# Assignment 3: Big Data Analytics
## Hierarchical Clustering on Cereals
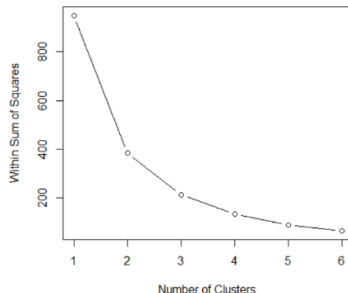## By: Valay Shah

## Purpose

The purpose of this report is to analyze the data file Cereals.csv which contains information on 77 different cereals and their characteristics. For the analysis, all cereals with missing values were removed, hierarchical clustering was applied to the data using Euclidean distance, and normalised measurements. Finally, the dendograms for single and complete linkages were compared and a recommendation for the number of clusters and cluster membership was determined.

## Removing all cereals with missing values

There are three rows which will be removed because of carbo, sugar, or potass columns containing blank cells. This will lead to the database going from 76 observations to 74 observations.

## Hierarchical clustering

Before we perform hierarchical clustering, non-numeric columns must be removed which means mfr and type columns will be removed from the database. These two variables are not entirely relevant to the database as the manufacturer of the cereal does not matter and almost all the cereal are the same type (type C). Using sapply will change all the numbers to standardized numbers and the distance between the variables will be normalized. Within sum of squares (wss) plot was used to help determine the number of clusters which is ideal for a lower variance between the points within the clusters.
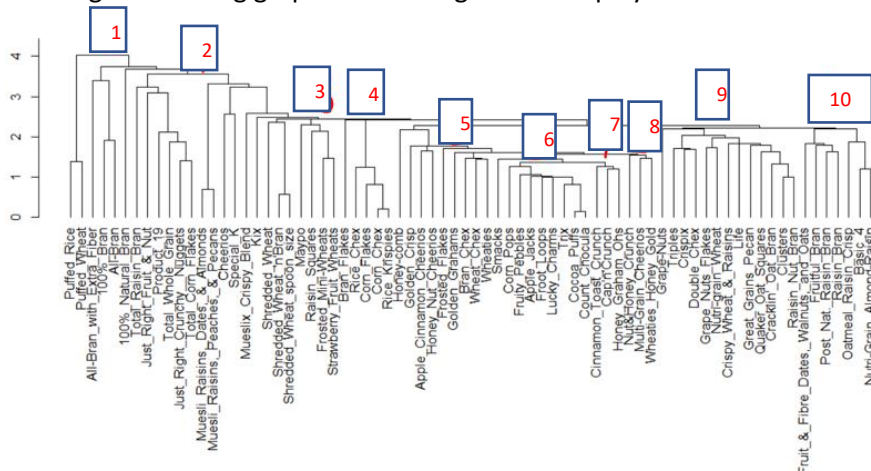


From the plot we can observe that 5 or 6 clusters are ideal for having a low sum of squares which is ideal due to low variance between the points within a cluster.
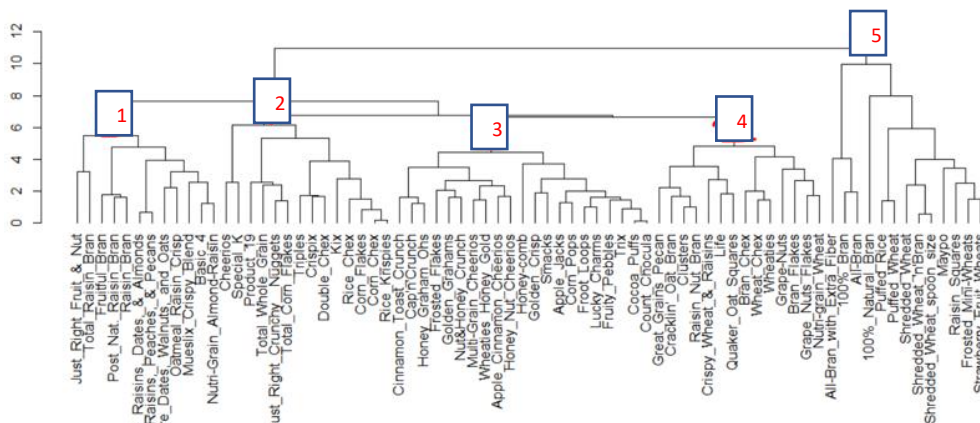
From this, we can move forward with using the single and complete clustering methods to analyze the dendograms.

## Single Clustering vs Complete Clustering

The single clustering graph with a hang of -6 is displayed below:

From this plot we can see that there are 10 distinct groups. Group 1 is comprised of nut bran type of cereals, group 2 contains crispy, wheat and bran type of cereals, group 3 has a healthier version of sugary/honey cereal, group 4 has unhealthy and sugary crunchy cereals, group 5 contains sugary cereals, group 6 has wheat, bran, and honey healthy cereals, group 7 has plain cereals, group 8 has frosty and fruity cereals, group 9 possesses cereals with nuts and fruits, and lastly group 10 is 100% natural and bran which are the healthiest type of cereals. From this dendogram, we can cut the tree into 10 clusters. The results of that show that most cereals fall in cluster 4 and there are not a lot of cereals in clusters 2,3 7, 8, and 10. For this reason, we can stick to a cutree of 6 clusters which once again shows that most cereals fall into cluster 4 but there are quite a few data points in clusters 5 and 6 as well. If we were to divide the above dendogram into 6 clusters, the clusters which are labeled 3,4,5,6 would likely all fall into cluster 4.



