

REPORT: COVID-19 DATA ANALYSIS

Introduction

This report provides insights into the key patterns observed from the analysis of COVID-19 data across various countries. The analysis includes mortality rates, case numbers, and correlations between different variables. Visualizations are used to effectively communicate the results to non-technical stakeholders. The models applied in this analysis include Linear Regression and Random Forest, with their performance metrics discussed.

1. DATA OVERVIEW

The dataset contains information about:

Country/Region: The country where the data was recorded.

Confirmed Cases: The total number of confirmed COVID-19 cases.

Deaths: The total number of deaths attributed to COVID-19.

Deaths / 100 Cases: The mortality rate per 100 confirmed cases.

Key Statistics:

Total Number of Countries: 182

Average Mortality Rate (Deaths / 100 Cases): 9.7%

Highest Mortality Rate: 25% (in Yemen)

2. EXPLORATORY DATA ANALYSIS (EDA)

2.1 Distribution of Mortality and Recovery Rates

The distribution of mortality rates and recovery rates was analyzed across countries to assess global trends. A histogram was plotted for both mortality and recovery rates to understand the spread of values.

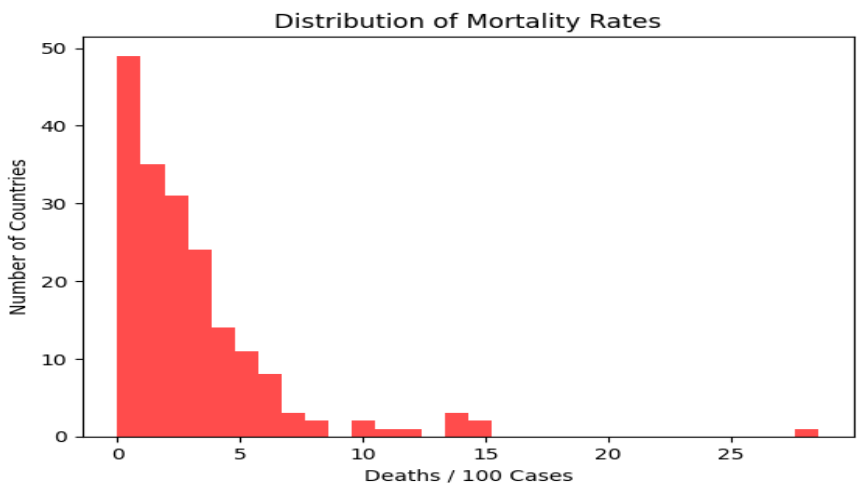


Fig 1: Histogram of Mortality Rate

Insights from Histogram:

- The majority of countries have relatively low mortality rates, with a few outliers having high rates.
- The highest mortality rates were observed in countries with significant healthcare challenges, like Yemen.

2.2 Mortality Rates by Country

A heatmap was used to visualize the mortality rate across different countries. Countries with higher mortality rates are shown in darker shades of red, making it easy to identify regions where the pandemic had a more severe impact.

