

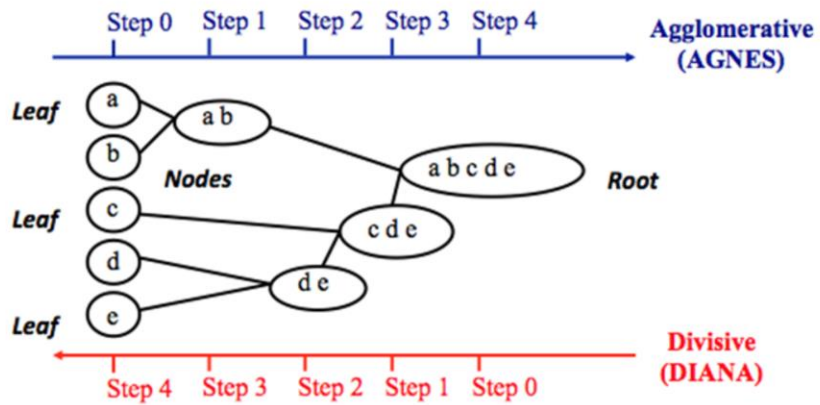
In this module, we discuss Agnes and Diana. The former provides a way to do Agglomerative Clustering, while the latter, Divisive Clustering.

AGNES & DIANA I

- Hierarchical clustering can be divided into two main types: **agglomerative** and **divisive**.
- **Agglomerative clustering:** It's also known as **AGNES** (Agglomerative Nesting). It works in a bottom-up manner. That is, each object is initially considered as a single-element cluster (leaf). At each step of the algorithm, the two clusters that are the most similar are combined into a new bigger cluster (nodes). This procedure is iterated until all points are member of just one single big cluster (root).
- **Divisive hierarchical clustering:** It's also known as **DIANA** (Divide Analysis) and it works in a top-down manner. The algorithm is an inverse order of AGNES. It begins with the root, in which all objects are included in a single cluster. At each step of iteration, the most heterogeneous cluster is divided into two. The process is iterated until all objects are in their own cluster.

As we discussed earlier, Hierarchical clustering can either be Agglomerative, i.e., a bottoms up approach, or a top-down approach, called Divisive. Here we will look at both approaches.

AGNES & DIANA II



This diagram depicts both approaches. While both approaches can lead to the same solution, the time complexity may vary depending on how the splitting (divisive), or merging (agglomerative) approaches are implemented.

AGNES : Pseudocode

- Basic algorithm is straightforward
 1. Compute the proximity matrix
 2. Let each data point be a cluster
 3. **Repeat**
 4. Merge the two closest clusters
 5. Update the proximity matrix
 6. **Until** only a single cluster remains
- Key operation is the computation of the **proximity** (or distance) of two clusters.
- Different approaches to defining the **distance** between clusters distinguish the different algorithms

The process is simple. Merge clusters based on their proximity to other clusters. Initially, all data points are in a cluster of size 1. Proximity is a measure of distance, and that can be defined in any fashion, including using well-known methods for distance.

Distance Measures

Distances can be defined in multiple ways, but in general, the following properties are required:

- Non-negative: $d_{ij} \geq 0$
- Self-proximity: $d_{ii} = 0$ (the distance from a record to itself is zero)
- Symmetry: $d_{ij} = d_{ji}$
- Triangle inequality: $d_{ij} \leq d_{ik} + d_{kj}$ (the distance between any pair cannot exceed the sum of distances between the other two pairs)

All distances follow some basic properties. These include, non-negativity, self proximity, symmetry, and the triangle inequality.

Distance Measures - Cont.

- **Euclidean Distance** $d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$
 - Scale dependent
 - Ignores relationship between variables
 - Sensitive to outliers
- **Manhattan Distance** $d_{ij} = \sum_{m=1}^p |x_{im} - x_{jm}|$
 - Useful in presence of outliers
- **Maximum Coordinate Distance**
 - Looks only at the measurement on which records i and j deviate the most

$$d_{ij} = \max_{m=1, 2, \dots, p} |x_{im} - x_{jm}|$$

While the Euclidean distance is the most commonly used one, other distance measures do exist. The City Block, or Manhattan Distance is also widely used.

