# Clustering Assignment

Prepared By:

Mohammad Aftab Alam

# Abstract

## Problem statement :-

HELP International is an international humanitarian NGO that is committed to fighting poverty and providing the people of backward countries with basic amenities and relief during the time of disasters and natural calamities. It runs a lot of operational projects from time to time along with advocacy drives to raise awareness as well as for funding purposes.

After the recent funding programmes, they have been able to raise around $ 10 million. Now the CEO of the NGO needs to decide how to use this money strategically and effectively. The significant issues that come while making this decision are mostly related to choosing the countries that are in the direst need of aid.

And this is where you come in as a data analyst. Your job is to categorise the countries using some socio-economic and health factors that determine the overall development of the country. Then you need to suggest the countries which the CEO needs to focus on the most. The datasets containing those socio-economic factors and the corresponding data dictionary

# Analysis Methodology

**Reading and Understanding of Data**
- Importing the .csv file
- Examine the data

**Data Quality**
- Missing value
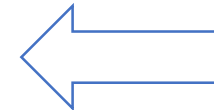- Duplicate data
- Spelling mistake checking

**Checking Outliers**
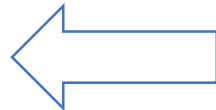- Removing the outlier where ever required as per understanding the problem statement.

**PCA (Principle Component Analysis)**
- To derive principal components.
- To check the variance ratios.
- Scree plot-plotting the cumulative variance against the number of components.
- Going ahead and doing dimensionality reduction using incremental PCA.
- Reducing the correlation to almost zero

**Scaling of the data**
- Standardizing all the continuous variables.

**Data Visualization**
- Visualizing few original data variables to look for any pattern or correlation.

**K means clustering**
- Identify the 'k' by silhouette analysis and sum of squared distances graph.
- Forming n –clusters on PCA modified data.
- Visualizing the clusters with various variables.
- Analyzing the clusters.
- Identifying the countries which requires aid.

**Hierarchical Clustering**
- Identify the 'n' via dendrogram.
- Forming n –clusters on PCA modified data.
- Visualizing the clusters with various variables.
- Analyzing the clusters.
- Identifying the countries which requires aid.

**Decision Making**
- Identifying the countries which requires aid by analyzing both K-means and Hierarchical Clustering results.

# Reading and Understanding of Data



- Importing the .csv file using pandas.

- Reading the .csv file.

- Checking Describe details

# Reading and Understanding of Data Cont...



- Checking size of the data.

- Checking infor all other details

# Data Quality



No Missing value found.

No Duplicate data found.

So we will move to the further step.

# Checking Outliers

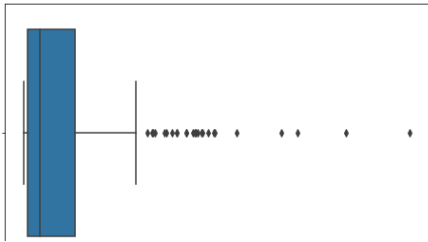## 'gdpp' variables

### 3. Checking Outliers

```
In [13]:  # Building boxplot for checking outliers

          def plot(col,x,y):
              plt.figure(figsize=(x,y))
              plt.subplot(1,1,1)
              sns.boxplot(col, data= country_details)
              return
```
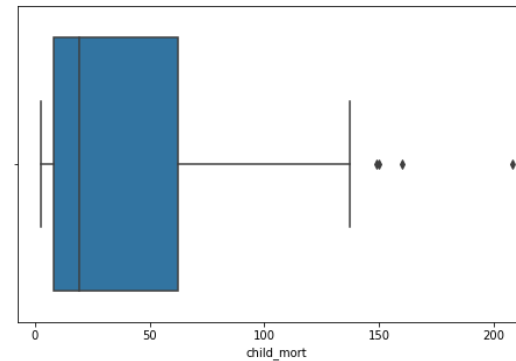
### 3.1.1 'gdpp'

```
In [14]:  plot('gdpp',8,5)
```



