

## HackBio Bio-Data Science Task 3

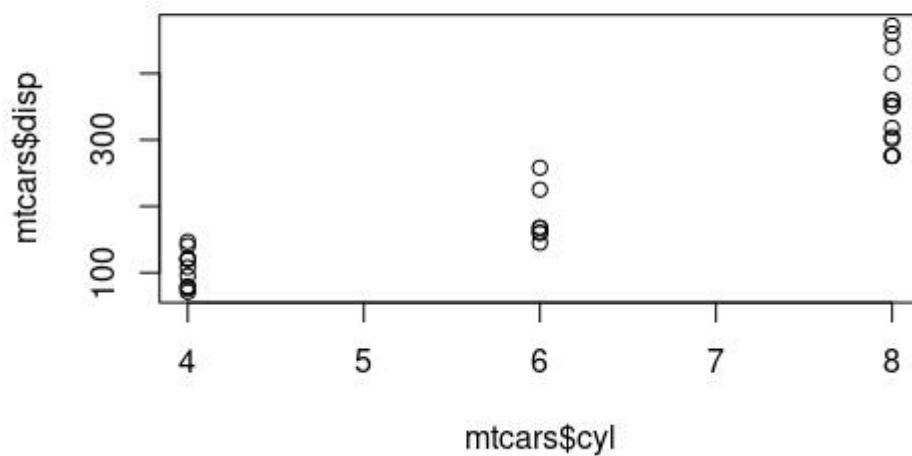
In this stage, I will continue to improve on my R programming skills. I have two tasks. First, I will perform K-means clustering and Hierarchical clustering on 'mtcars' built in data set in R' and secondly, the biological data set 'microbial\_stationary\_phase', that I performed pca on last week.

Consider the following explanations and plots.

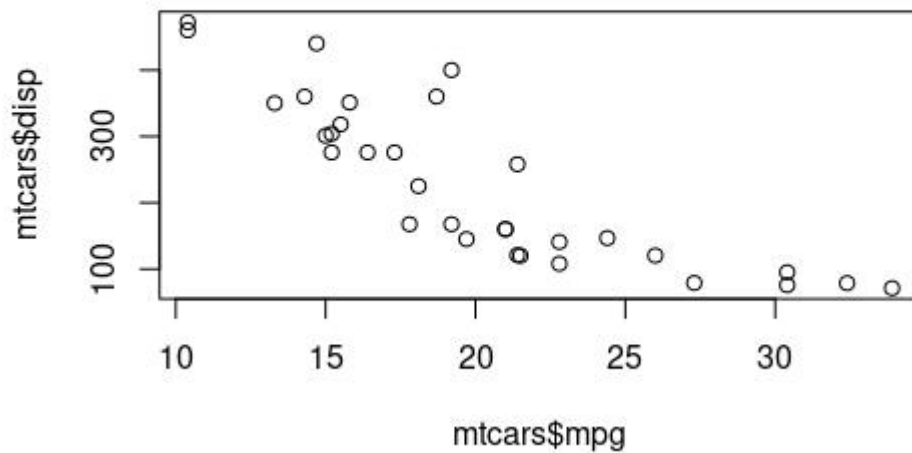
```
## Data set for this first analysis is 'mtcars'
# K-means clustering

# Our data is mtcars, an inbuilt data in R
data()
data(mtcars)
View(mtcars)

# For the purpose of visualization, I will like to use any two columns that correlate
cor(mtcars)
plot (mtcars$cyl, mtcars$disp)
```



```
plot (mtcars$mpg, mtcars$disp)
```

