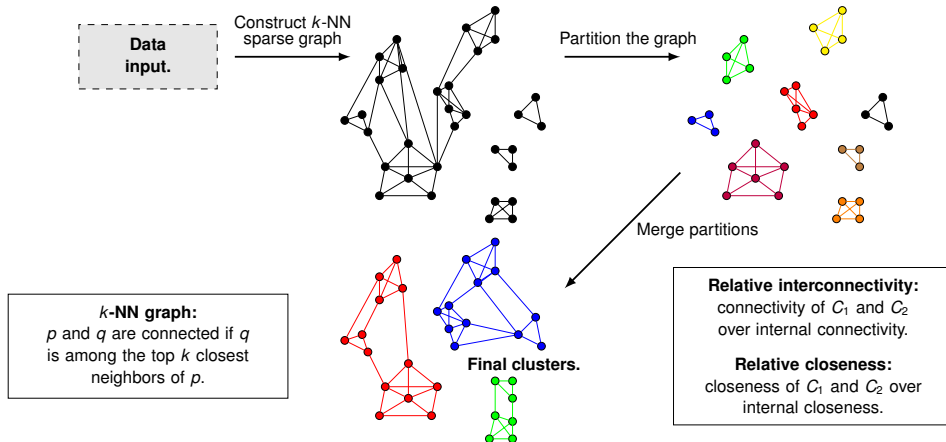
$$D = \sqrt{\frac{\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2}{n(n-1)}} = \sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}.$$



Phase 1:

- **For each point in the input:**
 - Find closest leaf-node entry.
 - Add point to leaf-node entry and update CF.
 - If $entry_diameter > max_diameter$, then split leaf node, and possibly parents.
 - Information about new point is passed toward the root of the tree.
- **Algorithm is $\mathcal{O}(n)$ and incremental.**
- **Concerns:**
 - Sensitive to insertion order of data points.
 - Since we fix the size of leaf nodes, clusters may not be so natural.
 - Clusters tend to be spherical given the radius and diameter measures.

- **Measures the similarity based on a dynamic model:**
 - Two clusters are merged only if the **interconnectivity** and **closeness (proximity)** between two clusters are high relative to the intraconnectivity of the clusters and the closeness of items within the clusters.
- **Graph-based, and a two-phase algorithm.**
 - Use a graph-partitioning algorithm:
 - Cluster objects into a large number of relatively small sub-clusters.
 - Use an agglomerative hierarchical clustering algorithm:
 - Find the genuine clusters by repeatedly combining these sub-clusters.



¹⁰G. Karypis et al., "Chameleon: Hierarchical clustering using dynamic modeling," *IEEE Computer*, vol. 32, no. 8, pp. 68–75, 1999. DOI: 10.1109/2.781637. [Online]. Available: <https://doi.org/10.1109/2.781637>.

Algorithmic hierarchical clustering.

- Nontrivial to choose a good distance measure.
- Hard to handle missing attribute values.
- Optimization goal not clear: heuristic, local search.

Probabilistic hierarchical clustering.

- Use probabilistic models to measure distances between clusters.
- **Generative model:**
 - Regard the set of data objects to be clustered as a sample of the underlying data-generation mechanism to be analyzed.
- Easy to understand, same efficiency as algorithmic agglomerative clustering method, can handle partially observed data.
- In practice, assume the generative models adopt **common distribution functions**, e.g., Gaussian distribution or Bernoulli distribution, governed by parameters.

- Given a set of 1-D points $X = \{x_1, \dots, x_n\}$ for clustering analysis and assuming they are generated by a Gaussian distribution:

$$N(\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mathbf{x} - \mu)^2}{2\sigma^2}\right).$$

- The probability that a point $x_i \in X$ is generated by the model:

$$P(x_i|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right).$$

- The likelihood that X is generated by the model:

$$L(N(\mu, \sigma)|X) := P(X|\mu, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right).$$

- The task of learning the generative model: find the parameters μ and σ , such that

$$N(\mu_0, \sigma_0) = \arg \max (L(N(\mu, \sigma)|X)).$$

- For a set of objects partitioned into m clusters C_1, \dots, C_m , the quality can be measured by:

$$Q(\{C_1, \dots, C_m\}) = \prod_{i=1}^m P(C_i),$$

where $P(C_i)$ is the maximum likelihood of C_i .

