

Clustering of countries

Siddharth Babu

Objective:

HELP International humanitarian NGO, committed to fight poverty and provide the people of backward countries with basic amenities and relief during the time of disasters and natural calamities. We run a lot of operational projects from time to time, along with advocacy, drives to raise awareness as well as for funding purposes.

Problem statement: During the recent funding programmes, we have been able to raise around \$ 10 million. As an analyst, we have to come up with the countries list that are in the direst need of aid.

Analysis methodology

Data collection and cleaning

Outlier analysis and removal

Visualizing the data

Scaling the data

PCA on the data

Hopkins Statistics

K means clustering

Hierarchical Clustering

Decision Making

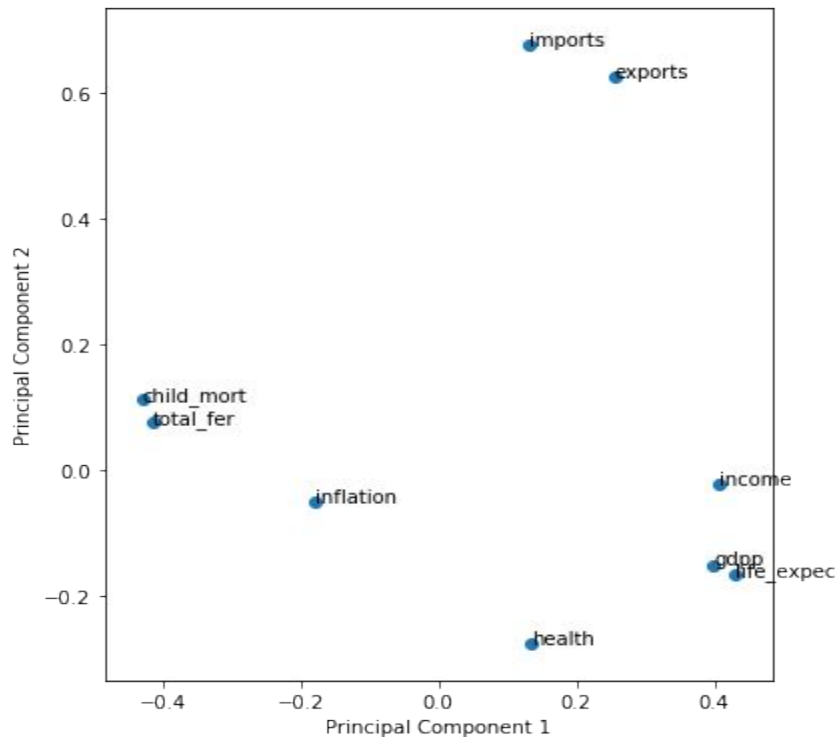
Correlation in data :

After data cleaning , we removed outlier from gdpp column because the country with high gdpp would not require any aid as there are already doing good.