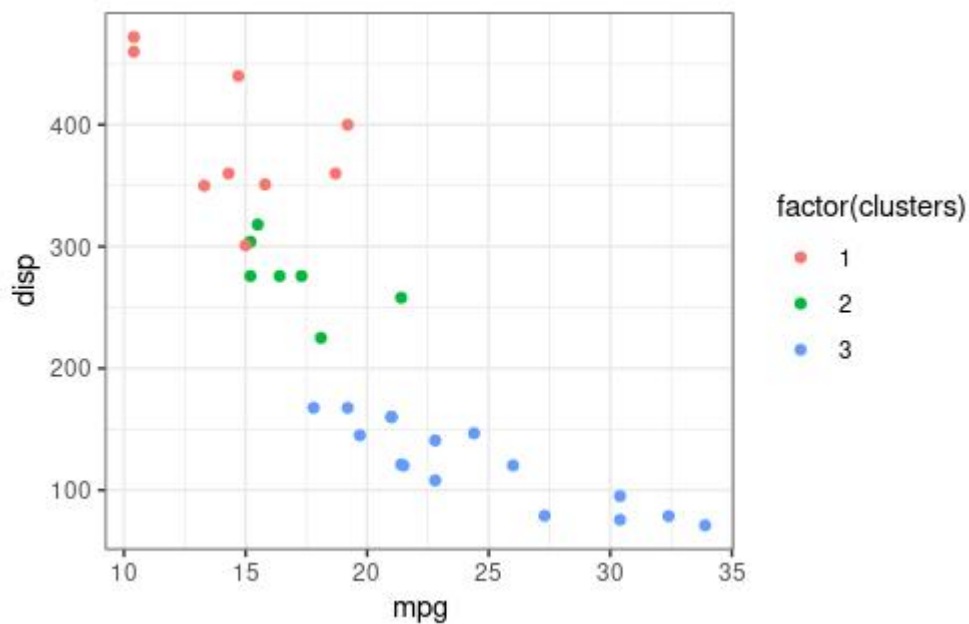


```
# So we can have same results
# I will set seed which is an arbitrary number. Since I want to get a consistent result
# from my K-Mean clustering (which is unsupervised). Picking '102' means that I have
# chosen to index to 102.
set.seed(102)

# I will now perform my KMEAN Clustering for now. I want to pick 3 centers
# (clusters)
mtcarsK3 <- kmeans (x = mtcars, centers = 3)

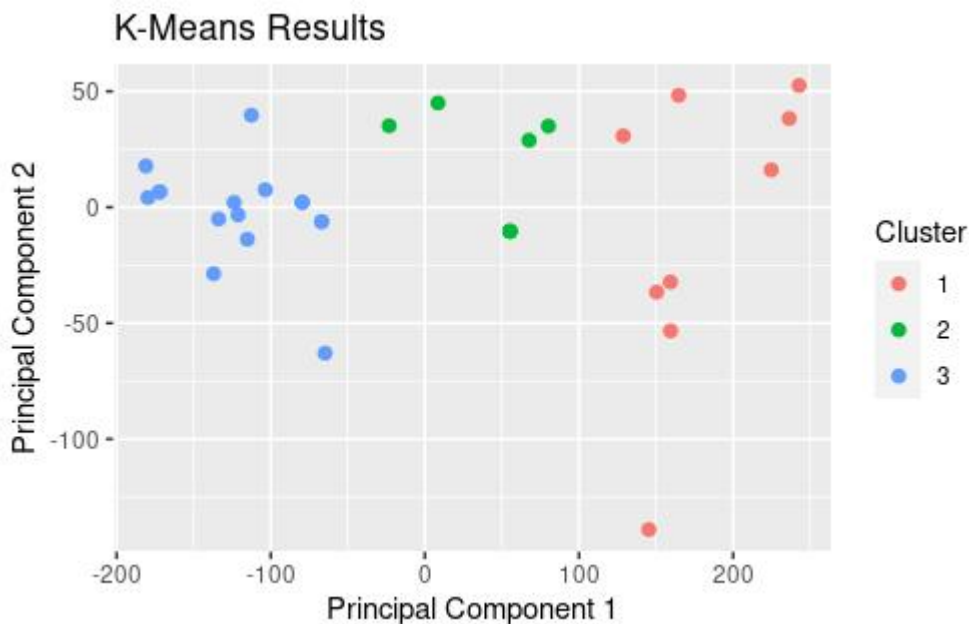
# Well, I can also add this cluster information to my dataset
mtcars$clusters <- c(mtcarsK3$cluster)

# I want to visualize this cluster information for each car
ggplot(mtcars, aes(x = mpg, y = disp, color = factor(clusters))) + geom_point() +
theme_bw()
```

