# MODULE 1. PEER REVIEW ASSIGNMENT-PREDICTIVE ANALYTICS & DATA MINING
## January 24th, 2021 - Aya Anisa Dwinidasari

## Deliverable 1

From the University data set given, after we adjust the raw data into the processed data as mentioned on the instruction, I conducted two methods for clustering the data. After I input the csv file into Rattle, I adjust the data into 100 partition, so it will gain 42 Seed. Then I conduct the two methods of clustering. There are:

| University | SAT | Top10 | Accept | SFRatio | Expenses | Grad |
|---|---|---|---|---|---|---|
| Harvard | 1400 | 91 | 14 | 11 | 39.525 | 97 |
| Princeton | 1375 | 91 | 14 | 8 | 30.22 | 95 |
| Yale | 1375 | 95 | 19 | 11 | 43.514 | 96 |
| Stanford | 1360 | 90 | 20 | 12 | 36.45 | 93 |
| MIT | 1380 | 94 | 30 | 10 | 34.87 | 91 |
| Duke | 1315 | 90 | 30 | 12 | 31.585 | 95 |
| Cal_Tech | 1415 | 100 | 25 | 6 | 63.575 | 81 |
| Dartmouth | 1340 | 89 | 23 | 10 | 32.162 | 95 |
| Brown | 1310 | 89 | 22 | 13 | 22.704 | 94 |
| Johns_Hopkins | 1305 | 75 | 44 | 7 | 58.691 | 87 |
| U_Chicago | 1290 | 75 | 50 | 13 | 38.38 | 87 |
| U_Penn | 1285 | 80 | 36 | 11 | 27.553 | 90 |
| Cornell | 1280 | 83 | 33 | 13 | 21.864 | 90 |
| Northwestern | 1260 | 85 | 39 | 11 | 28.052 | 89 |
| Columbia | 1310 | 76 | 24 | 12 | 31.51 | 88 |
| NotreDame | 1255 | 81 | 42 | 13 | 15.122 | 94 |
| U_Virginia | 1225 | 77 | 44 | 14 | 13.349 | 92 |
| Georgetown | 1255 | 74 | 24 | 12 | 20.126 | 92 |
| Carnegie_Mello | 1260 | 62 | 59 | 9 | 25.026 | 72 |
| U_Michigan | 1180 | 65 | 68 | 16 | 15.47 | 85 |
| UC_Berkeley | 1240 | 95 | 40 | 17 | 15.14 | 78 |
| U_Wisconsin | 1085 | 40 | 69 | 15 | 11.857 | 71 |
| Penn_State | 1081 | 38 | 54 | 18 | 10.185 | 80 |
| Purdue | 1005 | 28 | 90 | 19 | 9.066 | 69 |
| Texas_A&M | 1075 | 49 | 67 | 25 | 8.704 | 67 |

## Deliverable 2

1. Partitional Clustering (K-means Cluster Method)

From this method, first I iterate the cluster into 10 and I got the elbow plot in the figure below. The intersection between red line and blue line is located on the 3.6 horizontal axis, so we conclude that the number of clusters is approximately 3.6. After that, I input the 3.6 cluster without iteration so we got three within cluster sum of squares: **1.2948654, 0.9103000, and 0.8273532.**
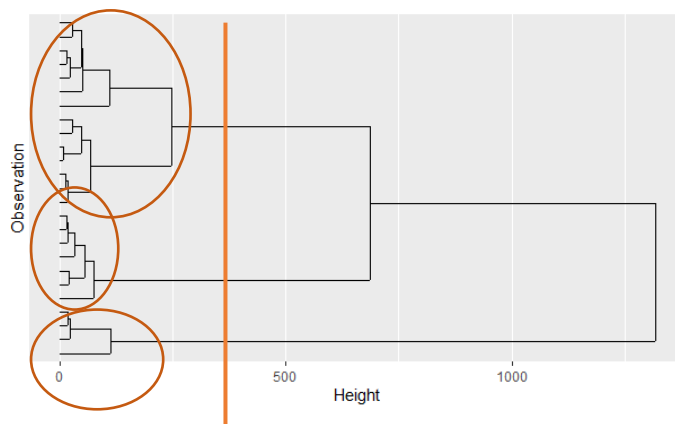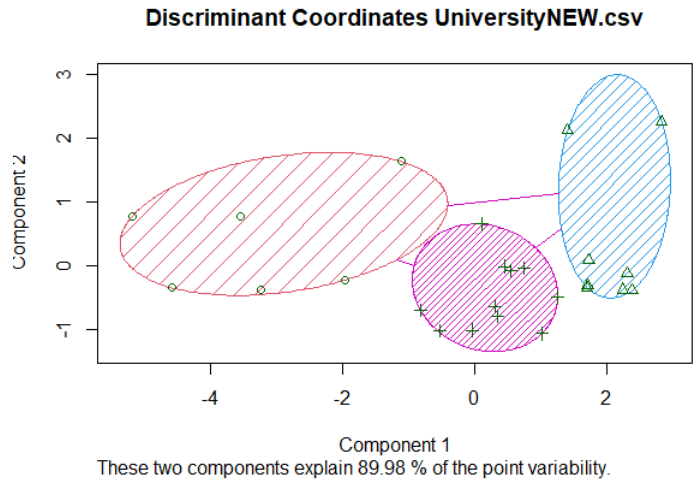




2. Hierarchical Method (dendogram method).

The second method I used is hierarchical method from discriminant and dendogram. From this method, it's shown that there are 3 clusters defined from the dendogram and from the K-means I got in the partitional methods on the step one. As we see from the discriminant coordinates with 89.98% point variability. The comparison between K-means clustering is there is an ambiguity in the dendogram, but after we compare the discriminant plot we can see clearly the cluster size so I can cut the dendogram cluster with vertical line to make the clustering clearer.

Discriminant Coordinates UniversityNEW.csv

These two components explain 89.98 % of the point variability.

### Deliverable 3

The challenge that I might find is the comparison between the two methods mightly give slight different result. As the matter of fact the K-means clustering in partitional methods need an insight of the data overall, but the dendogram methods can give a more adjusted result in clustering as well it can be helped by the discriminant plot. My finding is, from the University data set above, there are 3 cluster of "class" University from the distribution of the all variable (SAT, top 10, Accept, SF Ratio, Expenses)