Similarity Measures

- **Correlation Based Similarity**

- A popular similarity measure is the square of the Pearson correlation coefficient, r^2 . Such measures can be converted to distance measures as $d_{ij} = 1 - r_{ij}^2$

- **Mahalanobis Distance (Statistical Distance)**

- This metric has an advantage over the other metrics mentioned in that it takes into account the correlation between measurements. With this metric, measurements that are highly correlated with other measurements do not contribute as much as those that are uncorrelated or mildly correlated.
- S is the covariance matrix.

$$d_{ij} = \sqrt{(x_i - x_j)' S^{-1} (x_i - x_j)}$$

In the previous measures of distance, no consideration is given when values are correlated. It is possible to define distances that specifically account for this correlation. The Mahalanobis Distance is commonly used in Statistical Methods, and specifically accounts for the correlation between observations. These measures can be used as part of HC.

Distance Measures for Categorical Data

In the case of measurements with binary values, it is more intuitively appealing to use similarity measures than distance measures.

- Matching coefficient $(a + d) / n$
- Jacquard's coefficient $d / (b + c + d)$
 - This coefficient ignores zero matches. This is desirable when we do not want to consider two people to be similar simply because a large number of characteristics are absent in both.
- Mixed Data?

		Record <i>j</i>		
		0	1	
Record <i>i</i>	0	<i>a</i>	<i>b</i>	<i>a + b</i>
	1	<i>c</i>	<i>d</i>	<i>c + d</i>
		<i>a + c</i>	<i>b + d</i>	<i>n</i>

So far, the distances are all based on numerical measures. What about for categorical variables. In this case, we typically want to calculate a similarity index. The Matching coefficient and Jacquard's coefficients are two approaches for calculating the similarity index.

When the measurements are mixed (some continuous and some binary), a similarity coefficient suggested by Gower is very useful. Gower's similarity measure is a weighted average of the distances computed for each variable, after scaling each variable to a [0,1] scale.

Measuring Distance Between Two Clusters

- Minimum Distance: $\min(\text{distance}(A_i, B_j))$
- Maximum Distance: $\max(\text{distance}(A_i, B_j))$,
- Average Distance: $\text{Average}(\text{distance}(A_i, B_j))$
- Centroid Distance: $\text{distance}(\bar{X}_A, \bar{X}_B)$

$$\bar{x}_A = [(1/m \sum_{i=1}^m x_{1i}, \dots, 1/m \sum_{i=1}^m x_{pi})]$$

We define a cluster as a set of one or more records. How do we measure distance between clusters? The idea is to extend measures of distance between records into distances between clusters. Consider cluster A, which includes the m records A_1, A_2, \dots, A_m and cluster B, which includes n records B_1, B_2, \dots, B_n . The most widely used measures of distance between clusters are:

Minimum Distance: The distance between the pair of records A_i and B_j that are closest: $\min(\text{distance}(A_i, B_j))$

Maximum Distance: The distance between the pair of records A_i and B_j that are farthest: $\max(\text{distance}(A_i, B_j))$

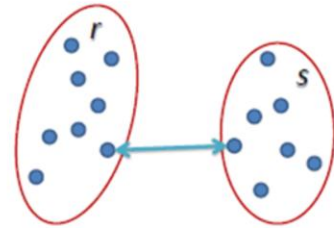
Average Distance: The average distance of all possible distances between records in one cluster and records in the other cluster: $\text{Average}(\text{distance}(A_i, B_j))$

Centroid Distance: The distance between the two cluster centroids. A cluster centroid is the vector of measurement averages across all the records in that cluster.

Each of these distances (minimum, maximum, average, and centroid distance) can be implemented in the hierarchical scheme as described as follows.

Single (Min) Linkage

- In single (min) linkage hierarchical clustering, the distance between two clusters is defined as the shortest distance between two points in each cluster.
- For example, the distance between clusters “r” and “s” to the left is equal to the length of the arrow between their two closest points.

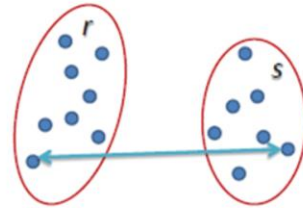


$$L(r, s) = \min(D(x_{ri}, x_{sj}))$$

In single linkage clustering, the distance measure that we use is the minimum distance (the distance between the nearest pair of records in the two clusters, one record in each cluster). This method has a tendency to cluster together at an early stage records that are distant from each other because of a chain of intermediate records in the same cluster. Such clusters have elongated sausage-like shapes when visualized as objects in space.

