# Clustering Assignment

Identify countries that are in dire need of financial aid

# Problem Statement

- Based on the data about socio-economic indicators for countries, we have to identify countries that are most vulnerable, and are in need of financial help in time of disaster and calamities .

- Objectively, we have to cluster similar countries and group, based on features – child mortality, exports, health, imports, income, inflation, life expectancy, total fertility and gdpp
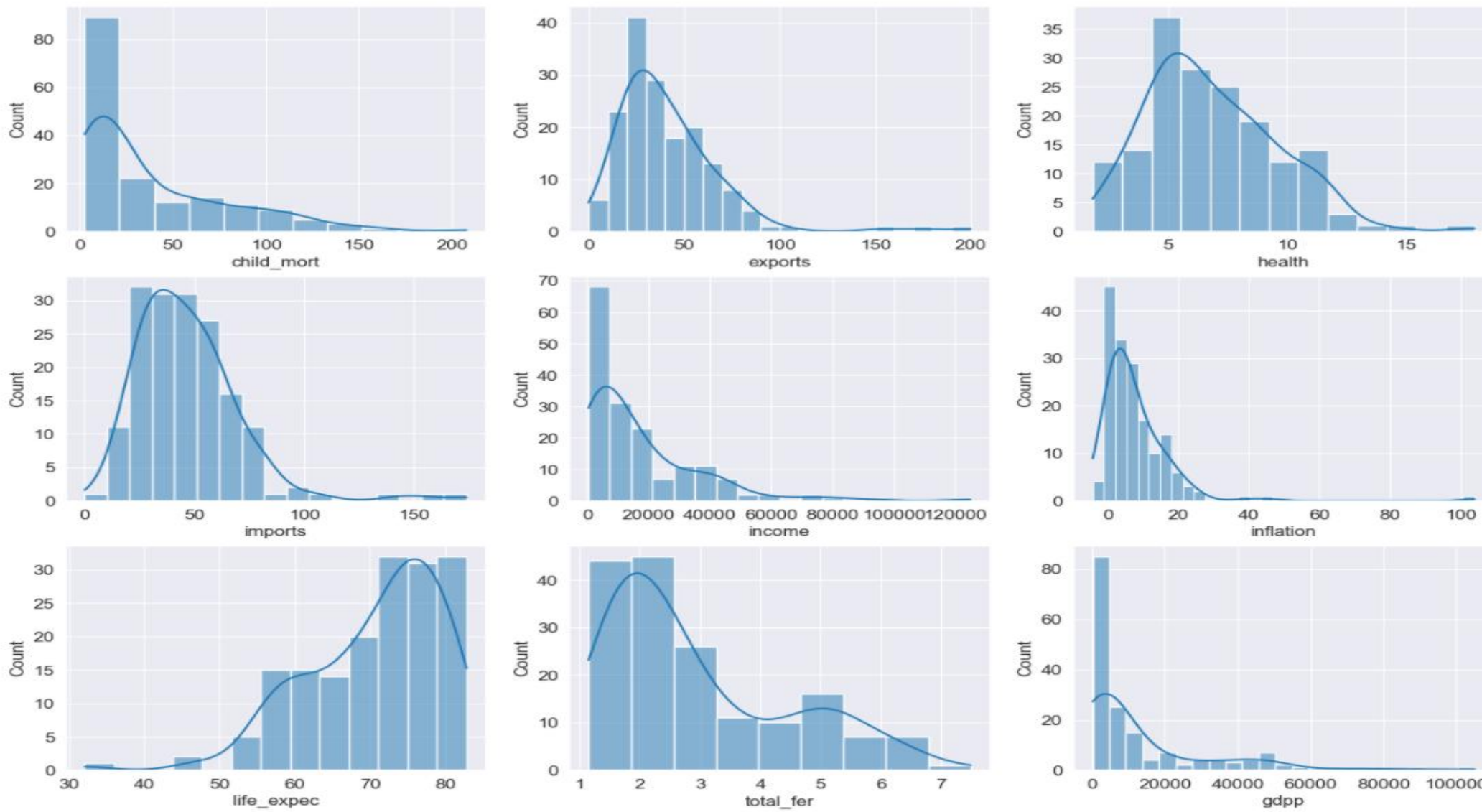
# Steps

