

# **Predicting Infant Mortality: A Global Analysis of 2023 Data**



**Aurora Degree & P.G College | M. Sc Applied Statistics**  
**Ruchitha | Shravanthika | Ramya | Madirai | Deepika**

# Abstract

**This study explores various machine learning models to predict infant mortality rates across different countries.**

**We evaluate multiple algorithms, including linear regression, random forests, support vector machines and more.**

# Content

<b>1. Introduction</b>	<b>4</b>
<b>2. Objectives</b>	<b>5</b>
<b>3. Literature Review</b>	<b>6</b>
<b>4. About data</b>	<b>8</b>
<b>5. Data Cleaning</b>	<b>10</b>
<b>6. EDA</b>	<b>12</b>
<b>7. Modelling</b>	<b>23</b>
<b>8. Result</b>	<b>29</b>
<b>9. Conclusion</b>	<b>30</b>
<b>10. Recommendation</b>	<b>32</b>

# INTRODUCTION

**The infant mortality rate (IMR) is often regarded as a barometer for overall welfare of a community or country.**

**It refers to the number of deaths per 1000 live births within the year**



# OBJECTIVES

## ● Objective 1

To identify the significant factors that influence the infant mortality globally.

## ● Objective 2

To develop the most accurate machine learning model for predicting infant mortality on a global scale

# LITERATURE REVIEW

1

In a study by [Zakir Hossain, Enamul Kabir, and Rumana Rois](#) (November 2021), machine learning algorithms were used to identify predictors of infant mortality in Bangladesh using data from the 2017–18 Demographic and Health Survey. Key features like age at first marriage, birth interval, and education were found significant. Among various ML models, random forest performed best with an accuracy of 89.3% and an AUC of 0.6613. The findings can guide policy-makers and public health interventions aimed at reducing infant mortality

2

[Leonardo Matsuno da Frota et al. \(March 2024\)](#) used machine learning to predict infant mortality in Brazil from 2.9 million Brazil's Unique Health System(SUS) data points. Survival Support Vector Machines and Extreme Gradient Boosting achieved high accuracy (c-index: 0.84 and 0.83). The Cox model also performed well (c-index: 0.83), highlighting machine learning's role in enhancing mortality predictions and informing health policy.

# About the Dataset

- **Data Source:**

The Dataset for this project is a secondary data taken from Kaggle repository.

<https://www.kaggle.com/datasets/nelgiriyeWithana/countries-of-the-world-2023/data>

- **Description:**

The Dataset contains 195 countries as rows and 27 features.

some of the features include Country, Population, Land Area, Agriculture Land, Co2- Emissions , CPI, Fertility Rate, Forested Area, Gasoline Price, GDP, Life expectancy, maternal mortality ratio, minimum wage, Density, Longitude, Latitude.



# Data overview

[illegible]



# Data Cleaning



# Handling Missing Data

dropped country with more than 90% missing values

Country	null_percentage
Eswatini	74.074
Vatican City	85.185
Monaco	55.556
Nauru	81.481
North Macedonia	77.778
Palestinian National Authority	92.593
Tuvalu	51.852

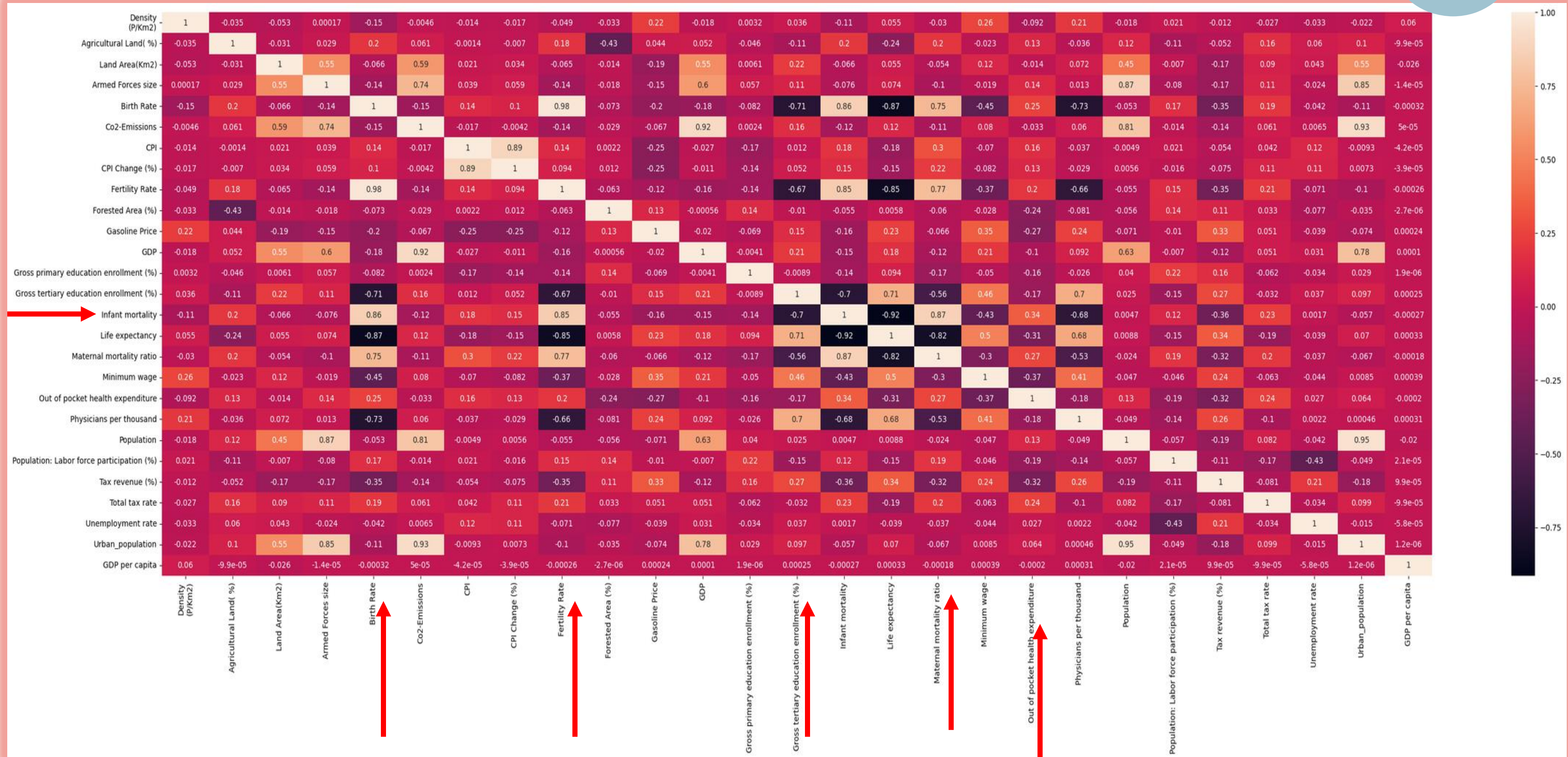
Minimum wage	20.745
Tax revenue (%)	10.638
Armed Forces size	9.043
Gasoline Price	7.447
Unemployment rate	6.383
Population: Labor force participation (%)	6.383
CPI	5.319
CPI Change (%)	4.787
Maternal mortality ratio	3.723
Total tax rate	2.660
Gross tertiary education enrollment (%)	2.660
Out of pocket health expenditure	2.128
Physicians per thousand	1.064
Infant mortality	0.532
Life expectancy	0.532
Gross primary education enrollment (%)	0.532
Co2-Emissions	0.532
Agricultural Land( %)	0.532
Forested Area (%)	0.532
Population	0.000
Country	0.000
Density\n(P/Km2)	0.000
GDP	0.000
Fertility Rate	0.000
Birth Rate	0.000
Land Area(Km2)	0.000
Urban_population	0.000

