

CLUSTERING AND PCA ASSIGNMENT

Sushma Subburayan

CLUSTERING AND PCA ASSIGNMENT

Problem Statement: Categorize the countries using some socio-economic and health factors that determine the overall development of the country and suggest at least five countries which are in direst need of aid.

DATA :

Country-data data set containing the socio economic factors of 167 countries

GOALS OF THE ANALYSIS:

To identify at least top five countries which are in the direst need of aid

- The following are the socio economic factors given for the countries
 - Name of the country
 - Child Mortality - Death of children under 5 years of age per 1000 live births exports
 - Exports of goods and services. Given as %age of the Total GDP
 - Health - Total health spending as %age of Total GDP
 - imports - Imports of goods and services. Given as %age of the Total GDP
 - Income - Net income per person
 - Inflation - The measurement of the annual growth rate of the Total GDP
 - Life Expectancy - The average number of years a new born child would live if the current mortality patterns are to remain the same
 - Total Fertility - The number of children that would be born to each woman if the current age-fertility rates remain the same.
 - GDPP - The GDP per capita. Calculated as the Total GDP divided by the total population.

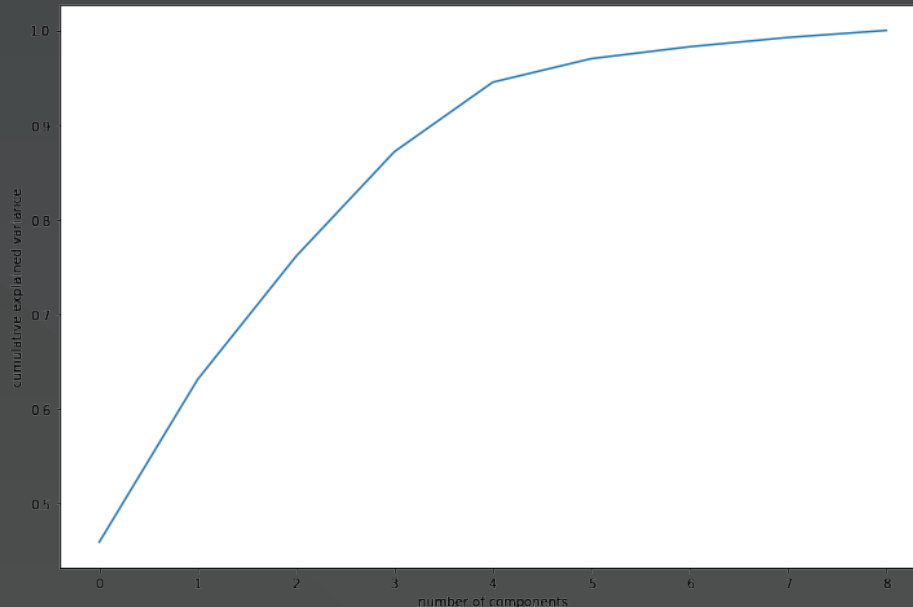
- 1 .Read and understood the given data set.
- 2. Cleaned the data
 - > a. Checked for null values
 - > b. Checked for duplicate records
- 3. Visualized the data
 - > a. Created pair plots for all the numeric variables in the data set
 - > b. Created box plots for all the categorical variables in the data set
- 4. Data Preparation
 - > a. Tried removing outliers(before PCA and after PCA),but as the outlier elimination is leading to loss of many countries and giving incorrect results ,went ahead without eliminating the outliers.
 - > b. Scaled all the numeric variables using standard scaler.

- 5. Hopkins Statistics

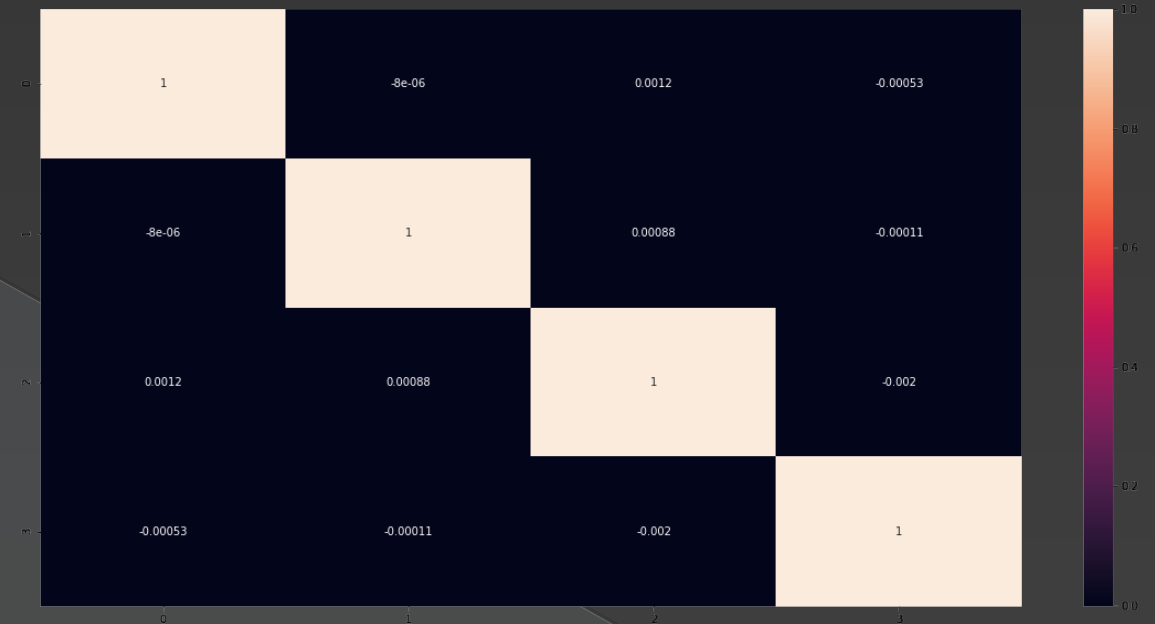
As the Hopkins statistics value was good, went ahead with PCA