- Understanding data and EDA
- Data processing
- Model – kmeans
- Model – kmeans (outlier removed)
- Model – hierarchical clustering (single and complete linkage)
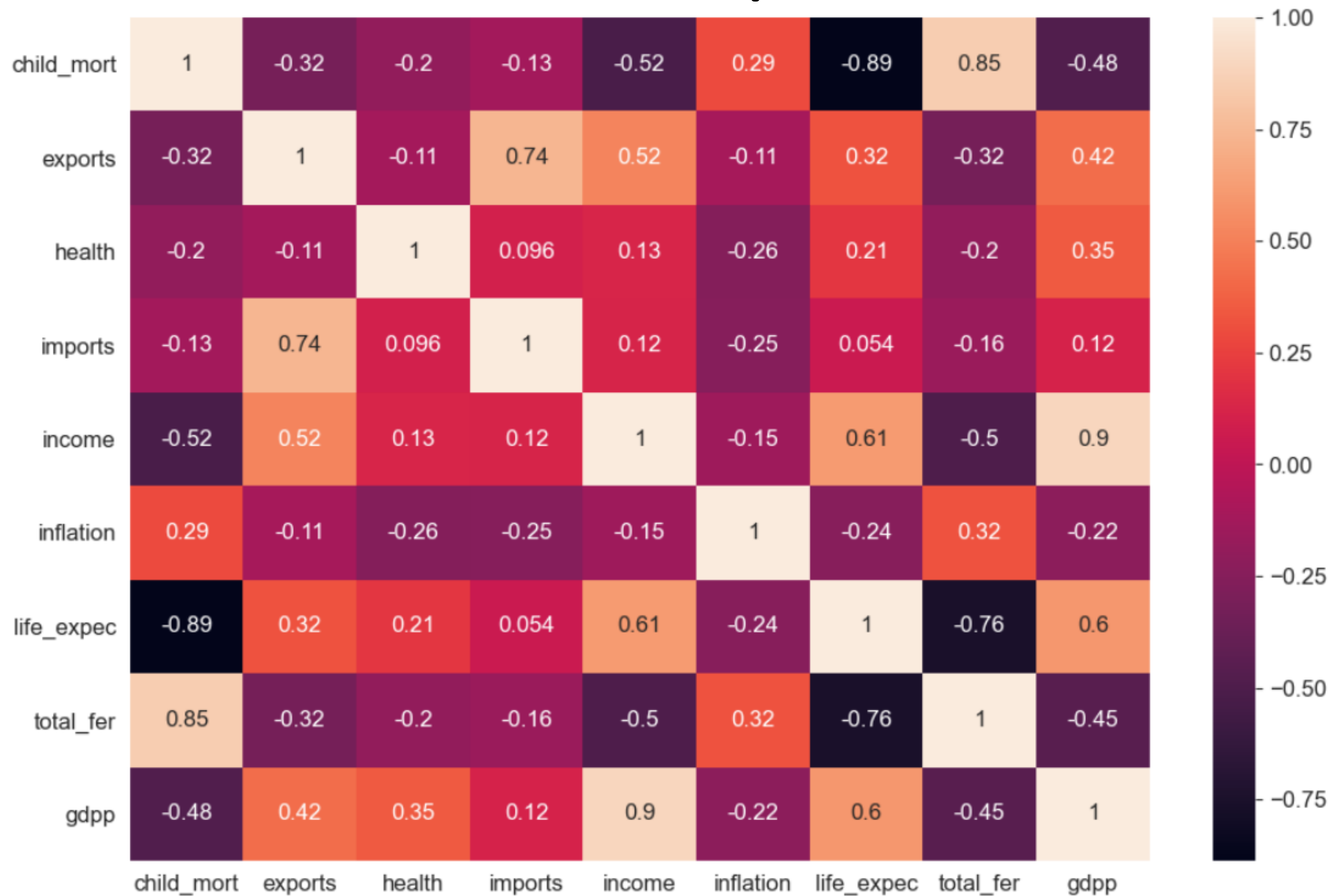- Model results

# EDA

- We have data about 167 countries.