Complete (Max) Linkage

- In complete (max) linkage hierarchical clustering, the distance between two clusters is defined as the longest distance between two points in each cluster.
- For example, the distance between clusters “r” and “s” to the left is equal to the length of the arrow between their two furthest points.



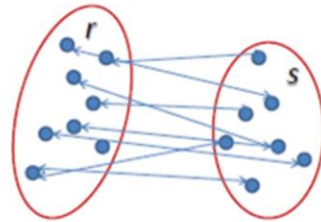
$$L(r, s) = \max(D(x_{ri}, x_{sj}))$$

In complete linkage clustering, the distance between two clusters is the maximum distance (between the farthest pair of records).

This method tends to produce clusters at the early stages with records that are within a narrow range of distances from each other. If we visualize them as objects in space, the records in such clusters would have roughly spherical shapes.

Average Linkage

- In average linkage hierarchical clustering, the distance between two clusters is defined as the average distance between each point in one cluster to every point in the other cluster.
- For example, the distance between clusters “r” and “s” to the left is equal to the average length each arrow between connecting the points of one cluster to the other.



$$L(r, s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$$

Average linkage clustering is based on the average distance between clusters (between all possible pairs of records).

Note that unlike average linkage, the results of the single and complete linkage methods depend only on the ordering of the inter-record distances. Linear transformations of the distances (and other transformations that do not change the ordering) do not affect the results.

Centroid Distance

The distance between the two cluster centroids. A cluster centroid is the vector of measurement averages across all the records in that cluster. For Cluster A, this is the vector $\bar{x}_A = [(1/m \sum_{i=1}^m x_{1i}, \dots, 1/m \sum_{i=1}^m x_{pi})]$.

The centroid distance between clusters A and B is then $\text{distance}(\bar{x}_A, \bar{x}_B)$.

Centroid linkage clustering is based on centroid distance, where clusters are represented by their mean values for each variable, which forms a vector of means. The distance between two clusters is the distance between these two vectors. In average linkage, each pairwise distance is calculated, and the average of all such distances is calculated. In contrast, in centroid distance clustering, just one distance is calculated: the distance between group means. This method is also called Unweighted Pair-Group Method using Centroids (UPGMC).

Ward's Method

Ward's method aims to minimize the total within-cluster variance. At each step the pair of clusters with minimum between-cluster distance are merged. In other words, it forms clusters in a manner that minimizes the loss associated with each cluster. At each step, the union of every possible cluster pair is considered and the two clusters whose merger results in minimum increase in information loss are combined. Here, information loss is defined by Ward in terms of an error sum-of-squares criterion (ESS).

Ref: <https://www.r-bloggers.com/how-to-perform-hierarchical-clustering-using-r/>

Ward's method is also agglomerative, in that it joins records and clusters together progressively to produce larger and larger clusters, but operates slightly differently from the general approach described previously. Ward's method considers the "loss of information" that occurs when records are clustered together. When each cluster has one record, there is no loss of information and all individual values remain available. When records are joined together and represented in clusters, information about an individual record is replaced by the information for the cluster to which it belongs. To measure loss of information, Ward's method employs a measure "error sum of squares" (ESS) that measures the difference between individual records and a group mean.

For example, consider the values (2, 6, 5, 6, 2, 2, 2, 2, 0, 0, 0) with a mean of 2.5. Their ESS is equal to

$$(2-2.5)^2 + (6-2.5)^2 + (5-2.5)^2 + \dots + (0-2.5)^2 = 50.5.$$

The loss of information associated with grouping the values into a single group is therefore 50.5. Now group the records into four groups: (0, 0, 0), (2, 2, 2, 2), (5), (6, 6). The loss of information is the sum of the ESS's for each group, which is 0 (each record in each group is equal to the mean for that group, so

the ESS for each group is 0). Thus, clustering the 10 records into 4 clusters results in no loss of information, and this would be the first step in Ward's method.

Ward's method tends to result in convex clusters that are of roughly equal size, which can be an important consideration in some applications (e.g., in establishing meaningful customer segments).

AGNES Versus DIANA

- Generally speaking, the output of AGNES and DIANA should be comparable.
- Both algorithms also have similar computation complexity
- AGNES is more commonly used though.
- Empirical results suggest that AGNES is good at identifying small clusters while DIANA is a better choice at identifying larger clusters.

This summarizes the choice of using Agnes or Diana. Both approaches produce similar results, though Agnes is usually more commonly used.