6. PCA

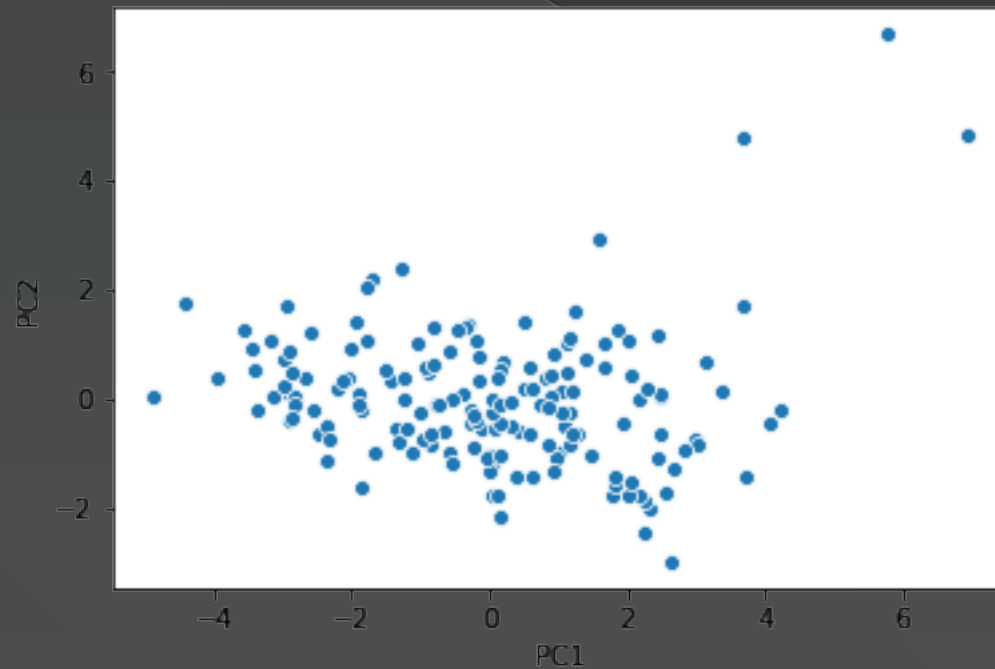
- > a. Performed PCA on the scaled data
- > b. Max variance was explained by the first 4 PCs (based on `pca.explained_variance_ratio_` and the scree plot)



- > c. Performed incremental PCA for efficiency with 4 components - four components were selected based on the scree plot
- > d. Created correlation matrix and heatmap for the principal components



e. Made a scatter plot to visualize the PCs

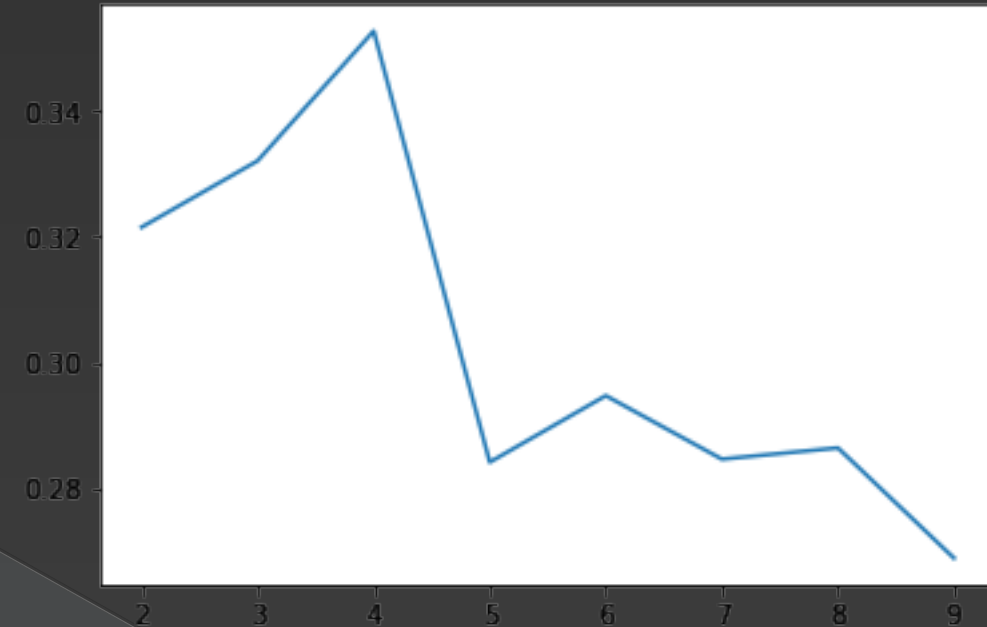
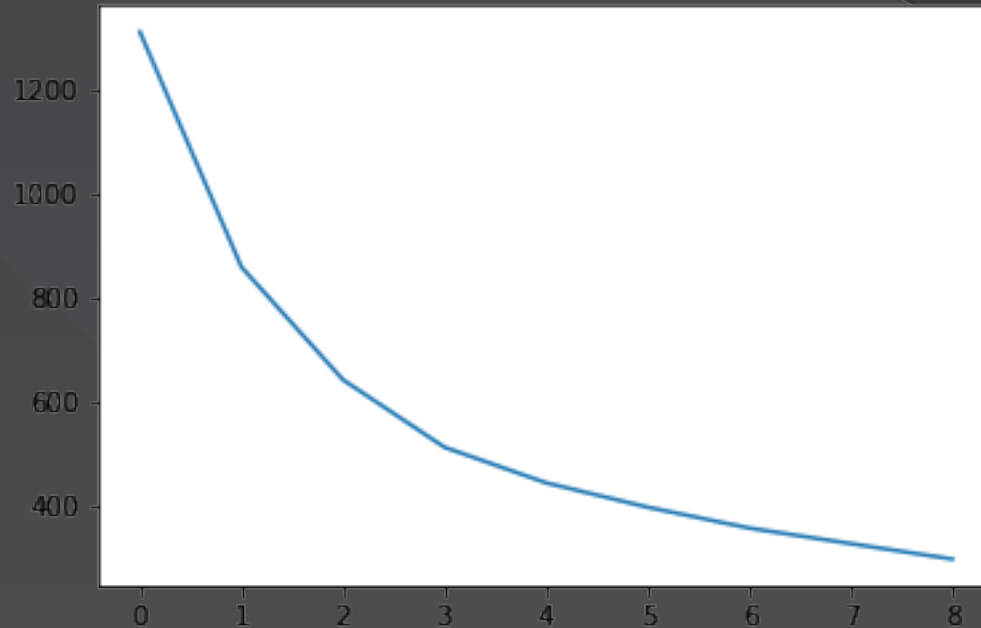


