# EAST WEST INSTITUTE OF TECHNOLOGY

## DEPARTMENT OF INFORMATION SCIENCE AND ENGINEERING

## INTERNSHIP PRESENTATION ON

## GDP ANALYSIS USING DATA SCIENCE

### Presented By

### Praveen R
### (1EW20IS056)

**EXTERNAL GUIDE**
**S.Romesh**
**Trainer, Technofly**

**INTERNAL GUIDE**
**Mrs. Veena N Iter**
**Asst. Prof**
**Dept. of ISE**

# TABLE OF CONTENTS

❖ ABSTRACT

❖ INTRODUCTION

❖ COMPANY PROFILE

❖ PROPOSED SYSTEM

❖ RESULTS

❖ CONCLUSION

❖ APPLICATIONS

❖ REFERENCES

# ABSTRACT

Gross Domestic Product (GDP) analysis using data science is a critical and multifaceted approach to understanding and evaluating the economic performance of a country or region. GDP serves as a fundamental metric for assessing the overall health and growth of an economy, and data science techniques have become indispensable in extracting valuable insights from the vast amount of data generated in modern economies. The abstract provides an overview of how data science is applied to GDP analysis, highlighting its significance, methodologies, and key takeaways. GDP is a comprehensive measure of a nation's economic activity, encompassing the total value of goods and services produced within a specific time frame. It serves as a crucial indicator for policymakers, businesses, and investors to make informed decisions. conclusion, GDP analysis with data science represents a powerful approach to understanding and managing economic performance. By harnessing the capabilities of data science, stakeholders can gain deeper insights, enhance decision-making, and adapt more effectively to the dynamic nature of modern economies.

# INTRODUCTION

❖ Gross Domestic Product (GDP) measures the economic performance of a country.

❖ Sum of all goods and services produced within a country in a specific time frame.

❖ It is one of the most widely used measures of a nation's economic performance and is essential for assessing and comparing the economic health and growth of different countries.

❖ Gross Domestic Product (GDP) is a key tool that guides investors, policymakers, and businesses in strategic decision-making.

❖ Data Science is a field of study that entails applying several scientific methods, algorithms, and processes to extract insights from massive amounts of data.

❖ In other words it is broadly defined as the capability of a machine to imitate intelligent human behaviour.

# ❖ COMPANY PROFILE

➢ Technofly was formed by professionals with formal qualifications and industrial experience in the fields of embedded systems, real-time software, process control and industrial electronics.

➢ The company is professionally managed and supported by qualified experienced specialists and consultants with experience in embedded systems - including hardware and software.

➢ Technofly Developed system software tools; these include C Compilers for micro-controllers and other supporting tools such as assembler, linker, simulator and Integrated Development Environment. Later Single Board Computers (SBCs) - were Developed and are still manufactured. Such hardware boards support a broad range of processors - including 8 bits, 16- and 32-bit processor. Since 2015, company also started offering design and development services

# ❖ PROPOSED SYSTEM

1. Data Collection

2. Data Cleaning

3. Exploratory Data Analysis (EDA)

4. Modeling and Forecasting

5. Interpretation and Reporting

# ➢ DATA COLLECTION

- Collecting data from a CSV file.

- Government economic reports.

- Population , Literacy rate , GDP ($ per capita)

- Birth and death rate

# ➢ DATA CLEANING

- Data cleaning is also referred to as data preparation, is a vital step that comprises reformatting the data, making data corrections, and merging data sets to enhance the data

- Normalizing data in data science refers to the process of rescaling or transforming your data so that it falls within a specific range or follows a particular distribution

*SNAPSHOTS*

# ➤ EXPLORATORY DATA ANALYSIS (EDA)

- <u>Bar plot</u> is a graphical representation of data using rectangular bars or columns. Bar plots are used to display categorical data and compare the values of different categories or groups.

- <u>Heatmap</u> is a graphical representation of data where individual values are represented as colors. Heatmaps are particularly useful for visualizing and exploring the patterns and relationships within large and complex datasets.