## 'child_mort' variables
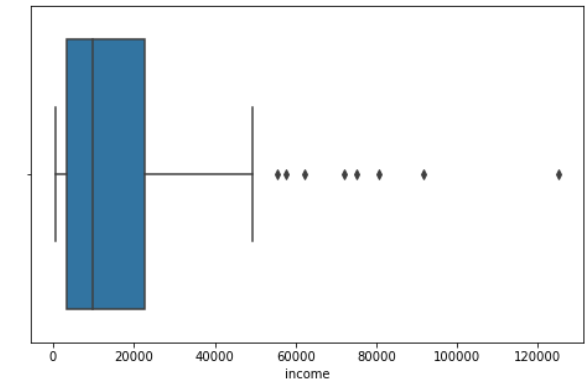
### 3.1.2 'child_mort'

```
In [15]:  plot('child_mort',8,5)
```



## 'income' variables
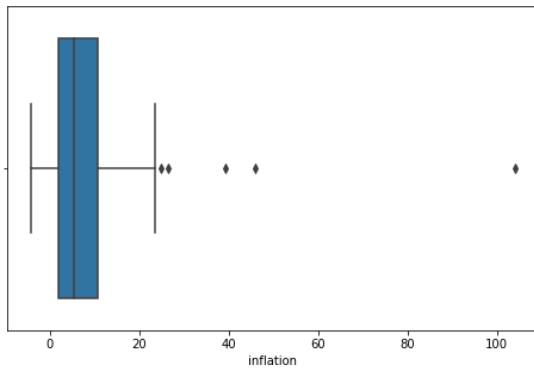
### 3.1.3 'income'

```
In [16]:  plot('income',8,5)
```



## 'inflation' variables

### 3.1.4 'inflation'

```
In [17]:  plot('inflation',8,5)
```
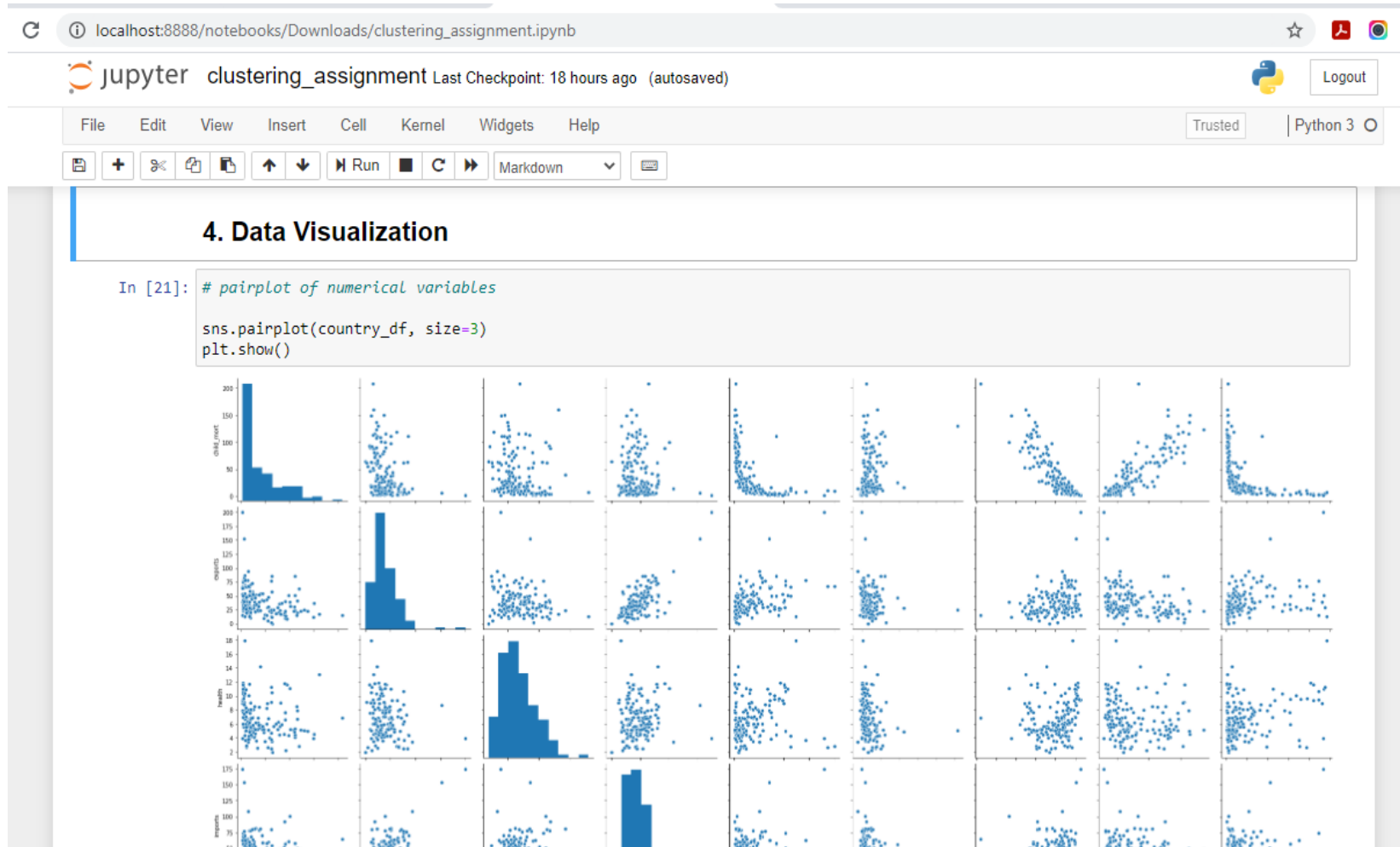