- Dataset is clean and no has no null values.

- Countries can be grouped in under-developed, developing and developed.

- Features like child mortality, inflation and total fertility are negative indicators (increase in them, is bad)

- All other are positive features, indicating progress of the country.
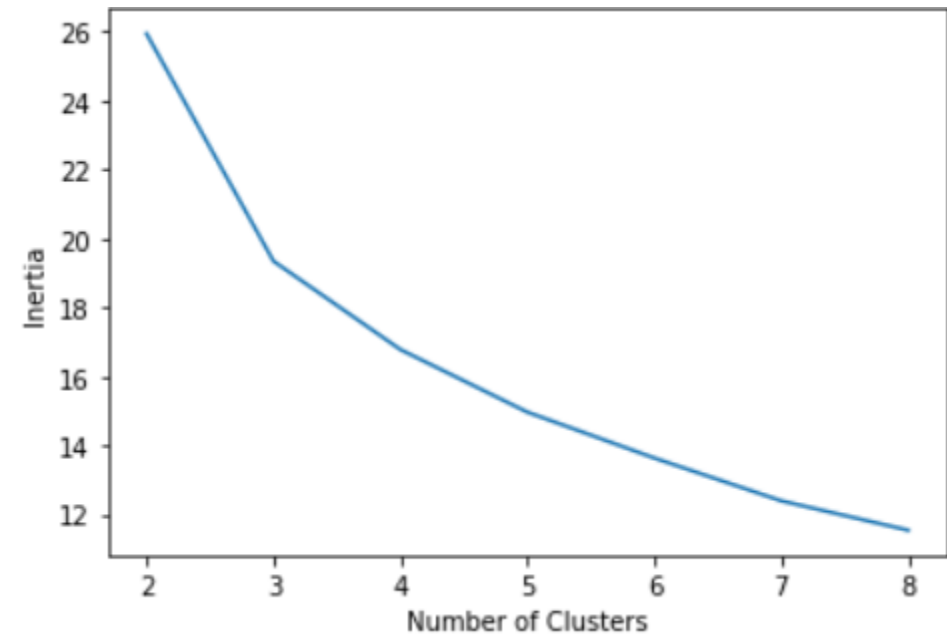
# Data Histograms