- <u>Correlation analysis</u> is a statistical method used to evaluate the strength and direction of the relationship between two or more variables. It measures the extent to which changes in one variable are associated with changes in another variable.
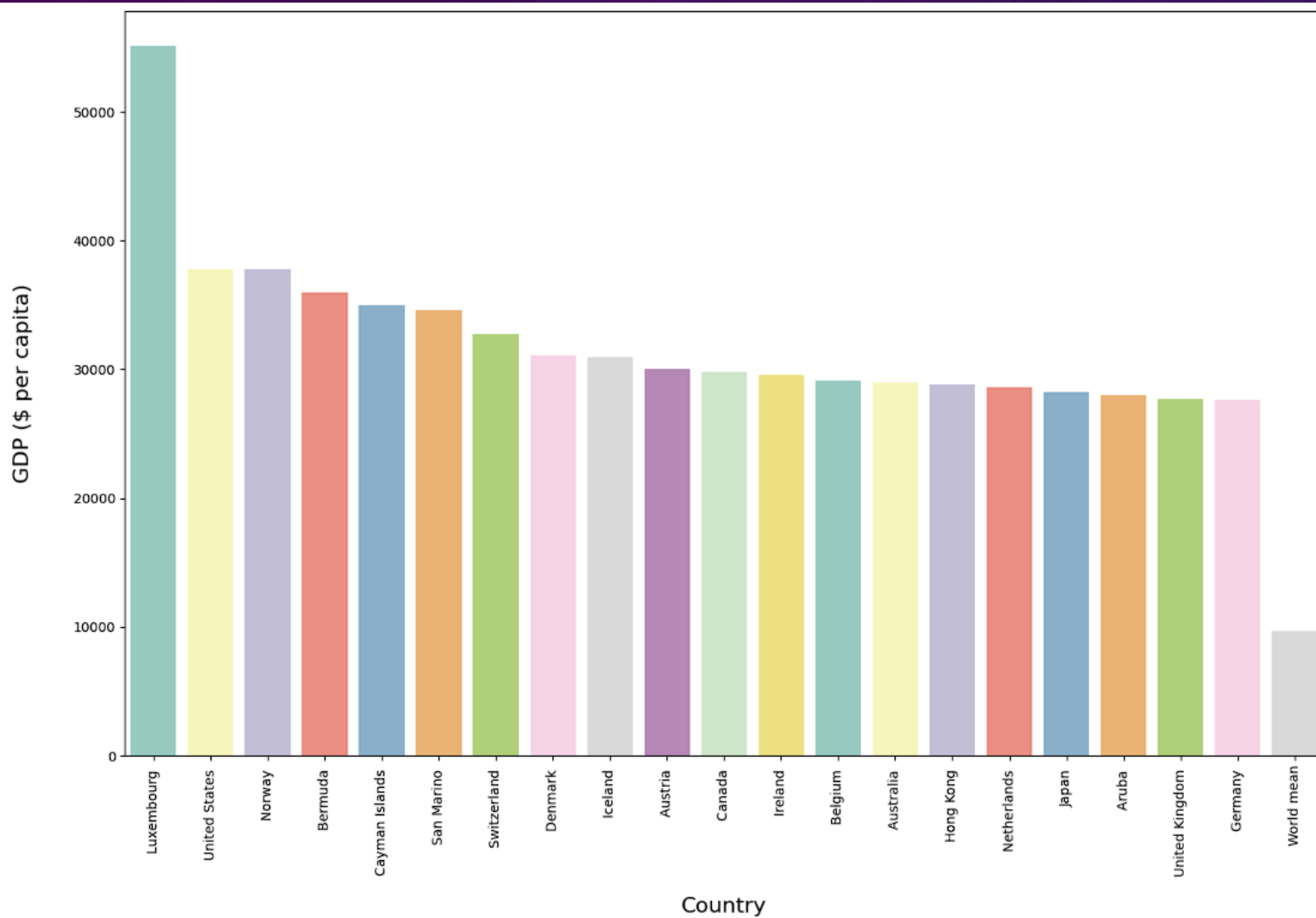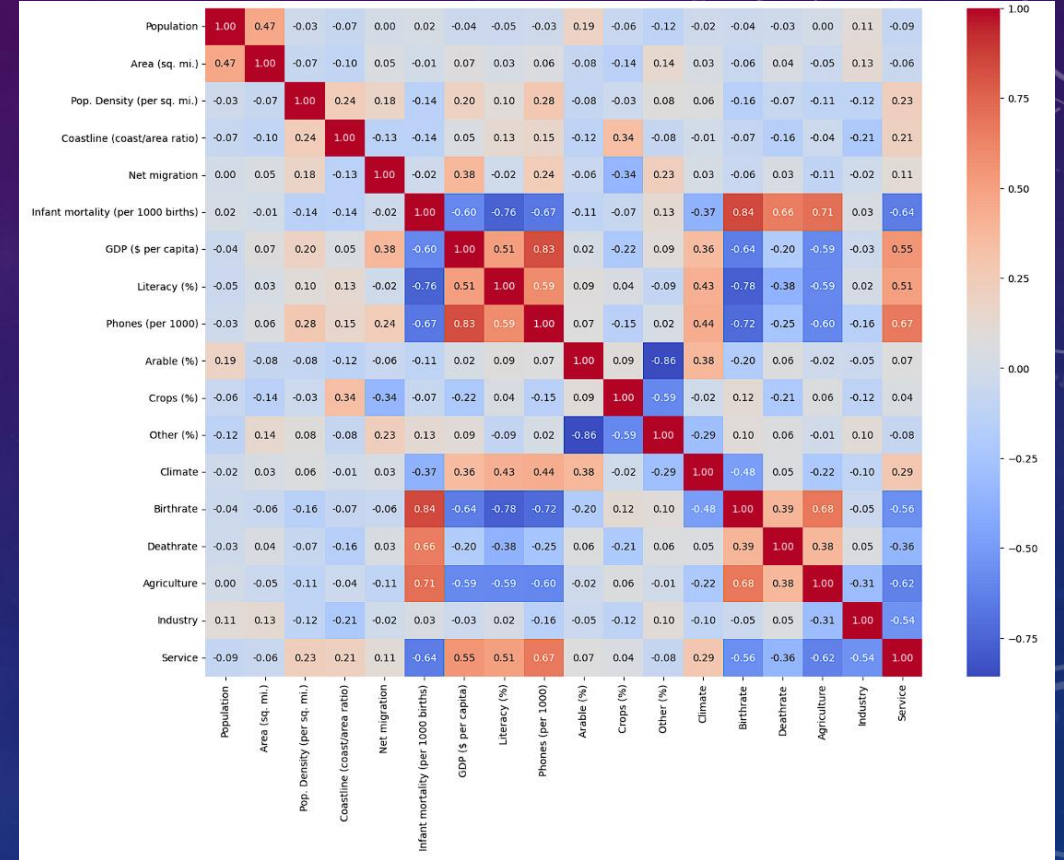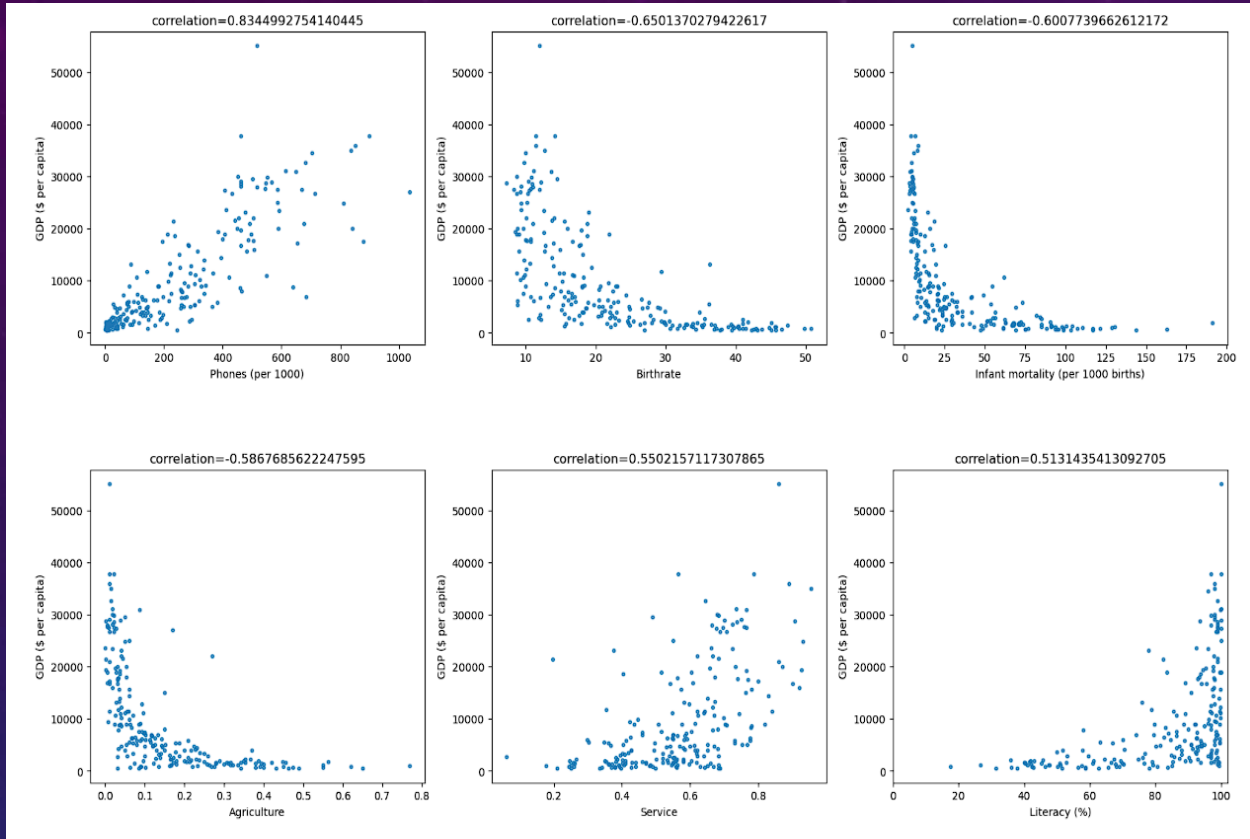
## ➢ PREDICTIVE MODELS