7. Hopkins Statistics

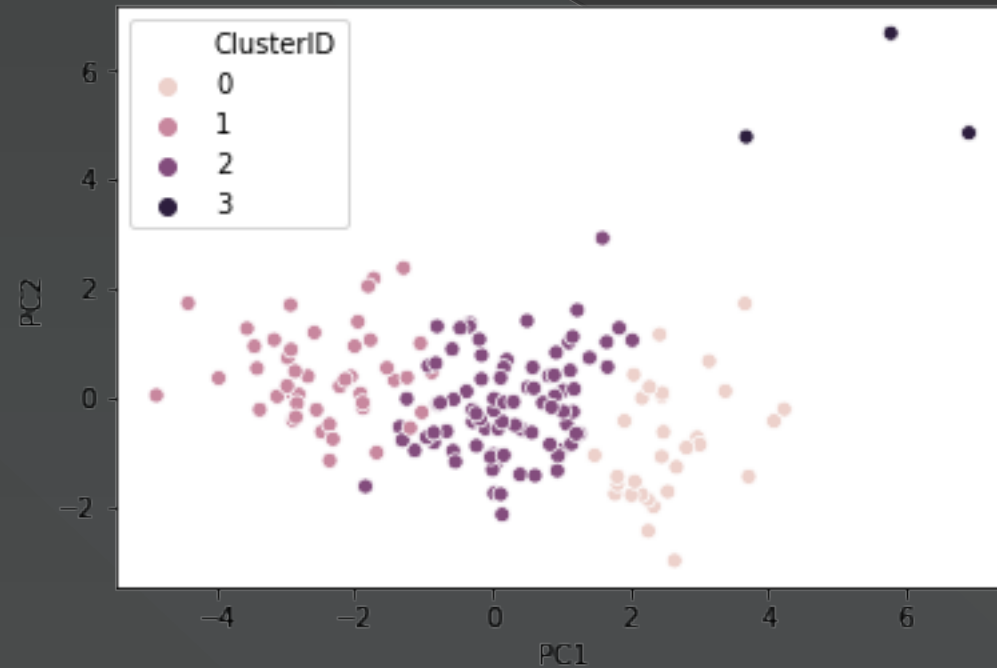
- As the Hopkins statistics for PCs was good, went ahead with K-Means Clustering

6. K-Means Clustering

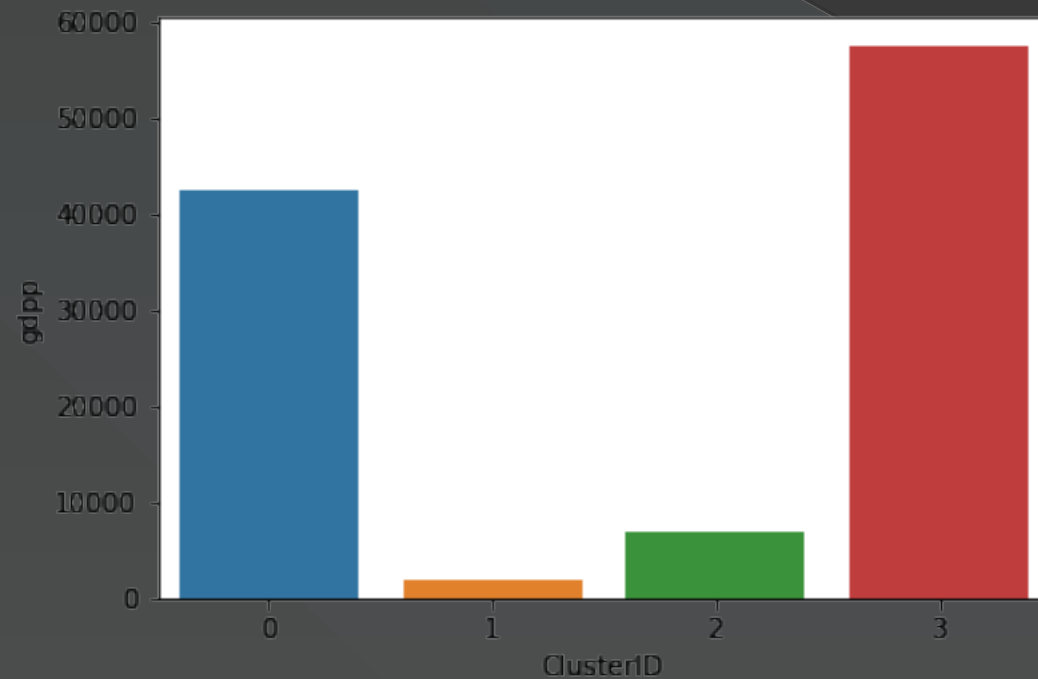
- a. Did silhouette score analysis and created elbow curve to identify the number of clusters

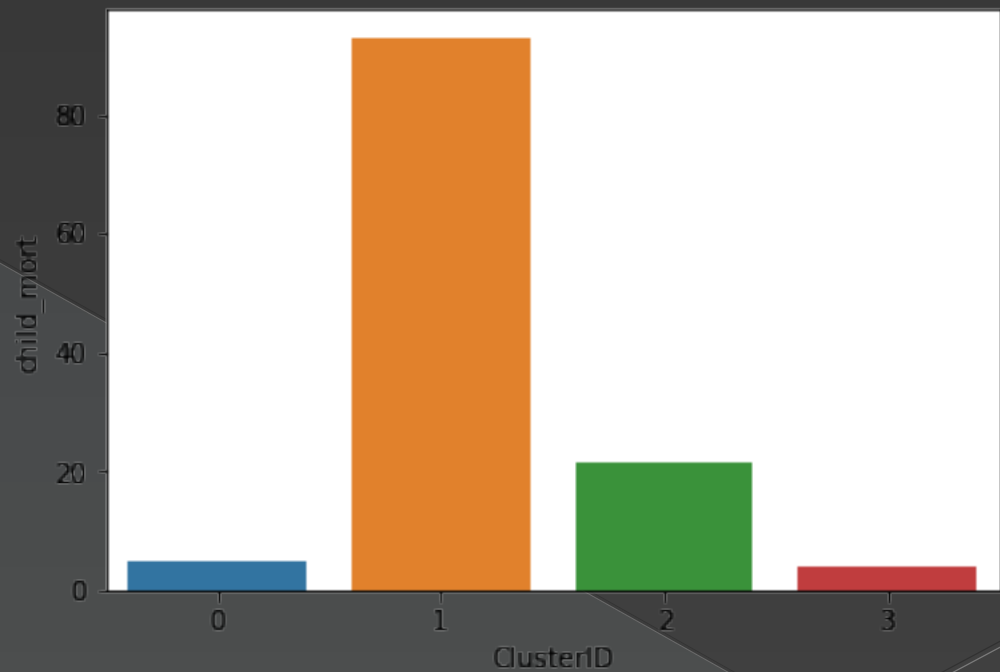
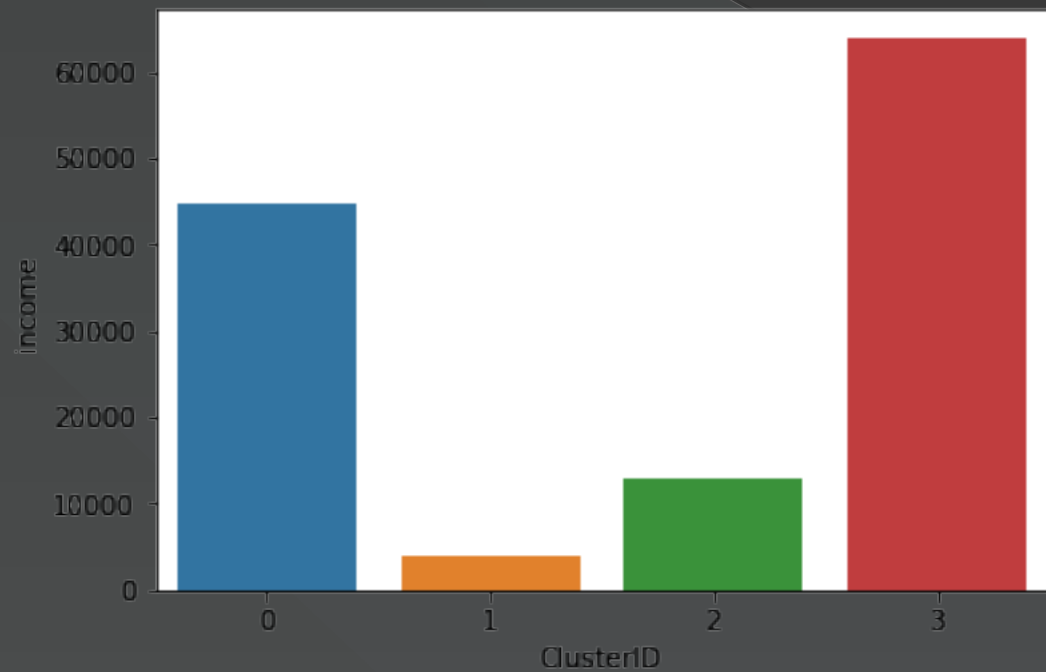