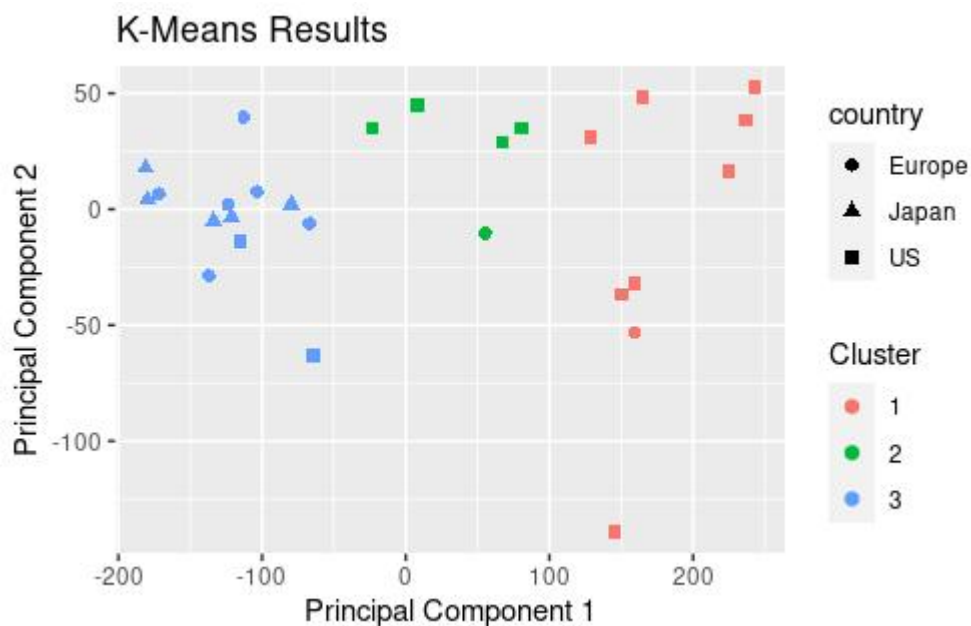


```
# I need to install "useful" for my K-Means Clustering
install.packages('useful')
library('useful')
```

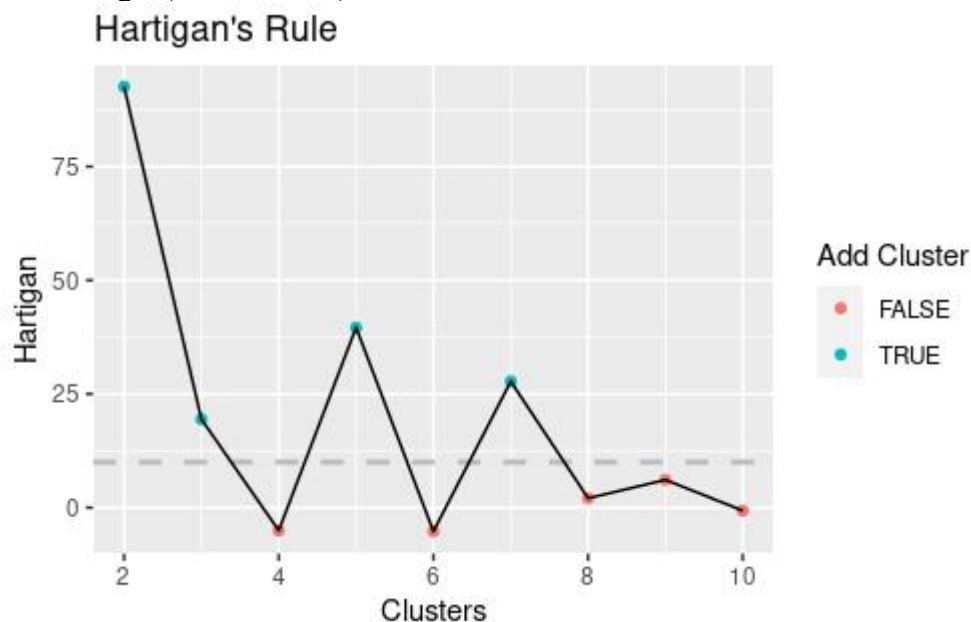
```
# I will create a new column for country
mtcars.country <- c(rep("Japan", 3), rep("US", 4), rep("Europe", 7), rep("US", 3),
"Europe", rep("Japan", 3), rep("US", 4), rep("Europe", 4), rep("US", 3))
mtcars$country <- c(mtcars.country)
# Let's now plot
plot(mtcarsK3, data = mtcars)
```



```
plot(mtcarsK3, data = mtcars, class = 'country')
```