**After seeing the outliers of some important variables we will conclude that :-**

- Variables 'gdpp', 'child_mort', 'income' and 'inflation' columns has high outliers.

- As of now, not remove the outliers as it may suits the business needs or a lot of countries are getting removed.

# Data Visualization
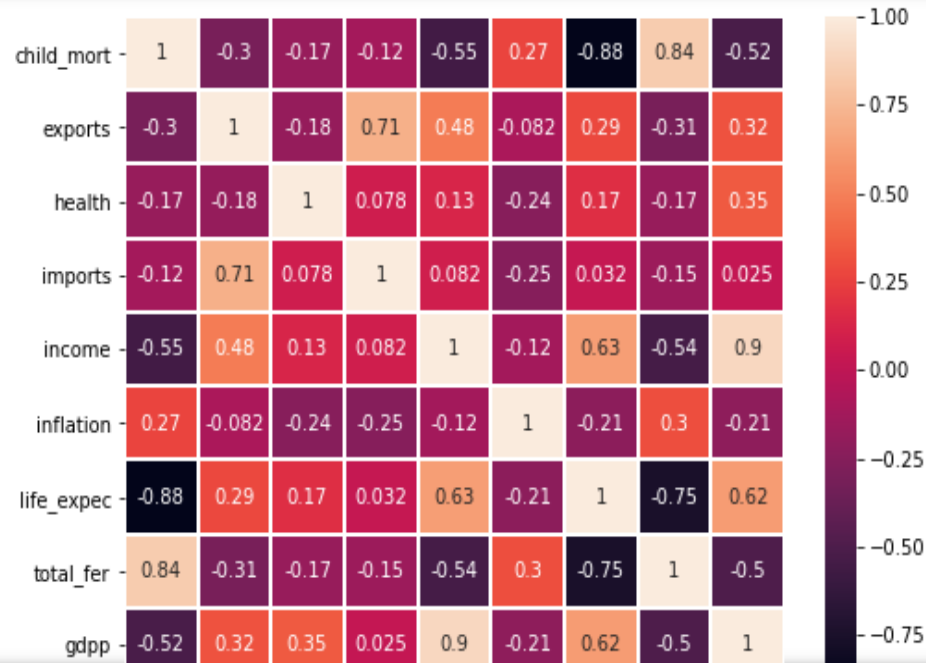
Visualization using Pairplot



- By seeing Pairplot, we are getting the clear result about the correlation between the variables.

