


Data Mining Fundamentals

Chapter 10. Cluster Analysis: Basic Concepts and Methods

Chapter 10. Cluster Analysis: Basic Concepts and Methods

- ❑ Cluster Analysis: An Introduction
- ❑ Partitioning Methods
- ❑ Hierarchical Methods 
- ❑ Density- and Grid-Based Methods
- ❑ Evaluation of Clustering
- ❑ Summary

Hierarchical Clustering Methods

- ❑ Basic Concepts of Hierarchical Algorithms
- ❑ Agglomerative Clustering Algorithms
- ❑ Divisive Clustering Algorithms
- ❑ Extensions to Hierarchical Clustering
- ❑ BIRCH: A Micro-Clustering-Based Approach
- ❑ CURE: Exploring Well-Scattered Representative Points
- ❑ CHAMELEON: Graph Partitioning on the KNN Graph of the Data
- ❑ Probabilistic Hierarchical Clustering

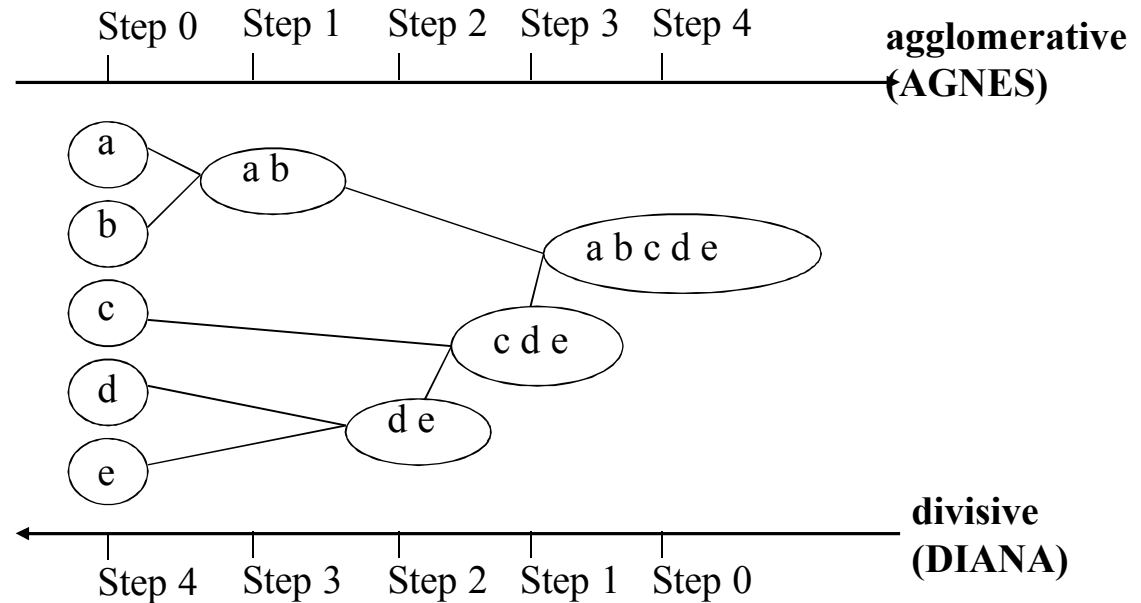
Hierarchical Clustering: Basic Concepts

□ Hierarchical clustering

- Generate a clustering hierarchy (drawn as a **dendrogram**)
- Not required to specify **K**, the number of clusters
- More deterministic
- No iterative refinement

