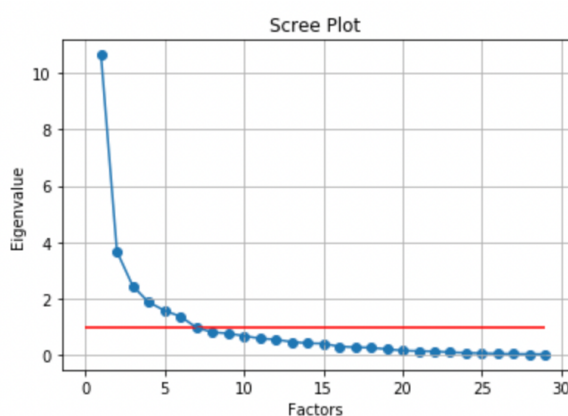# CLUSTERING ANALYSIS ON DIFFERENT COUNTRIES USING VARIOUS INDICATORS.

In this report, we are doing clustering analysis on different countries based on various indicators from the world bank data. First, we do some factor analysis to see what dimension these indicators vary and their representation. And then we do some clustering analysis based on different countries.

To avoid some potential bias caused by many indicators with significant difference, we have chosen growth indicators, instead of non-traditional ones.

By applying correlation matrix on all the indicators, we come up with the indicators which have positive relationship and negative correlation. In order to better understand the relationships between all the indicators, we have to apply factor analysis and reduce the number of indicators.

But the question still remains, how many factors to select? We can plot scree plot to determine the number of factors to select.


Scree Plot

The scree plot has number of factors on x-axis and the eigen value in the y-axis. The scree pot above shows that there could be 7 factors.

When we take a look at the heat map of all factors, we basically understand that correlation for variable and the factors. As we know that the 7 factors are representing 7 different aspects of the economy. However, some indicators represent less correlation and some high.
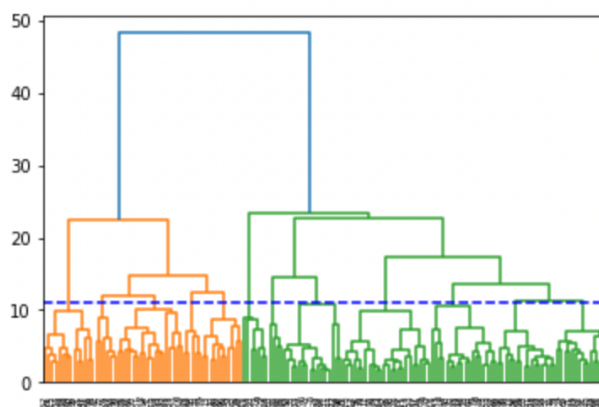
Now, we apply cluster analysis to classify the economies. We would apply most commonly used clustering method – hierarchical clustering with a bottom to top approach.

Before we apply hierarchical clustering, we have to standardize or normalize the indicators to the same scale. To prevent the scale difference between various indicators, we have to standardize the data.

After standardizing the data, we can apply clustering algorithms, by using the library, Agglomerative clustering.
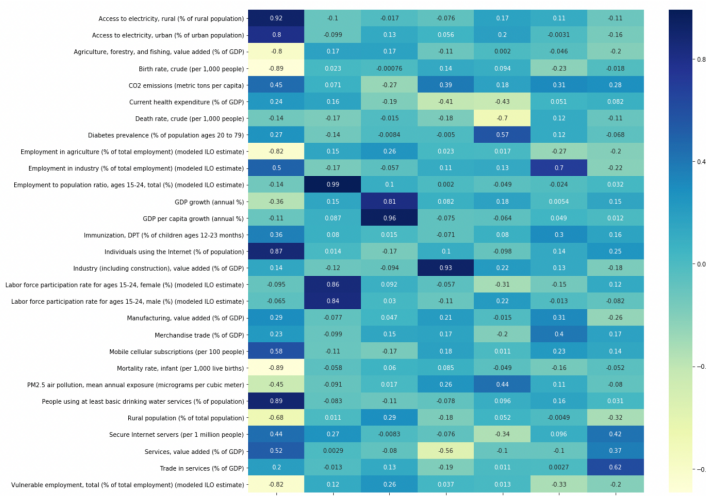
And to visualize the results of the clustering algorithm, we have to visualize it through the dendrogram, a tree like structure or diagram that represents a sequence of splits or merges.

The number of clusters finally formed is completely based on our requirement. The more the number of clusters, the more the information is detailed. If there are few clusters, then the economies may not be classified as per our judgment.



Based on the dendrogram above, there could be 12 clusters. Based on this, we can apply agglomerative clustering on the data set by 12 clusters and Euclidian distance and the linkage as ward's method.

After clustering many countries into 12 clusters, we perform a heat map, we observe that the 20 different variables have high correlation with only 7 factors as mentioned before.



Based on clusters characteristics, we further divided the entire clusters into 4 big categories – Most develop countries, developing, less developed and undeveloped economies.

Under most developed countries we see – US, UK and France. Similarly, we see countries from clusters 3,4,10 are least developed countries in economy.

It would be interesting if we would compare our clustering and classification with the world's bank data of high, upper middle, lower middle and low-income levels.