- Thus we can go through for Heatmap for more clear result.

# Data Visualization Cont...

## Visualization using Heatmap

```
In [72]:  # let draw a heatmap to understand the correlation better
          plt.figure(figsize=(8,6))
          sns.heatmap(country_df.corr(), annot=True, linewidth= 1)
          plt.show()
```



- ➢ By plotting heatmap, we can see high correlation between:
  - ▪ 'child_mort' and 'total_fer'.
  - ▪ 'exports' and 'imports'.
  - ▪ 'income' and 'gdpp'

- ➢ This will cause problem for the upcoming analysis, hence need to be removed but they have valuable information which we can't afford to loose.

- ➢ So, we will use PCA to overcome this multicollinearity.

- ➢ This will not only take care of multicollinearity but also will preserve the valuable information and also demensionality reduction.

# Scaling of the data



- We are using scaling of the data for standardizing all the continuous variables.

# PCA (Principle Component Analysis)



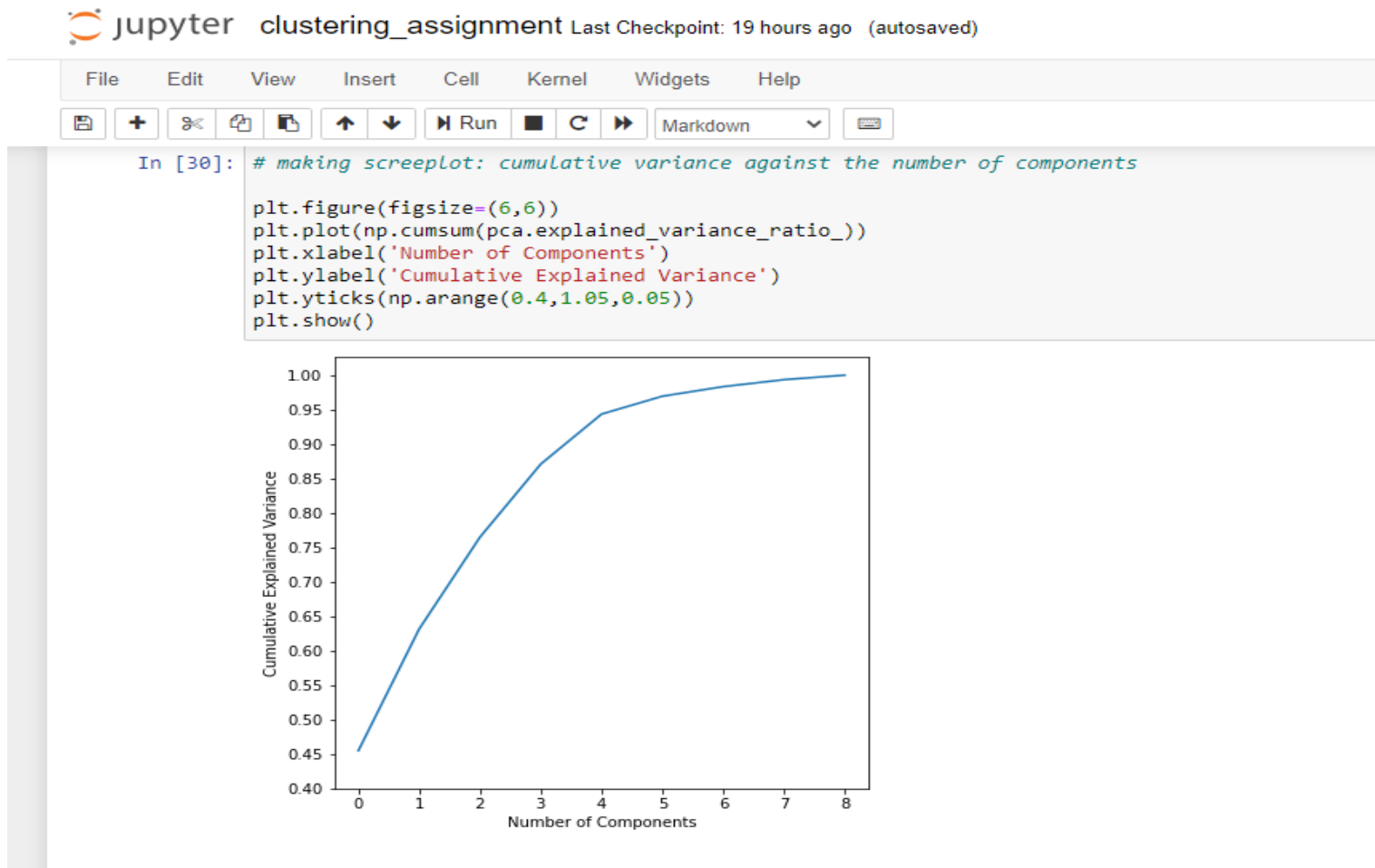Visualising the features by loaded along PC1 and PC2