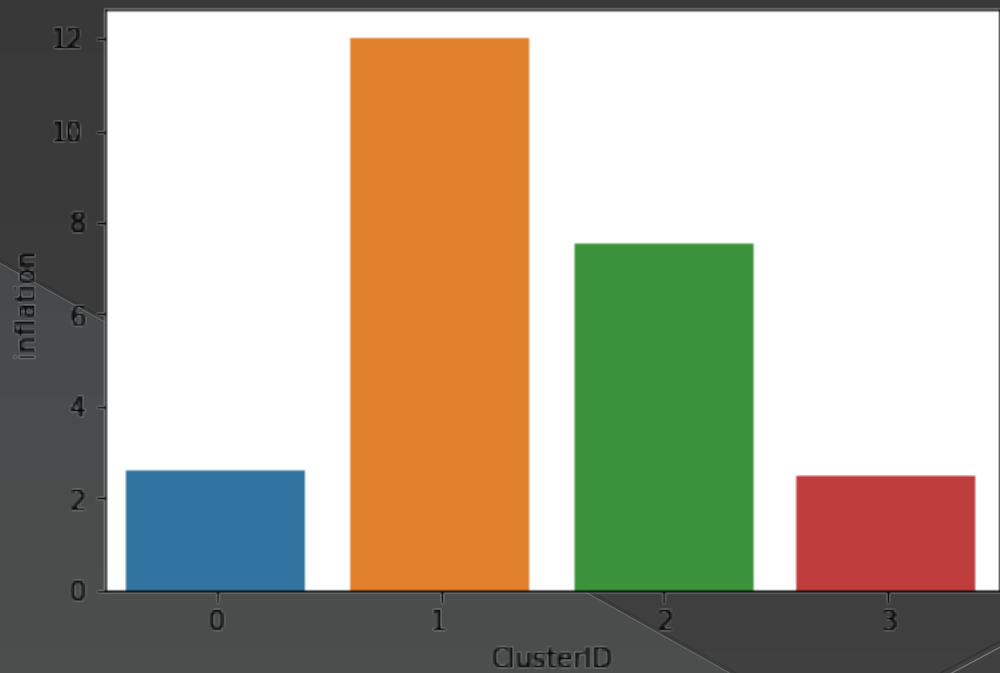
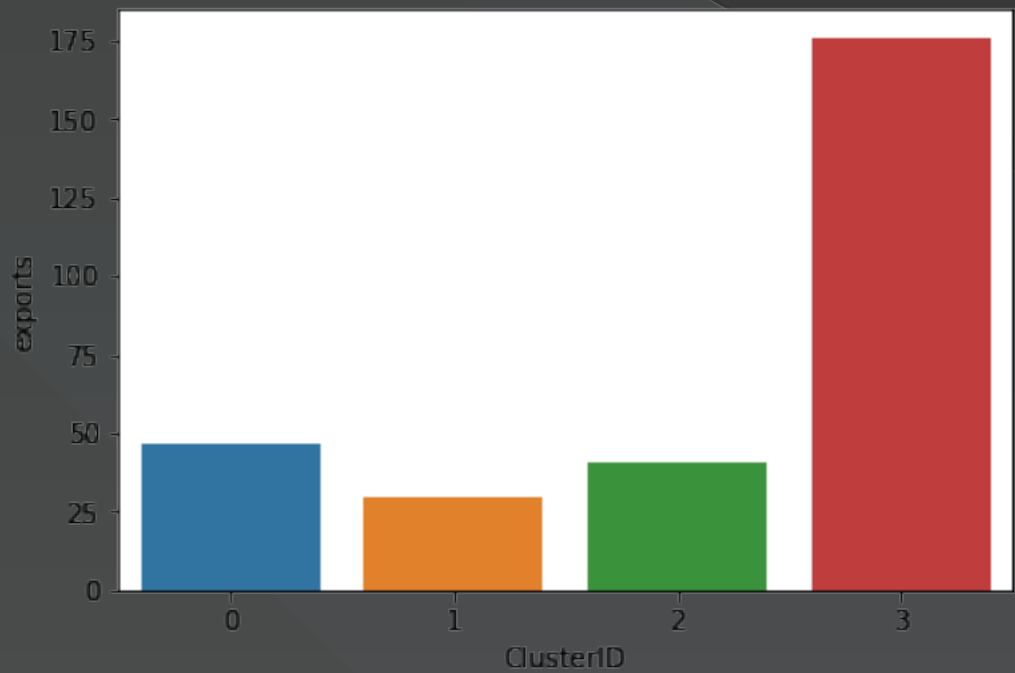
- > b. Performed K-Means clustering with four clusters
- > c. Analyzed the clusters formed