- Distance between clusters C_i and C_j can be computed as a proximity measure:

$$d(C_i, C_j) = -\log \frac{P(C_i \cup C_j)}{P(C_i)P(C_j)}.$$

Algorithm: Progressively merge points and clusters

- Input: $D = \{o_1, \dots, o_n\}$: a dataset containing n objects.
- Output: A hierarchy of clusters.
- Method:
 - Create cluster for each object $C_i = \{o_i\}$, for $1 \leq i \leq n$.
 - **For** $i = 1$ **to** n :
 - Find a pair of clusters C_i and C_j , such that
$$C_i, C_j = \arg \max_{i \neq j} \left(\log \left(\frac{P(C_i \cup C_j)}{P(C_i)P(C_j)} \right) \right).$$
 - **If** $\log \left(\frac{P(C_i \cup C_j)}{P(C_i)P(C_j)} \right) > 0$ **then** merge C_i and C_j ,
 - **Else stop.**

Density-based Methods

- **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:**
 - DBSCAN¹¹.
 - OPTICS¹².
 - DENCLUE¹³.
 - CLIQUE¹⁴ (more grid-based).

¹¹M. Ester et al., "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, Portland, Oregon, USA, E. Simoudis et al., Eds., AAAI Press, 1996, pp. 226–231.

¹²M. Ankerst et al., "OPTICS: ordering points to identify the clustering structure," in *SIGMOD 1999, Proceedings ACM SIGMOD International Conference on Management of Data, June 1-3, 1999, Philadelphia, Pennsylvania, USA*, A. Delis et al., Eds., ACM Press, 1999, pp. 49–60.

¹³A. Hinneburg and D. A. Keim, "An efficient approach to clustering in large multimedia databases with noise," in *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining (KDD-98)*, New York City, New York, USA, August 27-31, 1998, R. Agrawal et al., Eds., AAAI Press, 1998, pp. 58–65.

¹⁴R. Agrawal et al., "Automatic subspace clustering of high dimensional data for data mining applications," in *SIGMOD 1998, Proceedings ACM SIGMOD International Conference on Management of Data, June 2-4, 1998, Seattle, Washington, USA*, L. M. Haas and A. Tiwary, Eds., ACM Press, 1998, pp. 94–105.

- **Two parameters:**

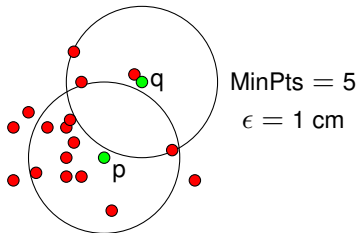
- ϵ : Maximum radius of a neighborhood.

- Defines the neighborhood of a point p : $N_\epsilon(p) := \{q \in D \mid d(p, q) \leq \epsilon\}$.
- Distance function is still needed.

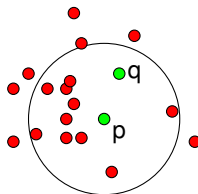
- **MinPts**: Minimum number of points in an ϵ -neighborhood of a point p .

- **Core-point condition:** $|N_\epsilon(p)| \geq \text{MinPts}$.

- Only these neighborhoods are considered.



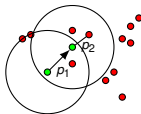
- A point q is said to be directly density-reachable from a core point p w.r.t. ϵ , MinPts, if q belongs to $N_\epsilon(p)$.