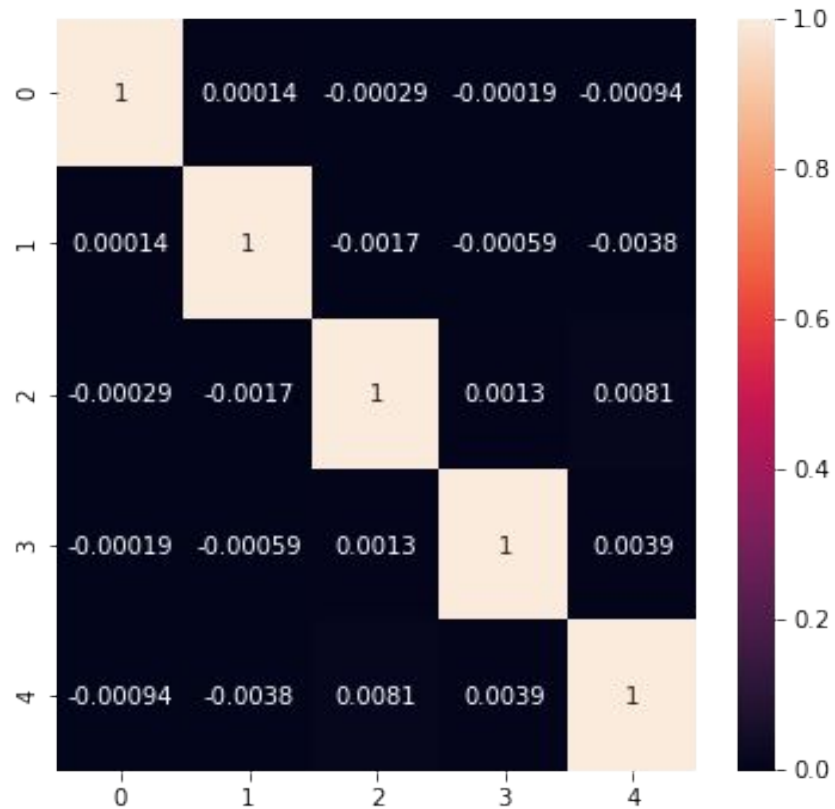
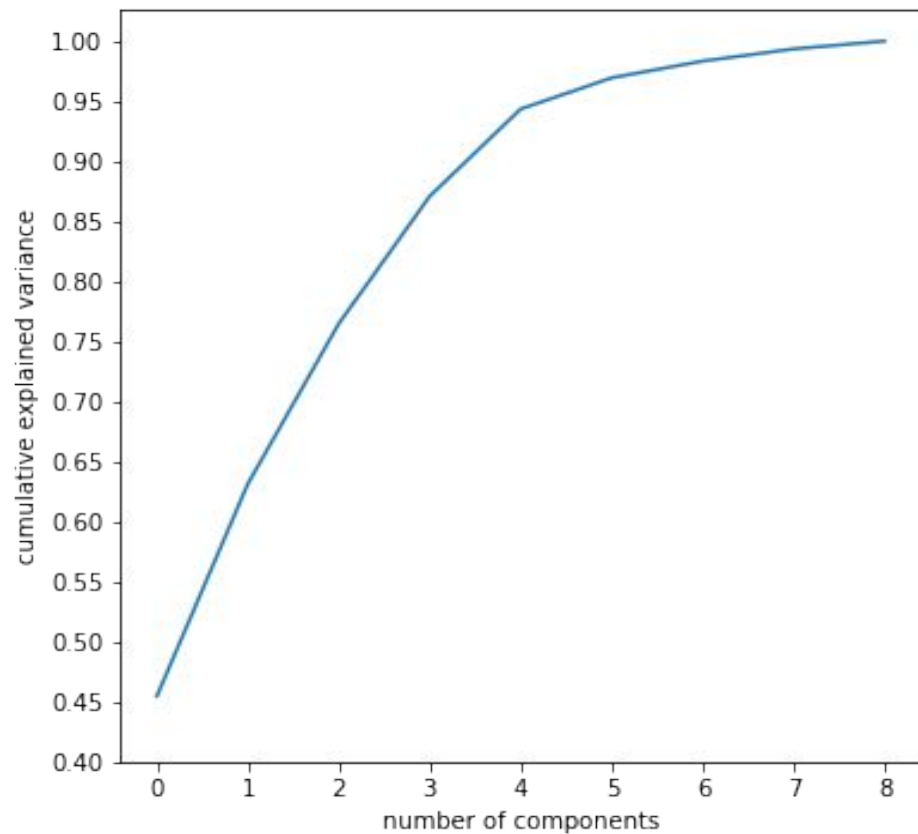
We did standardized scaling to standardize all parameters on cleaned, outlier removed data.

Looking at the heatmap, we see that few variables like (total fertility, child mortality) , (income , gdpp) and (imports and exports) have high correlation.