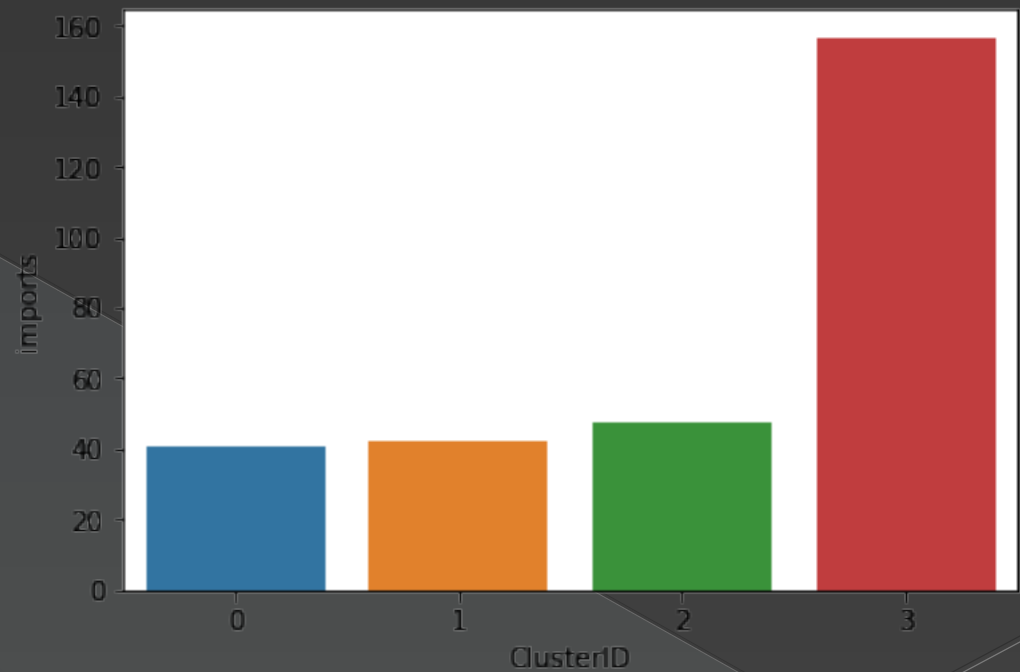
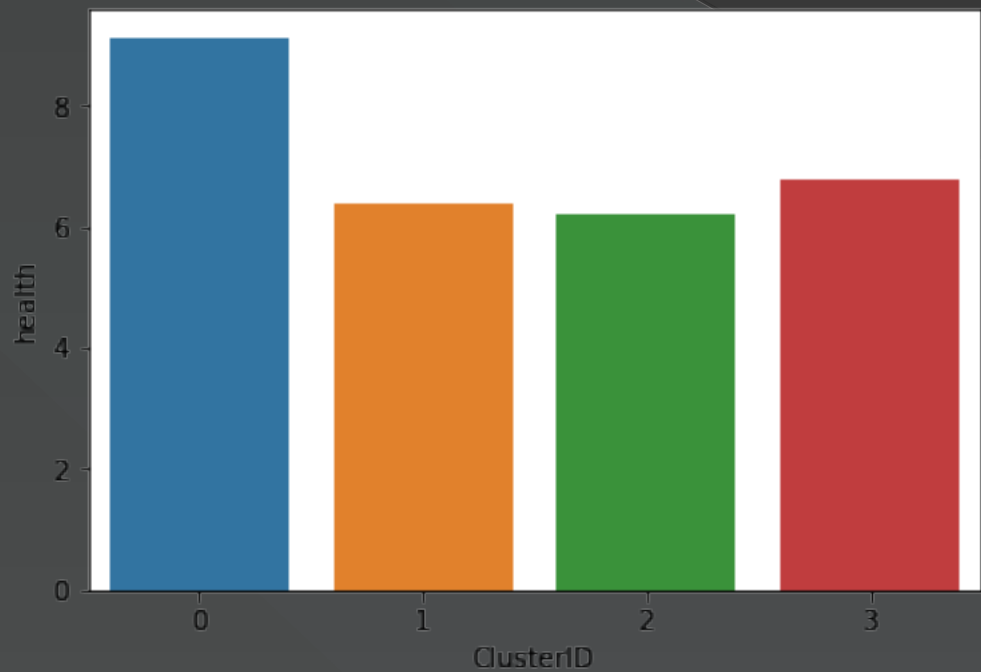


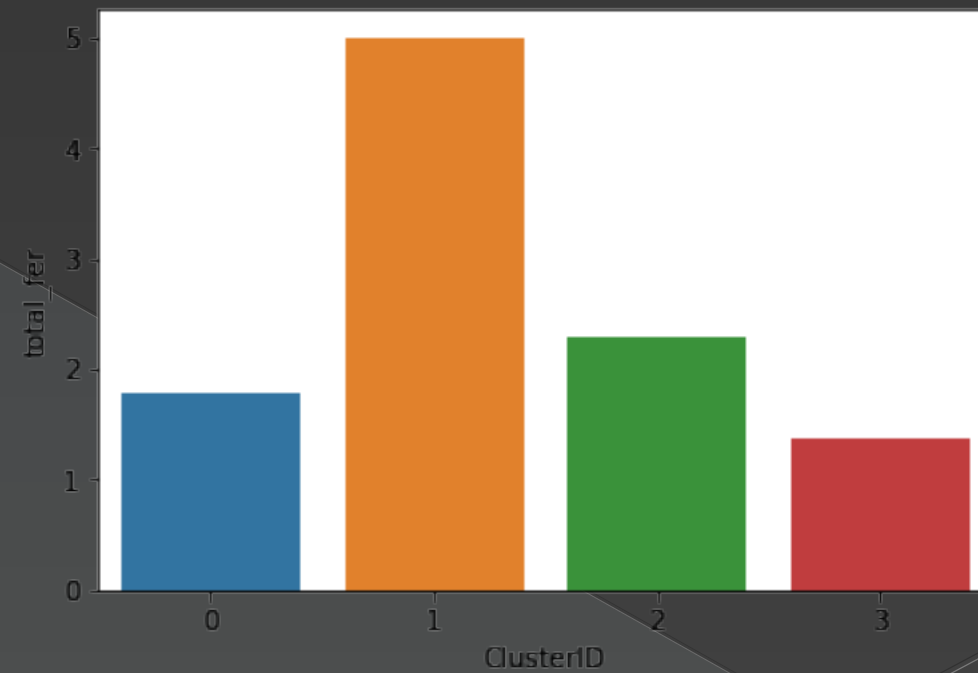
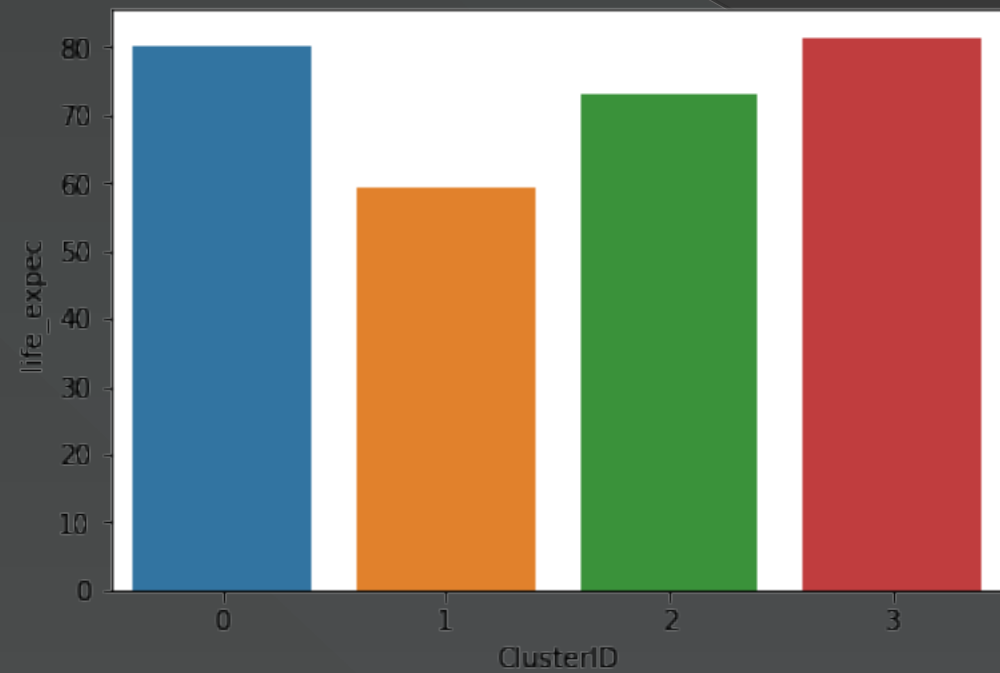
- > d. Plotted bar graphs between the clusters and the features and identified the clusters with countries in the direst need of aid (Clusters 1 and 2)