# Exploratory Data Analysis





# Correlation Matrix

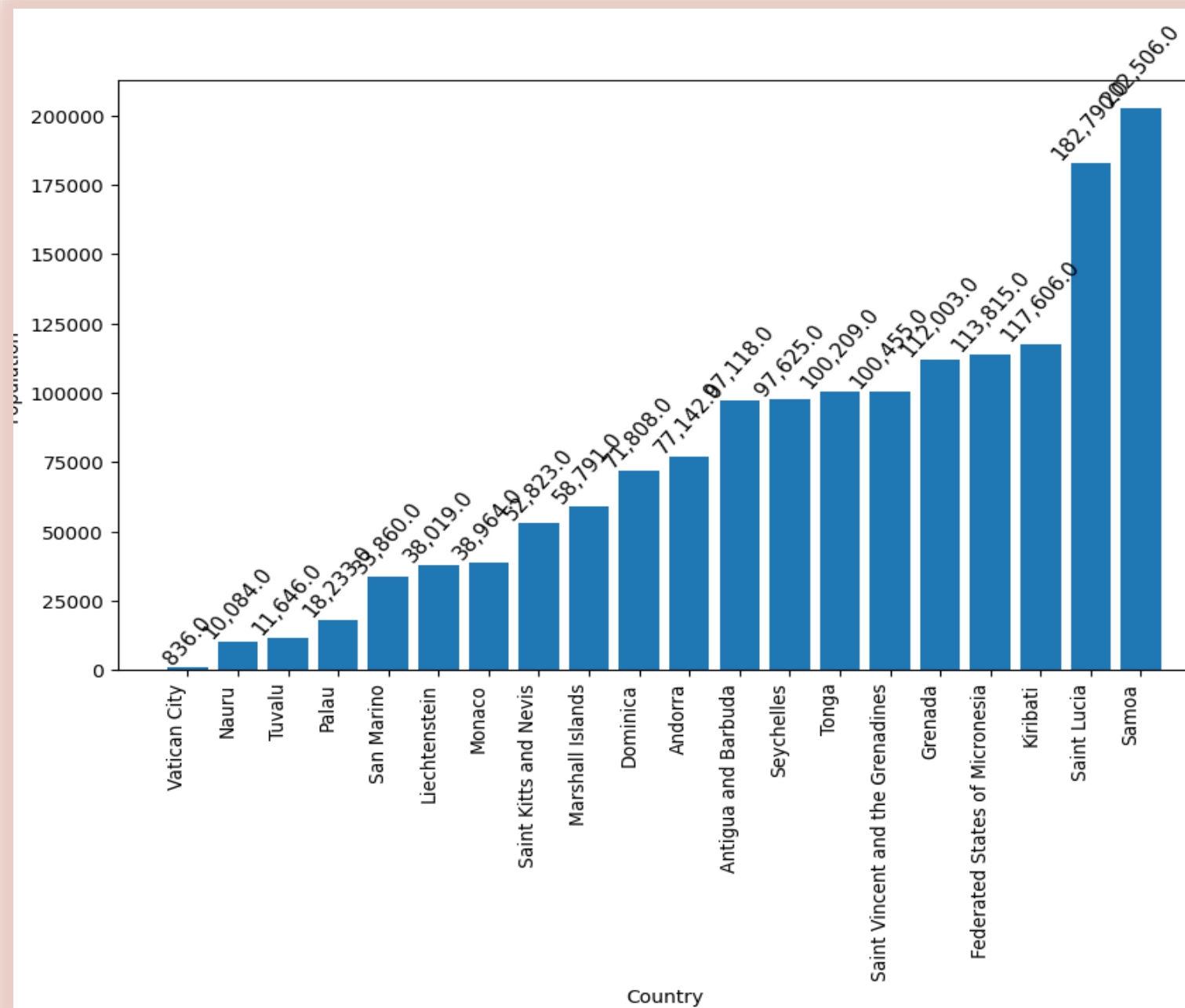




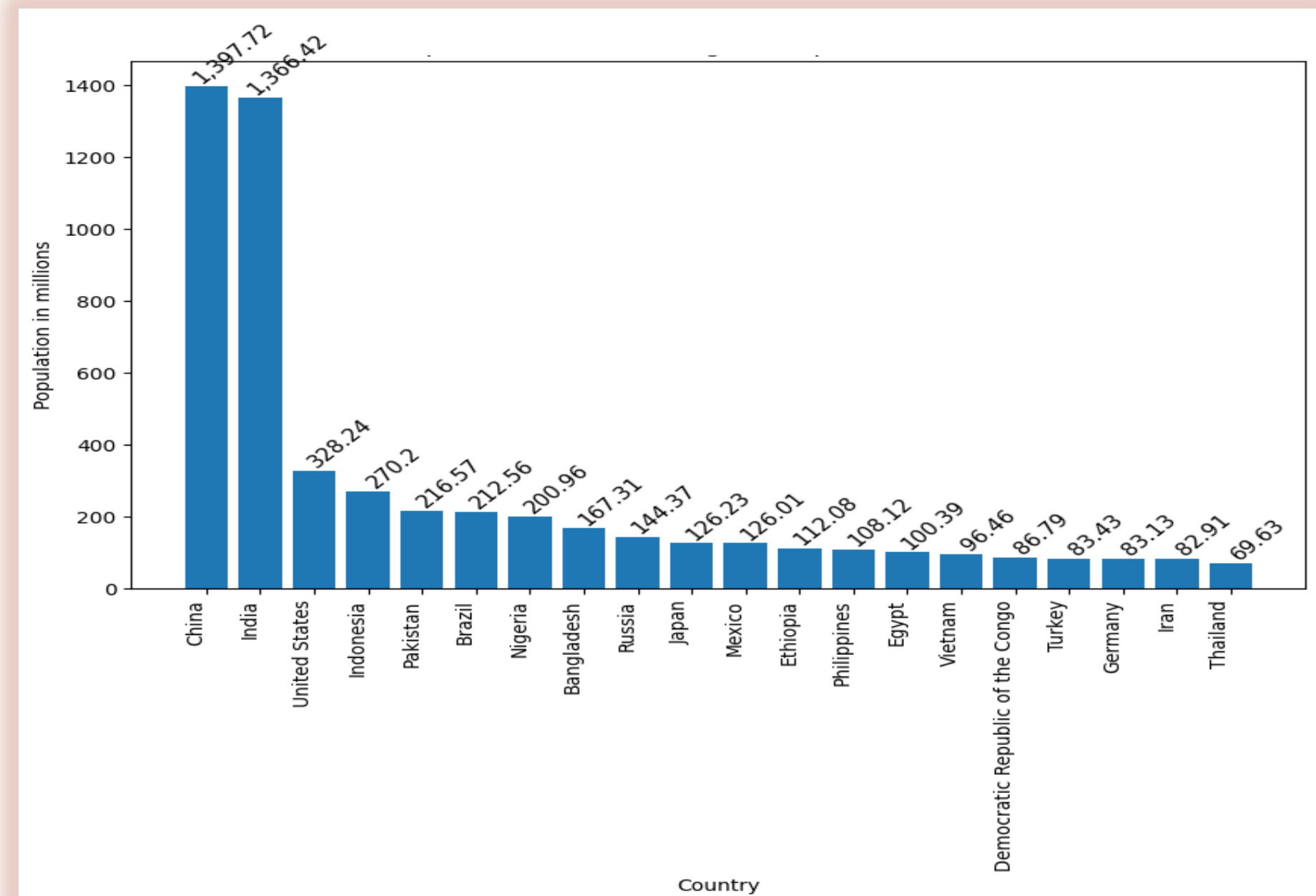
# Bar plot

## Countries vs Population

Bottom 20 Countries with lower Population

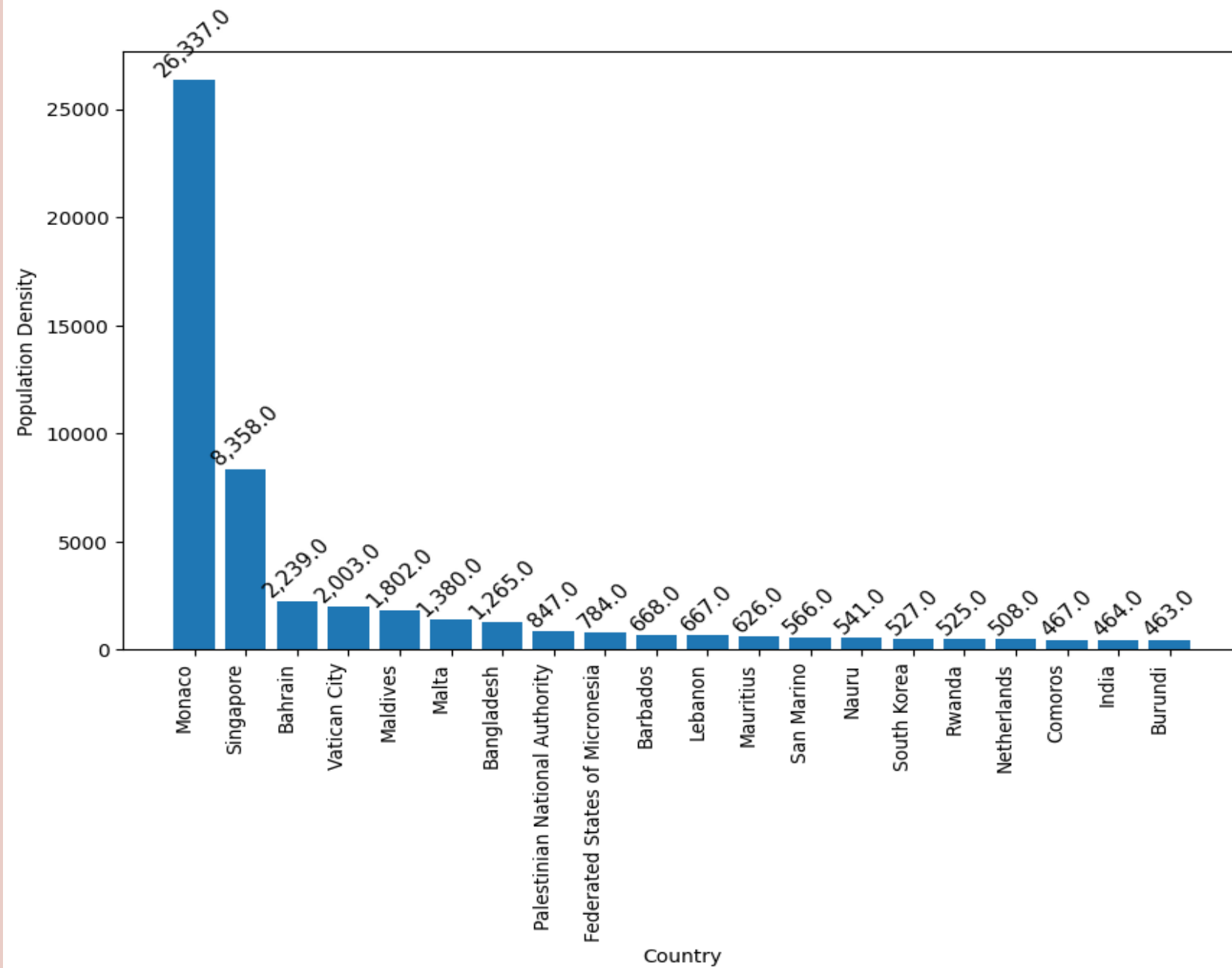


Top 20 Countries with High Population (in M)