- <u>Regression models</u> are a class of statistical models used in data analysis and machine learning to explore and quantify the relationship between one or more independent variables (predictors or features) and a dependent variable (the target or outcome). The primary goal of regression analysis is to understand how changes in the independent variables are associated with changes in the dependent variable.

- Machine learning algorithms are computational methods or procedures that enable computers and machines to learn from data and make predictions or decisions without being explicitly programmed to do so. These algorithms are a fundamental component of machine learning, a subfield of artificial intelligence (AI).

## ➢ **INTERPRETATION AND REPORTING**

- <u>Analysing model results</u> refers to the process of evaluating and interpreting the outcomes produced by a machine learning or statistical model. This step is crucial in understanding the model's performance, assessing its effectiveness, and gaining insights into the underlying patterns in the data.

- <u>Making forecasts</u> refers to the process of predicting future events, trends, or outcomes based on historical data, patterns, and models. Forecasting is a valuable tool used in various fields, including finance, economics, weather science, supply chain management, and many others, to make informed decisions and plans for the future

- <u>Preparing reports</u> involves the process of compiling, organizing, and presenting information in a structured and coherent manner to convey findings, analysis, or results to a specific audience

# RESULTS

```python
In [15]:  1  model = LinearRegression()
          2  model.fit(train_X, train_Y)
          3  train_pred_Y = model.predict(train_X)
          4  test_pred_Y = model.predict(test_X)
          5  train_pred_Y = pd.Series(train_pred_Y.clip(0, train_pred_Y.max()), index=train_Y.index)
          6  test_pred_Y = pd.Series(test_pred_Y.clip(0, test_pred_Y.max()), index=test_Y.index)