- We see that features like gdpp, life-expectancy and income are along the direction of PC1 and other features like total-fertility and child- mortality are along PC2 direction.

# PCA (Principle Component Analysis) Cont...



```
In [30]: # making screeplot: cumulative variance against the number of components

plt.figure(figsize=(6,6))
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')
plt.yticks(np.arange(0.4,1.05,0.05))
plt.show()
```

**Conclusion:**

- Around 95% of the cumulative variance is being explained by 5 components

Scree-plot, plotting b/w the cumulative variance against the number of components.
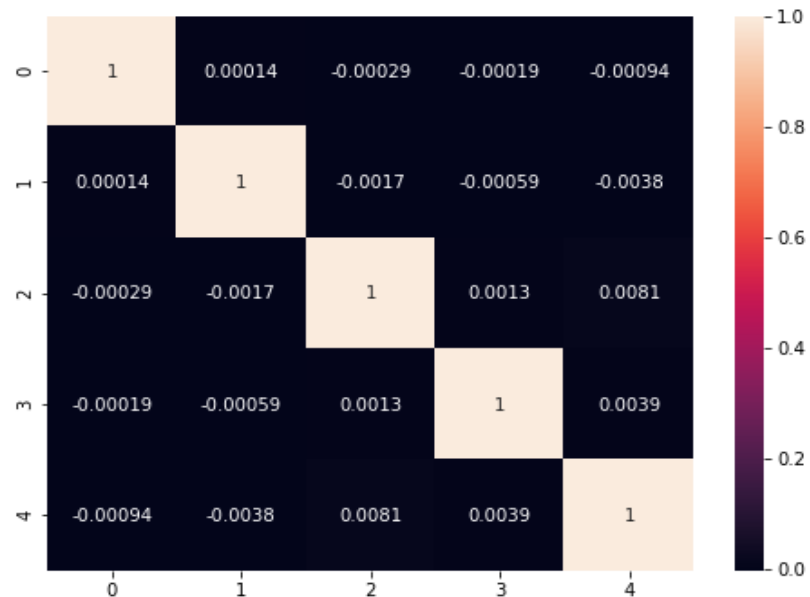
# PCA (Principle Component Analysis) Cont...