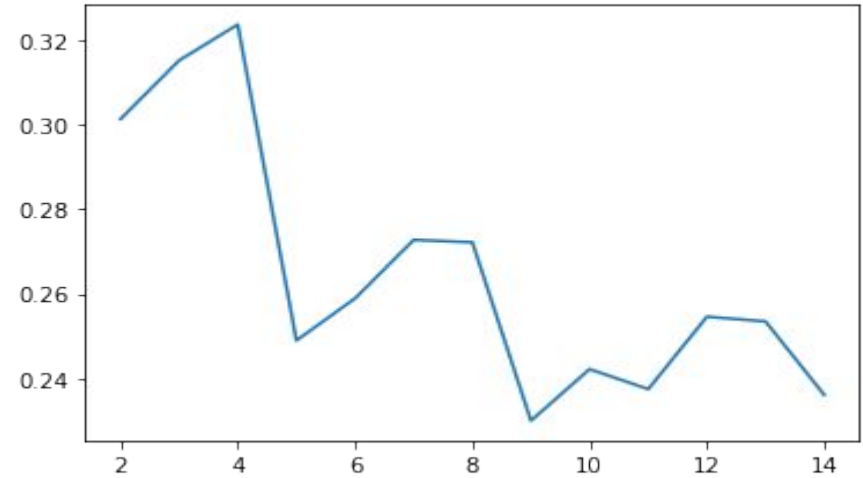
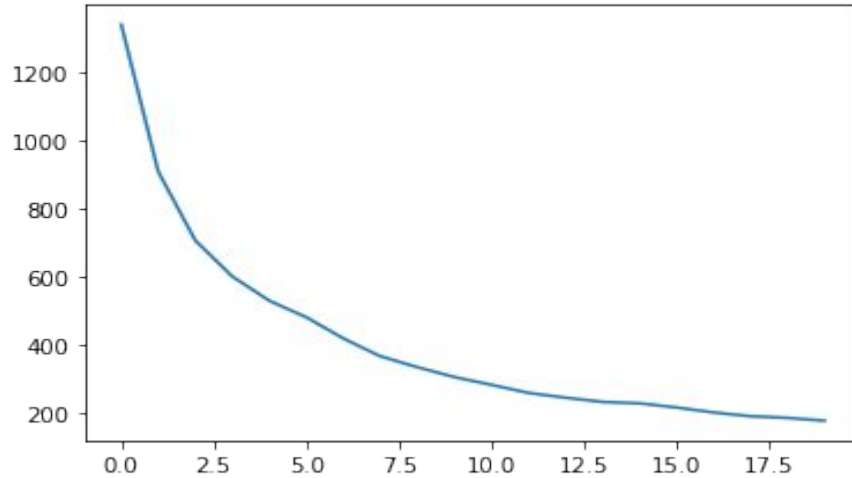
Principal component analysis