          7
          8  rmse_train = np.sqrt(mean_squared_error(train_pred_Y, train_Y))
          9  msle_train = mean_squared_log_error(train_pred_Y, train_Y)
         10  rmse_test = np.sqrt(mean_squared_error(test_pred_Y, test_Y))
         11  msle_test = mean_squared_log_error(test_pred_Y, test_Y)
         12
         13  print('rmse_train:',rmse_train,'msle_train:',msle_train)
         14  print('rmse_test:',rmse_test,'msle_test:',msle_test)
```

```
rmse_train: 4545.898405281641 msle_train: 4.844890347754545
rmse_test: 5631.797188276091 msle_test: 4.378515097409243
```

```python
In [16]:  1  model = RandomForestRegressor(n_estimators = 50,
          2                                max_depth = 6,
          3                                min_weight_fraction_leaf = 0.05,
          4                                max_features = 0.8,
          5                                random_state = 42)
          6  model.fit(train_X, train_Y)
          7  train_pred_Y = model.predict(train_X)
          8  test_pred_Y = model.predict(test_X)
          9  train_pred_Y = pd.Series(train_pred_Y.clip(0, train_pred_Y.max()), index=train_Y.index)
         10  test_pred_Y = pd.Series(test_pred_Y.clip(0, test_pred_Y.max()), index=test_Y.index)
         11
         12  rmse_train = np.sqrt(mean_squared_error(train_pred_Y, train_Y))
         13  msle_train = mean_squared_log_error(train_pred_Y, train_Y)
         14  rmse_test = np.sqrt(mean_squared_error(test_pred_Y, test_Y))
         15  msle_test = mean_squared_log_error(test_pred_Y, test_Y)
         16
         17  print('rmse_train:',rmse_train,'msle_train:',msle_train)
         18  print('rmse_test:',rmse_test,'msle_test:',msle_test)
```
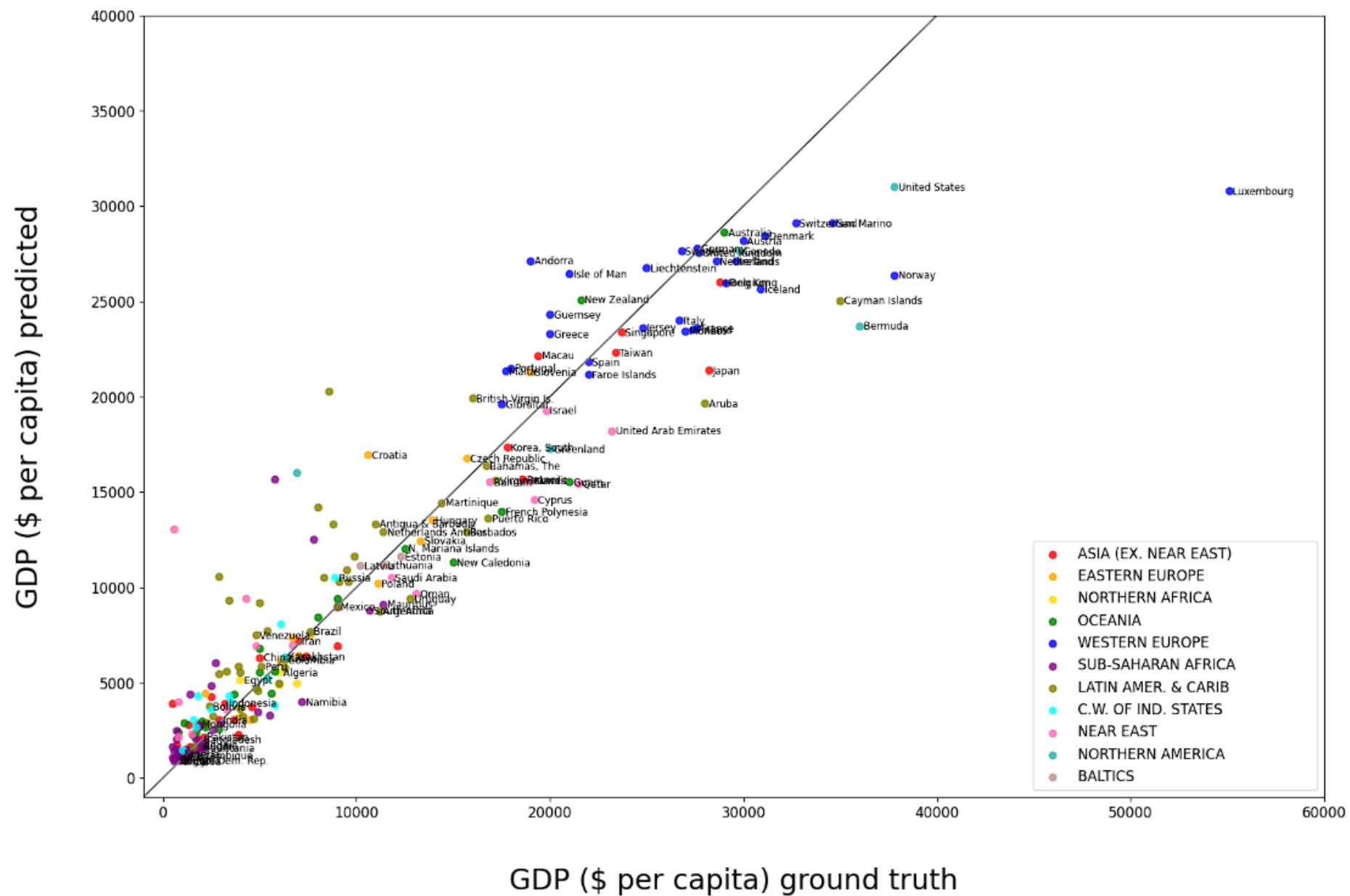
```
rmse_train: 3203.967879323867 msle_train: 0.1514357562756082
rmse_test: 3911.844310768375 msle_test: 0.3662583880957503
```
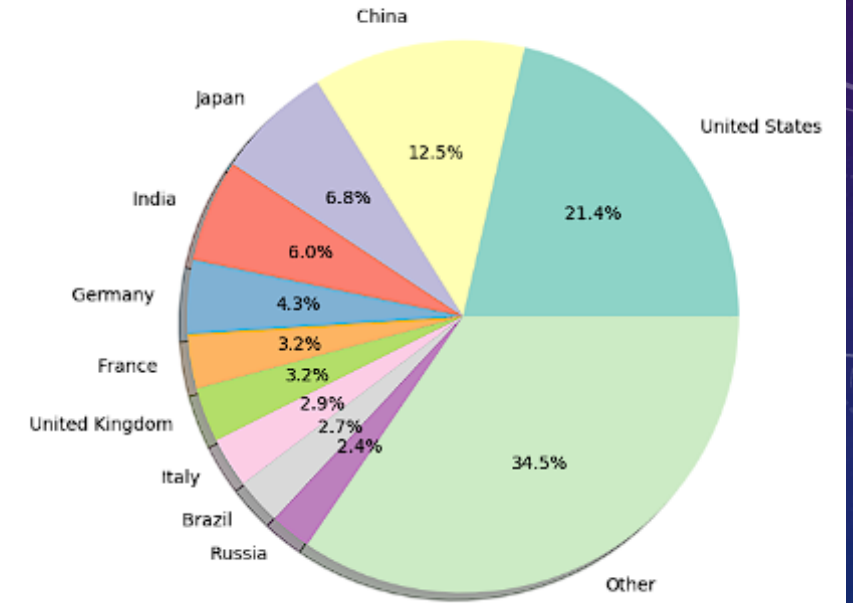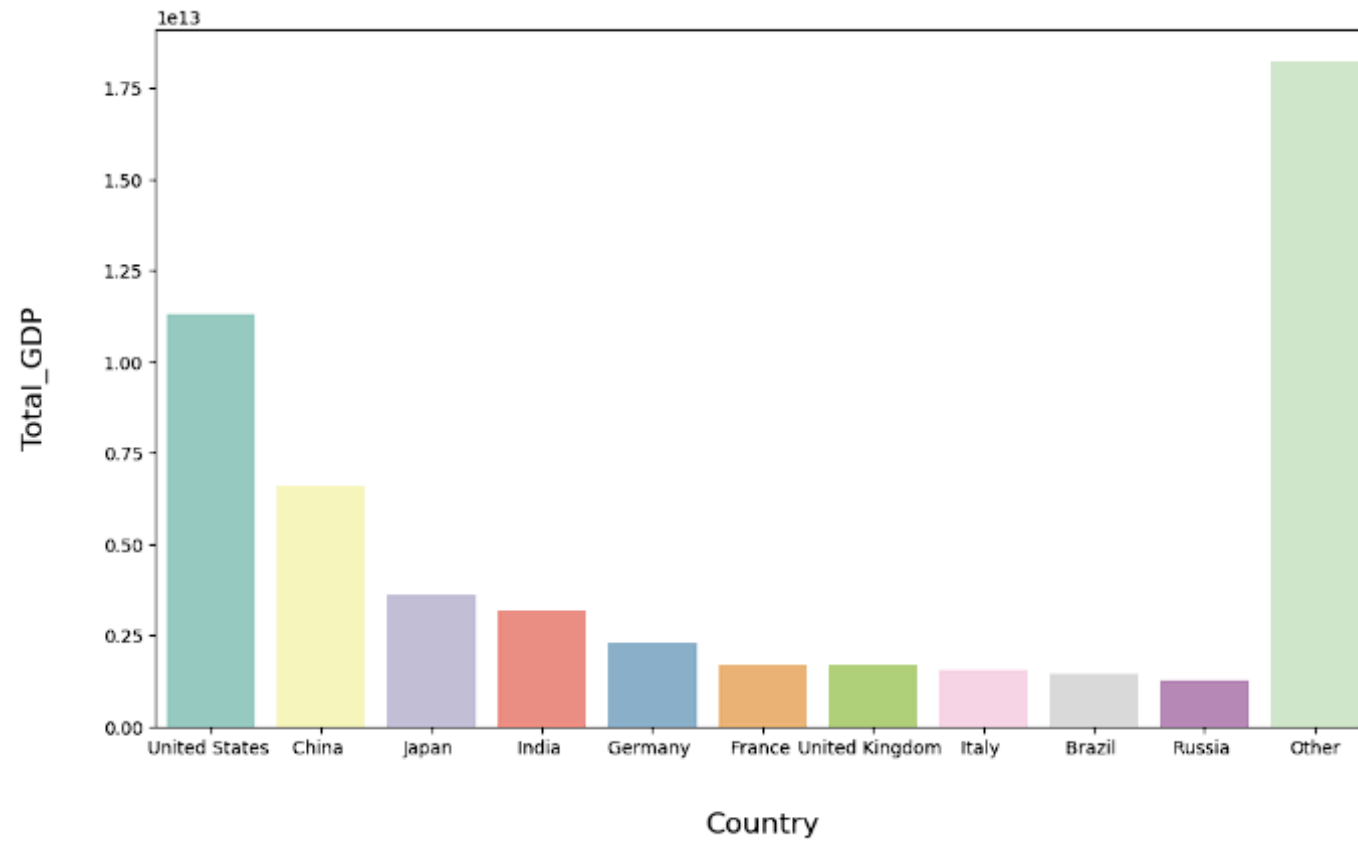
```python