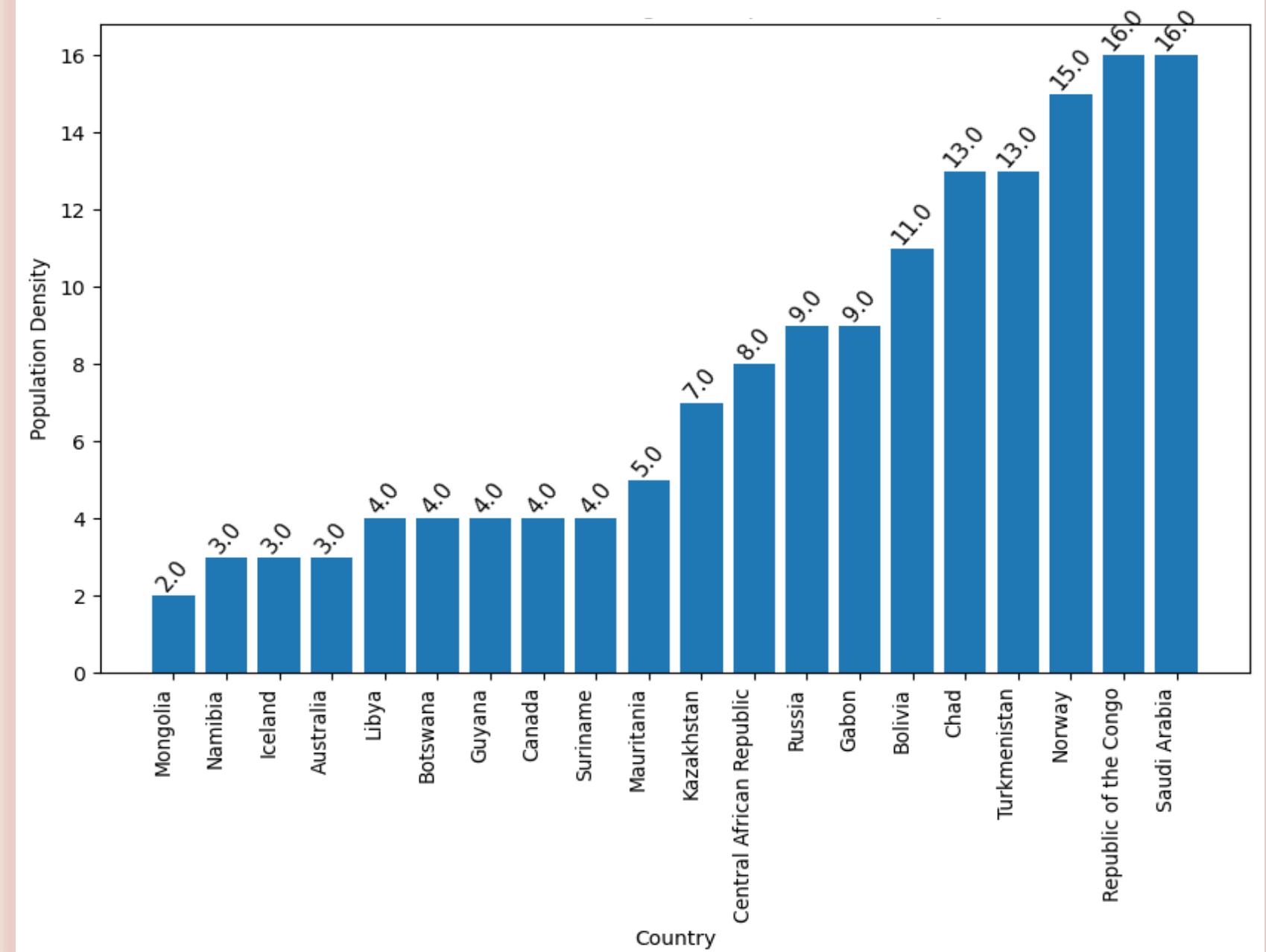


# Countries with Population Density

## Top 20 Countries with Highest Population Density

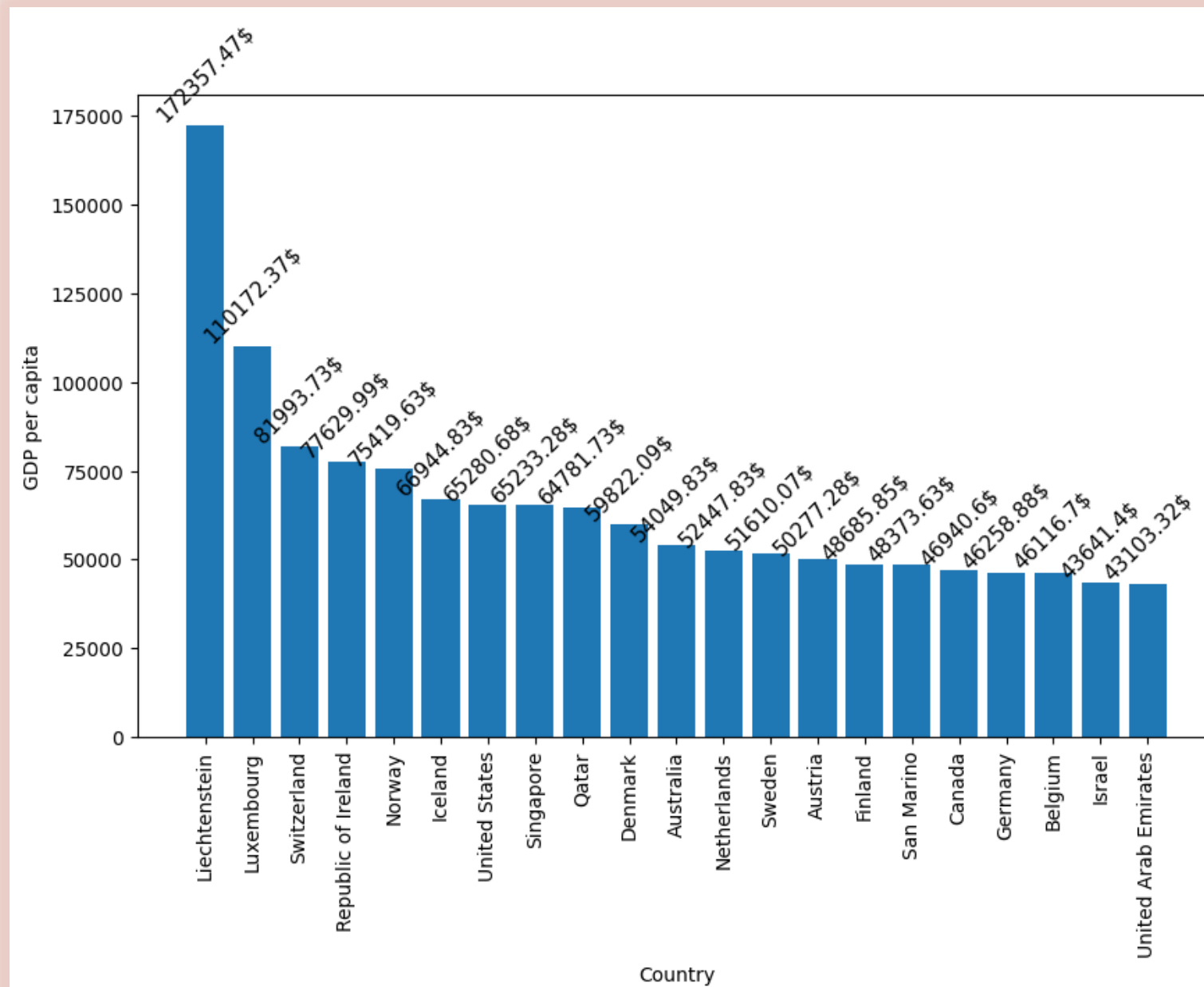


## Bottom 20 Countries with Lowest population Density

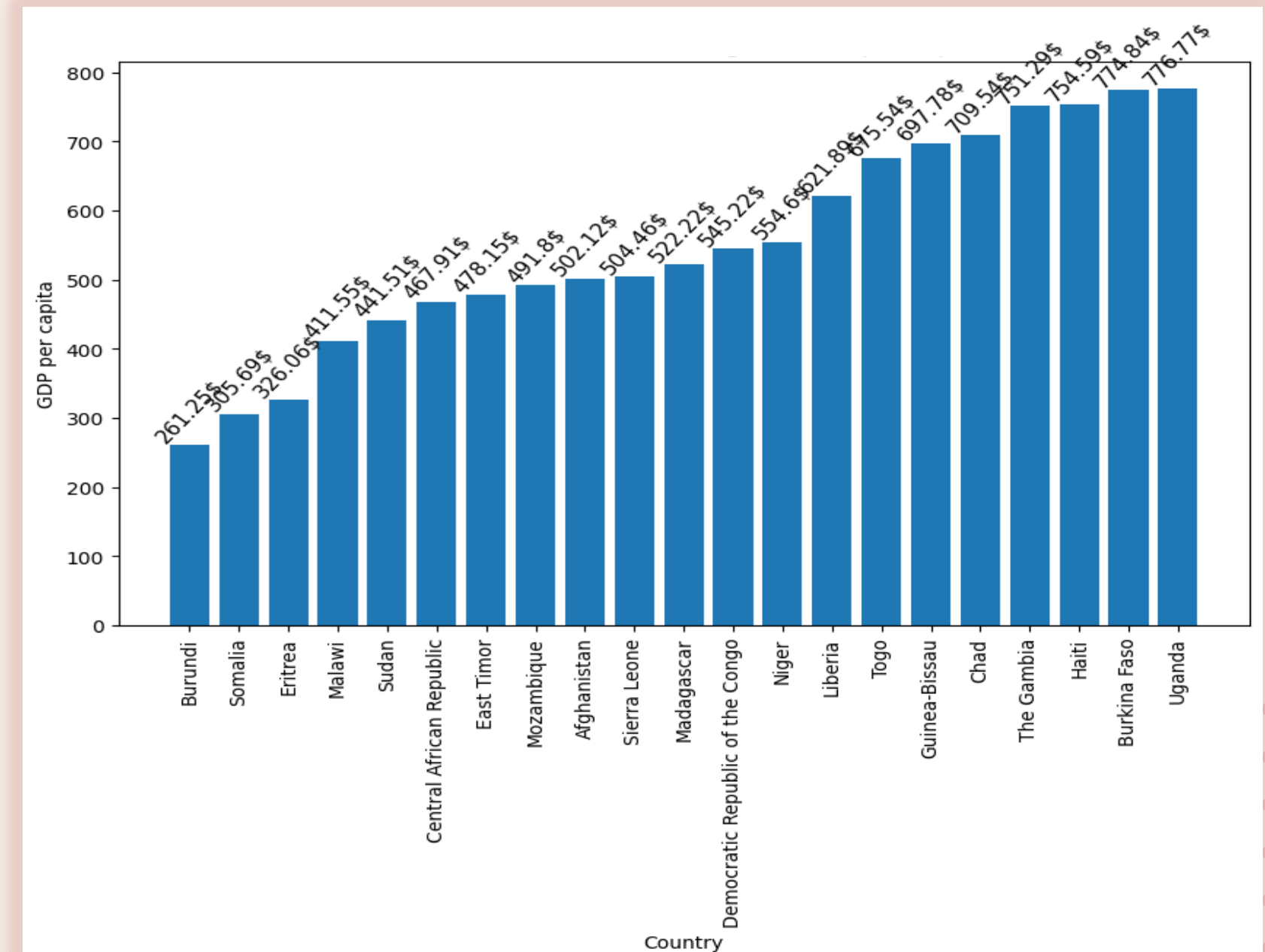


# Countries with GDP Per Capita

## Top 20 Countries with GDP per capita



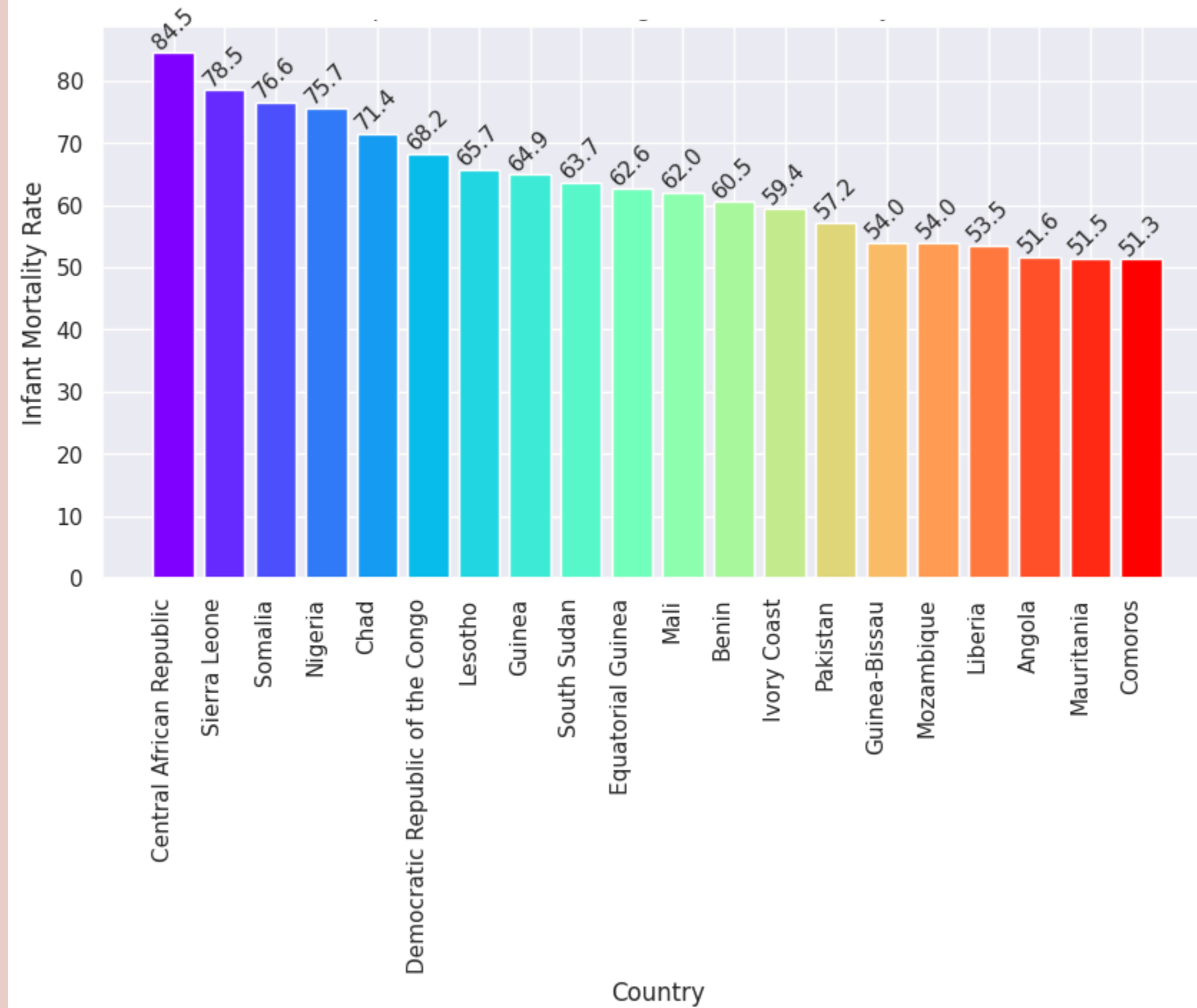
## Bottom 20 Countries with GDP per capita