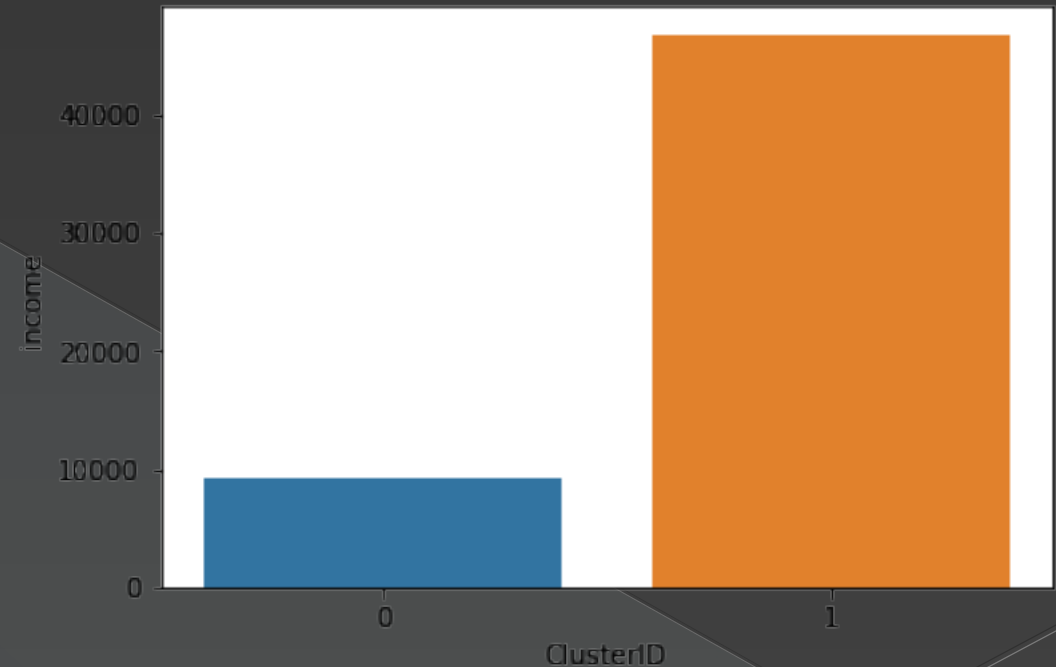


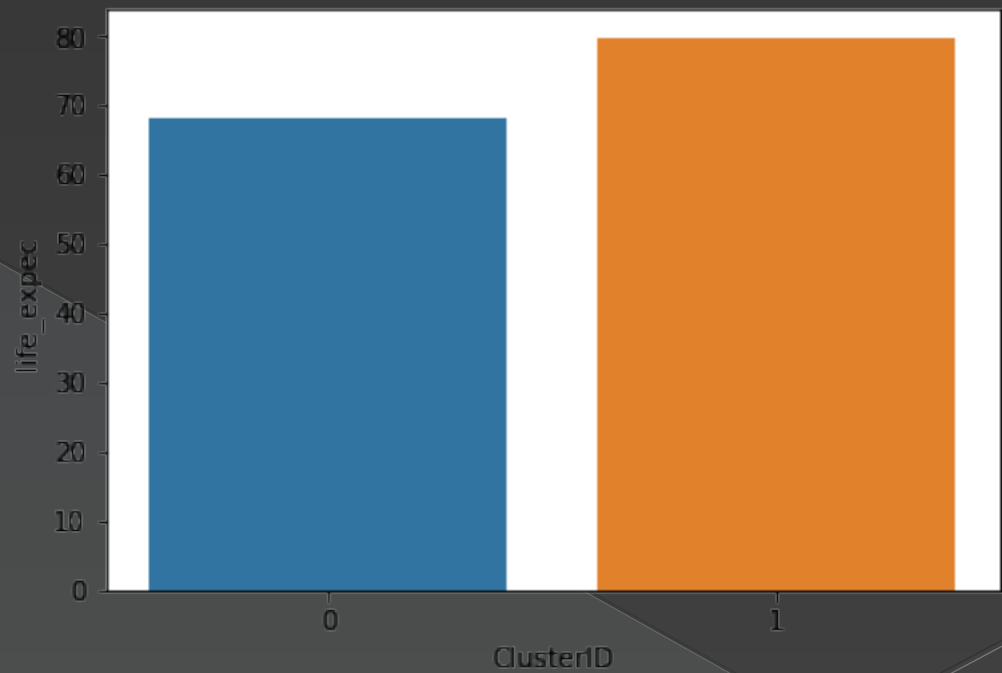
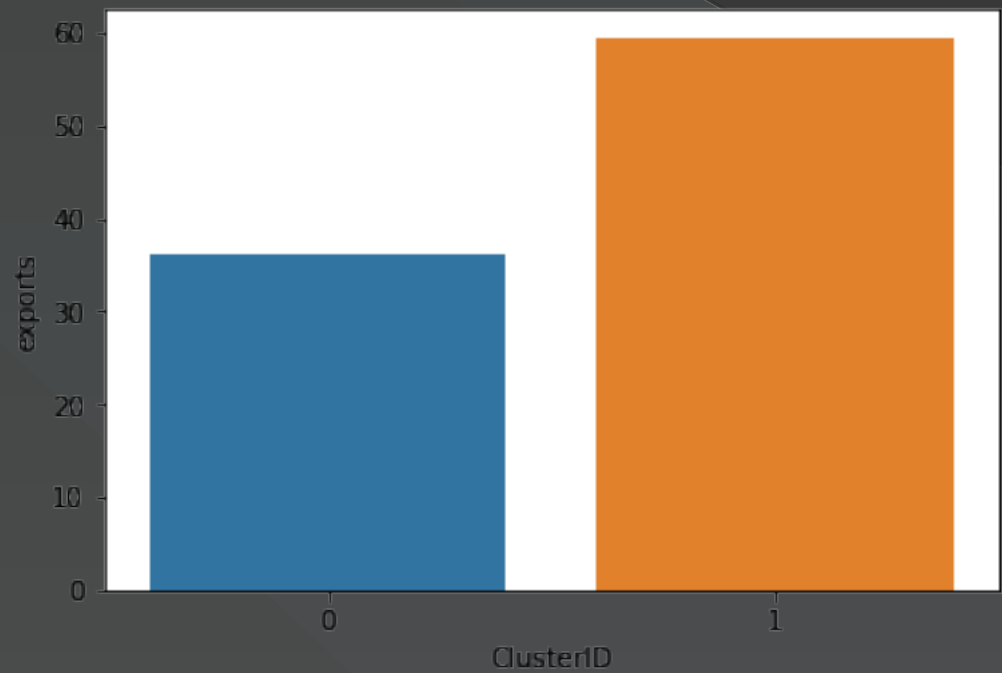


