

HackBio Bio-Data Science Task 3

In this stage, I will continue to improve on my R programming skills. I will perform K-means clustering and Hierarchical clustering on the biological data set 'microbial_stationary_phase', that I performed pca on last week.

Consider the following explanations and plots.

K-means clustering

1. Our second data set is 'microbial_stationary_phase.csv' from last week's task.

Import .csv file in R

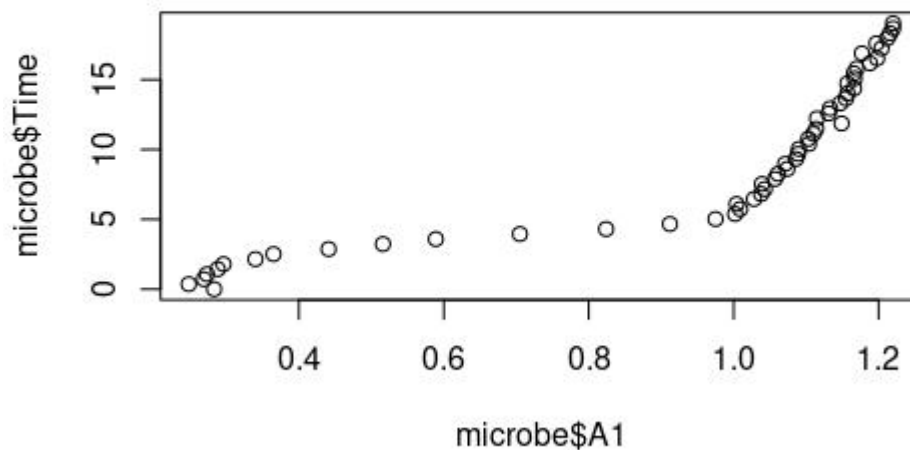
```
microbe <- read.csv(file.choose())
```

```
microbe
```

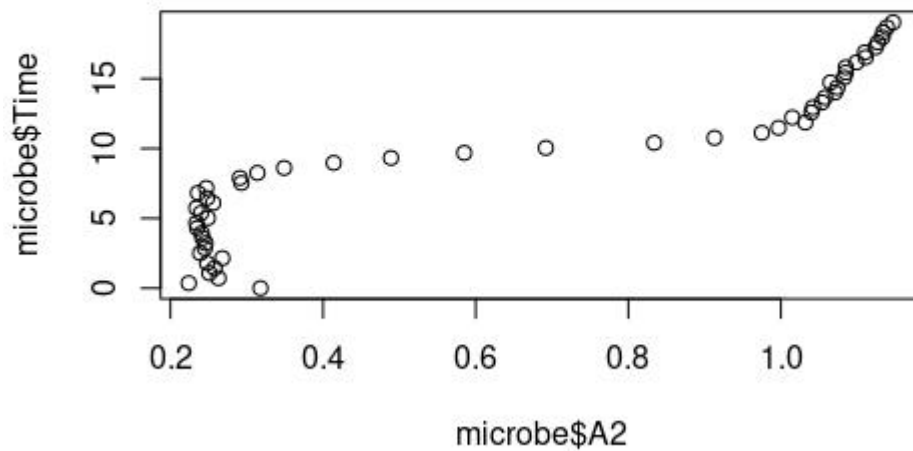
2. For the purpose of visualization, I will like to use any two columns that correlate

```
cor(microbe)
```

```
plot(microbe$A1, microbe$Time)
```