# Co-relation plots

# Kmeans (k – number of clusters)

- Scaled the data.
- Elbow, method used to identify optimal k
- Optimal k found = 3
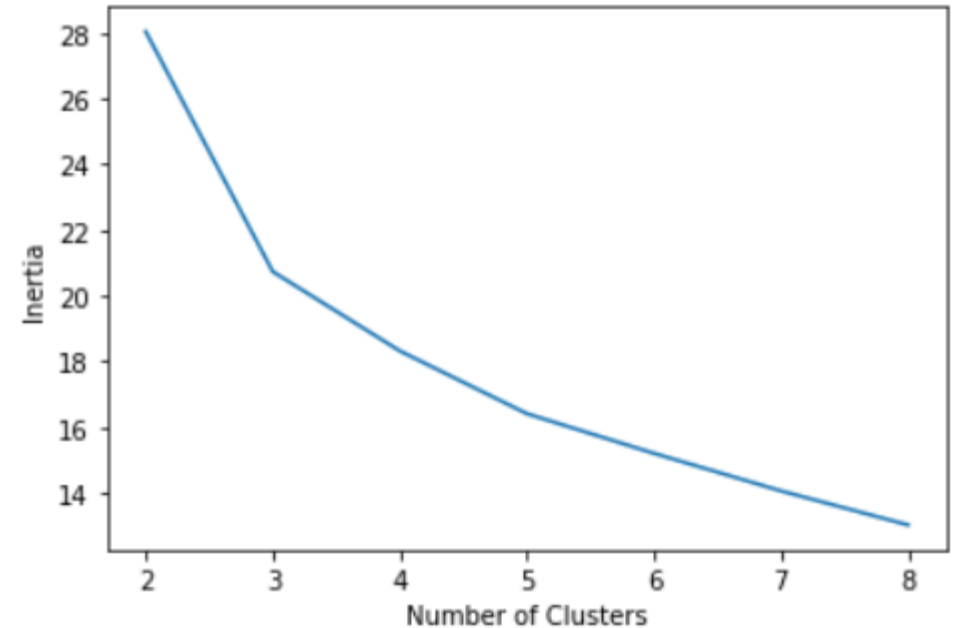- Inertia : Sum of square distances from cluster centre (SSD)



Finding optimal k by elbow method

# Outliers detection

- Outliers based on positive features.

- We have only 167 data points.

- Given our objective, we cannot discard outliers based on negative features.

- Outliers cleaned based on GDPp feature.

- Only countries with GDPp less than  - **Q3 + 1.5 x (Q3 – Q1)**
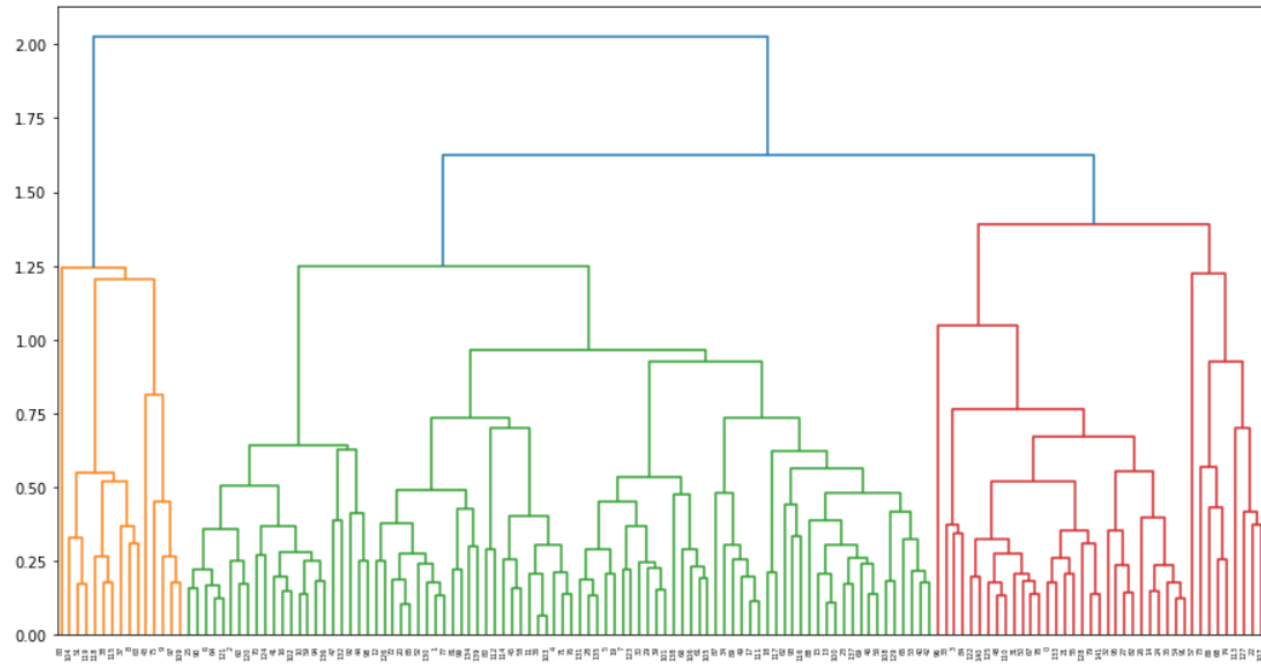
# kmeans – no_outliers

