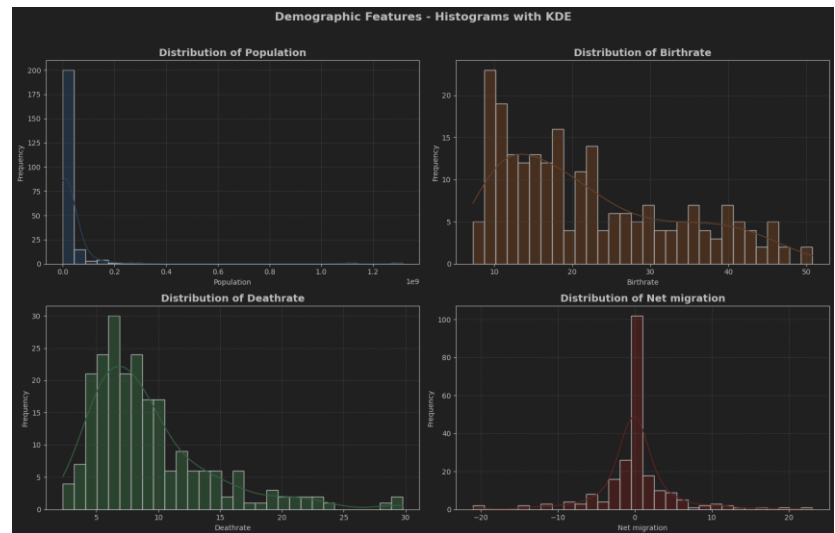# Data Collection and Preprocessing Phase