After doing dimensionality reduction via incremental PCA by taking 5 components, we see that the correlation in the data has almost reduced to zero.
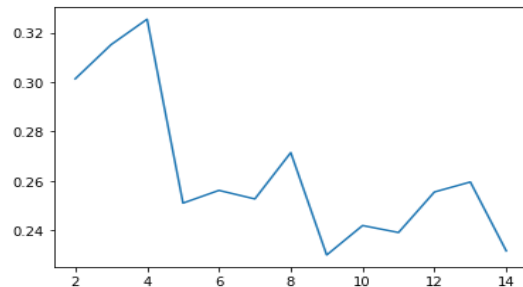
# K means clustering



Silhouette Analysis



Sum of squared distance

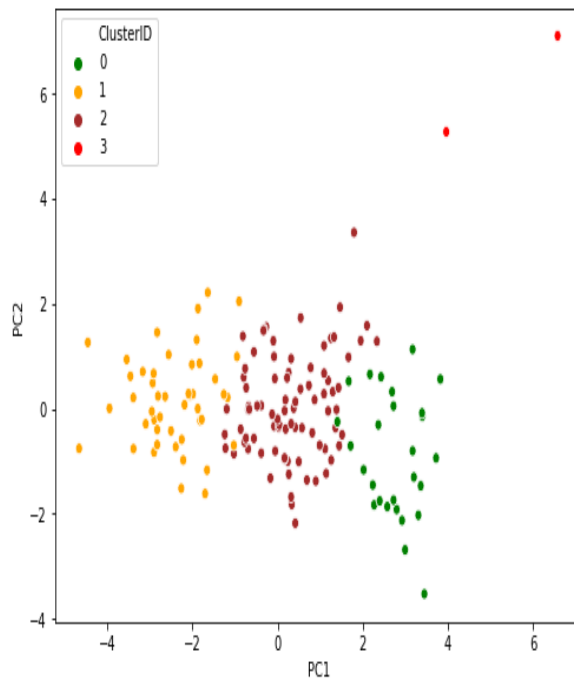- By looking silhouette analysis, we see the highest peak is at k=4 and in sum of squared distances graph, we see that the elbow is in the range of 3 to 5, so we are going ahead with k as 4.
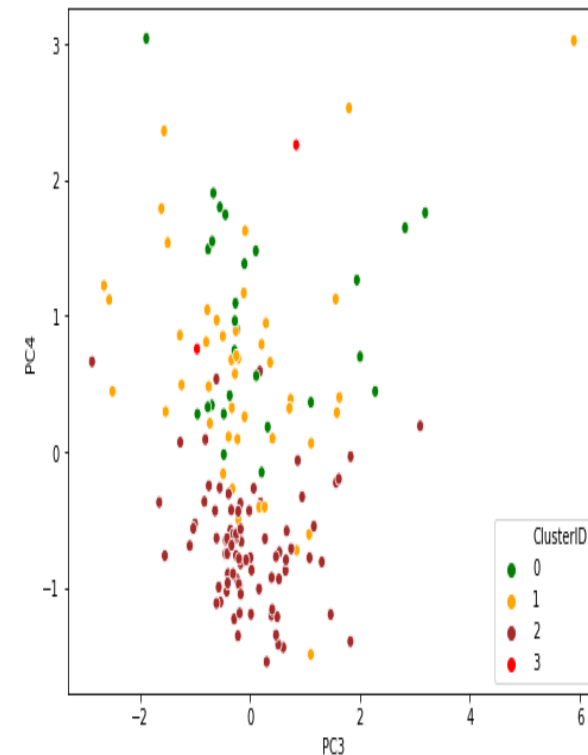
# K means clustering Cont...



Scatterplot Between PC1 and PC2 and Cluster-ID

Scatterplot Between PC3 and PC4 and Cluster-ID

Scatterplot b/w gdpp, child_mort and ClusterID

Scatterplot b/w income, gdpp and ClusterID

Scatterplot b/w child_mort, incomeand ClusterID

**Scatterplot Between Actual Variables and Cluster-ID**

# K means clustering Cont...

```
In [56]: plt.figure(figsize=(25,8))
         plt.subplot(1,3,1)
         plt.title('child_mort_mean', size=20)
         sns.barplot(x='ClusterID', y='child_mort_mean', data= country_analysis_df)
         plt.subplot(1,3,2)
         plt.title('gdpp_mean', size=20)
         sns.barplot(x='ClusterID', y='gdpp_mean', data= country_analysis_df)
         plt.subplot(1,3,3)
         plt.title('income_mean', size=20)
         sns.barplot(x='ClusterID', y='income_mean', data= country_analysis_df)
         plt.show()
```

- By seeing the graph, we can conclude that 'cluster-1' is our cluster of concern because :-

  - It has highest child_mort.
  - Lowest gdpp.
  - And lowest income.