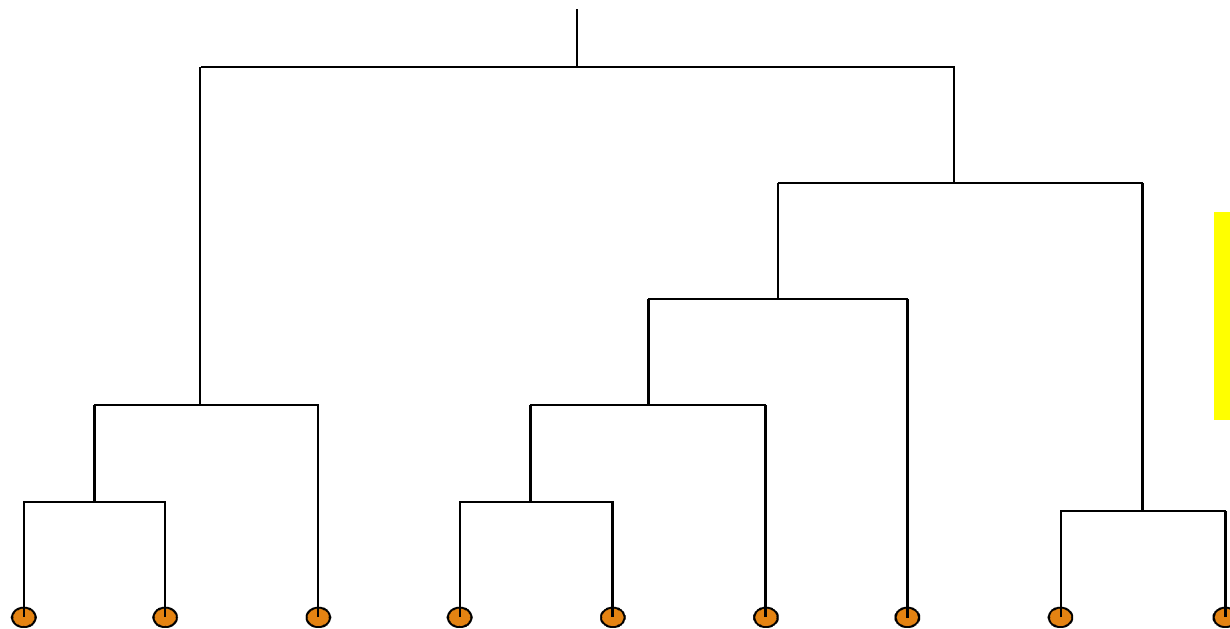
□ Two categories of algorithms:

- **Agglomerative**: Start with singleton clusters, continuously merge two clusters at a time to build a **bottom-up** hierarchy of clusters
- **Divisive**: Start with a huge macro-cluster, split it continuously into two groups, generating a **top-down** hierarchy of clusters



Dendrogram: Shows How Clusters are Merged

- ❑ Dendrogram: Decompose a set of data objects into a tree of clusters by multi-level nested partitioning
- ❑ A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster

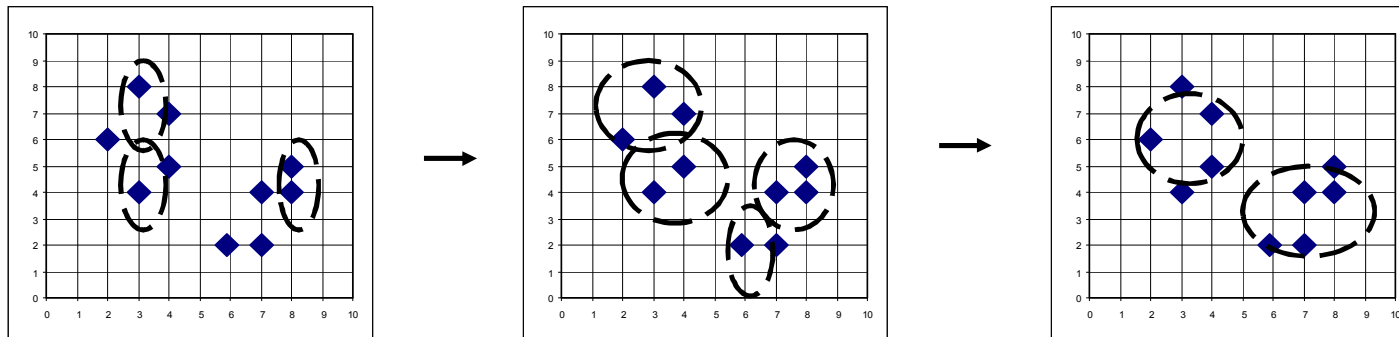


Hierarchical clustering
generates a dendrogram
(a hierarchy of clusters)

Agglomerative Clustering Algorithm

- AGNES (AGglomerative NESTing) (Kaufmann and Rousseeuw, 1990)

- Use the **single-link** method and the dissimilarity matrix
- Continuously merge nodes that have the least dissimilarity
- Eventually all nodes belong to the same cluster