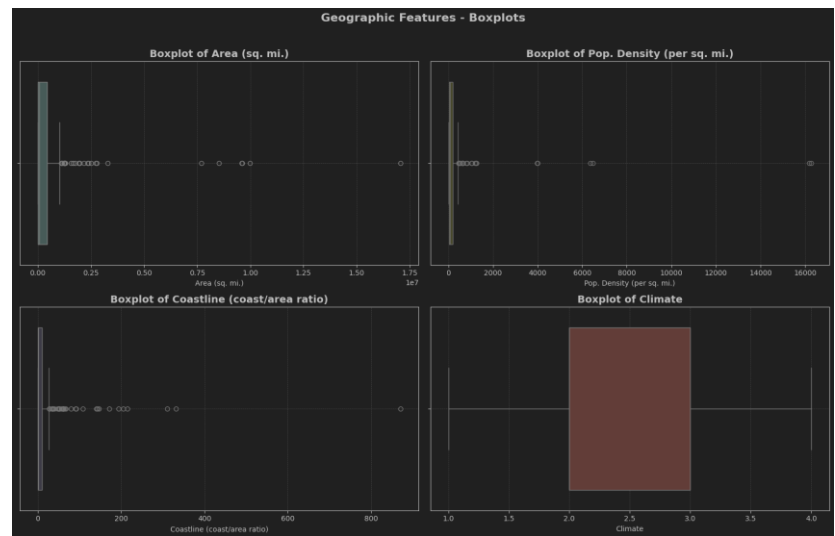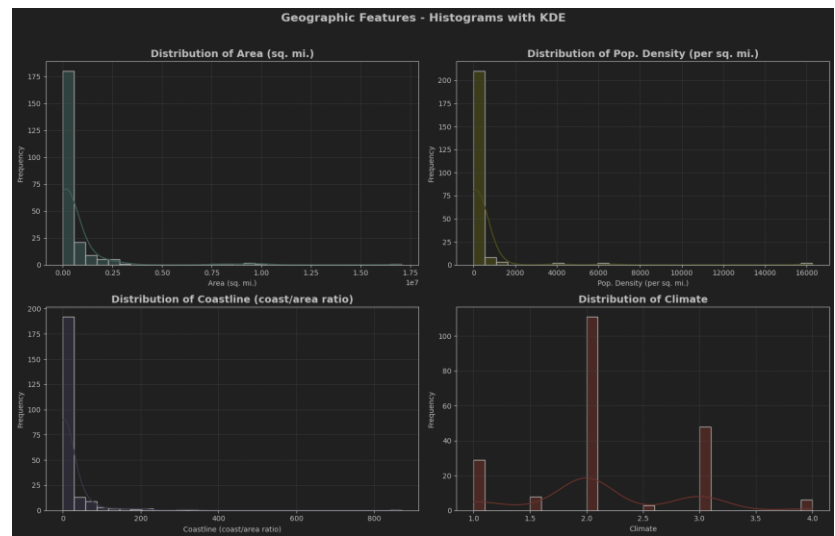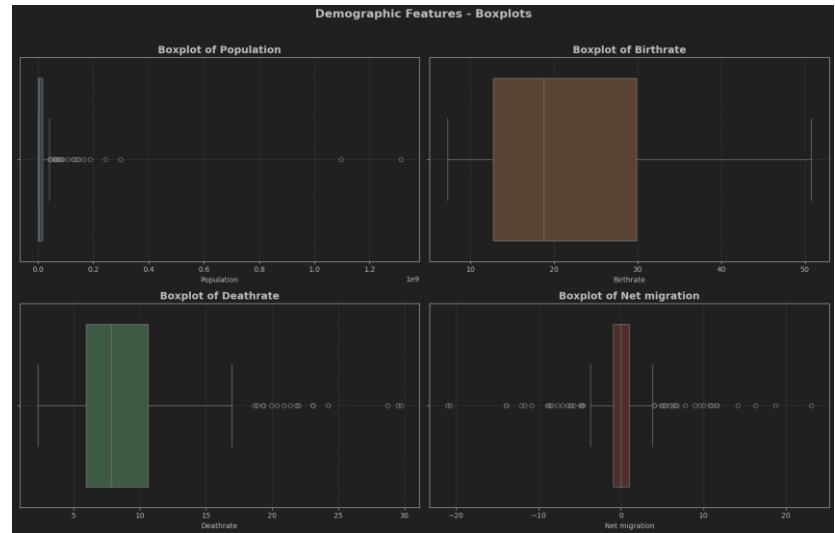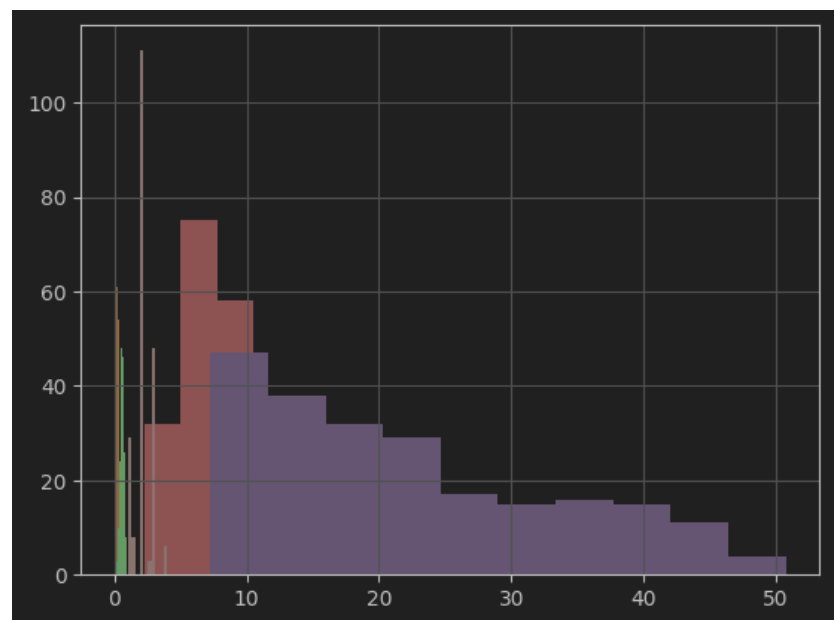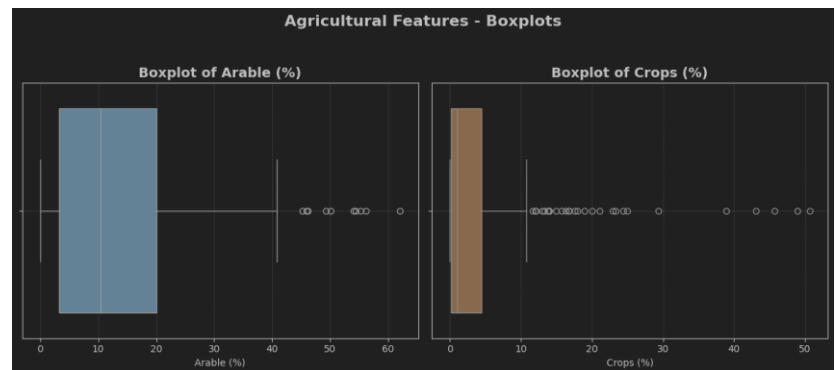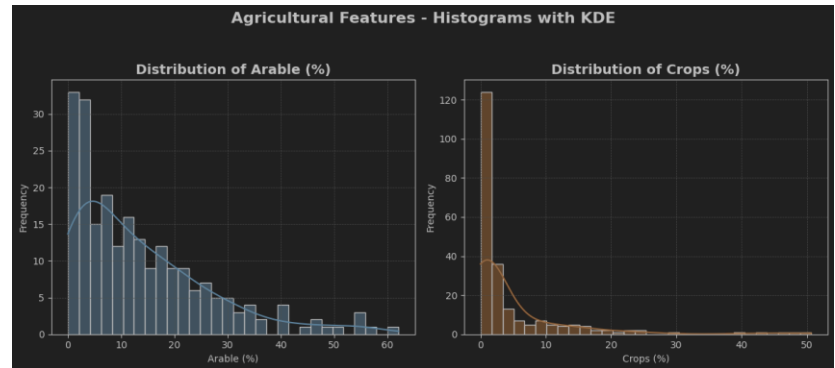| | |
|---|---|
| Date | 16 June 2025 |
| Team ID | SWTID1749653449 |
| Project Title | Economic Growth: A Machine Learning Approach to GDP per Capita Prediction |
| Maximum Marks | 6 Marks |