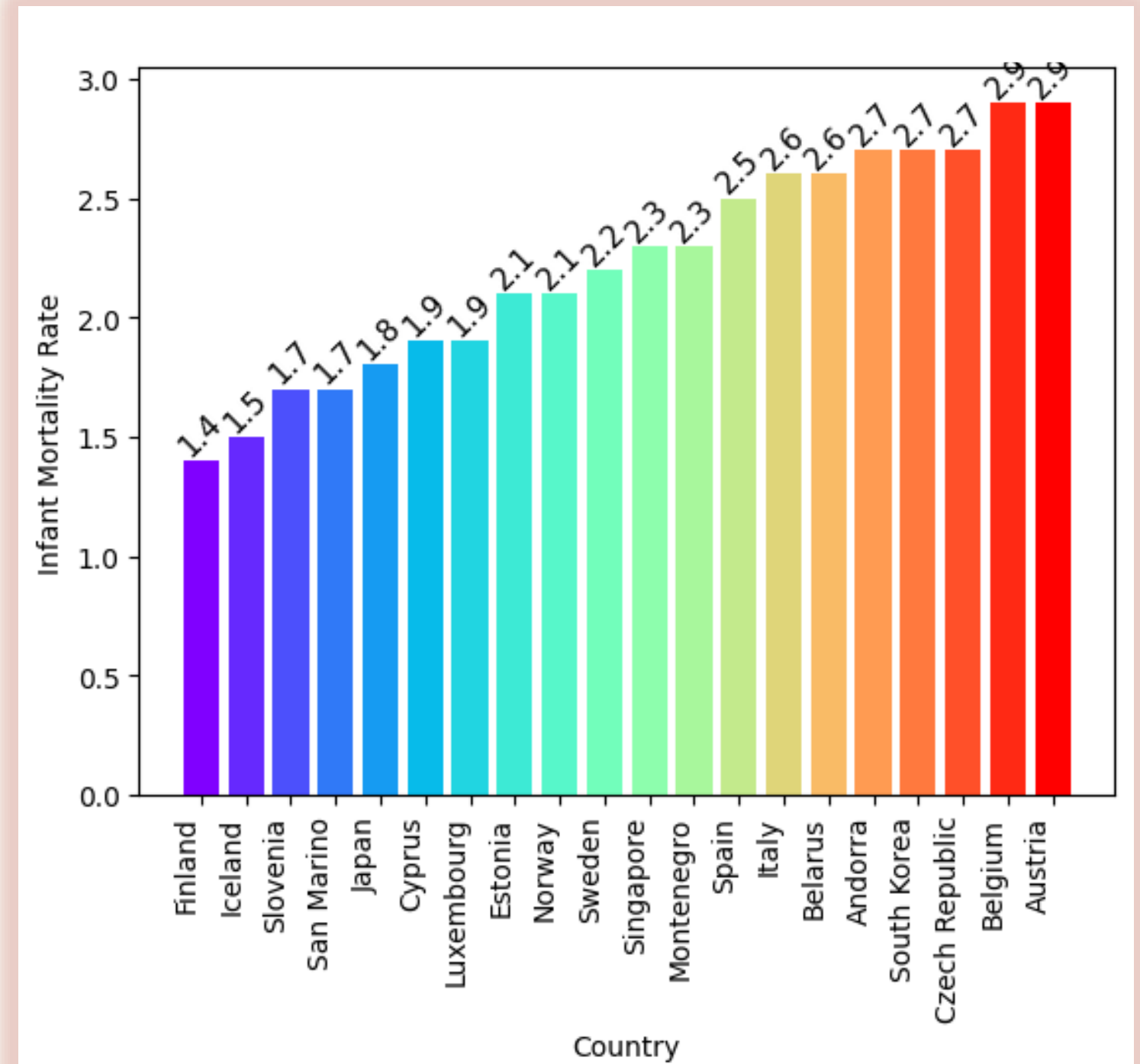


# Country with Infant Mortality

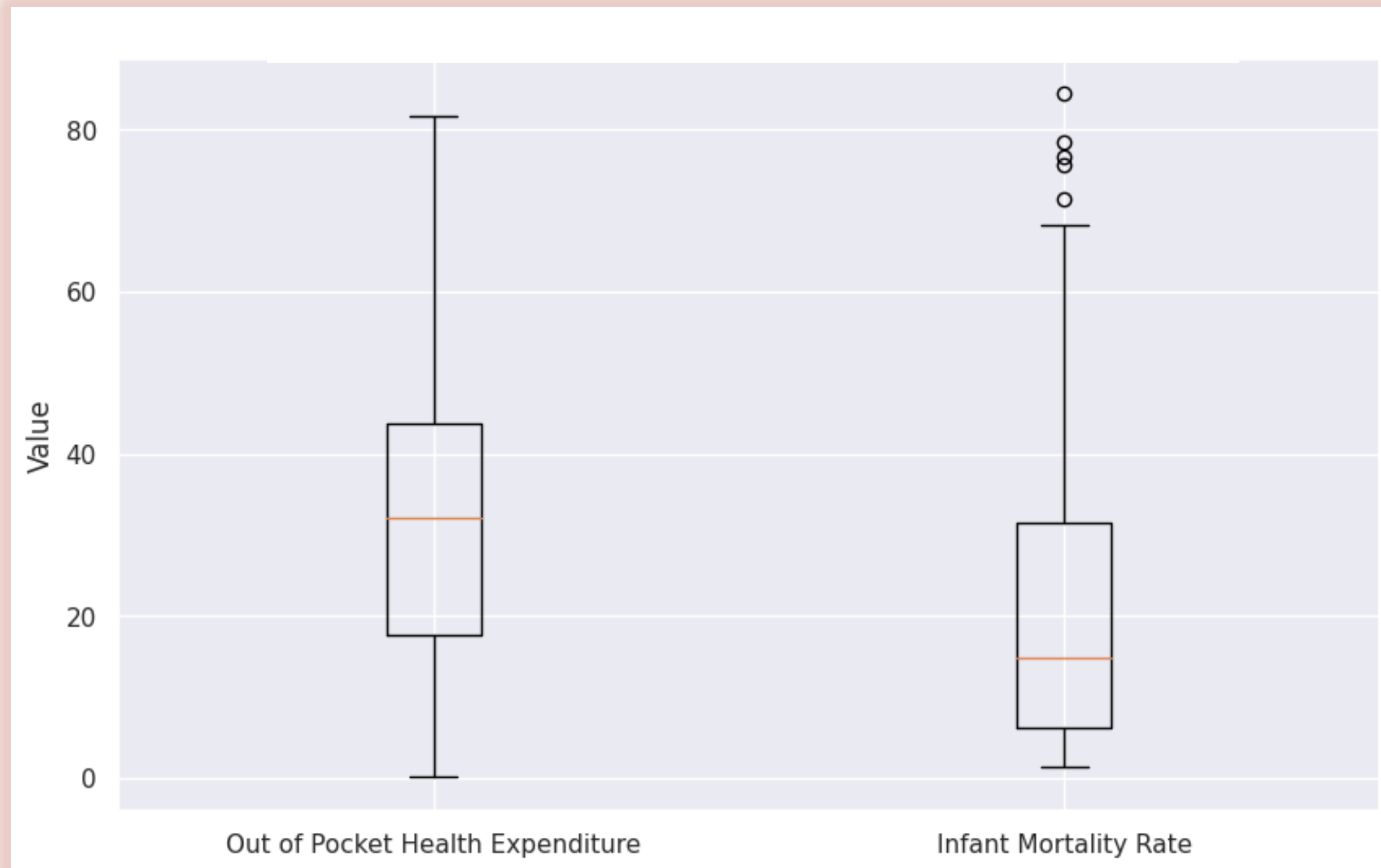
## Top 20 Countries with Highest IMR



## Bottom 20 Countries with lowest IMR



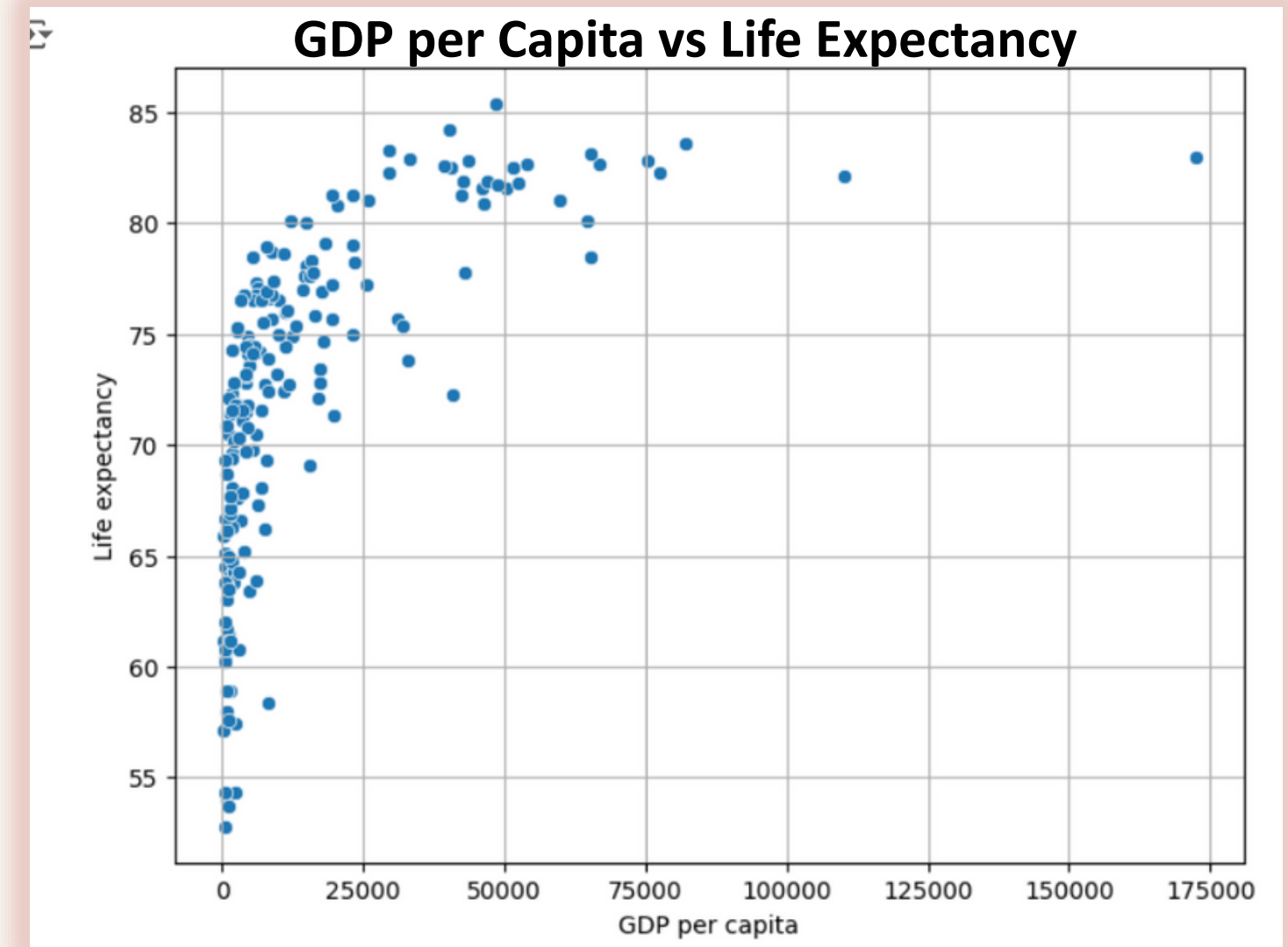
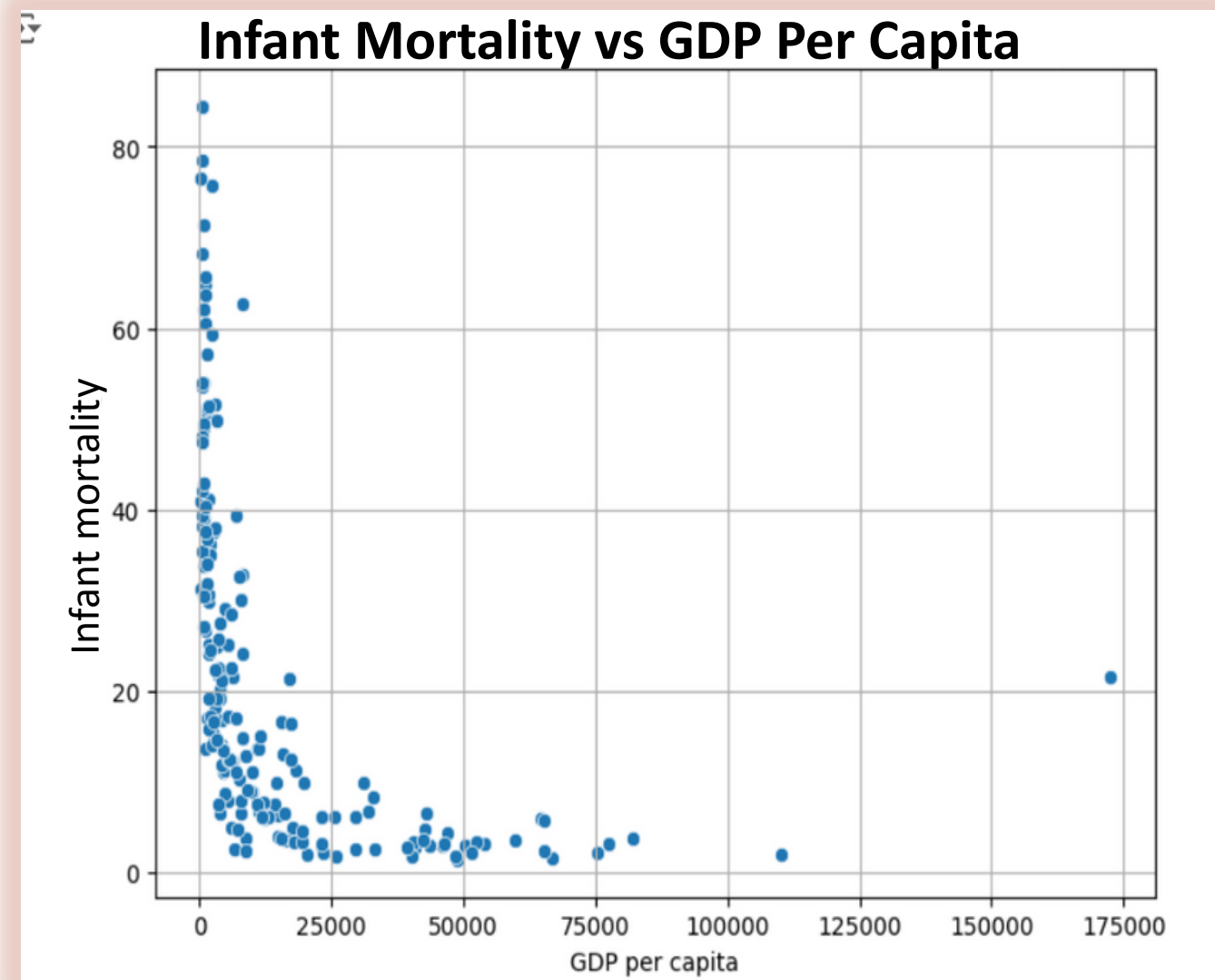
## Box Plot of Health Expenditure and Infant Mortality Rate



This means that in half of the countries, people spend 32% of their healthcare costs out of their own pockets.

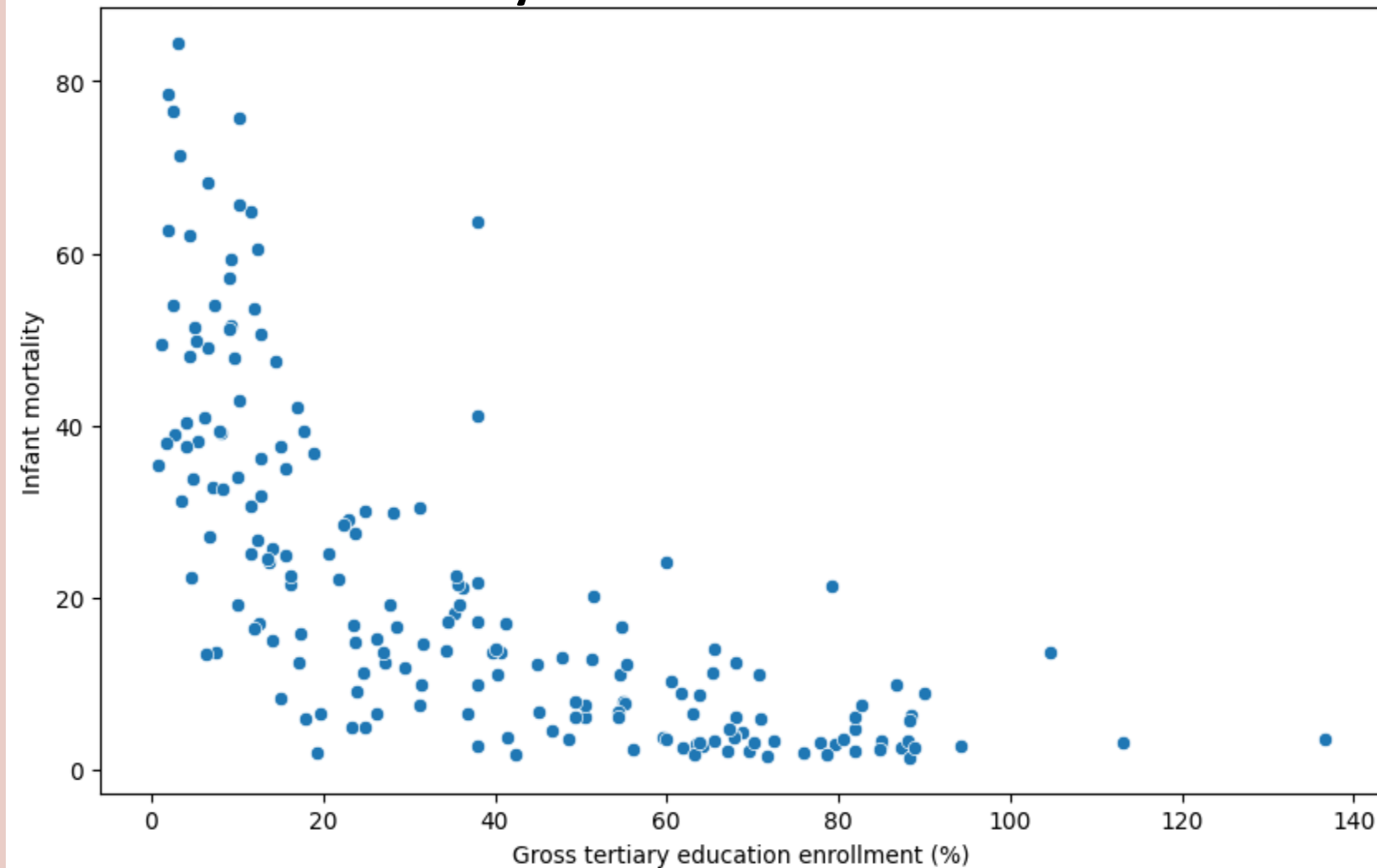
This indicates that in half of the countries, the infant mortality rate is 14 or lower.