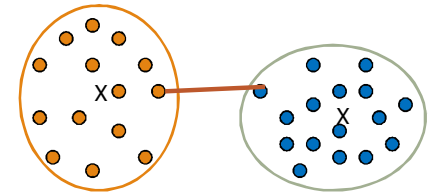
- Agglomerative clustering varies on different similarity measures among clusters

- Single link (nearest neighbor)
- Complete link (diameter)
- Average link (group average)
- Centroid link (centroid similarity)

Single Link vs. Complete Link in Hierarchical Clustering

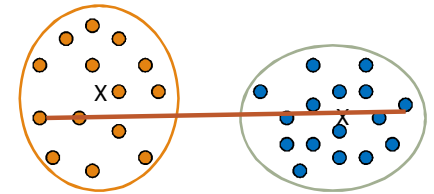
□ Single link (nearest neighbor)

- The similarity between two clusters is the similarity between their most similar (nearest neighbor) members
- Local similarity-based: Emphasizing more on close regions, ignoring the overall structure of the cluster
- Capable of clustering non-elliptical shaped group of objects
- Sensitive to noise and outliers



□ Complete link (diameter)

- The similarity between two clusters is the similarity between their most dissimilar members
- Merge two clusters to form one with the smallest diameter
- Nonlocal in behavior, obtaining compact shaped clusters
- Sensitive to outliers



Agglomerative Clustering: Average vs. Centroid Links

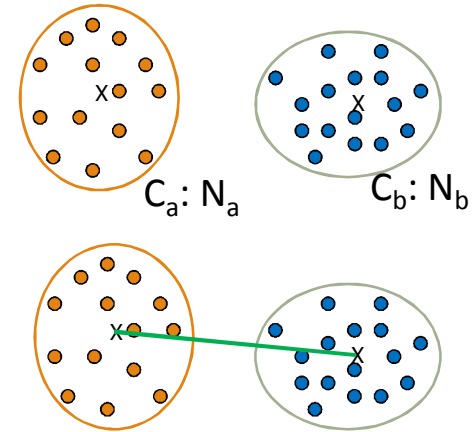
- Agglomerative clustering with **average link**

- **Average link:** The average distance between an element in one cluster and an element in the other (i.e., all pairs in two clusters)

- Expensive to compute

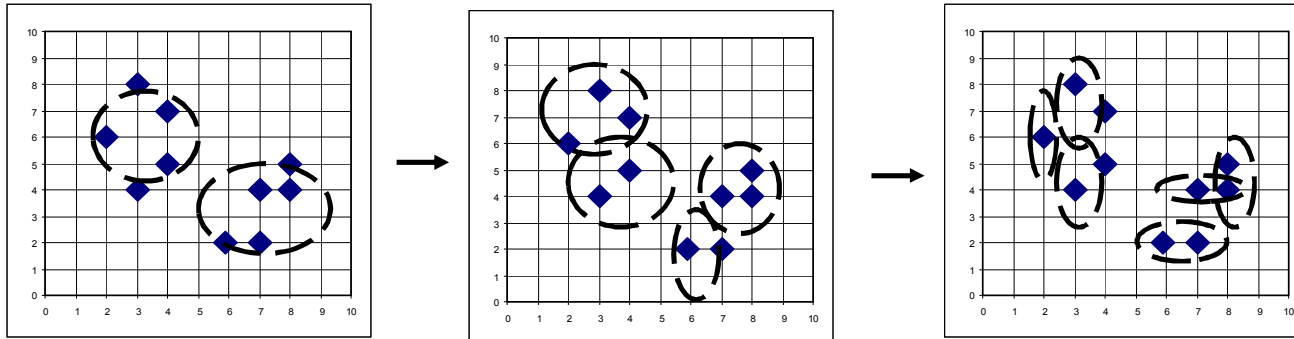
- Agglomerative clustering with **centroid link**

- **Centroid link:** The distance between the centroids of two clusters



Divisive Clustering

- DIANA (Divisive Analysis) (Kaufmann and Rousseeuw, 1990)
 - Implemented in some statistical analysis packages, e.g., Splus
- Inverse order of AGNES: Eventually each node forms a cluster on its own