MinPts = 5
 $\epsilon = 1 \text{ cm}$

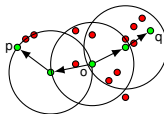
- **Density-reachable:**

- A point q is density-reachable from a point p w.r.t. ϵ , MinPts, if there is a **chain of (core) points** p_1, \dots, p_n , $p_1 = p$, $p_n = q$ such that p_{i+1} is directly density-reachable from p_i .

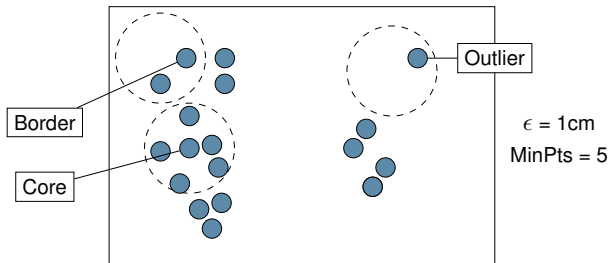


- **Density-connected:**

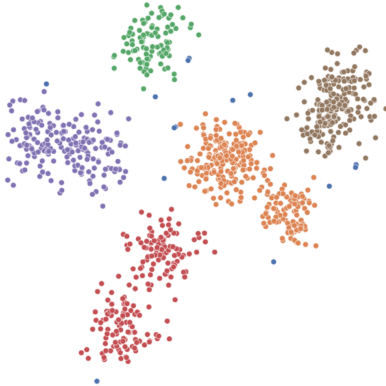
- A point p is density-connected to a point q w.r.t. ϵ , MinPts, if there is a (core) point o such that both p and q are density-reachable from o .



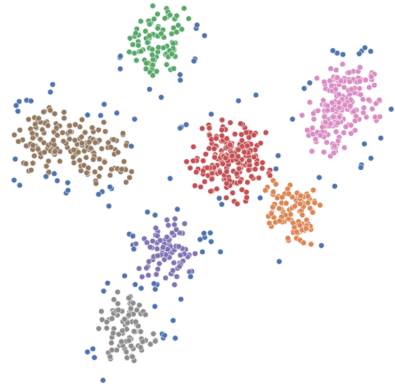
- **Relies on a density-based notion of cluster:**
 - A cluster is defined as a maximal set of density-connected points.
- **Discovers clusters of arbitrary shape in spatial databases with noise.**



- All objects in database D are marked unvisited.
- Randomly select unvisited object p and mark as visited.
- If ϵ -neighborhood of p contains less than MinPts objects:
 - Mark p as noise.
- Otherwise (that is, p is core point):
 - Create new cluster C for point p .
 - Add all objects in ϵ -neighborhood of p to candidate set N .
 - For each p' in N that does not yet belong to a cluster:
 - Add p' to C .
 - If p' is unvisited, mark as visited.
 - If p' is core point, add all objects in its ϵ -neighborhood to N .
 - Ends when N is empty, that is, C can no longer be expanded.
- Continue the process (randomly select next unvisited object in D) until all points have been visited.