# Scatter Plot

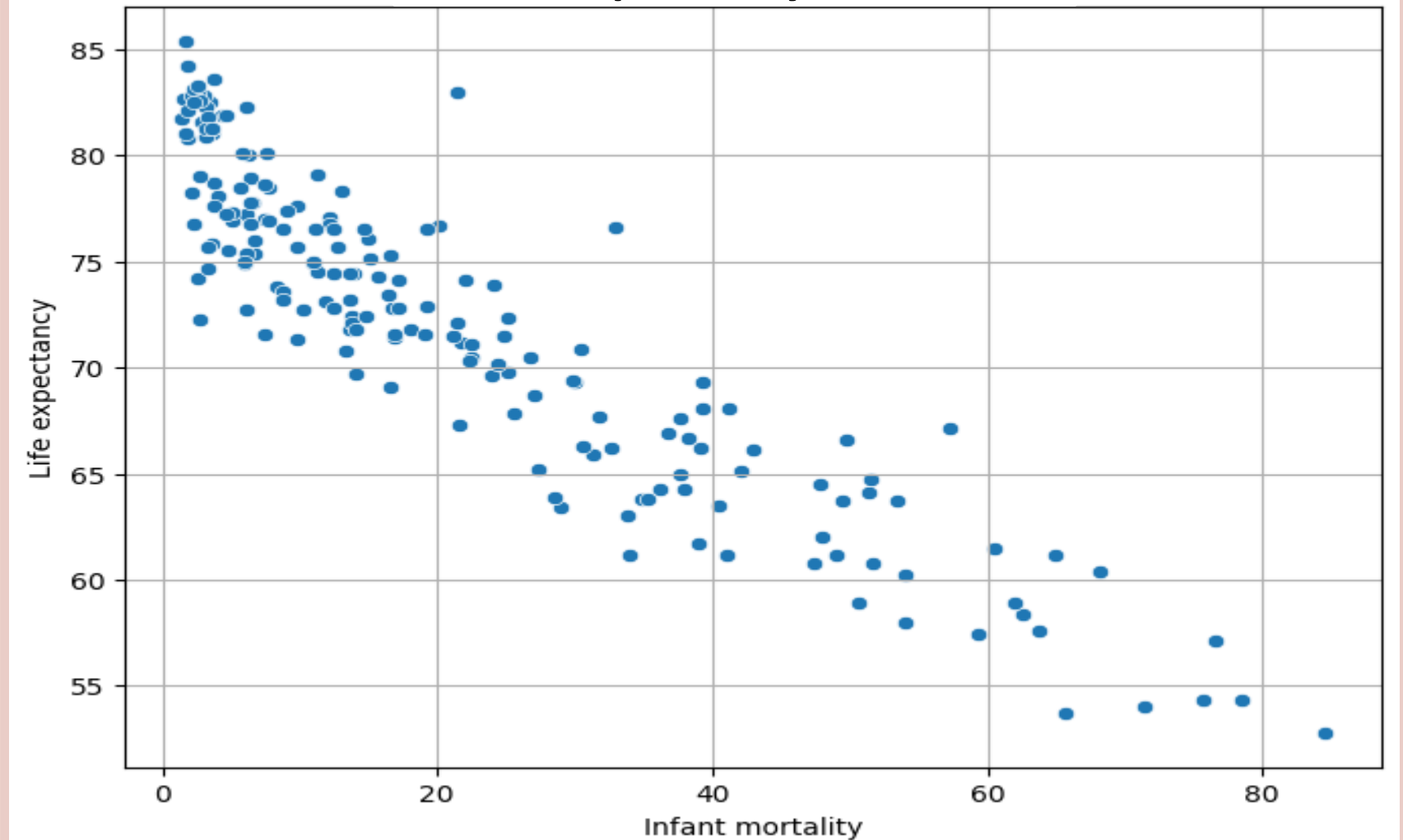


**Most of the countries have less than \$50,000 GDP per capita having Infant mortality is upto 80 years.**

**Most of the countries are concentrated with Life expectancy <80 and GDP per capita \$25,000.**

**Gross Tertiary Education enrolment vs IMR**

**The countries with higher education enrolment have lower infant mortality**  
**And higher infant mortality have lower education enrolment.**

**Life Expectancy vs IMR**

**The countries with higher life expectancy have lower infant mortality**  
**And higher infant mortality have lower life expectancy.**

# Multicollinearity check

20

Table 1

variables	VIF
Birth Rate	54.6
Urban_population	52.6
Fertility Rate	44.9
Population	28.7
Co2-Emissions	25.9
GDP	10.1
CPI	6
Armed Forces size	5.7
CPI Change (%)	5.6
Maternal mortality ratio	3.7
Physicians per thousand	3.2
Gross tertiary education enrollment (%)	3.1
GDP per capita	2.3
Land Area(Km2)	2.2
Minimum wage	2.1
Gasoline Price	1.8
Out of pocket health expenditure	1.7
Tax revenue (%)	1.6
Agricultural Land( %)	1.5
Life expectancy	1.5
Forested Area (%)	1.5
Population: Labor force participation (%)	1.5
Total tax rate	1.4
Unemployment rate	1.4
Gross primary education enrollment (%)	1.3
Density\n(P/Km2)	1.2

Table 2

variables	VIF
Urban_population	52.4
Population	28.6
Co2-Emissions	25.9
GDP	10.1
CPI	5.9
Armed Forces size	5.6
CPI Change (%)	5.5
Maternal mortality ratio	3.7
Fertility Rate	3.5
Gross tertiary education enrollment (%)	3
Physicians per thousand	2.8
GDP per capita	2.3
Minimum wage	2.1
Land Area(Km2)	2.1
Gasoline Price	1.7
Out of pocket health expenditure	1.6
Tax revenue (%)	1.6
Agricultural Land( %)	1.5
Life expectancy	1.5
Forested Area (%)	1.5
Population: Labor force participation (%)	1.5
Total tax rate	1.4
Unemployment rate	1.4
Gross primary education enrollment (%)	1.2
Density\n(P/Km2)	1.2

Table 3

variables	VIF
Co2-Emissions	16.9
GDP	10.1
Population	8.3
CPI	5.9
Armed Forces size	5.6
CPI Change (%)	5.5
Maternal mortality ratio	3.7
Fertility Rate	3.5
Gross tertiary education enrollment (%)	3
Physicians per thousand	2.8
GDP per capita	2.3
Minimum wage	2.1
Land Area(Km2)	2
Gasoline Price	1.7
Out of pocket health expenditure	1.6
Tax revenue (%)	1.6
Agricultural Land( %)	1.5
Life expectancy	1.5
Forested Area (%)	1.5
Population: Labor force participation (%)	1.5
Unemployment rate	1.4
Total tax rate	1.3
Gross primary education enrollment (%)	1.2
Density\n(P/Km2)	1.2

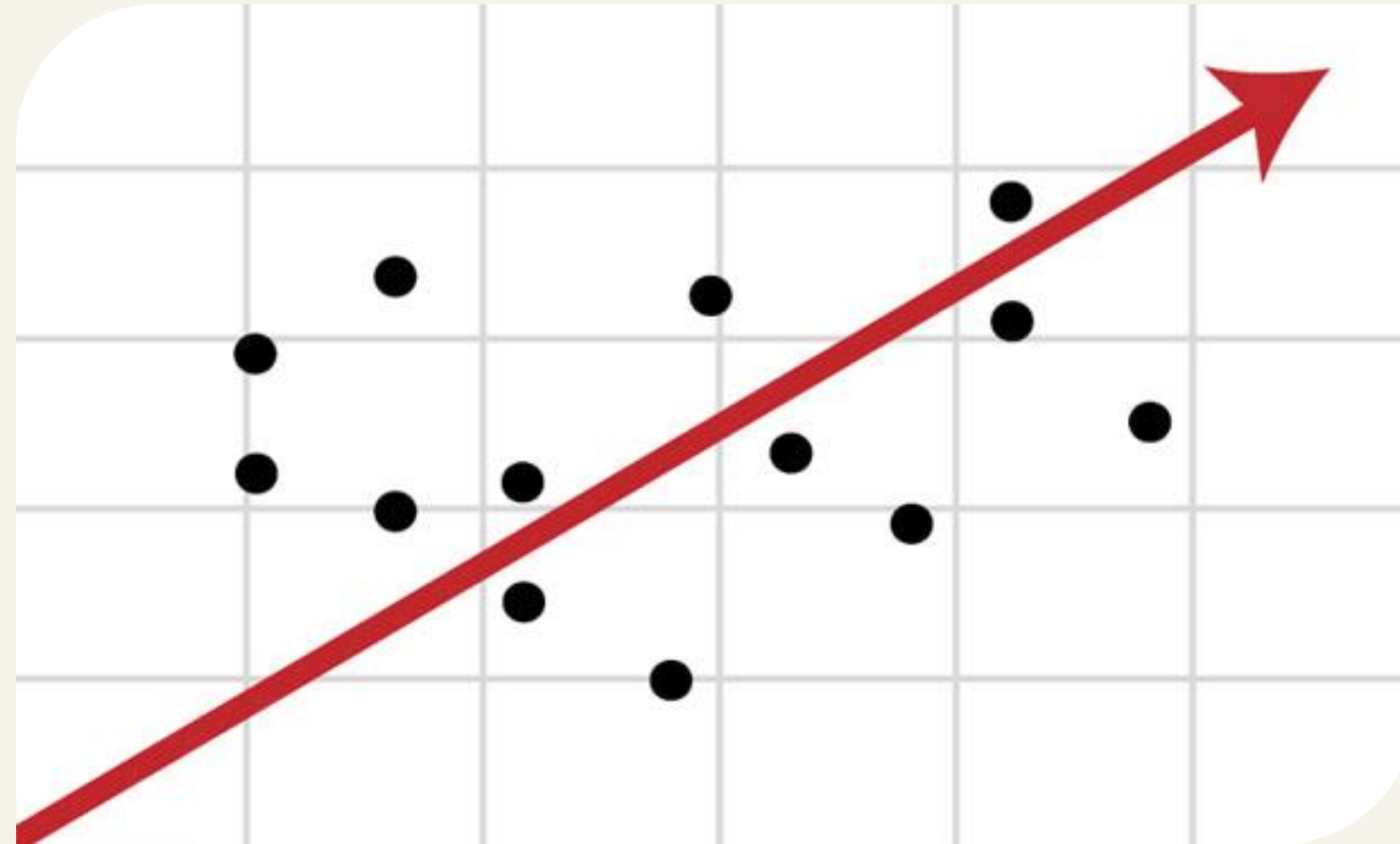
Table 4

variables	VIF
CPI	5.9
CPI Change (%)	5.5
Armed Forces size	5.4
Population	5.2
Maternal mortality ratio	3.6
Fertility Rate	3.5
Gross tertiary education enrollment (%)	3
Physicians per thousand	2.8
GDP	2.4
GDP per capita	2.2
Minimum wage	2
Land Area(Km2)	1.9
Gasoline Price	1.7
Tax revenue (%)	1.6
Agricultural Land( %)	1.5
Life expectancy	1.5
Out of pocket health expenditure	1.5
Forested Area (%)	1.5
Population: Labor force participation (%)	1.5
Unemployment rate	1.4
Total tax rate	1.3
Gross primary education enrollment (%)	1.2
Density\n(P/Km2)	1.2

The features like Birth Rate, Urban\_population, Co2-Emission were multicollinear in our dataset.