- Scale the data with MinMaxScaler()
- Elbow, method used to identify optimal k
- Optimal k found = 3
- Inertia : Sum of square distances from cluster centre (SSD)



Finding optimal k by elbow method

# Hierarchical Clustering

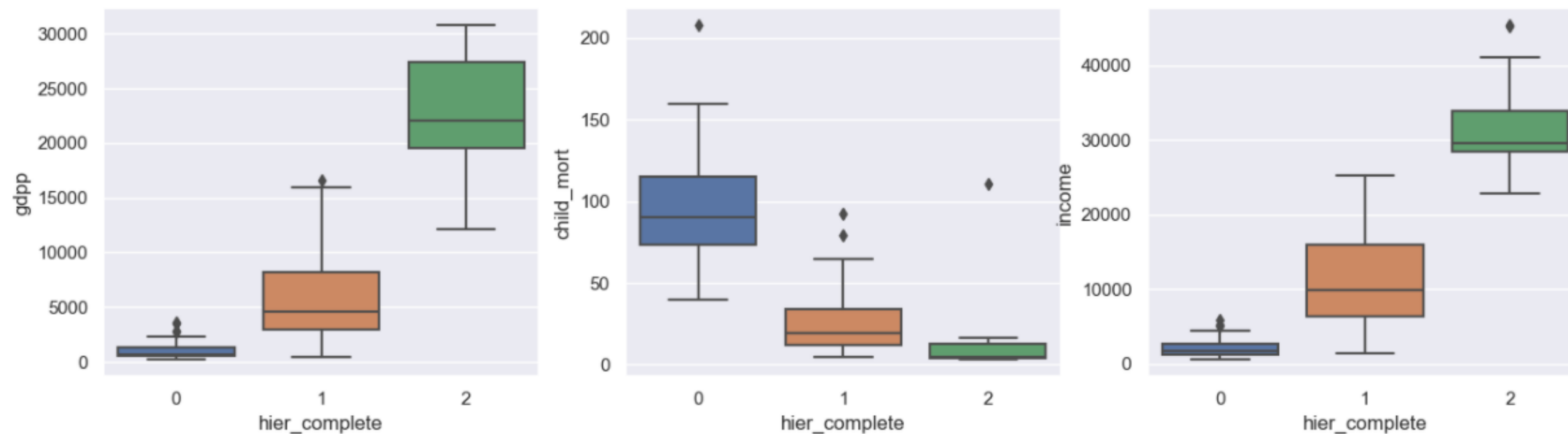• Solution with complete linkage (single linkage solution too skewed)
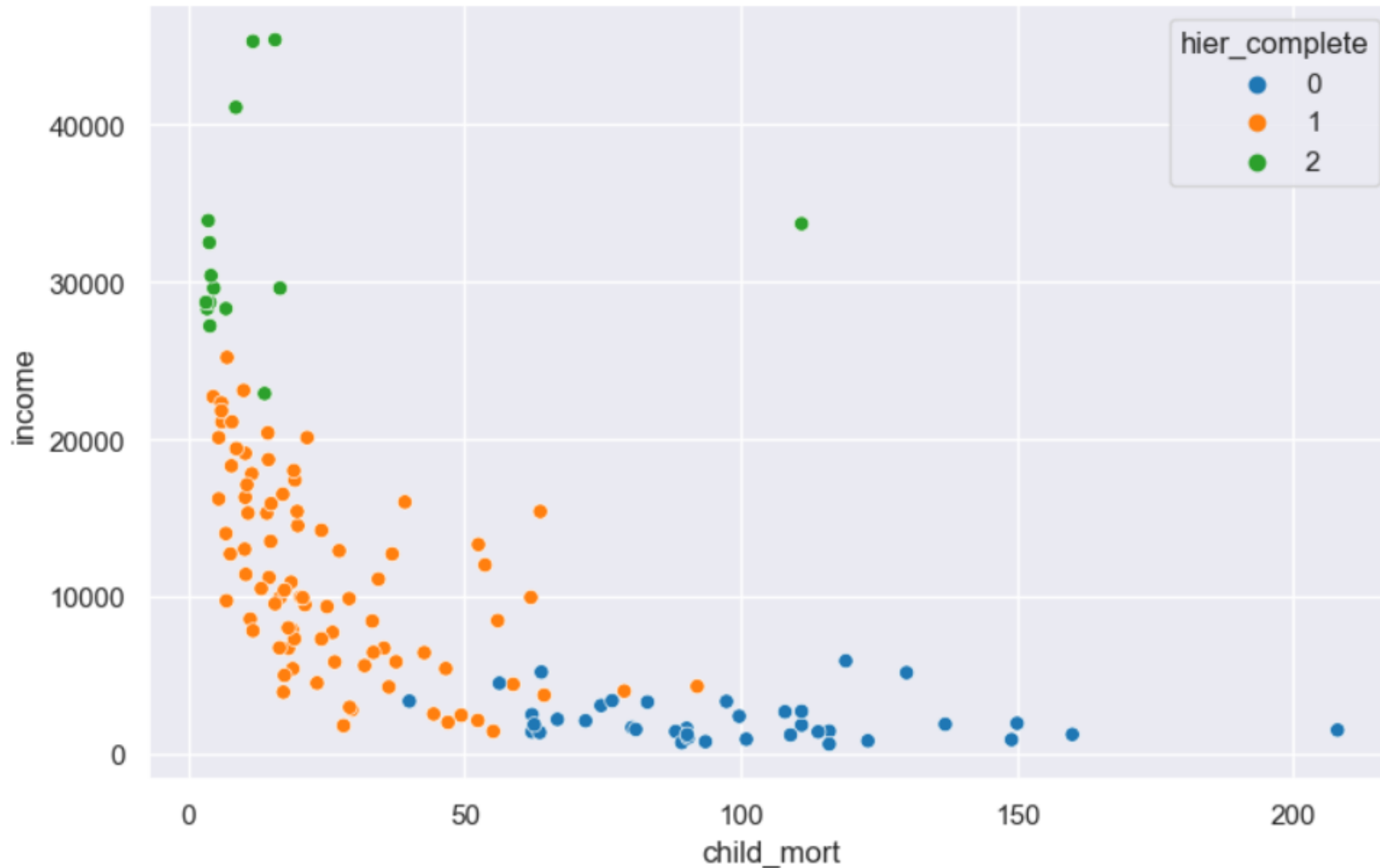
# Result and Visualizations-