$\epsilon = 1$ and $MinPts = 5$

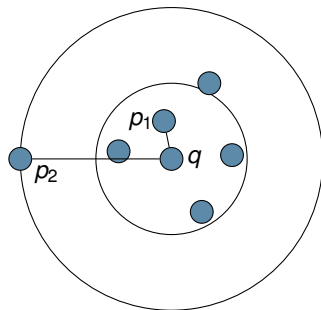


$\epsilon = 0.7$ and $MinPts = 10$

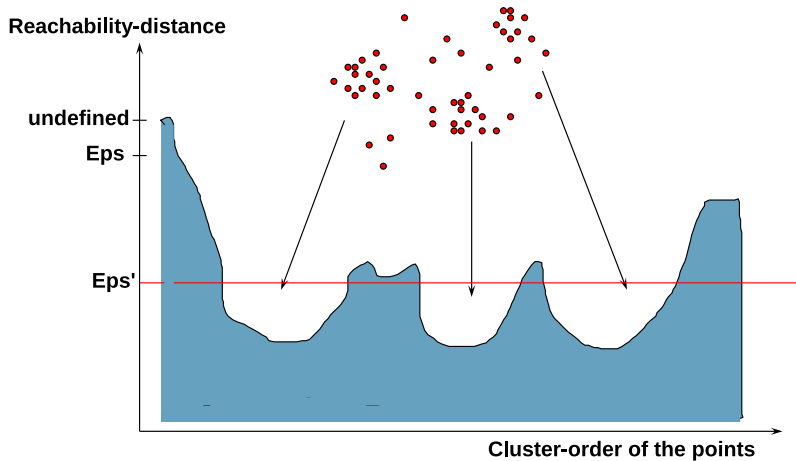
- **OPTICS: Ordering Points To Identify the Clustering Structure.**¹⁵
 - Avoid difficulty of finding the right parameters:
 - Only MinPts must be given, ϵ remains open.
 - Produces a special **order of the database** w.r.t. its density-based clustering structure.
 - This cluster-ordering contains info equivalent to the density-based 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 or using visualization techniques.
 - Process point in order: Select a point that is density-reachable w.r.t. the lowest ϵ value, so that clusters with higher density (lower ϵ) will be finished first.

¹⁵M. Ankerst et al., "OPTICS: ordering points to identify the clustering structure," in *SIGMOD 1999, Proceedings ACM SIGMOD International Conference on Management of Data, June 1-3, 1999, Philadelphia, Pennsylvania, USA*, A. Delis et al., Eds., ACM Press, 1999, pp. 49–60.

- **Core distance:**
 - of a point p .
 - Min ϵ s.t. p is core point.
- **Reachability Distance r :**
 - of a point p from q .
 - Minimum radius value that makes p directly density-reachable from q .
 - That is, $r(q, p) := \max\{\text{core-distance}(q), d(q, p)\}$.
- **Example:**
 - $\text{MinPts} = 5$, $\text{core_distance}(q) = 3\text{cm}$.
 - $d(q, p_1) = 2.8\text{cm} \implies r(q, p_1) = 3\text{cm}$.
 - $d(q, p_2) = 4\text{cm} \implies r(q, p_2) = 4\text{cm}$.



- **Maintains a list of OrderSeeds:**
 - Points sorted by **reachability-distance** from their resp. closest core points.
- **Begin with arbitrary point from input DB.**
- **For each point p under consideration:**
 - Retrieve the closest MinPts points of p , determine core-distance, set reachability-distance to undefined.
 - Write p to output.
 - If p is core point, for each point q in the ϵ -neighborhood of p :
 - Update reachability-distance from p .
 - Insert q into OrderSeeds (if q has not yet been processed).
 - Move to next object in OrderSeeds list with smallest reachability-distance (or input DB if OrderSeeds is empty).
- **Continue until input DB fully consumed (and OrderSeeds empty).**



- **DENsity-based CLUstEring**¹⁶

- **Using statistical density functions:**

- $f_{\text{Gaussian}}(x, y) = \exp\left(\frac{d(x, y)^2}{2\sigma^2}\right)$, influence of y on x .
- $f_{\text{Gaussian}}^D(x, x_i) = \sum_{i=1}^N \exp\left(\frac{d(x, x_i)^2}{2\sigma^2}\right)$, total influence on x .
- $\nabla_{x_i} f_{\text{Gaussian}}^D(x, x_i) = \sum_{i=1}^N (x_i - x) \cdot \exp\left(\frac{d(x, x_i)^2}{2\sigma^2}\right)$, gradient of x in direction x_i .

- **Major features:**

- Solid mathematical foundation.
- Good for data sets with large amounts of noise.
- Allows a compact mathematical description of arbitrarily shaped clusters in high-dimensional data sets.
- Significantly faster than existing algorithms (e.g., DBSCAN).
- But needs a large number of parameters.

¹⁶A. Hinneburg and D. A. Keim, "An efficient approach to clustering in large multimedia databases with noise," in *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining (KDD-98)*, New York City, New York, USA, August 27-31, 1998, R. Agrawal et al., Eds., AAAI Press, 1998, pp. 58–65.

