EuroVis 2016 Machine Learning Methods in Visualization for Big Data 6 June 2016, Groningen, the Netherlands

Hierarchical Clustering

Jaakko Peltonen^{1,2}

¹Aalto University, Department of Computer Science ²University of Tampere, School of Information Sciences

A single partitioning clustering - not enough?

A single partitioning clustering result e.g. from kmeans is just one of many clusterings for the data.

- How many clusters are needed? Maybe more than one level of detail is needed!
- One could run several clusterings (e.g. k-means) at different numbers of clusters - but results have no clear relationship to each other.
- **Hierarchical clustering** builds a tree of clusterings with changing number of clusters (detail level). Each level arises from the next by a simple **merge/split**.

Machine Learning Methods in Visualization for Big Data

Agglomerative hierarchical clustering

The most common hierarchical clustering is **agglomerative hierarchical clustering**.

- In a bottom-up strategy, at first each data item is its own cluster.
- Iteratively, two "closest" clusters are merged, until only one cluster remains.
- The reverse is called divisive clustering.

Agglomerative hierarchical clustering

Three types of agglomerative clustering:

• Single-linkage: distance of two clusters is the minimum between their members

$$d(c_{1,}c_{2}) = \min_{i \in c_{1}} \min_{j \in c_{2}} d(x_{i}, x_{j})$$

• Average-linkage (aka UPGMA): distance of two clusters is the average between their members

$$d(c_1, c_2) = E_{i \in c_1}[E_{j \in c_2}[d(x_i, x_j)]]$$

• Complete-linkage: distance of two clusters is the maximum between their members

$$d(c_{1,}c_{2}) = \max_{i \in c_{1}} \max_{j \in c_{2}} d(x_{i}, x_{j})$$

Machine Learning Methods in Visualization for Big Data

Agglomerative hierarchical clustering

Pseudocode: compute distances between all clusters, merge the two closest ones, repeat.

In each iteration, the number of active clusters decreases by 1.

A new data item can be assigned to a cluster by the same rules as the merging rules (minimum/average/maximum distance).

Example 50 clusters

average linkage



singlelinkage





Example 49 clusters

(merge in top left corner)

average linkage



singlelinkage





Example 48 clusters

(merge below middle)

average linkage



singlelinkage





Example 47 clusters

(merge at right)

average linkage



singlelinkage





Example 46 clusters

average linkage



singlelinkage





Example 45 clusters

average linkage



singlelinkage





Example 44 clusters

average linkage



singlelinkage





Example 43 clusters

average linkage



singlelinkage





Example 42 clusters

average linkage



singlelinkage





Example 41 clusters

average linkage



singlelinkage





Example 40 clusters

average linkage



singlelinkage





Example 39 clusters

average linkage



singlelinkage





Example 38 clusters

average linkage



singlelinkage





Example 37 clusters

average linkage



singlelinkage





Example 36 clusters

average linkage



singlelinkage





Example 35 clusters

average linkage



singlelinkage





Example 34 clusters

average linkage



singlelinkage





Example 33 clusters

average linkage



singlelinkage





Example 32 clusters

single-linkage has nonconvex Voronoi region near middle

average linkage



singlelinkage





Example 31 clusters

single-linkage has nonconvex Voronoi region near middle

average linkage



singlelinkage





Example 30 clusters

average linkage



singlelinkage





Example 29 clusters

average linkage



singlelinkage





Example 28 clusters

average linkage



singlelinkage





Example 27 clusters

average linkage



singlelinkage





Example 26 clusters

average linkage



singlelinkage





Example 25 clusters

average linkage



singlelinkage





Example 24 clusters

average linkage



singlelinkage





Example 23 clusters

subtrees within clusters start to be visible

average linkage



singlelinkage





Example 22 clusters

average linkage



singlelinkage





Example 21 clusters

average linkage



singlelinkage





Example 20 clusters

average linkage



singlelinkage





Example 19 clusters

single-linkage again has clearly nonconvex Voronoi region near middle

average linkage



singlelinkage





Example 18 clusters

average linkage



singlelinkage





Example 17 clusters

average linkage



singlelinkage





Example 16 clusters

average linkage



singlelinkage





Example 15 clusters

complete-linkage has points outside their cluster's Voronoi region

average linkage



singlelinkage





Example 14 clusters

average linkage



singlelinkage





Example 13 clusters

average linkage



singlelinkage





Example 12 clusters

average linkage



singlelinkage





Example 11 clusters

complete-linkage: merging two clusters shrank the Voronoi region

average linkage



singlelinkage





Example 10 clusters

average linkage



singlelinkage





Example 9 clusters

average linkage



singlelinkage

complete linkage



8

Example 8 clusters

average linkage



singlelinkage







singlelinkage

average linkage

complete linkage



-2





average linkage



singlelinkage





Example 5 clusters

average-linkage has points outside their cluster's Voronoi region

average linkage



singlelinkage







average linkage



singlelinkage





Example 3 clusters

nonconvex Voronoi regions for singlelinkage and complete-linkage

average linkage



singlelinkage







average linkage



singlelinkage





Example 1 cluster

trees show the different merging paths of the algorithms

average linkage



singlelinkage





Observations

- All three methods produce a hierarchical partitioning clustering but make different choices
- Single-linkage is more "greedy" and can yield large clusters early on
- Single-linkage iterations result in a union of two previous Voronoi regions. The other two methods change the shape of Voronoi regions in each merge.
- Average- and complete-linkage can yield Voronoi regions that do not contain all training points of that cluster.