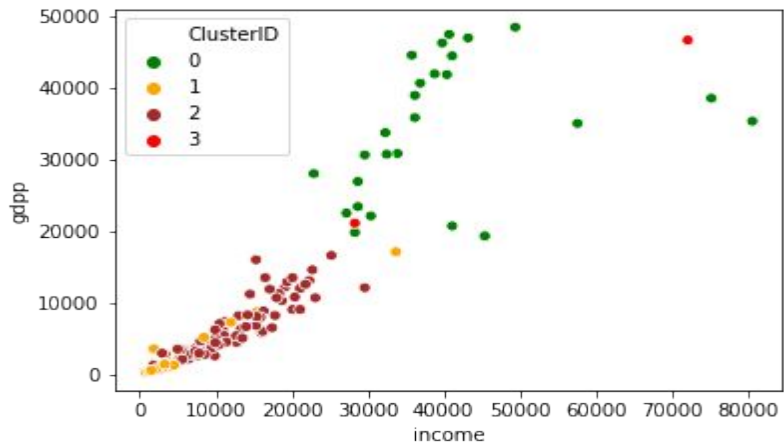
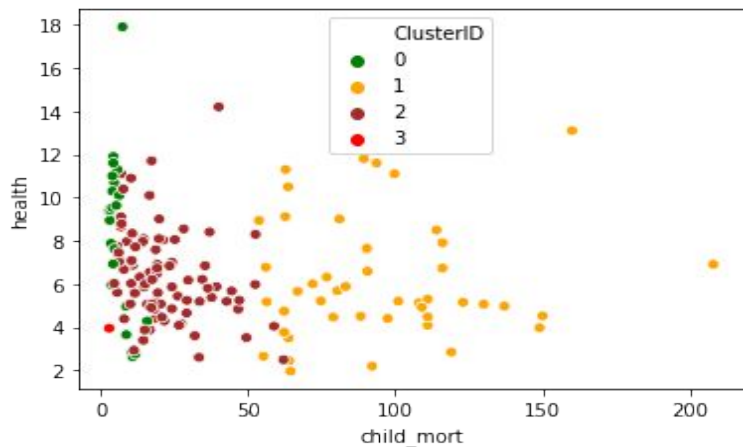
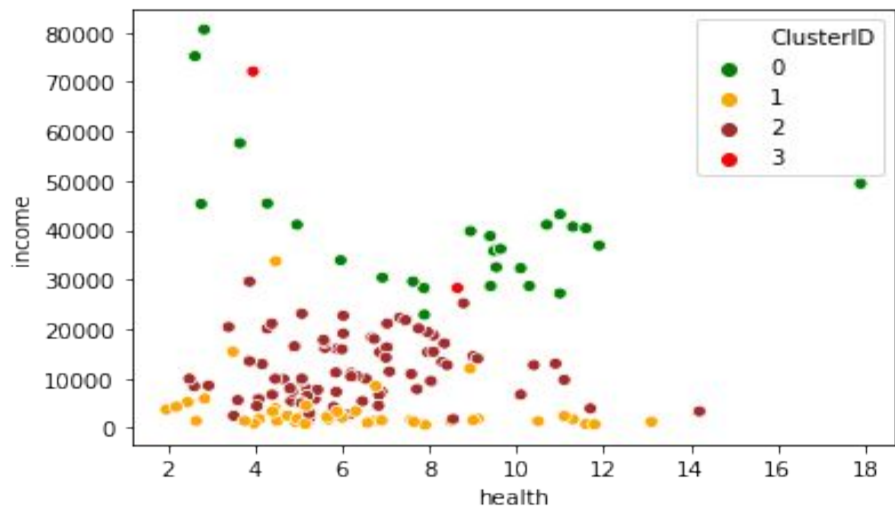
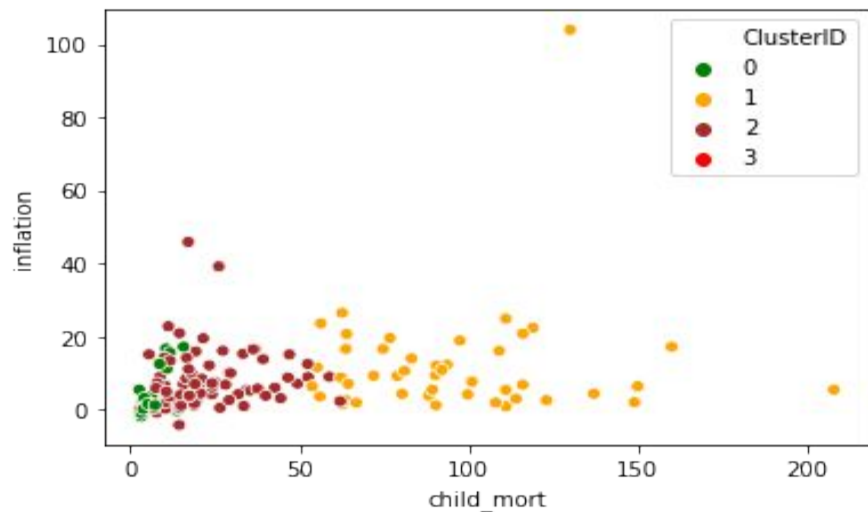
We see that features like gdpp, life expectancy and income are along the direction of PC1 and other features like total fertility and child mortality are along PC2 direction.