- **Density estimation:**
 - Estimation of an unobservable underlying probability density function based on a set of observed data.
 - Regarded as a sample.
- **Kernel density estimation:**
 - Treat an observed object as an indicator of high-probability density in the surrounding region.
 - Probability density at a point depends on the distances from this point to the observed objects.
- **Density attractor:**
 - Local maximum of the estimated density function.
 - Must be above threshold ξ .

- **Clusters:**
 - Can be determined mathematically by identifying density attractors.
 - That is, local maxima of the overall density function.
- **Center-defined clusters:**
 - Assign to each density attractor the points density-attracted to it.
- **Arbitrary shaped cluster:**
 - Merge density attractors that are connected through paths of high density ($>$ threshold).

Grid-based Methods

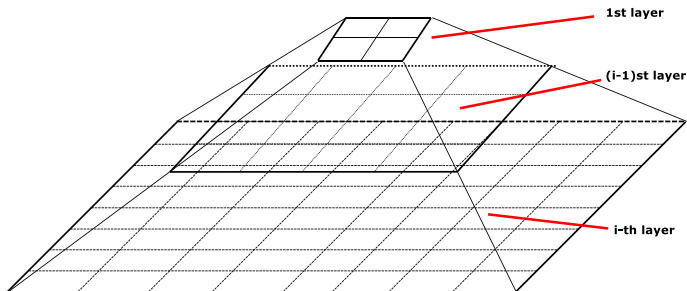
- Using multi-resolution grid data structure.
- Several interesting methods.
 - STING (a STatistical INformation Grid approach)¹⁷.
 - WaveCluster¹⁸:
 - A multi-resolution clustering approach using wavelet method.
 - Not part of this lecture.
 - CLIQUE¹⁹:
 - Both grid-based and subspace clustering.

¹⁷W. Wang et al., "STING: A statistical information grid approach to spatial data mining," in *VLDB'97, Proceedings of 23rd International Conference on Very Large Data Bases, August 25-29, 1997, Athens, Greece*, M. Jarke et al., Eds., Morgan Kaufmann, 1997, pp. 186–195. [Online]. Available: <http://www.vldb.org/conf/1997/P186.PDF>.

¹⁸G. Sheikholeslami et al., "Wavecluster: A multi-resolution clustering approach for very large spatial databases," in *VLDB'98, Proceedings of 24rd International Conference on Very Large Data Bases, August 24-27, 1998, New York City, New York, USA*, A. Gupta et al., Eds., Morgan Kaufmann, 1998, pp. 428–439. [Online]. Available: <http://www.vldb.org/conf/1998/p428.pdf>.

¹⁹R. Agrawal et al., "Automatic subspace clustering of high dimensional data for data mining applications," in *SIGMOD 1998, Proceedings ACM SIGMOD International Conference on Management of Data, June 2-4, 1998, Seattle, Washington, USA*, L. M. Haas and A. Tiwary, Eds., ACM Press, 1998, pp. 94–105.

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



- **Each cell at a higher level:**
 - Partitioned into a number of smaller cells at next lower level.
- **Statistical info of each cell:**
 - Calculated and stored beforehand and used to answer queries.
 - Count, plus for each attribute: mean, standard dev, min, max and type of distribution: normal, uniform, etc.
- **Parameters of higher-level cells:**
 - Can easily be calculated from parameters of lower-level cells.
- **Use a top-down approach:**
 - To answer spatial data queries (or: cluster definitions).
- **Start from a pre-selected layer:**
 - Typically with a small number of cells.
- **For each cell at the current level:**
 - Compute the confidence interval.