- Divisive clustering is a top-down approach
 - The process starts at the root with all the points as one cluster
 - It recursively splits the higher level clusters to build the dendrogram
 - Can be considered as a global approach
 - More efficient when compared with agglomerative clustering

More on Algorithm Design for Divisive Clustering

- Choosing which cluster to split
 - Check the sums of squared errors of the clusters and choose the one with the largest value
- Splitting criterion: Determining how to split
 - For categorical data, Gini-index can be used
- Handling the noise
 - Use a threshold to determine the termination criterion (do not generate clusters that are too small because they contain mainly noises)

Extensions to Hierarchical Clustering

- ❑ Major weaknesses of hierarchical clustering methods
 - ❑ Can never undo what was done previously
 - ❑ Do not scale well
 - ❑ Time complexity of at least $O(n^2)$, where n is the number of total objects
- ❑ Other hierarchical clustering algorithms
 - ❑ BIRCH (1996): Use CF-tree and incrementally adjust the quality of sub-clusters
 - ❑ CURE (1998): Represent a cluster using a set of well-scattered representative points
 - ❑ CHAMELEON (1999): Use graph partitioning methods on the K-nearest neighbor graph of the data

BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies)

- ❑ A multiphase clustering algorithm (Zhang, Ramakrishnan & Livny, SIGMOD'96)
- ❑ 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
- ❑ Key idea: Multi-level clustering
 - ❑ Low-level micro-clustering: Reduce complexity and increase scalability
 - ❑ High-level macro-clustering: Leave enough flexibility for high-level clustering
- ❑ *Scales linearly*: Find a good clustering with a single scan and improve the quality with a few additional scans

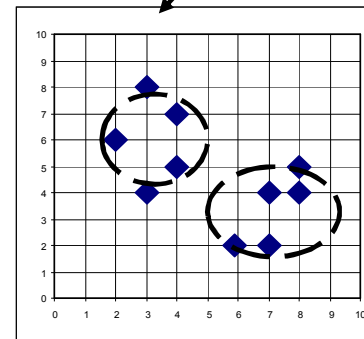
Clustering Feature Vector in BIRCH

□ Clustering Feature (CF): $CF = (N, LS, SS)$

□ N : Number of data points

□ LS : linear sum of N points: $\sum_{i=1}^N X_i$

□ SS : square sum of N points: $\sum_{i=1}^N X_i^2$



$CF = (5, (16,30),(54,190))$

(3,4)

(2,6)

(4,5)

(4,7)

(3,8)

□ Clustering feature:

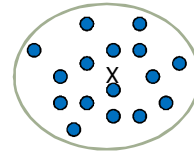
□ Summary of the statistics for a given sub-cluster: the 0-th, 1st, and 2nd moments of the sub-cluster from the statistical point of view

□ Registers crucial measurements for computing cluster and utilizes storage efficiently

Measures of Cluster: Centroid, Radius and Diameter

□ Centroid: \vec{x}_0

- the “middle” of a cluster
- n : number of points in a cluster
- \vec{x}_i is the i -th point in the cluster



$$\vec{x}_0 = \frac{\sum_i^n \vec{x}_i}{n} = \frac{LS}{N}$$

□ Radius: R

- Average distance from member objects to the centroid
- The square root of average distance from any point of the cluster to its centroid

$$R = \sqrt{\frac{\sum_i^n (\vec{x}_i - \vec{x}_0)^2}{n}} = \sqrt{\frac{nSS - 2LS^2 + nLS}{n^2}}$$

□ Diameter: D

- Average pairwise distance within a cluster
- The square root of average mean squared distance between all pairs of points in the cluster

$$D = \sqrt{\frac{\sum_i^n \sum_j^n (\vec{x}_i - \vec{x}_j)^2}{n(n-1)}} = \sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}$$