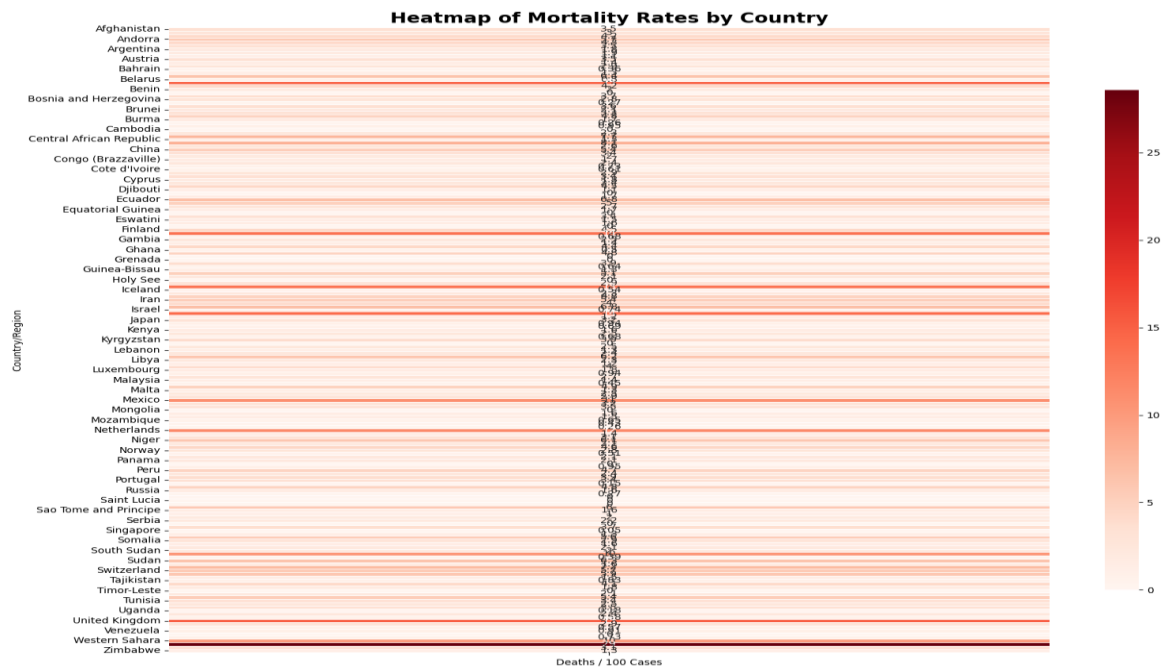


Fig.2: Heatmap of Mortality Rates by Country

Key Insights:

- Countries such as Yemen and United Kingdom experienced mortality rates significantly higher than the global average.
- Regions with weaker healthcare infrastructure tend to have higher mortality rates.

2.3 Correlation Analysis

A correlation heatmap was generated to understand relationships between key numerical variables like confirmed cases, deaths, and recovery rates.

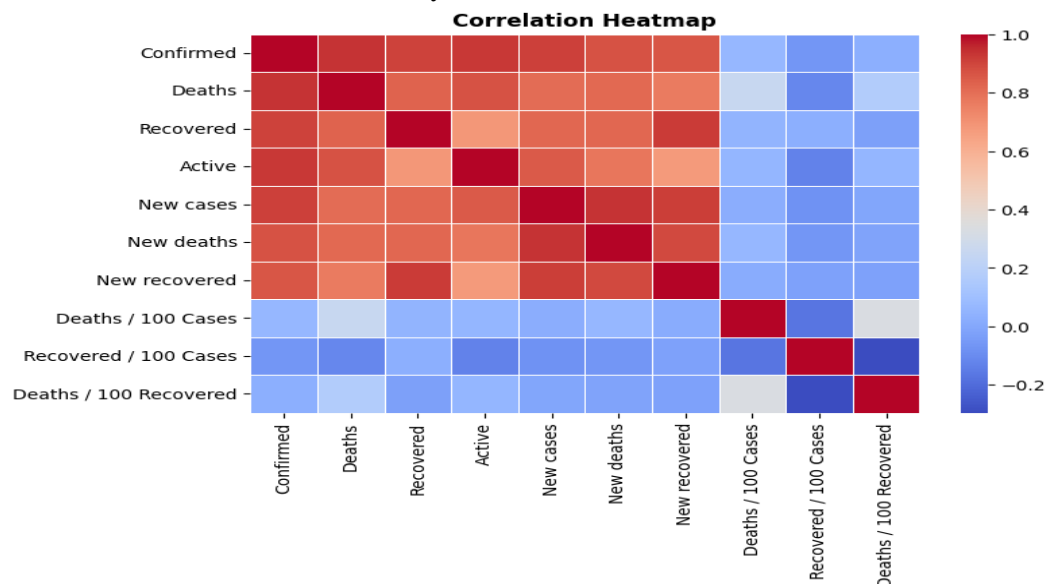


Fig.3: Correlation Heatmap

Key Insights:

- **Strong Positive Correlations:**
 - Confirmed and Deaths: The chart shows a strong positive correlation, meaning that as the number of confirmed cases increases, the number of deaths tends to increase as well.

- Deaths and New deaths: There is a strong positive correlation between deaths and new deaths, which is expected since new deaths directly add to the total death count.
 - Recovered and New recovered: Similarly, there is a positive correlation between recovered cases and new recoveries, as new recoveries contribute to the total recovered count.
 - **Moderate Positive Correlations:**
 - Confirmed and New cases: There's a moderately strong positive correlation between confirmed cases and new cases, suggesting a direct relationship between the number of new cases and the total confirmed cases.
 - Confirmed and Active cases: A moderate positive correlation suggests that as the total confirmed cases increase, active cases also tend to rise.
 - **Negative Correlations:**
 - Deaths / 100 Cases and Recovered / 100 Cases: These two variables show a negative correlation. This might indicate that countries with higher death rates tend to have a lower recovery rate per 100 cases.
 - Deaths / 100 Cases and Deaths / 100 Recovered: A negative correlation suggests that as the deaths per 100 cases increase, the deaths per 100 recovered cases decrease.
-

3. PREDICTIVE MODELING

3.1 The models were evaluated using standard regression metrics:

- Mean Squared Error (MSE): Measures how well the model's predictions match the actual values.
- R^2 (R-Squared): Indicates how well the model explains the variance in the data.