# DATA MODELLING

21





# Significance Test for Multiple Linear Regression 22

Ordinary Least squares

```

                                OLS Regression Results
=====
Dep. Variable:      Infant mortality    R-squared:                0.865
Model:              OLS                Adj. R-squared:          0.862
Method:             Least Squares       F-statistic:             294.2
Date:               Sun, 18 Aug 2024    Prob (F-statistic):      1.59e-78
Time:               12:25:20           Log-Likelihood:          -636.66
No. Observations:   188                AIC:                    1283.
Df Residuals:       183                BIC:                    1300.
Df Model:           4
Covariance Type:    nonrobust
=====
                                coef      std err          t      P>|t|      [0.025      0.975]
-----
const                0.9679         2.490         0.389     0.698     -3.945      5.881
Maternal mortality ratio  0.0411         0.004     11.164     0.000      0.034      0.048
Fertility Rate        5.1354         0.740         6.942     0.000      3.676      6.595
Physicians per thousand -2.1225         0.438        -4.846     0.000     -2.987     -1.258
Out of pocket health expenditure  0.1174         0.029         4.017     0.000      0.060      0.175
=====
Omnibus:            34.512    Durbin-Watson:           2.209
Prob(Omnibus):      0.000    Jarque-Bera (JB):        59.001
Skew:               0.941    Prob(JB):                1.54e-13
Kurtosis:           4.997    Cond. No.                1.37e+03
=====
```



# Multiple Linear Regression



$$Y(\text{Infant Mortality}) = 0.9679$$

$$+0.0411(\text{Maternal mortality ratio})$$

$$+5.1354(\text{Fertility Rate})$$

$$-2.1225(\text{Physicians per thousand})$$

$$+0.1174(\text{Out of pocket health expenditure})$$



Train-test	$R^2$	Train $R^2$
80-20	0.87	0.84
75-25	0.88	0.84
70-30	0.87	0.85
60-40	0.86	0.85

Here the Multiple Regression is the good fit with 88%  $R^2$

# Decision Tree

Train-test	max_features	min_samples_leaf	max_depth	$R^2$	Train $R^2$
80-20	4	3	4	0.73	0.85
75-25	4	3	4	0.74	0.80
70-30	4	3	4	0.69	0.86
60-40	4	3	4	0.73	0.87

Here the Decision Tree is **consistent** for different train-test ratio's as  $R^2$  is about 70%

# Random Forest

Train-test	max_features	min_samples_leaf	max_depth	$R^2$	Train $R^2$
80-20	4	3	4	0.83	0.94
75-25	4	3	4	0.85	0.94
70-30	4	3	4	0.86	0.95
60-40	4	3	4	0.84	0.95

Here the Random Forest Tree is **consistent** for different train-test ratio's as  $R^2$  is in around 83% to 86%.

# KNN

Train-test	n	$R^2$	Train $R^2$
80-20	10	0.80	0.85
75-25	10	0.82	0.84
70-30	10	0.80	0.85
60-40	10	0.80	0.86

Here the KNN is **consistent** for different train-test ratio's as  $R^2$  is about 80%.

# XGBoosting

Train-test	max_depth	n_estimators	sub_sample	$R^2$	Train $R^2$
80-20	4	100	0.031	0.71	0.78
75-25	4	100	0.031	0.64	0.65
70-30	4	100	0.031	0.53	0.70
60-40	4	100	0.031	0.55	0.76


Here the XGBoosting is **not consistent** for different train-test ratio's as  $R^2$  is ranging between 50% - 70%.

# Bagging

train-test	n_estimator	max_samples	max_features	R^2	MSE
80-20	99	0.8	4	0.87	33.37
75-25	99	0.8	4	0.86	44.92
70-30	99	0.8	4	0.86	48.88
60-40	99	0.8	4	0.86	47.49

Here the Bagging is consistent for different train-test ratio's as  $R^2$  is ranging is around 86%.

# Results



Models/splits	80-20	75-25	70-30	60-40
Multiple Linear Regression	0.87	0.88	0.87	0.86
Random Forest	0.85	0.87	0.86	0.87
Decision Tree	0.73	0.74	0.69	0.72
KNN	0.80	0.82	0.80	0.80
SVR	0.73	0.74	0.75	0.73
Bagging	0.87	0.86	0.86	0.86
XGBoosting	0.71	0.64	0.53	0.55

Here the Multiple linear is consistent for different train-test ratio's as  $R^2$  is 88%.



# Conclusion

- On comparing different machine learning algorithms the Multiple Linear Regression is the best model with best fit of 88%.

- **From Multiple Linear Regression:**

Having more doctors and longer life expectancy for countries helps reduce infant deaths, while having more children per family and higher health costs tend to increase it.

# Linear model

$Y(\text{Infant Mortality}) = 0.9679$

$+0.0411(\text{Maternal mortality ratio})$

$+5.1354(\text{Fertility Rate})$

$-2.1225(\text{Physicians per thousand})$

$+0.1174(\text{Out of pocket health expenditure})$

# Recommendation

- Number of doctors(Physicians) of the country can reduce the infant mortality.
- In order to reduce the infant mortality in the country, government's must focus on improving the numbers of doctors and improving the Health care infrastructure.
- Promote Family Planning and Education, a need for policies that promote family planning, reproductive health education, and women's empowerment, aiming to reduce high fertility rates



# Thank you

**Ruchitha | Shravanthika | Ramya | Madirai | Deepika**



---

# Contributions

34

Ruchitha

Worked on Exploratory Data Analysis:

Shravanthika

Worked on Data Cleaning, and imputing methods

Ramya

Worked on Multicollinearity

Madirai

Worked on data modeling

Deepika

Worked on data modeling

# Appendix

# Appendix

# Data Description:

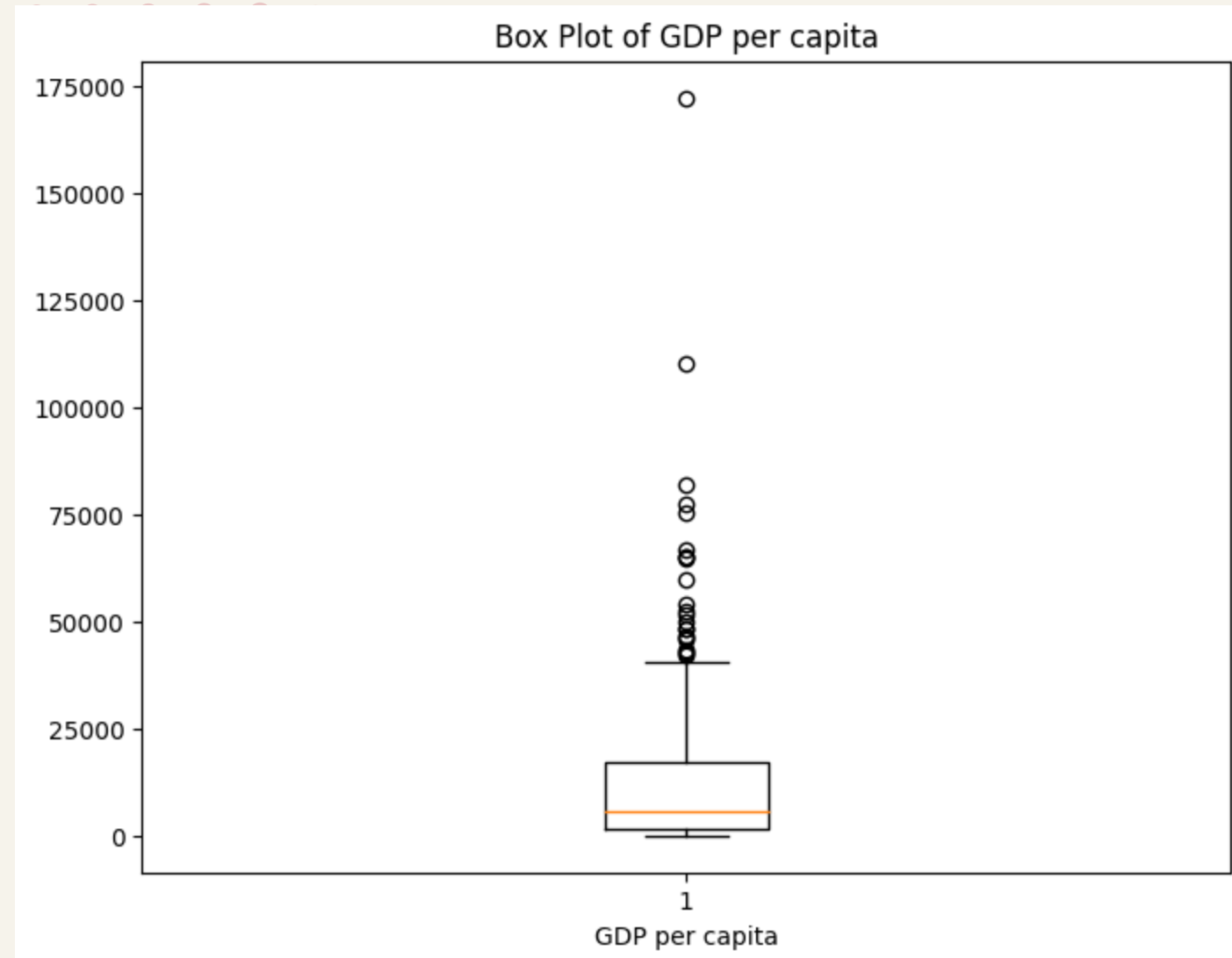
Variable Name	Description
Country	Name of the country.
Density (P/Km²)	Population density measured in persons per square kilometer.
Agricultural Land (%)	Percentage of total land area that is used for agriculture.
Land Area (Km²)	Total land area of the country in square kilometers.
Armed Forces Size	Number of personnel in the country's armed forces.
Birth Rate	Number of live births per 1,000 people in a year.
CO2 Emissions	Amount of carbon dioxide emissions in metric tons.
CPI	Consumer Price Index, a measure of the average change in prices over time.
CPI Change (%)	Percentage change in the Consumer Price Index over a specified period.
Fertility Rate	Average number of children born to a woman during her lifetime.
Forested Area (%)	Percentage of the total land area covered by forests.
Gasoline Price	Price of gasoline per liter.
GDP	Gross Domestic Product, total monetary value of all goods and services produced in a country.
Gross Primary Education Enrollment (%)	Percentage of children of official primary school age enrolled in primary school.
Gross Tertiary Education Enrollment (%)	Percentage of individuals of official tertiary education age enrolled in tertiary education.
Infant Mortality	Number of deaths of infants under one year old per 1,000 live births.
Life Expectancy	Average number of years a person is expected to live.
Maternal Mortality Ratio	Number of maternal deaths per 100,000 live births.
Minimum Wage	Lowest legal wage that can be paid to workers.
Out of Pocket Health Expenditure	Percentage of health expenses paid directly by individuals rather than covered by insurance.
Physicians per Thousand	Number of physicians per 1,000 people in the population.
Population	Total number of people living in the country.
Population: Labor Force Participation (%)	Percentage of the working-age population that is part of the labor force.
Tax Revenue (%)	Government tax revenue as a percentage of GDP.
Total Tax Rate	Total tax burden on businesses as a percentage of commercial profits.
Unemployment Rate	Percentage of the labor force that is unemployed and actively seeking employment.
Urban Population	Percentage of the total population living in urban areas.
GDP per Capita	GDP divided by the total population, representing average economic output per person.



```
1 # Get pairs of variables with correlation greater than 0.5
2 corr_matrix = data1.iloc[:,1:].corr()
3 high_corr_pairs = []
4 for i in range(len(corr_matrix.columns)):
5     for j in range(i + 1, len(corr_matrix.columns)):
6         if abs(corr_matrix.iloc[i, j]) > 0.5:
7             high_corr_pairs.append((corr_matrix.columns[i], corr_matrix.columns[j]))
8
9 print("Pairs of variables with correlation greater than 0.5:")
10 for pair in high_corr_pairs:
11     print(pair)
12
```

```
Pairs of variables with correlation greater than 0.5:
('Land Area(Km2)', 'Armed Forces size')
('Land Area(Km2)', 'Co2-Emissions')
('Land Area(Km2)', 'GDP')
('Land Area(Km2)', 'Urban_population')
('Armed Forces size', 'Co2-Emissions')
('Armed Forces size', 'GDP')
('Armed Forces size', 'Population')
('Armed Forces size', 'Urban_population')
('Birth Rate', 'Fertility Rate')
('Birth Rate', 'Gross tertiary education enrollment (%)')
('Birth Rate', 'Infant mortality')
('Birth Rate', 'Life expectancy')
('Birth Rate', 'Maternal mortality ratio')
('Birth Rate', 'Physicians per thousand')
('Birth Rate', 'GDP per capita')
('Co2-Emissions', 'GDP')
('Co2-Emissions', 'Population')
('Co2-Emissions', 'Urban_population')
```