The CF Tree Structure in BIRCH

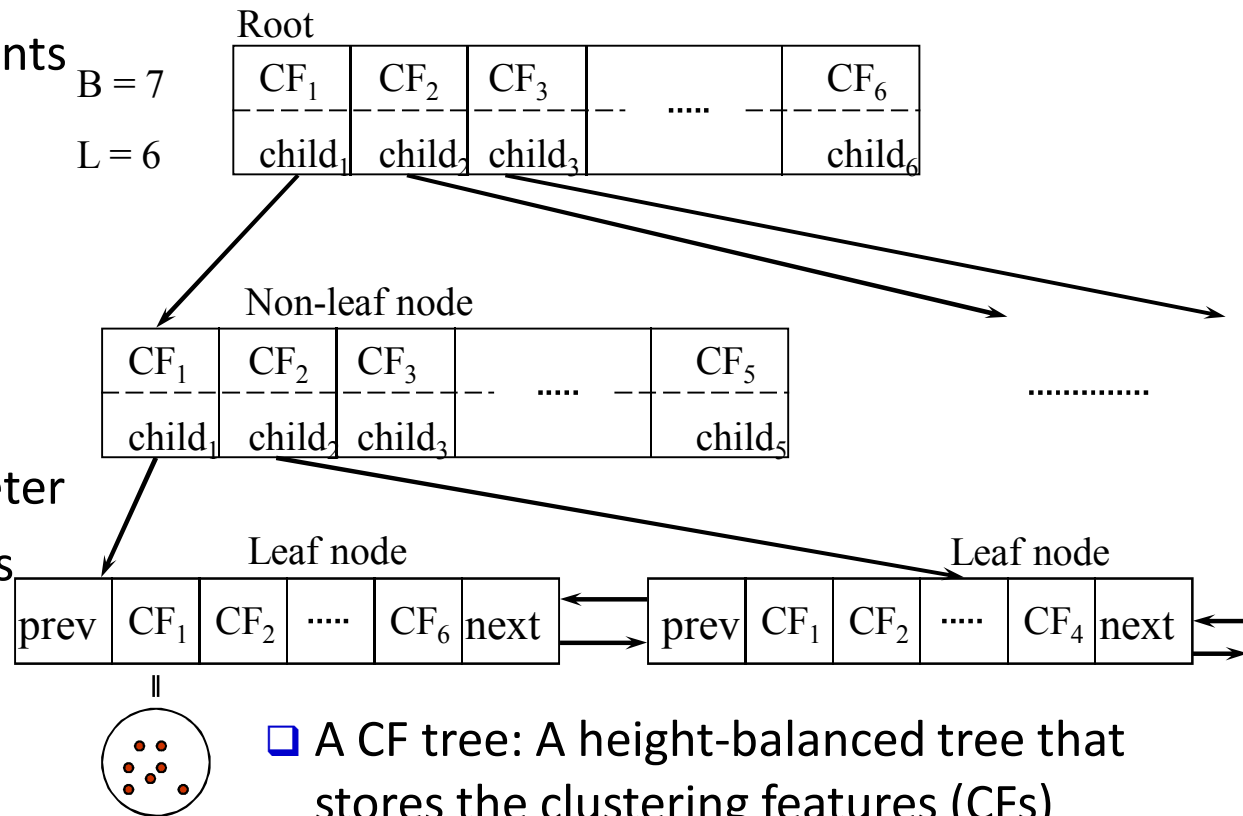
- Incremental insertion of new points (similar to B+-tree)

- For each point in the input

- Find closest leaf entry
- Add point to leaf entry and update CF
- If entry diameter $>$ max_diameter
 - split leaf, and possibly parents

- A CF tree has two parameters

- Branching factor: Maximum number of children
- Maximum diameter of sub-clusters stored at the leaf nodes



- A CF tree: A height-balanced tree that stores the clustering features (CFs)
- The non-leaf nodes store sums of the CFs of their children

BIRCH: A Scalable and Flexible Clustering Method

- ❑ An integration of agglomerative clustering with other (flexible) clustering methods
 - ❑ Low-level micro-clustering
 - ❑ Exploring CP-feature and BIRCH tree structure
 - ❑ Preserving the inherent clustering structure of the data
 - ❑ Higher-level macro-clustering
 - ❑ Provide sufficient flexibility for integration with other clustering methods
- ❑ Impact to many other clustering methods and applications
- ❑ Concerns
 - ❑ Sensitive to insertion order of data points
 - ❑ Due to the fixed size of leaf nodes, clusters may not be so natural
 - ❑ Clusters tend to be spherical given the radius and diameter measures