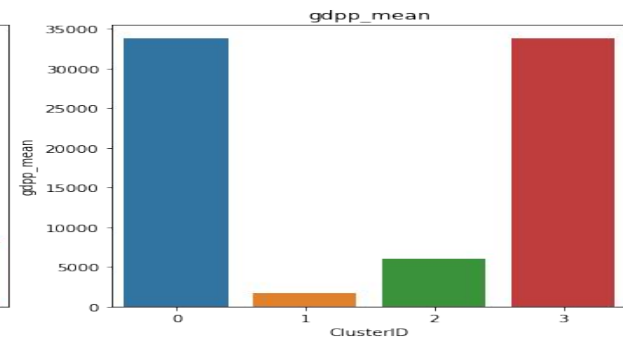
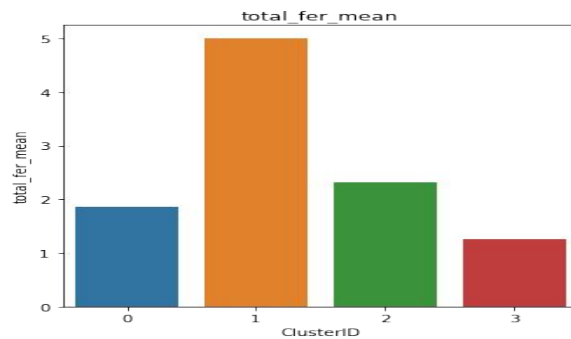
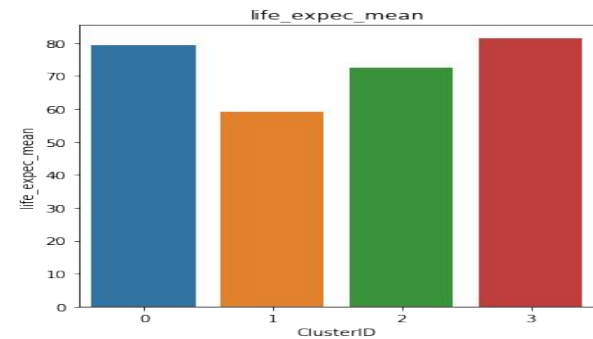
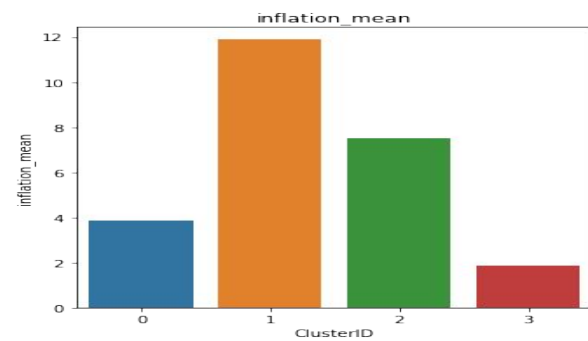
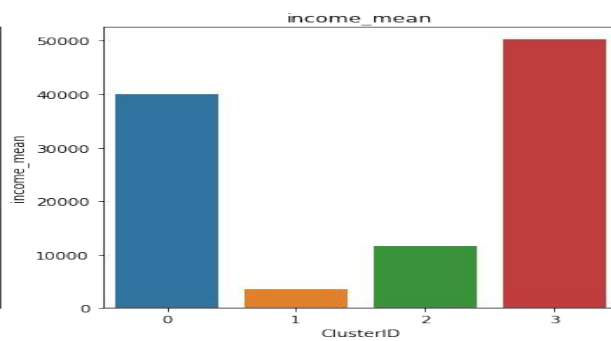