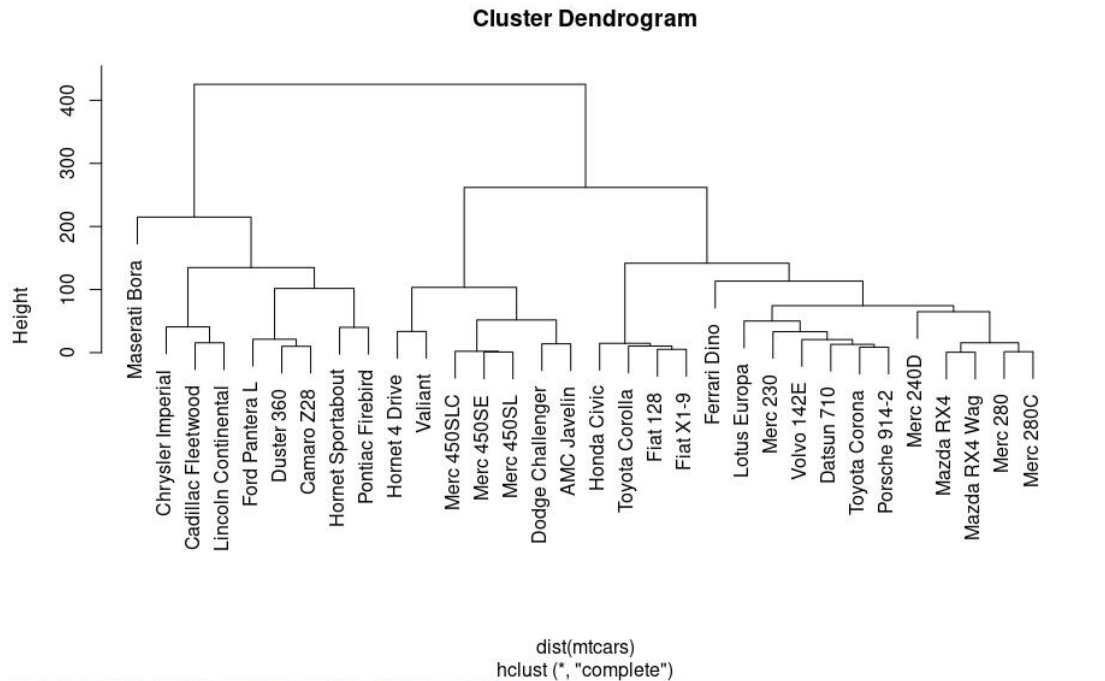
```
# Choosing the right number of clusters: I will set my maximum cluster to 10.
# NB: I need to re-run my data set since the columns of 'country' and 'cluster' which I
# previously added will not all our code to run because it doesn't require alphabets.
mtcarsBEST <- FitKMeans(mtcars, max.clusters=10)
mtcarsBEST
# Once I got a FALSE, I stopped counting for clusters
PlotHartigan(mtcarsBEST)
```



# I will now plot again, this time using the number of clusters that I determined. Since I determined 3 clusters just like previously, there is no need to repeat all over.

```
# Hierarchical Clustering is used to cluster clusters into clusters
# Let's see how to implement it: I will start by calculating the distance between the
# rows. Next, I will pick a method. In this case I will be using the "complete method".
# "hclust and dist" are inbuilt functions in R.
```

