# Data Collection and Preprocessing Phase
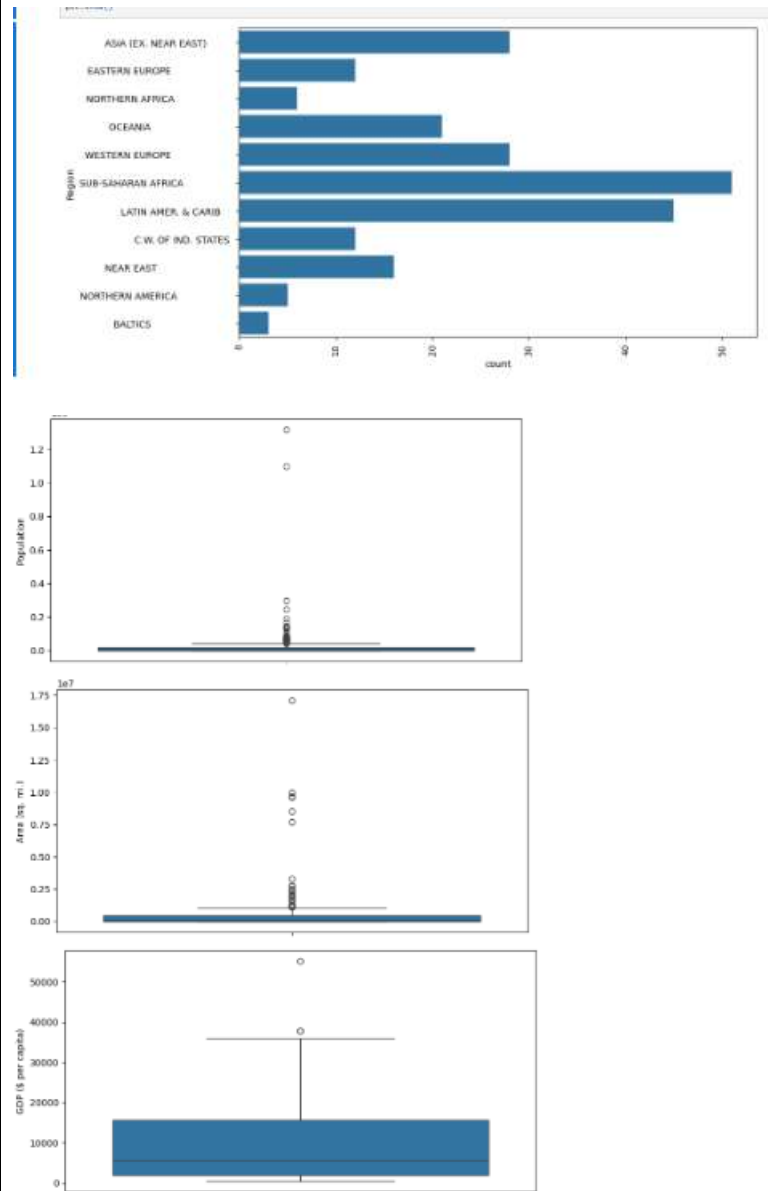
| Date | 2 July 2024 |
|---|---|
| Team ID | SWTID1720086522 |
| Project Title | Forecasting Economic Prosperity: Leveraging Machine Learning For GDP Per Capita Prediction |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview | **Dimension:**<br>227 rows X 20 columns<br>**Descriptivestatistics:**<br> |

| | |
|---|---|
| Univariate Analysis |  |
| Bivariate Analysis | - |

| | |
|---|---|
| Multivariate Analysis |  |
| Outliers and Anomalies | - |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |
| Handling Missing Data |  |
| Data Transformation |  |

| Feature Engineering | `df.drop(['Literacy (%)','Net migration','Population','Area (sq. mi.)','Coastline (coast/area ratio)'], axis=1,inplace=True)` |
|---|---|
| Save Processed Data | ```import pickle
pickle.dump(rf,open('models.pkl','wb'))``` |