

# 1 Introduction to hierarchical clustering

The idea behind hierarchical clustering is to capture the different abstraction layers in a hierarchy, i.e. moving from the broadest levels of differences down to finer distinctions. The idea is to not only speed up the retrieval process (which drops from a linear amount of comparisons to logarithmic) but also to increase the accuracy of retrieval. The standard approach is to cluster at different levels of dissimilarity (generally by using Euclidean distance) and build a tree from this. We explore an alternative approach by generalizing  $k$ -means to a hypersphere, though many of the paradigms of hierarchical clustering also apply.

## 2 Divisive and agglomerative clustering

Divisive clustering[?] is also known as the “top down” method. This takes in a data set as one large cluster. Then as it moves down the hierarchy, it recursively splits the data until it reaches the leaves of the tree, i.e. each observation is its own singleton cluster.

Agglomerative clustering can be thought to do the opposite of divisive clustering and is known as a “bottom up” strategy. It takes in a data set and it initially looks at each observation as its own singleton cluster. Then, based on its linkage criterion, it will merge recursively to form new clusters until all the data is in one large cluster.

A dendrogram[?] can be used in both cases to record and visualize the hierarchical nature of the clusters and their distance to each other.

## 3 Cluster Dissimilarity Measurements

In forming clusters, a measurement must be made to determine which clusters split or merge in regards to divisive or agglomerative clustering, respectively. There are many different ways to measure dissimilarity but this paper will explore single linkage, complete linkage, and centroid linkage. Each of the measurements below refer to ways in measuring distances between  $k$  number of clusters  $c_{i=1\dots k}$  in relationship with specific points  $x_{i=1\dots k}$  in those clusters. Because different linkage patterns will result in different clusters being formed, it is an important parameter to vary in hierarchical clustering.

### 3.1 Single Linkage

Single linkage takes the distance between the closest observations, or minimum distance, in each cluster and merges the those with the shortest distance. This works well for data sets that desire to cluster together long chained data points. Let  $D(c_i, c_j)$  be the distance between two centroids for  $i, j \in \{1, 2, \dots, k\}$ , and  $D(x_i, x_j)$  be the distance between two data points in  $i, j \in \{1, 2, \dots, n\}$ . Then single linkage defines the following relationship:

$$D(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2) \quad (1)$$

### 3.2 Complete Linkage

Complete linkage is similar to single linkage. It takes the distance between the farthest observations, or the maximum distance, in each cluster. The clusters that have the shortest measured distance are merged together. This measurement is ideal for data that is spherically organized relative to each other. The relationship is defined as

$$D(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2) \quad (2)$$

### 3.3 Centroid Linkage

Using centroid linkage involves taking the average of all points in a cluster and assigning that value as the mean centroid. Then the centroids with the shortest distances are merged together. Centroid linkage is defined as

$$D(c_1, c_2) = D\left(\left(\frac{1}{|c_1|} \sum_{x \in c_1} \vec{x}\right), \left(\frac{1}{|c_2|} \sum_{x \in c_2} \vec{x}\right)\right) \quad (3)$$

## 4 Clustering implementations

We developed both hierarchical divisive and agglomerative clustering functions to explore which approach clustered more effectively and efficiently (i.e. algorithm time complexity). In our particular application, instead of flat

clustering, we had to cluster over the curved surface of a sphere. In order to do this, we utilized the SPKmeans function [?]. We also chose our linkage criterion to be based on centroid linkage.

Essentially, we developed the divisive clustering by following the intuition mentioned above. For our agglomerative clustering implementation, we did a slightly different variation hoping its run time would be slightly better than the traditional method. In our code, intuitively, we initiated each point as its own cluster with its own centroid. Then we recursively merged two clusters together based on centroid linkage to form a new cluster with its own cluster centroid each time. This process stops when there is a single point to represent the entire data's average centroid.

We tested with a relatively small, three dimensional data set of 150 points. In both divisive and agglomerative, we input different cluster values into SPKmeans which determines how often the function gets called. In the divisive method, we asked SPKmeans to return two clusters. This means the number of times we call SPKmeans is based on how many nodes we have which is  $m = 2n - 1$  where  $n$  is the number of data points. For the agglomerative case, we asked SPKmeans to return the  $n/2$  clusters. This means means we had to call SPKmeans at every level which is  $l = \lceil \log_2 n \rceil$ .

As a benchmark, we also used the MATLAB `linkage` function to construct a pairwise-distance based tree (merging nodes together by single-linkage) and compared the results.

## 4.1 Implementation details

In  $k$ -means, there is a choice of initial mean initialization. This affects the initial distribution of means and therefore potentially the number of iterations until convergence, but also may be expensive computationally. For instance, a good setting of the initial means will reduce the likelihood of empty clusters found, which is dealt with by splitting the largest cluster and is therefore computationally expensive.

For the agglomerative algorithm, it makes sense to spend time initializing the means (we chose to maximize the distance between the initial means), since the number of calls to SPKMeans is small and the number of clusters for each call is high. For the divisive algorithm, it doesn't make sense to incur this expense when SPKMeans is called many times with only the branching-factor number of clusters (generally a small constant). This optimization and lack thereof, respectively, reduces the running time for both algorithms substantially.