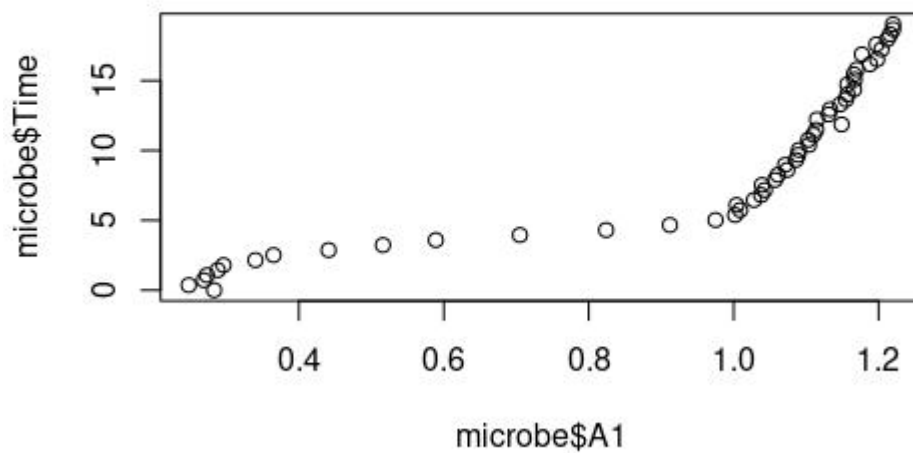
```
hcmtcars <- hclust(d=dist(mtcars), method='complete')
plot(hcmtcars)
```



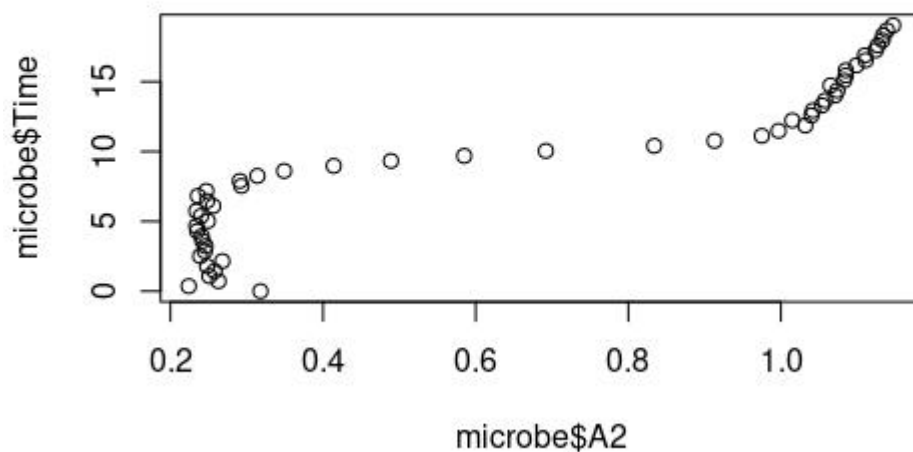
```
# Second K-means clustering
# Our second data set is 'microbial_stationary_phase.csv' from last week's task

# Import .csv file in R
microbe <- read.csv(file.choose())
microbe

# For the purpose of visualization, I will like to use any two columns that correlate
cor(microbe)
plot (microbe$A1, microbe$Time)
```



```
plot (microbe$A2, microbe$Time)
```