**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.
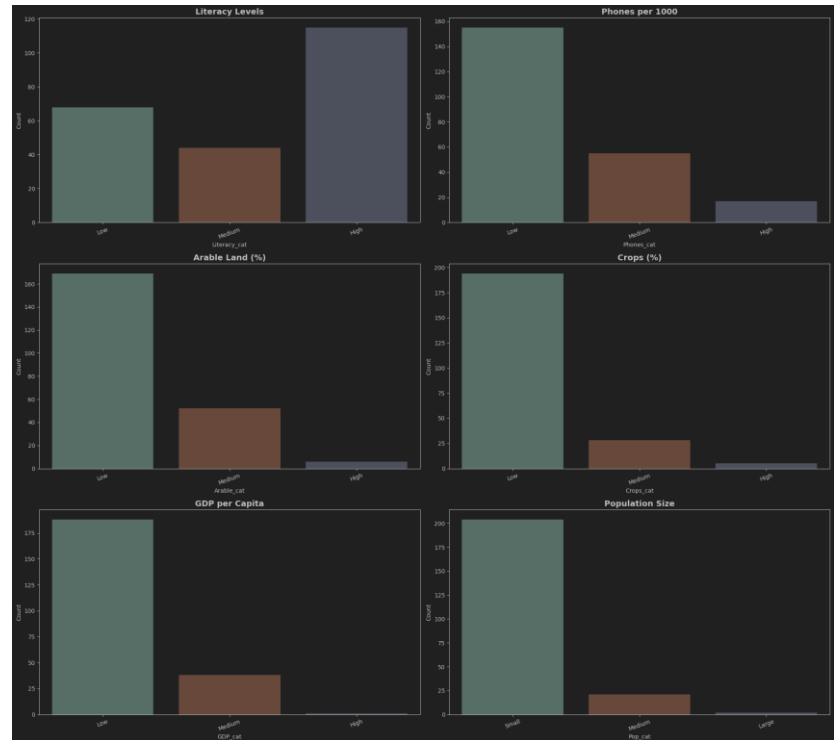
| Section | Description |
|---|---|
| Data Overview | ```
# Understanding the data
df.shape

(227, 20)
``` <br> Dimensions: 227 rows x 20 Columns <br> Descriptive Stats: <br><br> (descriptive statistics table image) |
| Univariate Analysis | |

Histograms with KDE for Economic Features



Demographic Features - Histograms with KDE

Demographic Features - Boxplots



Geographic Features - Histograms with KDE



Geographic Features - Boxplots

Agricultural Features - Histograms with KDE



Agricultural Features - Boxplots



| Bivariate Analysis | |

Pair Plot:

Pairplot of Selected Features

.

| | |
|---|---|
| Multivariate Analysis |  |

Correlation Heatmap of Selected Features

Outliers and Anomalies

**Data Preprocessing Code Screenshots**

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



```python
# Dropping unnecessary columns for model building
df.drop([
    "Country",                          # Unique identifier
    "Other (%)",                        # Highly correlated
    "Infant mortality (per 1000 births)",  # Highly correlated
    "Literacy_cat", "Phones_cat",       # Used only for visualization
    "Arable_cat", "Crops_cat",
    "GDP_cat", "Pop_cat"
], axis=1, inplace=True)
```

```python
# Handling missing values
df.isnull().sum()
```

```
Region                        0
Population                    0
Area (sq. mi.)                0
Pop. Density (per sq. mi.)    0
Coastline (coast/area ratio)  0
Net migration                 3
GDP ($ per capita)            1
Literacy (%)                 18
Phones (per 1000)             4
Arable (%)                    2
Crops (%)                     2
Climate                      22
Birthrate                     3
Deathrate                     4
Agriculture                  15
```

```python
# Dropping uneccessary values
df = df.dropna(subset=['Net migration'])
df = df.dropna(subset=['GDP ($ per capita)'])
df = df.dropna(subset=['Phones (per 1000)'])
df = df.dropna(subset=['Arable (%)'])
df = df.dropna(subset=['Crops (%)'])
df = df.dropna(subset=['Birthrate'])
df = df.dropna(subset=['Deathrate'])
```

| Data Transformation | ```
df['Literacy (%)'] = df['Literacy (%)'].fillna(df['Literacy (%)'].mean())
df['Climate'] = df['Climate'].fillna(df['Climate'].mean())
df['Service'] = df['Service'].fillna(df['Service'].mean())
df['Agriculture'] = df['Agriculture'].fillna(df['Agriculture'].mean())
df['Industry'] = df['Industry'].fillna(df['Industry'].median())
``` |
|---|---|
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | __ Jupyter Notebook was Used __ |