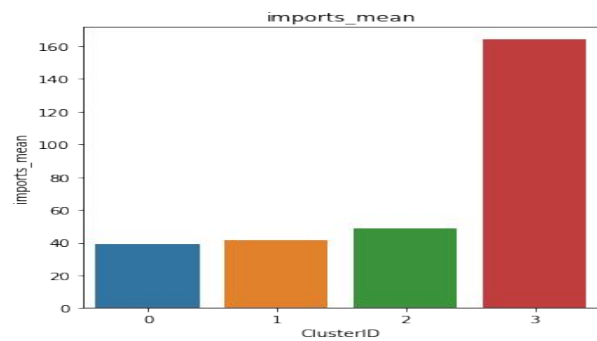
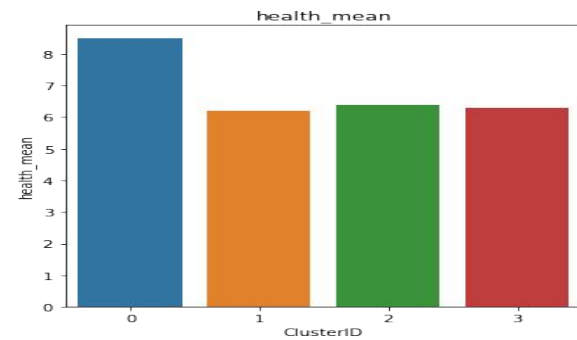
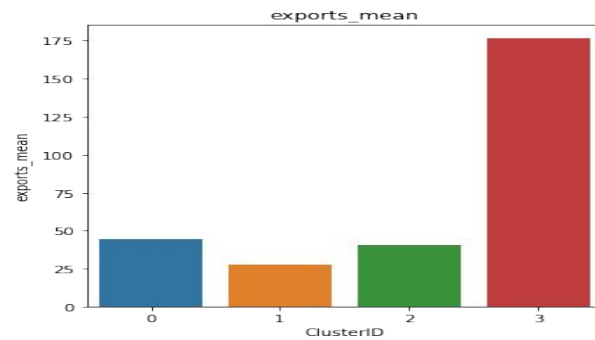
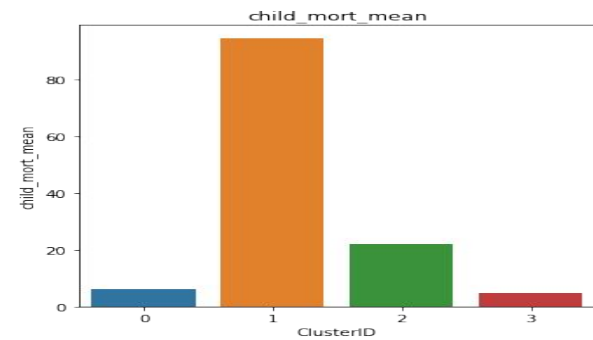