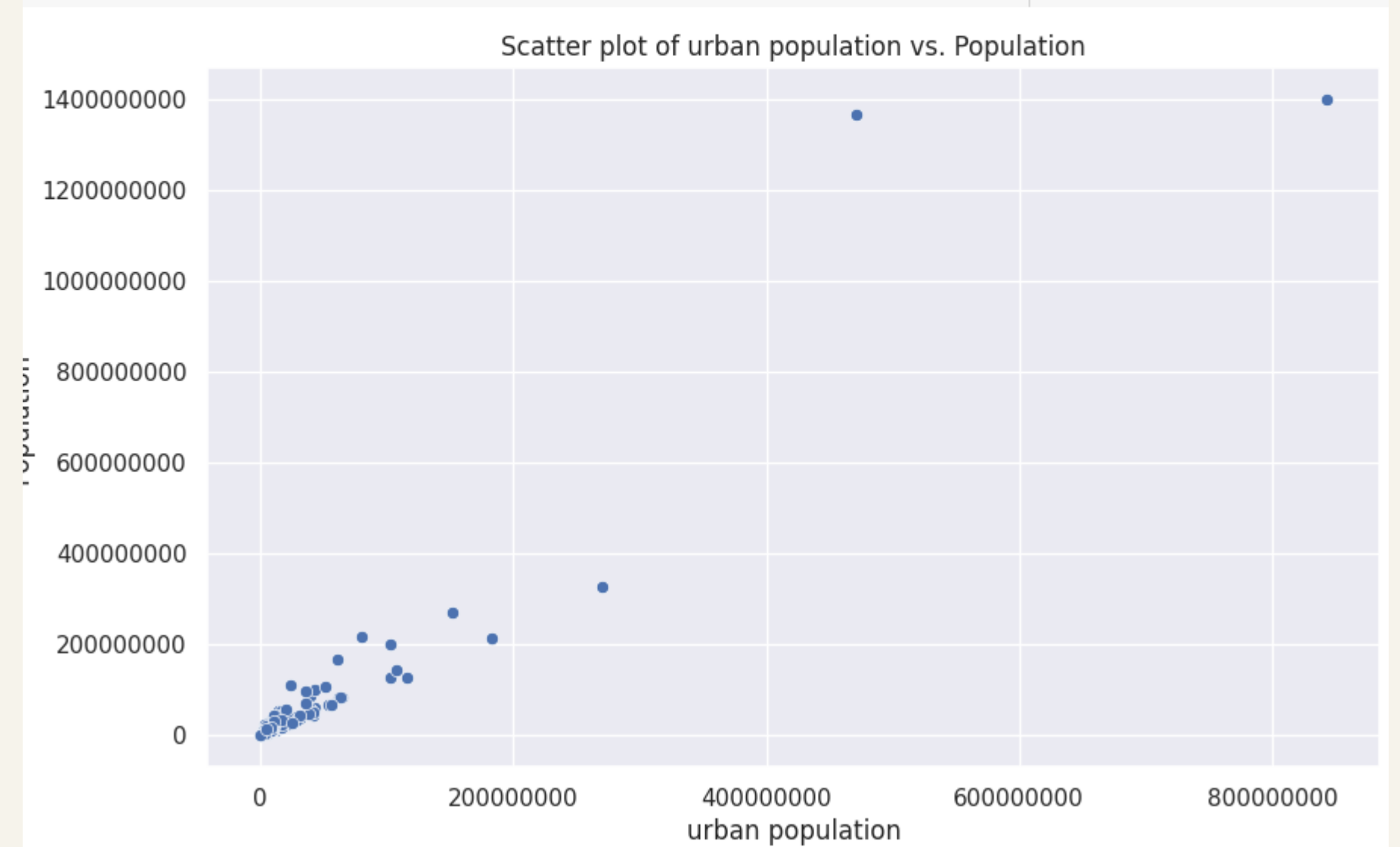
**Pair of Features which are highly Correlated with each other**

**Box Plot of GDP per Capita**

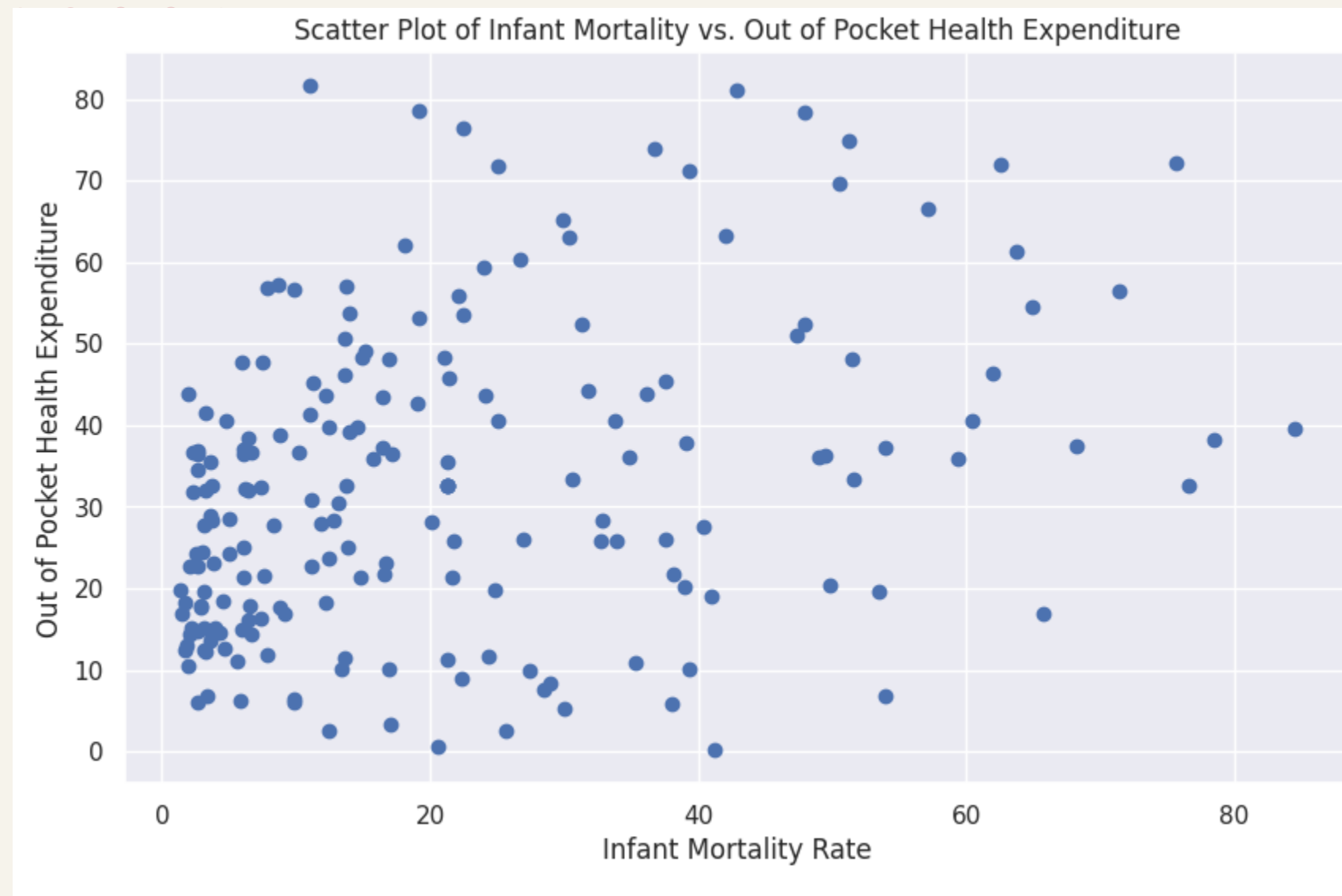
```
count 188.00  
mean 14987.56  
std 22377.72  
min 261.25  
25% 1887.26  
50% 5978.25  
75% 17498.87  
max 172357.47
```

**Statistics of GDP per capita****Scatter Plot between Urban Population and Population**

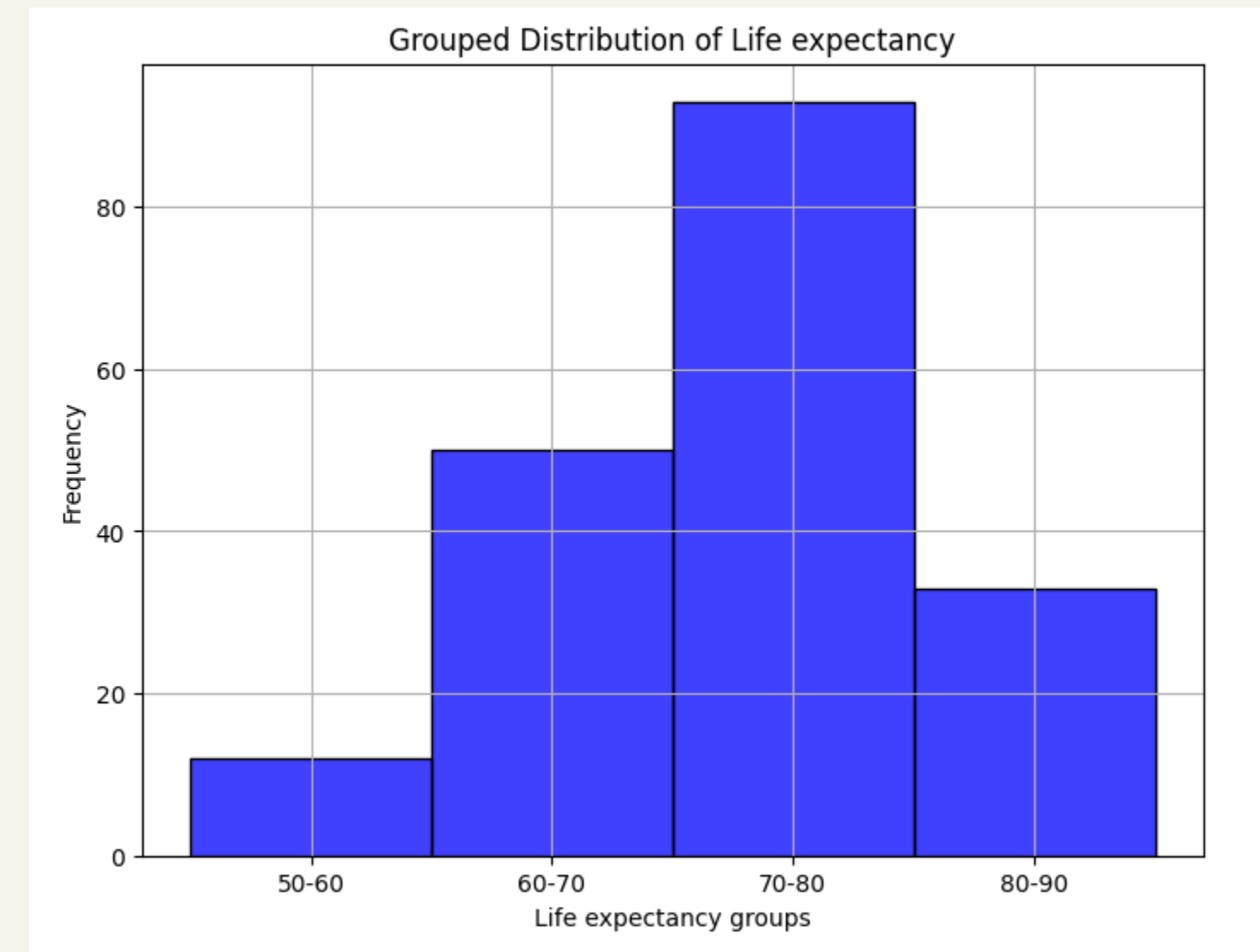
```
plt.figure(figsize=(10, 6))  
sns.scatterplot(x='Urban_population', y='Population', data=data1)  
plt.title('Scatter plot of urban population vs. Population')  
plt.ticklabel_format(style='plain', axis='y')  
plt.ticklabel_format(style='plain', axis='x')  
plt.xlabel('urban population')  
plt.ylabel('Population')  
plt.show()
```



## Scatter Plot of IMR and Pocket health Expenditure



## Grouped Life Expectancy into Group



## Code for Variance Inflation Factor (VIF)

```

1 # Import library for VIF
2 from statsmodels.stats.outliers_influence import variance_inflation_factor
3
4 def calc_vif(X):
5
6     # Calculating VIF
7     vif = pd.DataFrame()
8     vif["variables"] = X.columns
9     vif["VIF"] = [variance_inflation_factor(X.values, i).round(1) for i in range(X.shape[1])]
10
11     return(vif)
12
13 vif_df = calc_vif(X)
14 vif_df.sort_values(by='VIF',ascending=False)

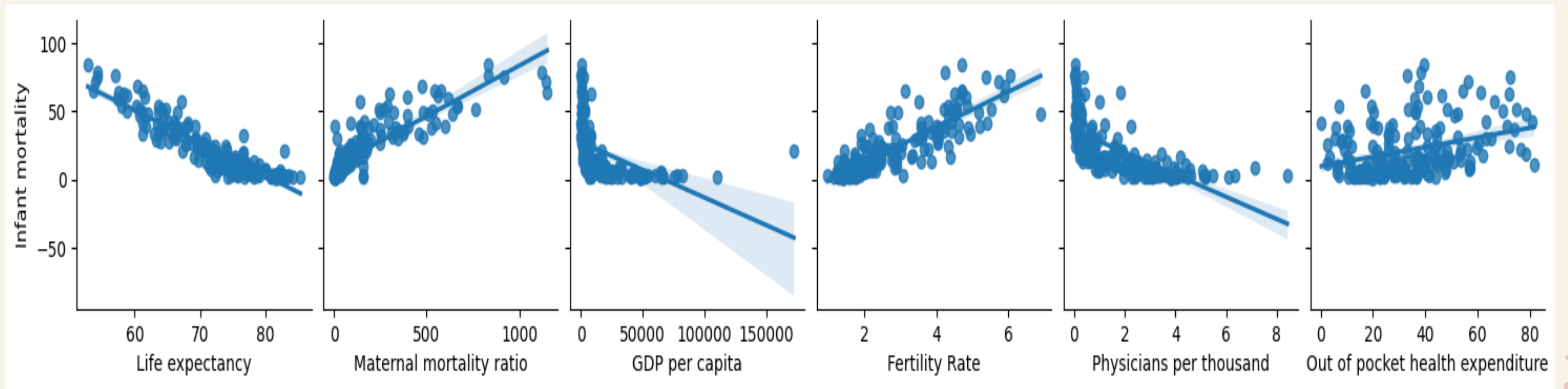
```

	variables	VIF
4	Birth Rate	54.6
24	Urban_population	52.6
8	Fertility Rate	44.9
19	Population	28.7
5	Co2-Emissions	25.9
11	GDP	10.1
6	CPI	6.0
3	Armed Forces size	5.7
7	CPI Change (%)	5.6
15	Maternal mortality ratio	3.7
18	Physicians per thousand	3.2
13	Gross tertiary education enrollment (%)	3.1
25	GDP per capita	2.3
2	Land Area(Km2)	2.2
16	Minimum wage	2.1
10	Gasoline Price	1.8
17	Out of pocket health expenditure	1.7
21	Tax revenue (%)	1.6
1	Agricultural Land( %)	1.5
14	Life expectancy	1.5