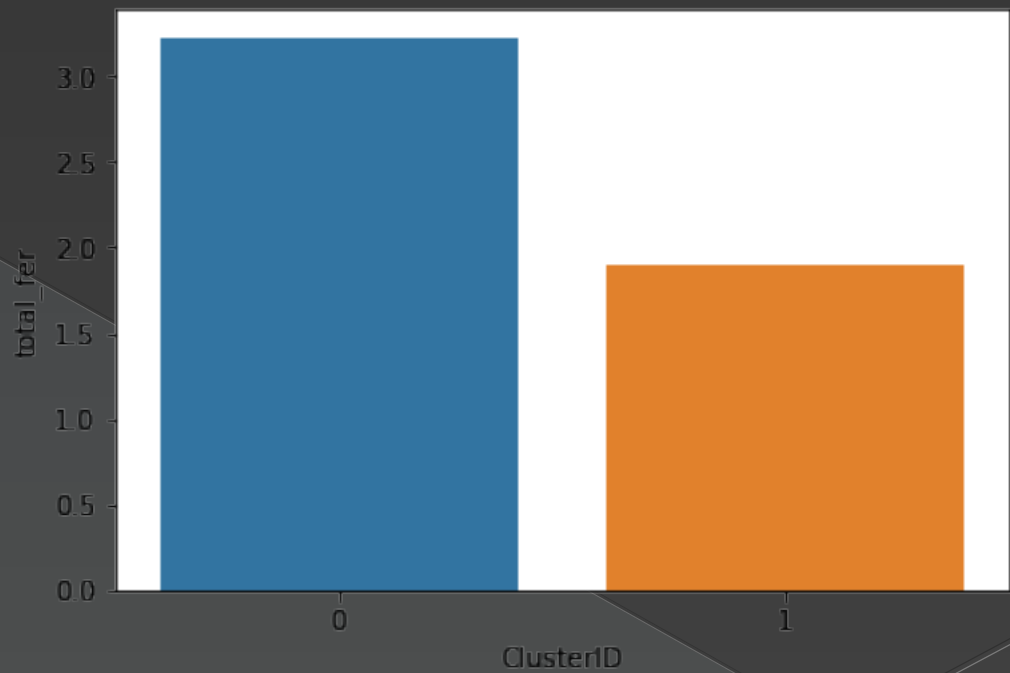
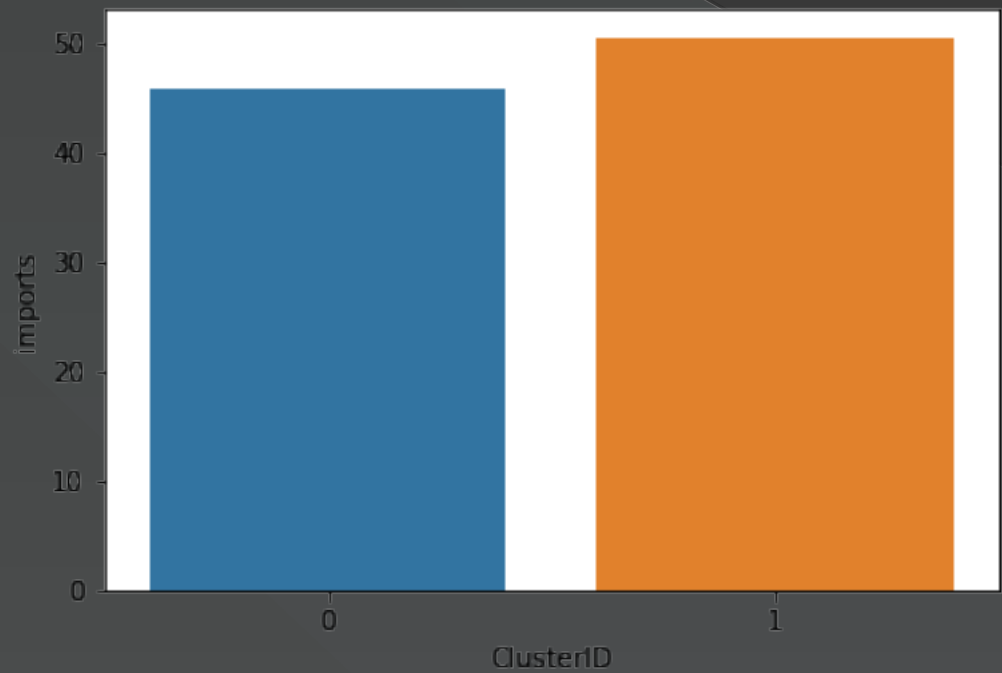