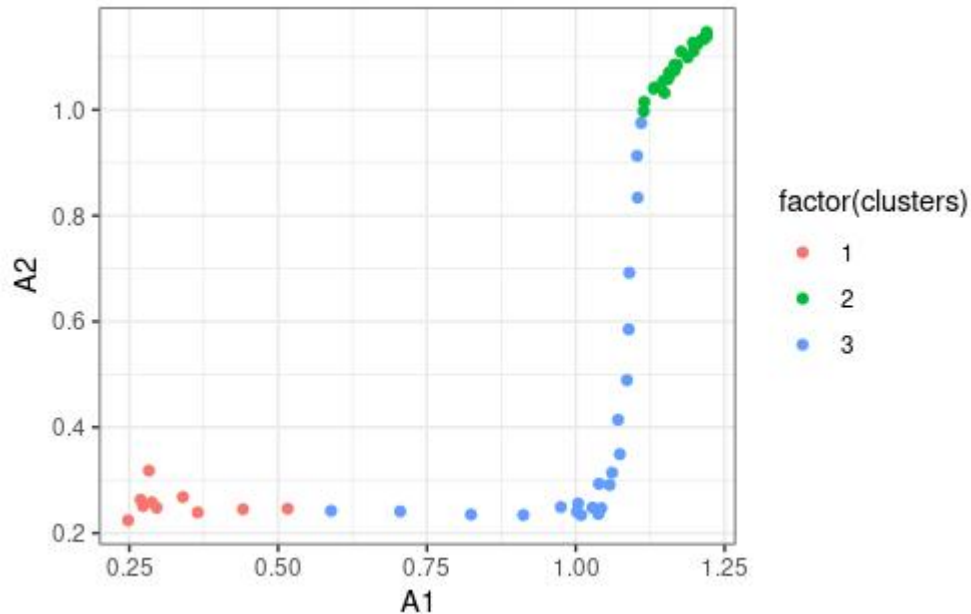


```
# So we can have same results
# I will set seed which is an arbitrary number. Since I want to get a consistent result
# from my K-Mean clustering (which is unsupervised). Picking '102' means that I have
# chosen to index to 102.
set.seed(102)
```

```
# I will now perform my KMEAN Clustering for now. I want to pick 3 centers
# (clusters)
microbeK3 <- kmeans (x = microbe, centers = 3)
```

```
# Well, I can also add this cluster information to my dataset
microbe$clusters <- c(microbeK3$cluster)
```

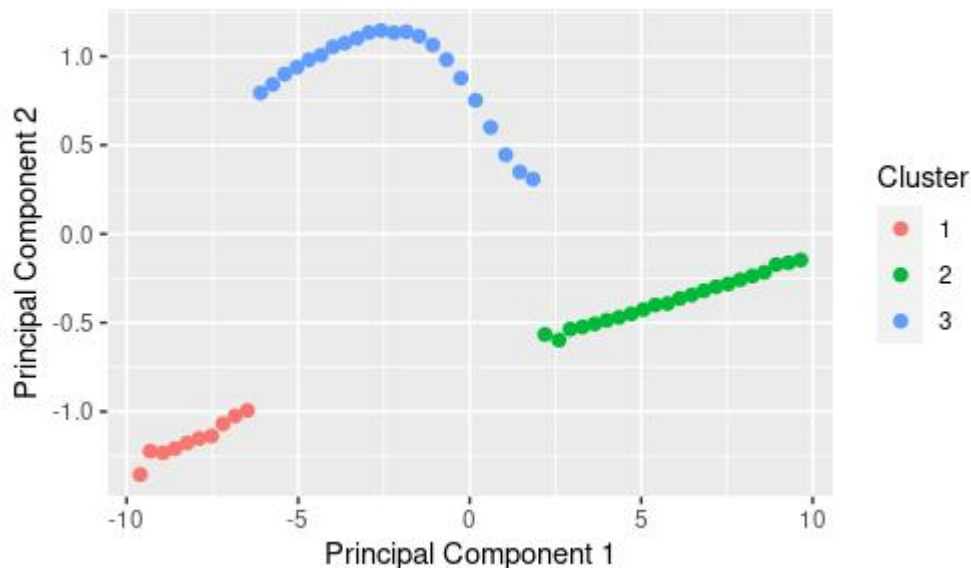
```
# I want to visualize this cluster information for each microbe
ggplot(microbe, aes(x = A1, y = A2, color = factor(clusters))) + geom_point() +
theme_bw()
```