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



3. 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)
```

4. I will now perform my KMEAN Clustering for now. I want to pick 3 centers (clusters)

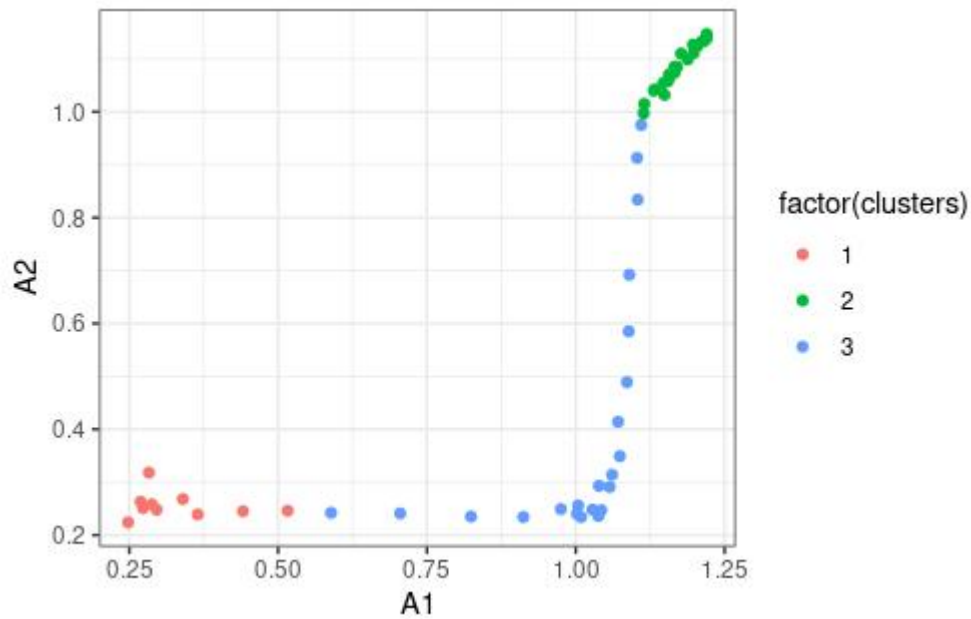
```
microbeK3 <- kmeans (x = microbe, centers = 3)
```

5. Well, I can also add this cluster information to my dataset

```
microbe$clusters <- c(microbeK3$cluster)
```

6. I want to visualize this cluster information for each microbe

```
ggplot(microbe, aes(x = A1, y = A2, color = factor(clusters))) + geom_point() +  
theme_bw()
```

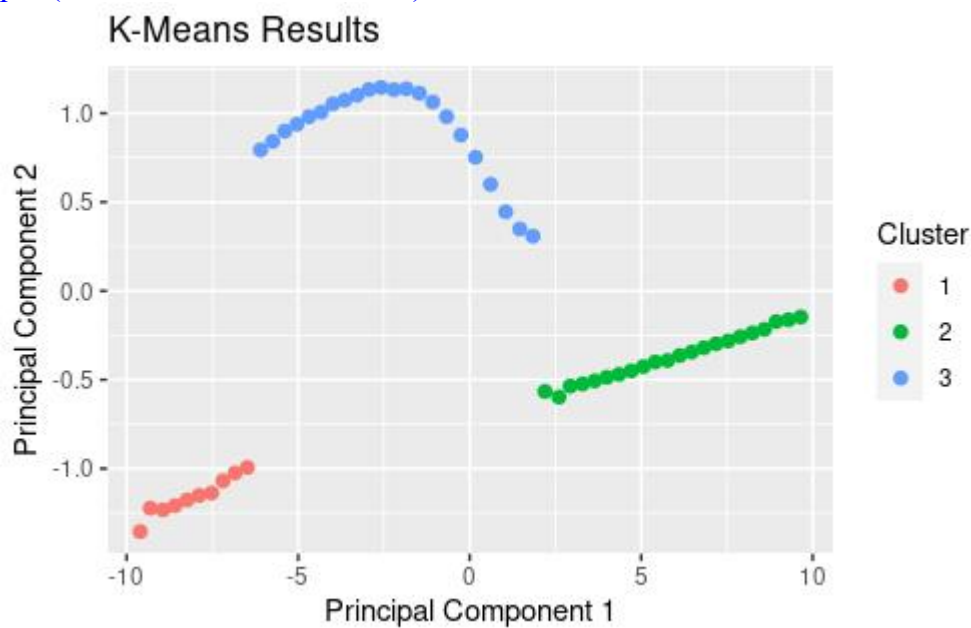


7. Install useful

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

8. Let's plot