K-means clustering

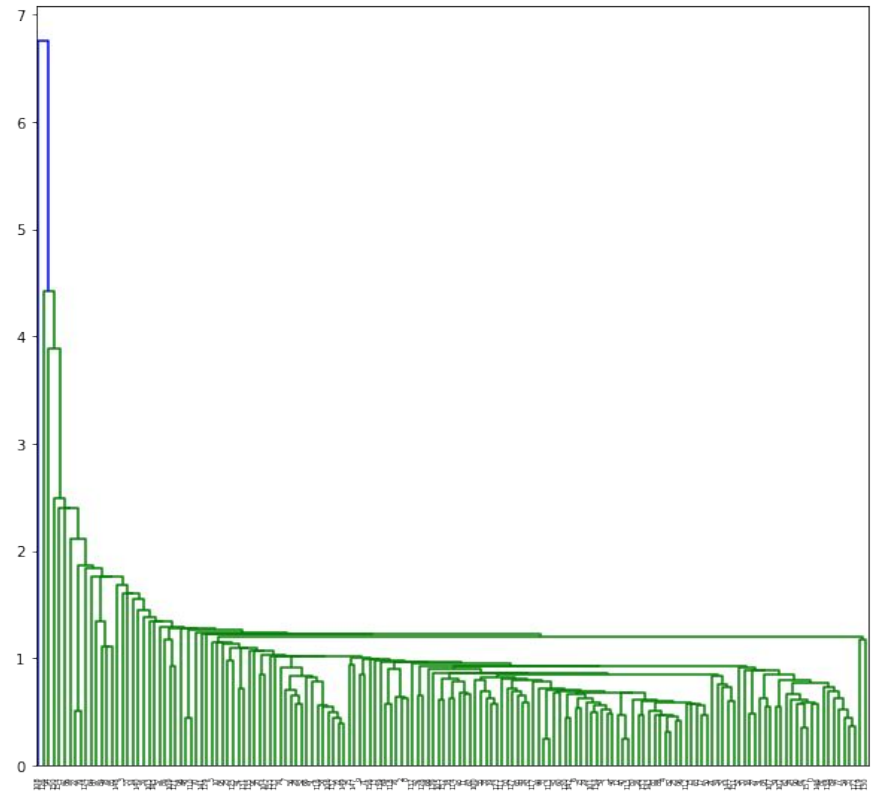
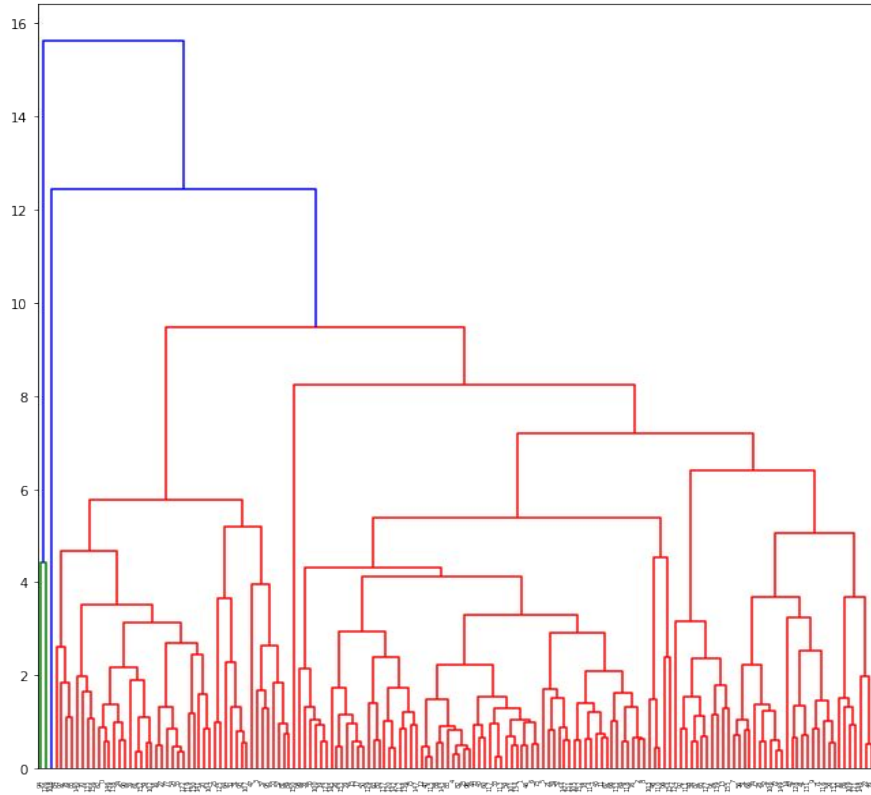


By looking silhouette analysis, we see the highest peak is at $k = 4$ and in sum of squared distances graph, we see that the elbow is in the range of 3 to 5, so let us take k as 4.





Hierarchical clustering



summary

As by both K means and Hierarchical clustering method - we have got same countries which requires aid. The following are the countries which are in direst need of aid by considering socio – economic factor into consideration:

	country	child_mort	exports	health	imports	income	inflation	life_expec	total_fer	gdpp	ClusterID
25	Burundi	93.6	8.92	11.60	39.2	764	12.30	57.7	6.26	231	0
85	Liberia	89.3	19.10	11.80	92.6	700	5.47	60.8	5.02	327	0
36	Congo, Dem. Rep.	116.0	41.10	7.91	49.6	609	20.80	57.5	6.54	334	0
107	Niger	123.0	22.20	5.16	49.1	814	2.55	58.8	7.49	348	0
125	Sierra Leone	160.0	16.80	13.10	34.5	1220	17.20	55.0	5.20	399	0
89	Madagascar	62.2	25.00	3.77	43.0	1390	8.79	60.8	4.60	413	0
102	Mozambique	101.0	31.50	5.21	46.2	918	7.64	54.5	5.56	419	0
30	Central African Republic	149.0	11.80	3.98	26.5	888	2.01	47.5	5.21	446	0
90	Malawi	90.5	22.80	6.59	34.9	1030	12.10	53.1	5.31	459	0
141	Togo	90.3	40.20	7.65	57.3	1210	1.18	58.7	4.87	488	0