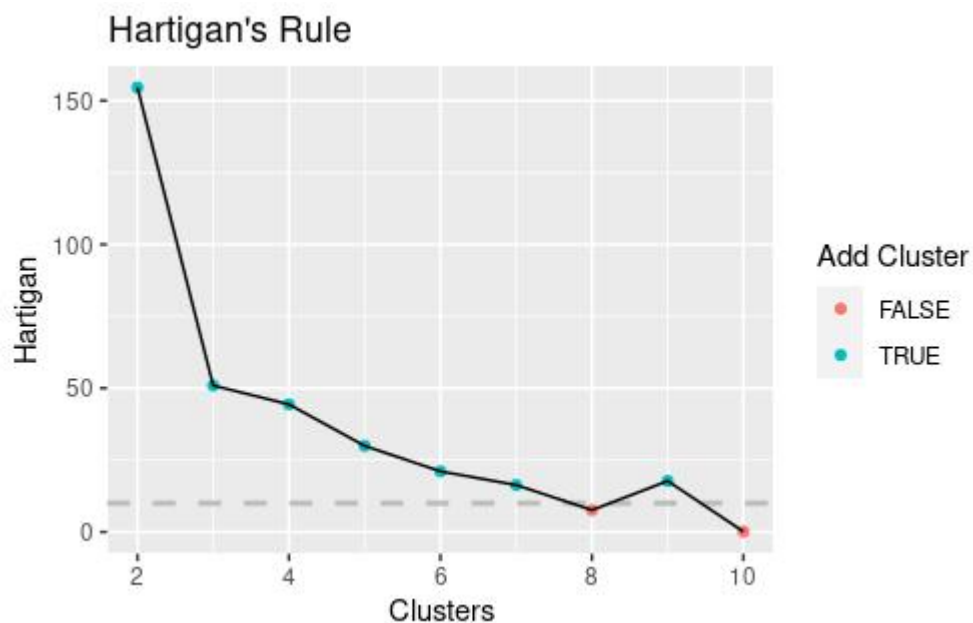
```
# Install useful
install.packages('useful')
library('useful')

# Let's plot
plot(microbeK3, data = microbe)
```

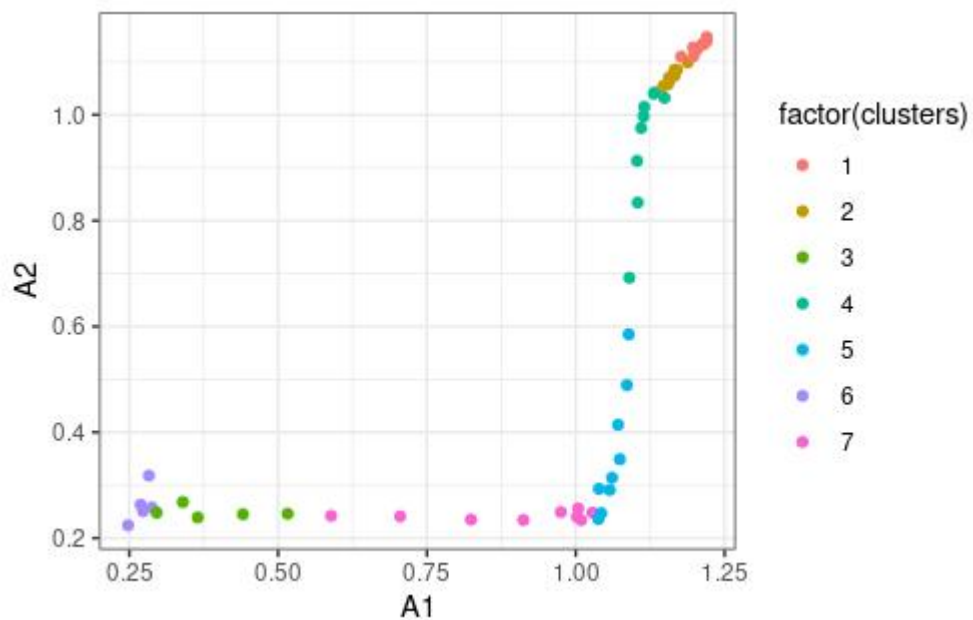
### K-Means Results



```
# Choosing the right number of clusters: I will set my maximum cluster to 10.
microbeBEST <- FitKMeans(microbe, max.clusters=10)
# Once I got a FALSE, I stopped counting for clusters
PlotHartigan(microbeBEST)
```



# I will now plot again, this time using the number of clusters that I determined



# Hierarchical Clustering is used to cluster clusters into clusters

# Let's see how to implement it: I will start by calculating the distance between the rows. Next, I will pick a method. In this case I will be using the "complete method". "hclust and dist" are inbuilt functions in R.

```
hcmicrobe <- hclust(d=dist(microbe), method='complete')
plot(hcmicrobe)
```