In [17]:  1  plt.figure(figsize=(18,12))
          2
          3  train_test_Y = train_Y.append(test_Y)
          4  train_test_pred_Y = train_pred_Y.append(test_pred_Y)
          5
          6  data_shuffled = data.loc[train_test_Y.index]
          7  label = data_shuffled['Country']
          8
          9  colors = {'ASIA (EX. NEAR EAST)          ':'red',
         10           'EASTERN EUROPE                      ':'orange',
         11           'NORTHERN AFRICA                     ':'gold',
         12           'OCEANIA                             ':'green',
         13           'WESTERN EUROPE                      ':'blue',
         14           'SUB-SAHARAN AFRICA                  ':'purple',
         15           'LATIN AMER. & CARIB     ':'olive',
         16           'C.W. OF IND. STATES ':'cyan',
         17           'NEAR EAST                           ':'hotpink',
         18           'NORTHERN AMERICA                    ':'lightseagreen',
```

```
In [6]: import pandas as pd
        data = pd.read_csv('world.csv',decimal=',')
        print('number of missing data:')
        print(data.isnull().sum())
        data.describe(include='all')
```

```
number of missing data:
Country                              0
Region                               0
Population                           0
Area (sq. mi.)                       0
Pop. Density (per sq. mi.)           0
Coastline (coast/area ratio)         0
Net migration                        3
Infant mortality (per 1000 births)   3
GDP ($ per capita)                   1
Literacy (%)                        18
Phones (per 1000)                    4
Arable (%)                           2
Crops (%)                            2
Other (%)                            2
Climate                             22
Birthrate                            3
Deathrate                            4
Agriculture                         15
Industry                            16
Service                             15
dtype: int64
```