- **Cells labeled relevant or not relevant.**
 - At the specified confidence level.
 - Irrelevant cells removed from further consideration.
- **When finished examining the current layer, proceed to the next lower level.**
 - Only look at cells that are children of relevant cells.
- **Repeat this process until bottom layer is reached.**
- **Find regions (clusters) that satisfy the density specified.**
 - Breadth-first search, at bottom layer.
 - Examine cells within a certain distance from center of current cell (often just the neighbors).
 - If average density within this small area is greater than density specified, mark area and put relevant cells just examined into queue.
 - Examine next cell from queue and repeat procedure, until end of queue.
 - Then one region has been identified.

- **Advantages:**

- Query-independent, easy to parallelize, incremental update.
- $\mathcal{O}(K)$, 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.

- Using **subspaces** (lower-dimensional) of a high-dimensional data space that allow better clustering than original space
- **CLIQUE can be considered as both density-based and grid-based.**
 - Partitions each dimension into non-overlapping intervals, thereby partitioning the entire data space into **cells**.
 - Uses density threshold to identify **dense** cells.
 - A cell is dense, if the number of data points mapped to it exceeds the density threshold.

- **Monotonicity of dense cells w.r.t. dimensionality:**

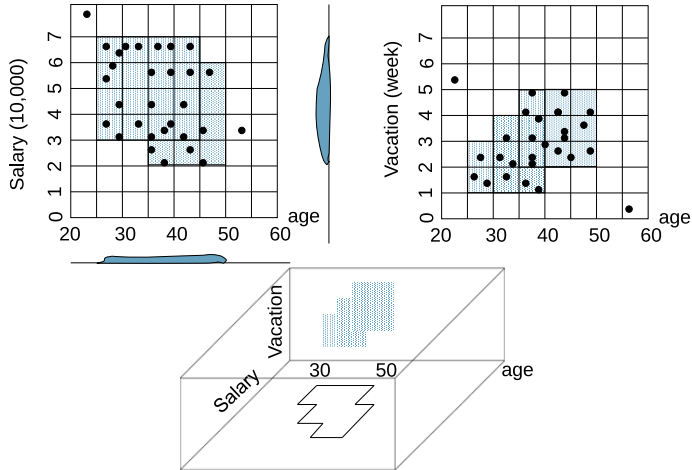
- Based on the a priori property used in frequent-pattern and association-rule mining.
- k -dimensional cell c can have at least l points only, if every $(k - 1)$ -dimensional projection of c (which is a cell in a $(k - 1)$ -dimensional subspace) has at least l points, too.

- **Clustering step 1:**

- Partition each dimension into intervals, identify intervals containing at least l points.
- Iteratively join k -dimensional dense cells c_1 and c_2 in subspaces (D_{i1}, \dots, D_{ik}) and (D_{j1}, \dots, D_{jk}) with $D_{i1} = D_{j1}$ and $D_{i2} = D_{j2}$ and $\dots D_{i(k-1)} = D_{j(k-1)}$ and c_1 and c_2 share the same intervals to those dimensions.
- New $(k + 1)$ -dimensional candidate cell c in space $(D_{i1}, \dots, D_{ik}, D_{jk})$ tested for density.
- Iteration terminates when no more candidate cells can be generated or no candidate cells are dense.

- **Clustering step 2:**

- Use dense cells in each subspace to assemble clusters.
- Apply Minimum Description Length (MDL) principle to use the maximal regions to cover connected dense cells.
 - **Maximal region:** hyper rectangle where every cell falling into the regions is dense, and region cannot be extended further in any dimension.
- Simple greedy approach:
 - Start with arbitrary dense cell.
 - Find maximum region covering that cell.
 - Work on remaining dense cells that have not yet been covered.



- **Strength:**

- Automatically finds subspaces of the highest dimensionality such that high-density clusters exist in those subspaces.
- Insensitive to the order of records in input and does not presume any canonical data distribution.
- Scales linearly with the size of input and has good scalability as the number of dimensions in the data is increased.

- **Weaknesses:**

- Dependent on proper grid size and density threshold.
- Accuracy of clustering result may be degraded at the expense of simplicity of the method.