<u>Complete Cluster Plot</u>

From the complete cluster plot above, it paints a much clearer picture for which group go together than single cluster plot. In this plot, we can clearly see 5 overarching distinct groups: Group 1 – Brans and wheats, Group 2: Flakes, Chex and Raisins, Group 3 – Sugary cereals, Group 4 – Plain cereals, Group 5 – Fruits, nuts, and dry fruit cereals. This complete cluster graph is much easier to read than the single cluster graph as it groups cereals with maximum distance from one another apart from each other as opposed to the two closest pairings we find in the single cluster graph. The complete linkage dendogram tells a much clearer story by breaking the clusters in 5 distinct groups and one can clearly determine healthy cereal groups and unhealthy groups.

**Recommendation**

As noted above, the complete linkage dendogram shows a much clearer picture than the single linkage dendogram and the groups can be easily identified and sorted into 5 distinct cereal groups. As a result, the complete linkage dendogram's results are what I would recommend. We should cut the data into 5 clusters using cutree which gives the cluster membership as follows:

## Code:

```r
# read in data
cereals.df <- read.csv("C:/Users/valay/OneDrive/Desktop/Ivey/Winter Term 2021/Big Da

#remove any data that is na or blank

cereals.df = cereals.df[-which(is.na(cereals.df$carbo)),]
cereals.df = cereals.df[-which(is.na(cereals.df$potass)),]
cereals.df

# Make row names the names of the utilities, remove numerical labeling of rows
row.names(cereals.df) <- cereals.df[,1]

#Before we can normalize the variables, all variables must have only numeric values
#So we can remove the first two columns since they are not numerical
#the minus 1 removes a column in the cereals file
cereals.df <- cereals.df[,-1]
cereals.df
#removes mfr column
cereals.df <- cereals.df[,-1]
cereals.df
#removes type column
cereals.df <- cereals.df[,-1]
cereals.df
```

```r
26  # normalize input variables
27  # sapply changes everything to standardized numbers
28  cereals.df.norm <- sapply(cereals.df, scale)
29
30  # add row names to cereals.df
31  row.names(cereals.df.norm) <- row.names(cereals.df)
32  #you can only run the above line once; can't put row names once row names have already been esta
33
34  # compute normalized distance based on all variables
35  d.norm <- dist(cereals.df.norm, method = "euclidean")
36
37
38
39  # in hclust() set argument method =
40  # to "ward.D", "single", "complete", "average", "median", or "centroid"
41
42  # plot an empty scatter plot
43  km<-kmeans(cereals.df.norm, 6)
44  #the 6 means 6 clusters
45
46  plot(c(0), xaxt = 'n', ylab = "", type = "l", ylim =
47       c(min(km$centers), max(km$centers)), xlim = c(0, 8))
48
49  # label x-axes
50  axis(1, at = c(1:13), labels = names(cereals.df))
51
52  # plot centroids
53  for (i in c(1:6))
54    lines(km$centers[i,], lty = i, lwd = 2, col = ifelse(i %in% c(1, 3, 5),"black", "dark grey"))
55
56  # name clusters
57  text(x = 0.5, y = km$centers[, 1], labels = paste("Cluster", c(1:6)))
58
59  #averages within clusters for each variable
60  km$centers
61
```

```r
62  #distance between points within clusters
63  km$withinss
64
65  #how many data points are within each cluster is the size of the cluster
66  km$size
67
68  #1-6 clusters
69  #wss is within sum of squares meaning the distance between points within each cluster
70  #the higher the number of clusters, the lower the variance b/w points within clusters
71  wss <- numeric(6)
72
73  #mean of distance between points within a cluster vs number of clusters
74  for (k in 1:6) wss[k] <- mean(kmeans(cereals.df.norm, centers = k, nstart = 25)$withinss)
75  plot(1:6, wss, type="b",xlab="Number of Clusters", ylab="Within Sum of Squares")
76
77
78
79
80  #SINGLE CLUSTER
81  hc1 <- hclust(d.norm, method = "single")
82  hc1
83
84  plot(hc1, hang = -1, ann = FALSE)
85
86  #the graph (dendrogram) shows us the shortest distance NY is to anything else is 3.6
87  #closest Idaho and Puget are to anything else is 2.3
88
89  cutree(hc1,k=6)
90  #this shows just 6 clusters; number of clusters is 6 clusters
91
92  cutree (hc1, k=10)
93  #this shows 10 clusters
94
95  #h will be cut the tree based on a number
96  #if h is 3 then it clusters everything below 3 as one cluster and for everything above 3 will have
97  #own seperate cluster
```

```r
103  #COMPLETE CLUSTER
104  hc2 <- hclust(d.norm, method = "complete")
105  hc2
106
107  plot(hc2, hang = -1, ann = FALSE)
108
109  #the graph (dendrogram) shows us the shortest distance NY is to anything else is 3.6
110  #closest Idaho and Puget are to anything else is 2.3
111
112  cutree(hc2,k=2)
113  #this shows just 2 clusters; number of clusters is 2 clusters
114
115  cutree (hc2, k=5)
116  #this shows 6 clusters
117
118  #h will be cut the tree based on a number
119  #if h is 3 then it clusters everything below 3 as one cluster and for everything above 3 will have
120  #own seperate cluster
```