Out[6]:

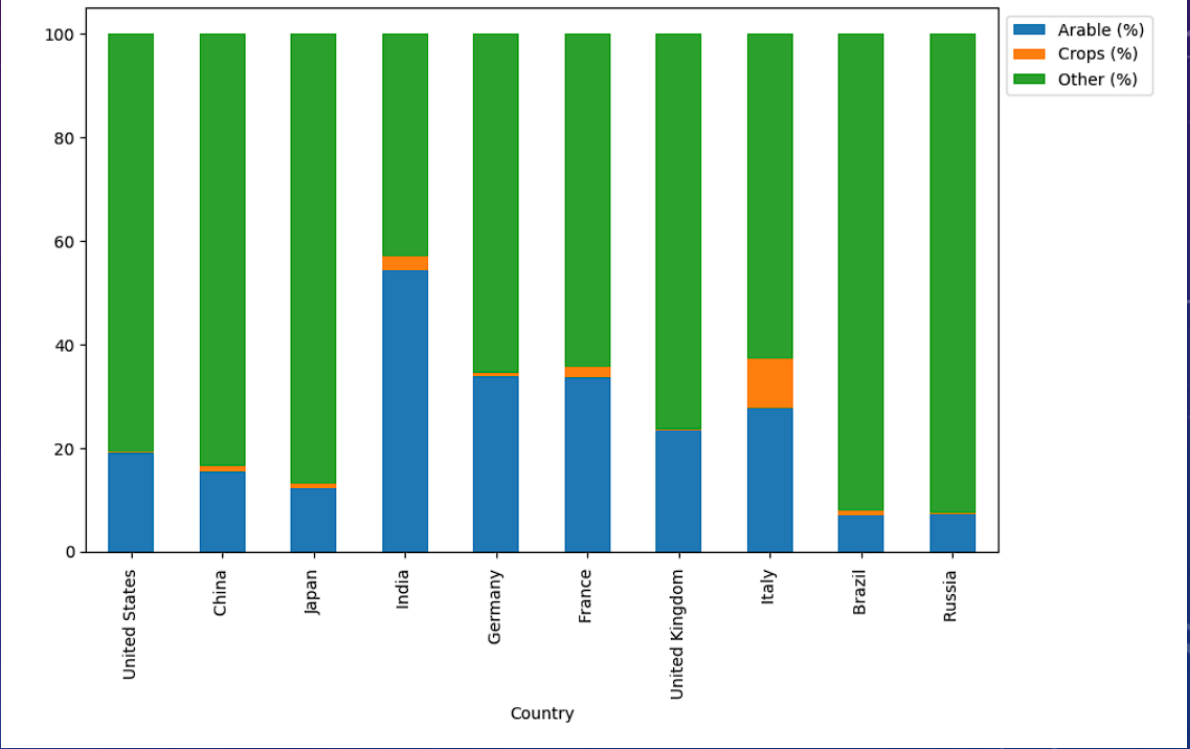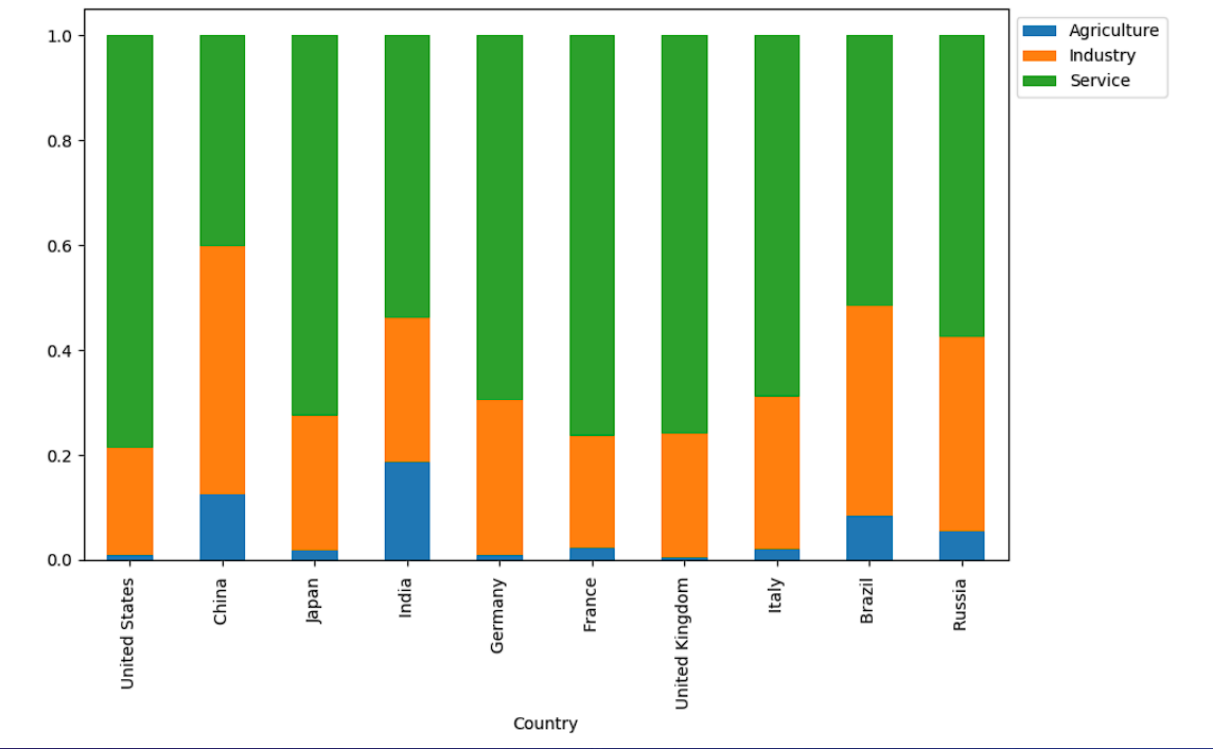| | Country | Region | Population | Area (sq. mi.) | Pop. Density (per sq. mi.) | Coastline (coast/area ratio) | Net migration | Infant mortality (per 1000 births) | GDP ($ per capita) | Literacy (%) | Phones (per 1000) | Arable ( |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 227 | 227 | 2.270000e+02 | 2.270000e+02 | 227.000000 | 227.000000 | 224.000000 | 224.000000 | 226.000000 | 209.000000 | 223.000000 | 225.0000 |
| unique | 227 | 11 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | N |
| top | Afghanistan | SUB-SAHARAN AFRICA | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | N |
| freq | 1 | 51 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | N |
| mean | NaN | NaN | 2.874028e+07 | 5.982270e+05 | 379.047137 | 21.165330 | 0.038125 | 35.506964 | 9689.823009 | 82.838278 | 236.061435 | 13.797 |
| std | NaN | NaN | 1.178913e+08 | 1.790282e+06 | 1660.185825 | 72.286863 | 4.889269 | 35.389899 | 10049.138513 | 19.722173 | 227.991829 | 13.0404 |
| min | NaN | NaN | 7.026000e+03 | 2.000000e+00 | 0.000000 | 0.000000 | -20.990000 | 2.290000 | 500.000000 | 17.600000 | 0.200000 | 0.0000 |
| 25% | NaN | NaN | 4.376240e+05 | 4.647500e+03 | 29.150000 | 0.100000 | -0.927500 | 8.150000 | 1900.000000 | 70.600000 | 37.800000 | 3.2200 |
| 50% | NaN | NaN | 4.786994e+06 | 8.660000e+04 | 78.800000 | 0.730000 | 0.000000 | 21.000000 | 5550.000000 | 92.500000 | 176.200000 | 10.4200 |
| 75% | NaN | NaN | 1.749777e+07 | 4.418110e+05 | 190.150000 | 10.345000 | 0.997500 | 55.705000 | 15700.000000 | 98.000000 | 389.650000 | 20.0000 |
| max | NaN | NaN | 1.313974e+09 | 1.707520e+07 | 16271.500000 | 870.660000 | 23.060000 | 191.190000 | 55100.000000 | 100.000000 | 1035.600000 | 62.1100 |

```
In [4]:    1  data.groupby('Region')[['GDP ($ per capita)','Literacy (%)','Agriculture']].median()
```