Evaluation of Clustering

- **Assess if non-random structure** exists in the data by measuring the probability that the data is generated by a uniform data distribution.
- **Data with random structure:**
 - Points uniformly distributed in data space.
- **Clustering may return clusters, but:**
 - Artificial partitioning.
 - Meaningless.

- **Data set:**

- Let $X = \{x_i \mid i = 1, \dots, n\}$ be a collection of n patterns in a d -dimensional space such that $x_i = (x_{i1}, x_{i2}, \dots, x_{id})$.

- **Random sample of data space:**

- Let $Y = \{y_j \mid j = 1, \dots, m\}$ be m sampling points placed at random in the d -dimensional space, with $m \ll n$.

- **Two types of distances defined:**

- u_j as minimum distance from y_j to its nearest pattern in X and
- w_j as minimum distance from a randomly selected pattern in X to its nearest neighbor in X .

- **The Hopkins statistic in d dimensions is defined as:**

$$H = \frac{\sum_{j=1}^m u_j}{\sum_{j=1}^m u_j + \sum_{j=1}^m w_j}.$$

- Compares nearest-neighbor distribution of randomly selected locations (points) to that for randomly selected patterns.
- Under the null hypothesis, H_0 , of uniform distribution:
 - Distances from sampling points to nearest patterns should, on the average, **be the same** as the interpattern nearest-neighbor distances, implying randomness.
 - H should be about 0.5.
- When patterns are aggregated or clustered:
 - Distances from sampling points to nearest patterns should, on the average, **larger** as the interpattern nearest-neighbor distances.
 - H should be larger than 0.5.
 - Almost equal to 1.0 for very well clustered data.

- **Empirical method:**

- # of clusters $\approx \sqrt{\frac{n}{2}}$ for a dataset of n points.

- **Elbow method:**

- Use the turning point in the curve of sum of within-cluster variance w.r.t. the # of clusters.

- **Cross-validation method:**

- Divide a given data set into m parts.
- Use $m - 1$ parts to obtain a clustering model.
- Use the remaining part to test the quality of the clustering.
E.g., for each point in the test set, find the closest centroid, and use the sum of squared distances between all points in the test set and the closest centroids to measure how well the model fits the test set.
- For any $k > 0$, repeat it m times, compare the overall quality measure w.r.t. different k 's, and find # of clusters that fits the data the best.

Two methods:

1. **Extrinsic: supervised, i.e., the ground truth is available.**
 - Compare a clustering against the ground truth using certain clustering quality measure.
 - Ex. BCubed precision and recall metrics.
2. **Intrinsic: unsupervised, i.e., the ground truth is unavailable.**
 - Evaluate the goodness of a clustering by considering how well the clusters are separated, and how compact the clusters are.
 - Ex. silhouette coefficient.

- **Clustering-quality measure:** $Q(C, C_*)$.
 - For a clustering C given the ground truth C_* .
- **Q is good, if it satisfies the following four essential criteria:**
 - Cluster homogeneity:
 - The purer, the better.
 - Cluster completeness:
 - Should assign objects that belong to the same category in the ground truth to the same cluster.
 - Rag bag:
 - Putting a heterogeneous object into a pure cluster should be penalized more than putting it into a rag bag (i.e., "miscellaneous" or "other" category).
 - Small cluster preservation:
 - Splitting a small category into pieces is more harmful than splitting a large category into pieces.


Summary

- **Cluster analysis:**
 - Groups objects based on their similarity and has wide applications.
- **Measure of similarity:**
 - Can be computed for various types of data.
- **Clustering algorithms can be categorized into:**
 - Partitioning methods (k -means and k -medoids).
 - Hierarchical methods (BIRCH and CHAMELEON; probabilistic hierarchical clustering).
 - Density-based methods (DBSCAN, OPTICS, and DENCLUE).
 - Grid-based methods (STING, CLIQUE).
 - Model-based methods.
- **Quality of clustering results.**

Any questions about this chapter?

Ask them now or ask them later in our forum:



 https://www.studon.fau.de/studon/goto.php?target=lcodes_OLYeD79h