## Code for Ordinary Least Square

```
1 x_nomulti_colinearity=X.drop(['Urban_population', 'Birth Rate', 'Co2-Emissions'],axis=
2
3 import statsmodels.api as sm
4
5 x_train_constant = sm.add_constant(x_nomulti_colinearity)
6
7 model = sm.OLS(y, x_train_constant).fit()
8
9 # Print the summary to get p-values
10 print(model.summary())
11 print(model.pvalues.sort_values(ascending=False).round(4))
```

Pair plot for Significant Features



# Multiple Linear Regression

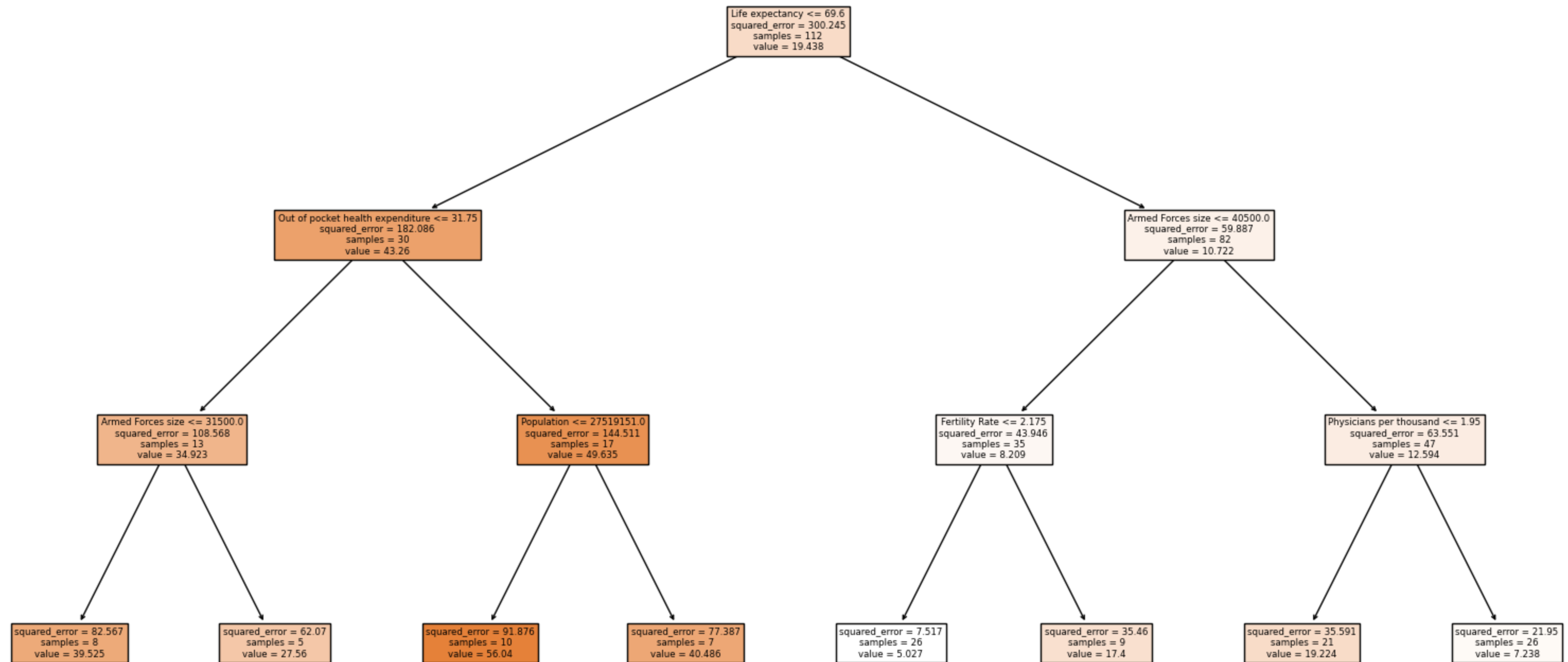
```
1 import statsmodels.api as sm
2 import statsmodels.formula.api as smf
3 import sklearn as sk
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6 from sklearn.metrics import mean_squared_error, r2_score
7 print("FOR DIFFERENT TRAIN TEST SPLITS: MULTIPLEREGRESSION")
8 for i in [0.2,0.25,0.3,0.4]:
9     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=i, random_state=0)
10
11     # instantiate and fit
12     lm2 = LinearRegression()
13     lm2.fit(X_train, y_train)
14     y_pred = lm2.predict(X_train)
15
16     r2 = r2_score(y_train, y_pred)
17     print("-----\n")
18     print(f"R^2 Score for the training set {(1-i)*100}% is : {r2.round(4)}")
19
20     y_pred = lm2.predict(X_test)
21     r2 = r2_score(y_test, y_pred)
22     print(f"R^2 Score for the testing {(i)*100}% set is : {r2.round(4)}")
23     print("Mean squared error:",(metrics.mean_squared_error(y_test, y_pred)))
24
25
```

# Decision Tree

```
1 from sklearn.tree import DecisionTreeRegressor, plot_tree
2 from sklearn.metrics import mean_squared_error, r2_score
3 from sklearn.model_selection import GridSearchCV
4
5 for i in [x_nomulti_colinearity]:
6     X_train, X_test, y_train, y_test = train_test_split(i, y, test_size=0.2, random_state=0)
7     regressor = DecisionTreeRegressor(max_features=4,min_samples_leaf=3,max_depth=3)
8     regressor.fit(X_train, y_train)
9     y_pred = regressor.predict(X_test)
10    print("\nTrain Mean Squared Error: ", mean_squared_error(y_train, regressor.predict(X_train)).round(3))
11    print("Test Mean Squared Error: ", mean_squared_error(y_test, y_pred).round(3))
12    print("\nTrain R^2 Score: ", r2_score(y_train, regressor.predict(X_train)).round(3))
13    print("Test R^2 Score: ", r2_score(y_test, y_pred).round(3))
14
```



# Decision Tree Diagram



# Random Forest

46

```
1 for i in [0.2,0.25,0.3,0.4]:
2     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=i, random_state=0)
3
4     # Creating the Random Forest Regressor model
5     model = RandomForestRegressor(n_estimators=100,max_depth=3,min_samples_leaf=2,min_samples_split=2,max_features=4)
6
7     # Training the model on the training data
8     model.fit(X_train, y_train)
9
10    # Predicting on the test data
11    y_pred = model.predict(X_train)
12
13
14    r2 = r2_score(y_train, y_pred)
15    print("-----\n")
16    print(f"R^2 Score for the training set {(1-i)*100}% is : {r2.round(4)}")
17
18    y_pred = model.predict(X_test)
19    r2 = r2_score(y_test, y_pred)
20    print(f"R^2 Score for the testing {(i)*100}% set is : {r2.round(4)}")
21    print("Mean squared error:",(metrics.mean_squared_error(y_test, y_pred)))
22
23
```

```
1 from sklearn.neighbors import KNeighborsRegressor
2 from sklearn.metrics import accuracy_score, r2_score
3 from sklearn.preprocessing import StandardScaler
4
5 scaler = StandardScaler()
6 y_array = y.to_numpy()
7 y_resaped = y_array.reshape(-1, 1)
8
9
10 y_scaled = pd.DataFrame(scaler.fit_transform(y_resaped), columns=['Infant mortality'])
11
12 for i in [x_nomulti_colinearity]:
13     # Fit the scaler to the data and transform it
14     df_standardized = pd.DataFrame(scaler.fit_transform(i), columns=i.columns)
15     X_train, X_test, y_train, y_test = train_test_split(df_standardized, y_scaled, test_size=0.2, random_state=0)
16     for j in range(1,20):
17         model=KNeighborsRegressor(n_neighbors=j)
18         model.fit(X_train,y_train)
19         y_pred=model.predict(X_test)
20         print(f'-----N={j}-----')
21         print("Train Mean Squared Error: ", mean_squared_error(y_train, model.predict(X_train)).round(3))
22         print("Test Mean Squared Error: ", mean_squared_error(y_test, y_pred).round(3))
23         print()
24         print("Train R^2 Score: ", r2_score(y_train, model.predict(X_train)).round(3))
25         print("Test R^2 Score: ", r2_score(y_test, y_pred).round(3))
26
```

# Support Vector Machine

48

```
1 from sklearn.svm import SVR
2
3
4 for i in [0.2,0.25,0.3,0.4]:
5     X_train, X_test, y_train, y_test = train_test_split(df_standardized, y_scaled, test_size=i, random_state=0)
6
7     model=SVR(C=0.89)
8     model.fit(X_train,y_train)
9
10    y_pred = model.predict(X_train)
11    r2 = r2_score(y_train, y_pred)
12    print("-----\n")
13    print(f"R^2 Score for the training set {(1-i)*100}% is : {r2.round(4)}")
14
15    y_pred = model.predict(X_test)
16    r2 = r2_score(y_test, y_pred)
17    print(f"R^2 Score for the testing {(i)*100}% set is : {r2.round(4)}")
18    print("Mean squared error:",(metrics.mean_squared_error(y_test, y_pred)))
19
```

# XGBossting

49

```
1 from sklearn.model_selection import train_test_split
2 from xgboost import XGBRegressor
3 from sklearn.metrics import r2_score
4
5
6 for i in [0.2,0.25,0.3,0.4]:
7     X_train, X_test, y_train, y_test = train_test_split(x_nomulti_colinearity, y, test_size=i, random_state=0)
8
9     # Creating the XGBoost Regressor model
10    model = XGBRegressor(max_depth=4,
11                          n_estimators=100,
12                          subsample=0.0310,
13
14                          random_state=0)
15
16    # Training the model on the training data
17    model.fit(X_train, y_train)
18
19    # Predicting on the test data
20
21
22    y_pred = model.predict(X_train)
23    r2 = r2_score(y_train, y_pred)
24    print("-----\n")
25    print(f"R^2 Score for the training set {(1-i)*100}% is : {r2.round(4)}")
26
27    y_pred = model.predict(X_test)
28    r2 = r2_score(y_test, y_pred)
29    print(f"R^2 Score for the testing {(i)*100}% set is : {r2.round(4)}")
30    print("Mean squared error:",(metrics.mean_squared_error(y_test, y_pred)).round(3))
```

# Bagging

50

```
1 from sklearn.ensemble import BaggingRegressor
2 from sklearn.tree import DecisionTreeRegressor
3
4
5 X=x_nomulti_colinearity=data1.drop(['Country','Urban_population','Birth Rate','Co2-Emissions'],axis=1)
6 y=data1['Infant mortality']
7 for i in [0.2,0.25,0.3,0.4]:
8     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=i, random_state=0)
9
10    # Creating the Bagging Regressor model with Decision Tree as base estimator
11    model = BaggingRegressor(
12        base_estimator=DecisionTreeRegressor(),
13        n_estimators=99, # Number of trees in the ensemble
14        random_state=42,
15        max_samples=0.8, # Fraction of samples to use for each tree
16        max_features=4
17    )
18
19    # Training the model on the training data
20    model.fit(X_train, y_train)
21
22    # Predicting on the test data
23    y_pred = model.predict(X_test)
24
25    # Calculating the R2 score
26    r2 = r2_score(y_test, y_pred)
27
28    # Output the R2 score
29    print(f'\nR2 for test ration {i*100}% score: {r2:.4f}')
30    print("Mean squared error:",(metrics.mean_squared_error(y_test, y_pred)))
```



# Thank you

**Ruchitha | Shravanthika | Ramya | Madirai | Deepika**




---






# 80-20 split



Models/splits	R <sup>2</sup>
Multiple Regression	0.927
Random Forest	0.853
Dession Tree	0.731
KNN	0.803
SVR	0.735
Bagging	0.870
XGBoosting	0.712

Here for 80-20 train-test split Multiple Regression is the best model as compared to the other models.


# 75-25 split



Models/splits	R <sup>2</sup>
Multiple Regression	0.933
Random Forest	0.873
Dession Tree	0.744
KNN	0.820
SVR	0.749
Bagging	0.862
XGBoosting	0.647

Here for 75-25 train-test split Multiple Regression is the best model as compared to the other models.

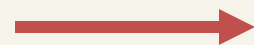
# 70-30 split



Models/splits	R <sup>2</sup>
Multiple Regression	0.926
Random Forest	0.864
Dession Tree	0.693
KNN	0.807
SVR	0.758
Bagging	0.866
XGBoosting	0.539

Here for 70-30 train-test split Multiple Regression is the best model as compared to the other models.

# 60-40 split



Models/splits	R <sup>2</sup>
Multiple Regression	0.913
Random Forest	0.878
Dession Tree	0.729
KNN	0.802
SVR	0.736
Bagging	0.869
XGBoosting	0.558

Here for 60-40 train-test split Multiple Regression is the best model as compared to the other models.