# upGrad
*#LifeKoKaroLift*

# Clustering: Pre-Assignment Session

**Course :** Data Science

**Lecture On :** Pre-Assignment

**Instructor :** Sumit Shukla

upGrad

**Assignment Explanation**

**upGrad**

*[handwritten: → Cluster → match sus → all ch → data pur → 2,3 → kmean ( randon → 1,0]*
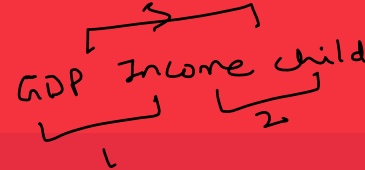
**Step to procedure in the Assignment**

*[handwritten: % of GDPP]*

Let's first understand the problem statement:

Identify top countries that are direst need of aid. 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.

| | country | child_mort | exports | health | imports | income | inflation | life_expec | total_fer | gdpp |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 90.2 | 10.0 | 7.58 | 44.9 | 1610 | 9.44 | 56.2 | 5.82 | 553 |
| 1 | Albania | 16.6 | 28.0 | 6.55 | 48.6 | 9930 | 4.49 | 76.3 | 1.65 | 4090 |
| 2 | Algeria | 27.3 | 38.4 | 4.17 | 31.4 | 12900 | 16.10 | 76.5 | 2.89 | 4460 |
| 3 | Angola | 119.0 | 62.3 | 2.85 | 42.9 | 5900 | 22.40 | 60.1 | 6.16 | 3530 |
| 4 | Antigua and Barbuda | 10.3 | 45.5 | 6.03 | 58.9 | 19100 | 1.44 | 76.8 | 2.13 | 12200 |

**upGrad**

*[handwritten: GDP Income child, with annotations 1 and 2]*

## Step to procedure in the Assignment

*[handwritten: Export, Import and health → %ouge of GDPP ↓ convert them to acutal values; 80+ → Good Rang, 60-80 → Avg., Below → Bad]*

1. Data Understanding.
   a. Hint: Don't forget to read the data description properly.

1. Perform Clustering.  *[handwritten: EDA]*
   a. Data preparation for clustering.
      i. Outlier treatment
      ii. Hopkins check
   b. Clustering
      i. K-MEANS
         1. Run K-Means and choose K using both Elbow and Silhouette score
         2. Run K-Means with the chosen K
         3. Visualise the clusters
         4. Clustering profiling using "gdpp, child_mort and income"
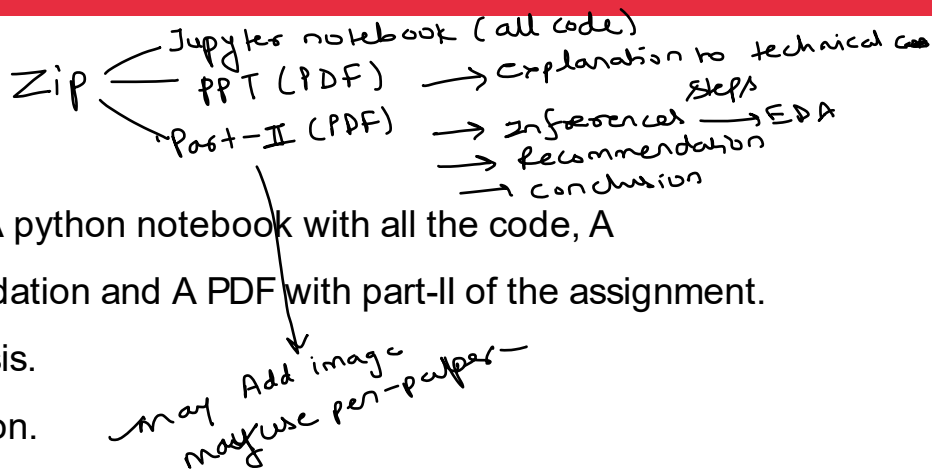
**Step to procedure in the Assignment**

1. Perform Clustering.
   a. Clustering
      i. Hierarchical Clustering
         1. Use both Single and Complete linkage
         2. Choose one method based on the results
         3. Visualise the clusters
         4. Clustering profiling using "gdpp, child_mort and income"

1. Country Identification
   a. Based on the analysis, choose the countries that are in need for the aid.
   b. Choose the countries based on some socio-economic and health factors.

*using both k-mean & hierarcs*

*Find out top-5 countries that are in need of the AID.*

Zip
- Jupyter notebook (all code)
- PPT (PDF) → Explanation to technical co...
  steps
- Part-II (PDF) → Inferences → EDA
  → Recommendation
  → Conclusion

- You need to comment your code properly.

- You need to submit 3 files zipped as one file. A python notebook with all the code, A

  PPT(Converted as PDF) with all the recommendation and A PDF with part-II of the assignment.

- PPT should cover main points from your analysis.

- Choose K wisely as this reflects the final solution.

- Mention your assumptions in your notebook, If taken any.

- For any help regarding assignment, use Discussion Forum.

- For any help regarding coding error, use online portals such as StackOverflow.

May Add image
Mayfuse pen-paper—

**upGrad**

*#LifeKoKaroLift*

# Thank You!

23/05/19