Out[4]:

| Region | GDP ($ per capita) | Literacy (%) | Agriculture |
|---|---|---|---|
| ASIA (EX. NEAR EAST) | 3450.0 | 90.60 | 0.1610 |

Out[4]:

| Region | GDP ($ per capita) | Literacy (%) | Agriculture |
|---|---|---|---|
| ASIA (EX. NEAR EAST) | 3450.0 | 90.60 | 0.1610 |
| BALTICS | 11400.0 | 99.80 | 0.0400 |
| C.W. OF IND. STATES | 3450.0 | 99.05 | 0.1980 |
| EASTERN EUROPE | 9100.0 | 98.60 | 0.0815 |
| LATIN AMER. & CARIB | 6300.0 | 94.05 | 0.0700 |
| NEAR EAST | 9250.0 | 83.00 | 0.0350 |
| NORTHERN AFRICA | 6000.0 | 70.00 | 0.1320 |
| NORTHERN AMERICA | 29800.0 | 97.50 | 0.0100 |
| OCEANIA | 5000.0 | 95.00 | 0.1505 |
| SUB-SAHARAN AFRICA | 1300.0 | 62.95 | 0.2760 |
| WESTERN EUROPE | 27200.0 | 99.00 | 0.0220 |

# ❖ CONCLUSION

In conclusion, our journey through GDP analysis using data science has provided us with valuable insights into the economic performance of our nation. We've seen how data science techniques can transform raw economic data into actionable information for policymakers, businesses, and individuals. The world is always changing, and we can keep using data science to understand it better. Our hope is that this helps everyone, from the government to businesses and individuals, make smarter decisions for a brighter economic future.

# ❖ APPLICATIONS

- **Economic Policy Formulation:** Governments use GDP analysis to design economic policies. By understanding the current economic situation and predicting future trends, policymakers can make informed decisions about taxation, spending, and interest rates to promote economic growth and stability.

- **Business Strategy:** Companies use GDP data to make strategic decisions. For instance, retailers may use GDP forecasts to plan inventory levels, and investment firms may adjust their portfolios based on economic outlooks

- **International Trade:** Exporters and importers use GDP data to identify markets with growth potential. A country with a rising GDP may represent an attractive export destination

- **Entrepreneurship Opportunities:** Entrepreneurs can spot business opportunities by analyzing GDP trends. For example, during economic expansion, there may be increased demand for new services or products.

# ❖ REFERENCES

➢ https://www.w3schools.com/datascience/

➢ https://www.investopedia.com/terms/p/per-capita-gdp.asp

➢ https://www.guru99.com/data-science-tutorial.html

THANK YOU