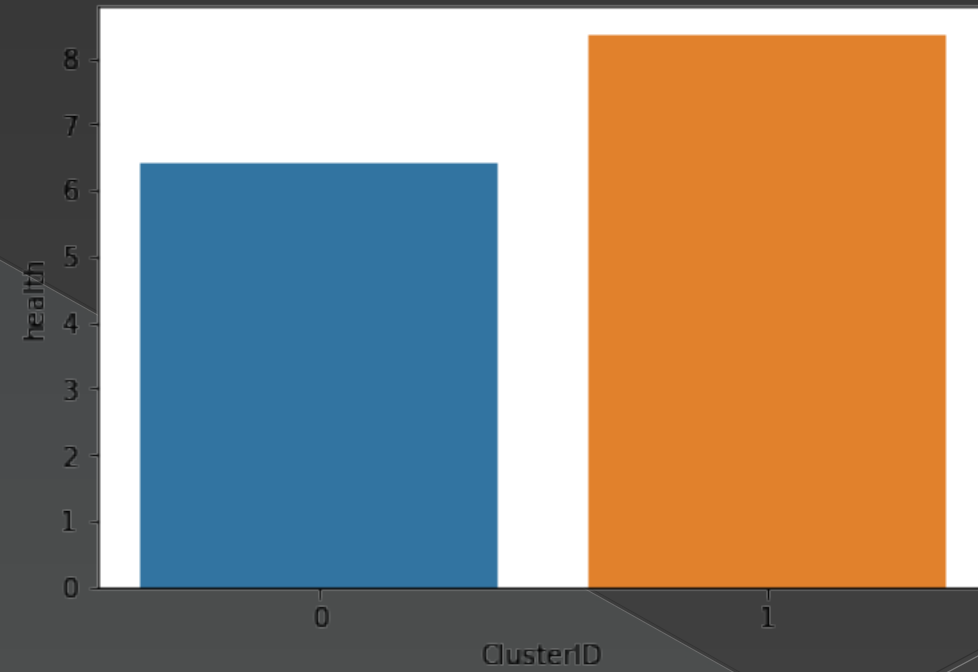
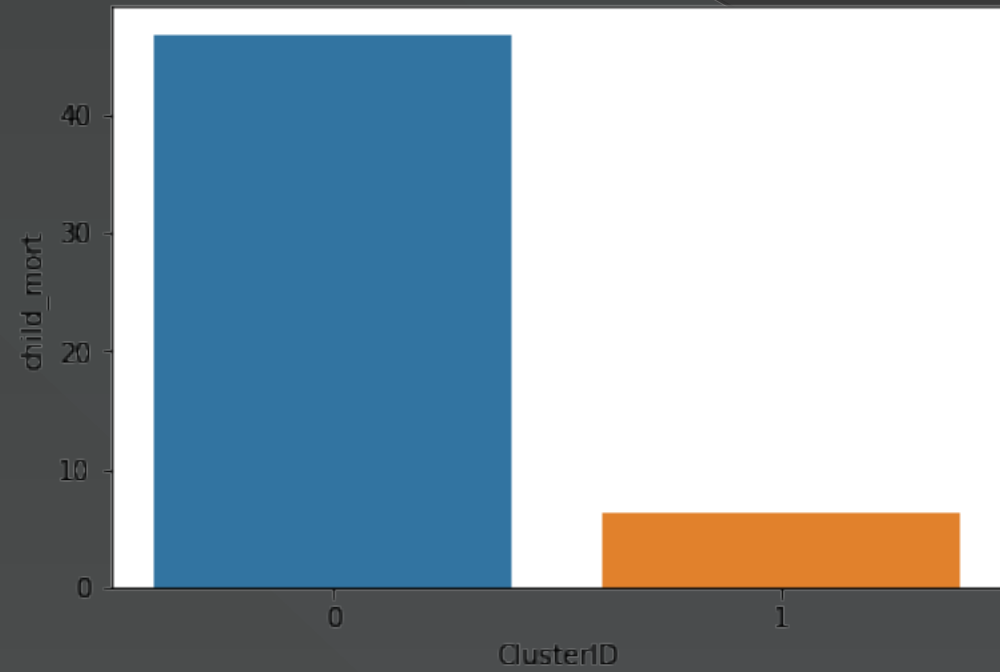


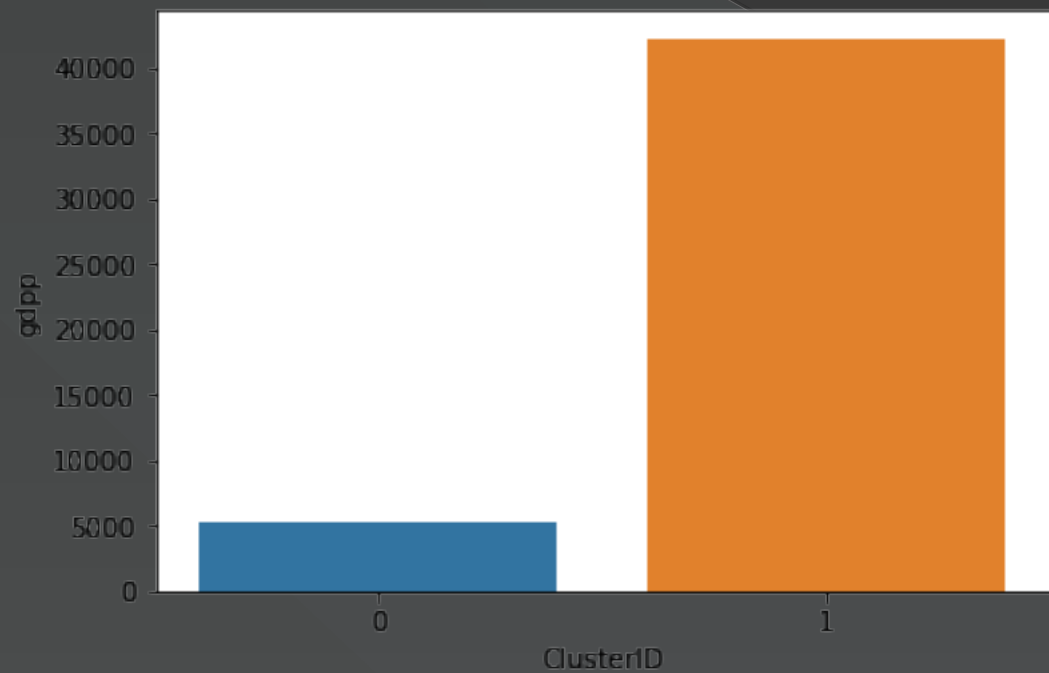
- > b. Analyzed the clusters formed
- > c. Plotted bar graphs between the clusters and the features and identified the clusters with countries in the direst need of aid (Clusters 0)











- > d. Identified top five countries from the cluster which are in direst need of the aid
- 8. The list of top five countries which are in the direst need of aid are:
 - > 1. Congo, Dem. Rep.
 - > 2. Burundi
 - > 3. Liberia
 - > 4. Niger
 - > 5. Central African Republic

THANK YOU