**Analysing the Clusters by using these three variables- 'gdpp', 'child_mort', 'income'**

As per K-Means clustering, the country which are direst need of aid are :
- Burundi
- Liberia
- Congo, Dem, Rep.
- Niger
- Sierra Leone
- Madagascar
- Mozambique
- Central African Republic
- Malawi
- Eritrea

# Hierarchical Clustering

```
dendrogram(mergings)
plt.show()
```



- By single method Hierarchical clustering the things are not clear.
- so now we go for complete method hierarchical clustering.

Single Method Hierarchical Clustering

# Hierarchical Clustering Cont...



Complete Method Hierarchical Clustering

Scatterplot b/w gdpp, child_mort and ClusterID

Scatterplot b/w child_mort , income and ClusterID

Scatterplot b/w child_mort , income and ClusterID

**Scatterplot Between original datapoint and cluster-id**

# Hierarchical Clustering Cont...

```
In [68]:  plt.figure(figsize=(25,8))
          plt.subplot(1,3,1)
          plt.title('child_mort_mean_hc', size=20)
          sns.barplot(x='ClusterID', y='child_mort_mean_hc', data= country_analysis_df_hc)
          plt.subplot(1,3,2)
          plt.title('gdpp_mean', size=20)
          sns.barplot(x='ClusterID', y='gdpp_mean_hc', data= country_analysis_df_hc)
          plt.subplot(1,3,3)
          plt.title('income_mean', size=20)
          sns.barplot(x='ClusterID', y='income_mean_hc', data= country_analysis_df_hc)
          plt.show()
```



- By seeing above graph, we can conclude that here are two clusters that is 'cluster- 3' and 'cluster- 0' for our cluster of concern but we leave 'cluster- 3' as it has only 1 country and we use 'cluster- 0' because:-

  - Cluster-0 has highest child_mort.
  - Lowest gdpp.
  - And lowest income.