39 countries were common out of all models, that were clustered as under-developed.

Afghanistan, Angola, Benin, Burkina Faso, Burundi, Cameroon, Central African Republic, Chad, Comoros, Congo, Dem. Rep., Congo, Rep., Cote d'Ivoire, Gambia, Ghana, Guinea, Guinea-Bissau, Haiti, Kenya, Kiribati, Lesotho, Liberia, Madagascar, Malawi, Mali, Mauritania, Micronesia, Fed. Sts., Mozambique, Niger, Nigeria, Rwanda, Senegal, Sierra Leone, Sudan, Tanzania, Timor-Leste, Togo, Uganda, Yemen, Zambia
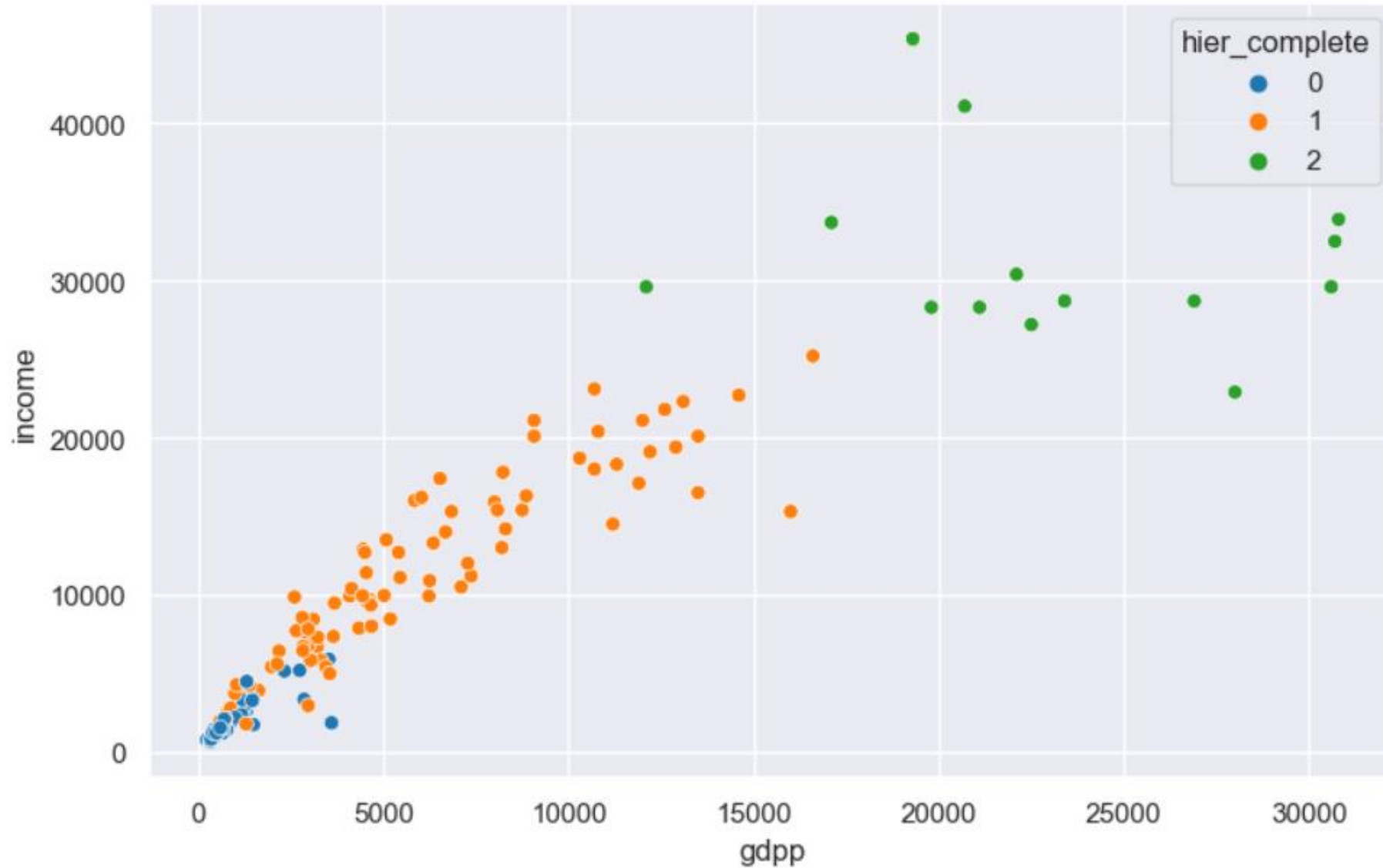
# Cluster 0 - under-developed countries

# Income vs child_mort (cluster 0 – underdeveloped)

# Income vs GDPp (cluster 0 – underdeveloped)

# GDPp vs child mort (cluster 0 – underdeveloped)