```
plot(microbeK3, data = microbe)
```

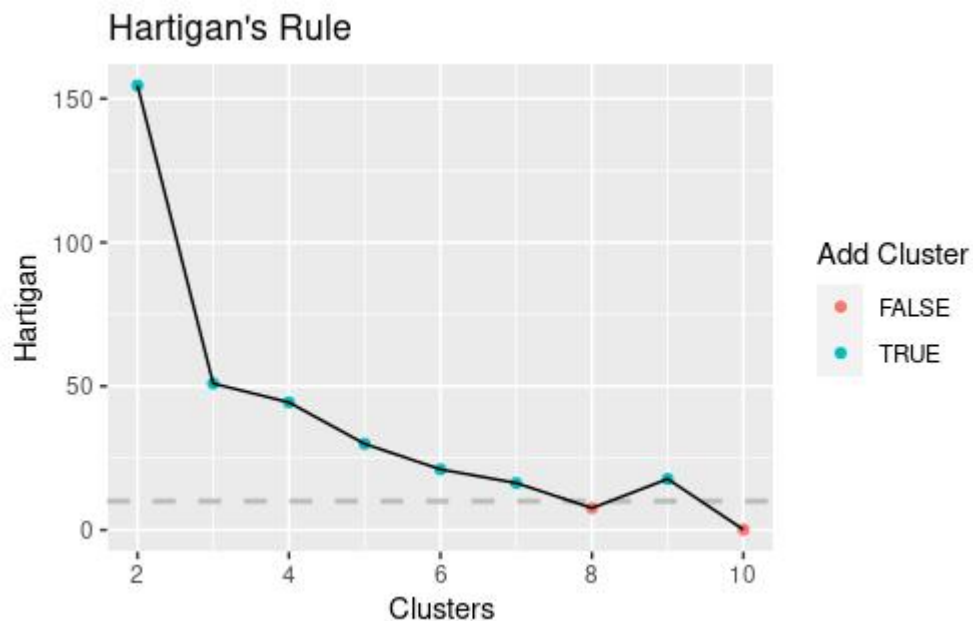


9. Choosing the right number of clusters: I will set my maximum cluster to 10.

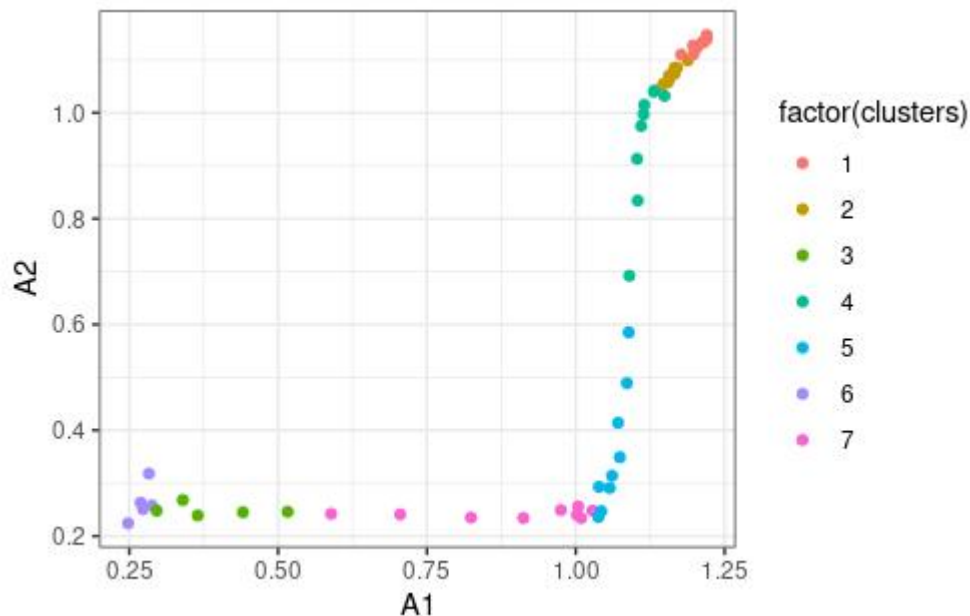
```
microbeBEST <- FitKMeans(microbe, max.clusters=10)
```

10. Once I got a FALSE, I stopped counting for clusters

```
PlotHartigan(microbeBEST)
```



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



12. 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)
```