**Analysing the Hierarchical-Clusters by using these three variables- 'gdpp', 'child_mort', 'income'**

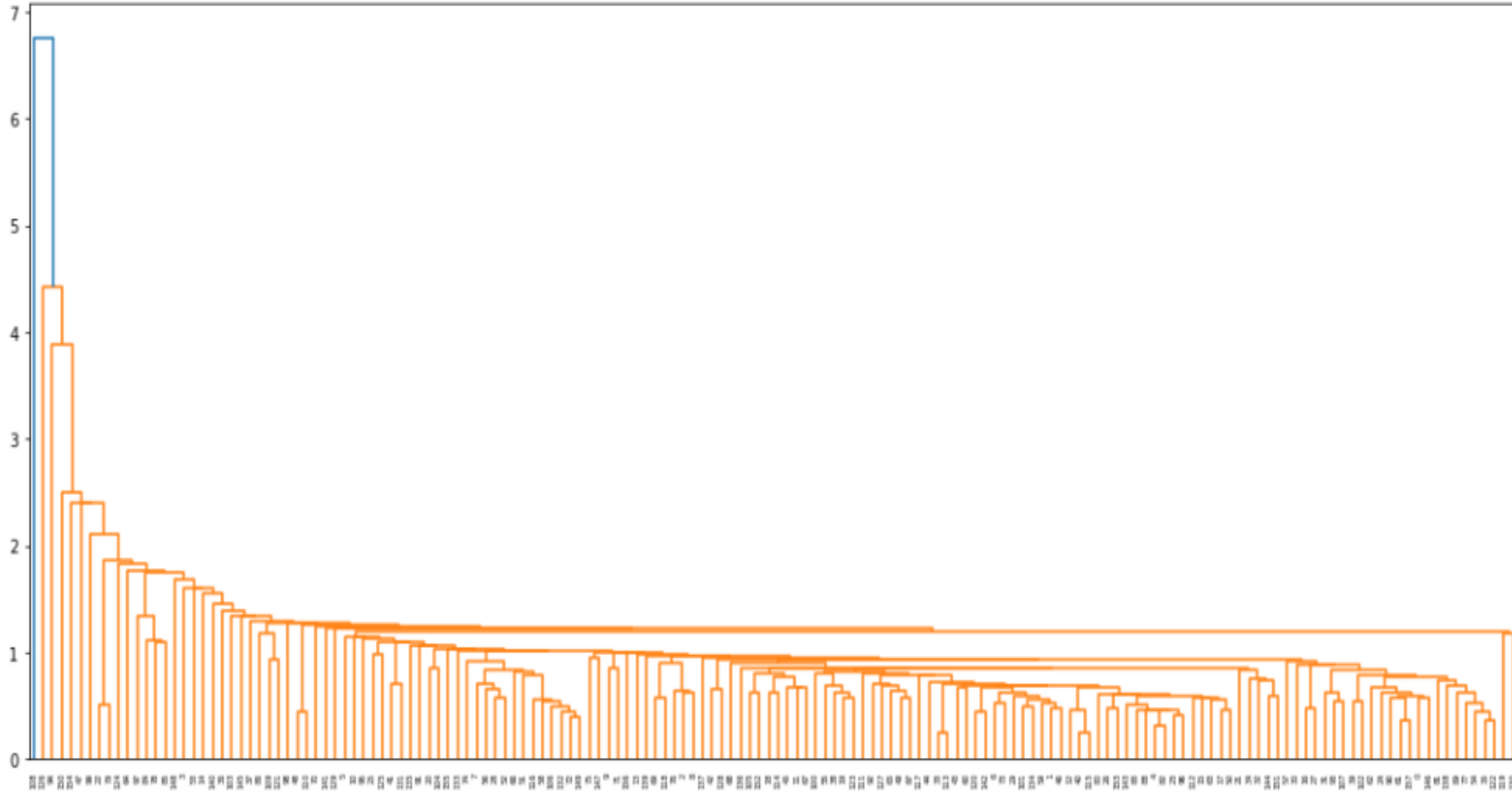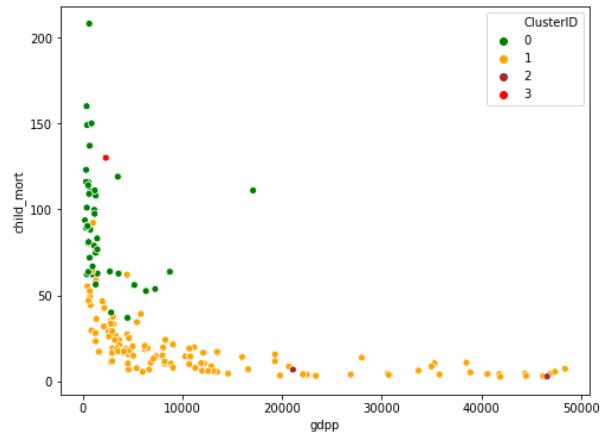**As per Hierarchical clustering, the country which are direst need of aid are :**

- Burundi
- Liberia
- Congo, Dem, Rep.
- Niger
- Sierra Leone
- Madagascar
- Mozambique
- Central African Republic
- Malawi
- Togo

# Decision Making

We got same countries by both K-Means and Hierarchical clustering. Therefore following are the countries ehic are direst need of aid by considering socio-economic factor into consideration:-

- Burundi
- Liberia
- Congo, Dem, Rep.
- Niger
- Sierra Leone
- Madagascar
- Mozambique
- Central African Republic
- Malawi