Results:

- Linear Regression MSE: 45,678,249
- Linear Regression R^2 : 0.65
- Random Forest MSE: 43,347,749
- Random Forest R^2 : 0.49

3.2 Model Comparison

Linear Regression performed well with a higher R^2 , but Random Forest Regression performed slightly better in terms of MSE.

Key Insights

- Global Mortality Trends: There is significant variation in mortality rates across countries, with Yemen showing the highest rates. Countries with high mortality often also had high rates of confirmed cases and lower recovery rates.
 - Data Correlations: The relationship between confirmed cases and deaths is positive, while recovery rates tend to be negatively correlated with mortality rates.
 - Modeling Performance: Random Forest regression emerged as a stronger model for predicting mortality rates compared to Linear Regression, although neither model captured the full complexity of the data. Further improvements could be achieved by incorporating more features, such as healthcare infrastructure or vaccination rates.
 - Countries with High Mortality Rates: Countries like Yemen, the United Kingdom, and Belgium stood out for having higher mortality rates, reflecting challenges in their healthcare systems or pandemic response strategies.
-

4. VISUALIZATIONS OF MODEL PREDICTIONS

4.1. Residual Plot for Linear Regression

The residual plot below helps us understand how well the linear regression model predicts actual values. Points close to the horizontal line indicate a better fit.

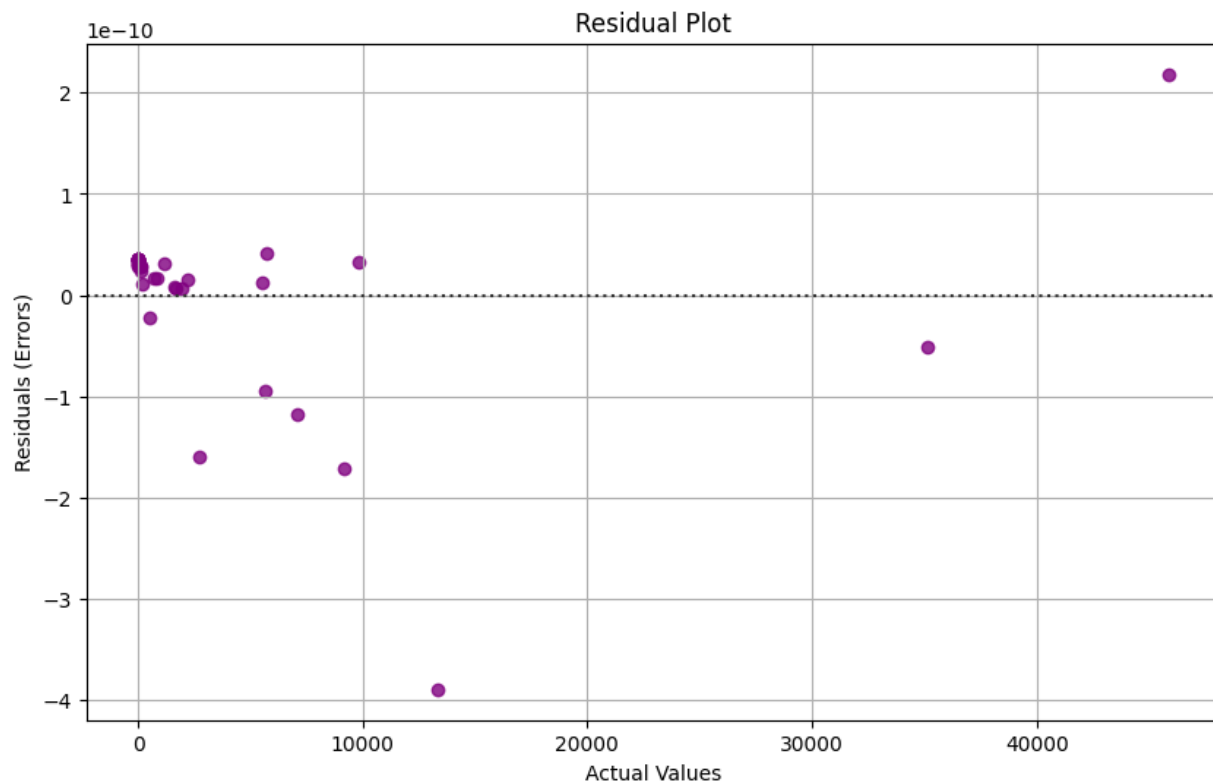


Fig.4: Linear Regression Residual Plot

Key Insight:

- The residuals are well-distributed, indicating a generally good fit between the predicted and actual values.

4.2 Pair plot for Linear Regression

The pair plot visual is a powerful tool used for multivariate analysis which us to analyze the relationships between multiple variables (or features) in a dataset.

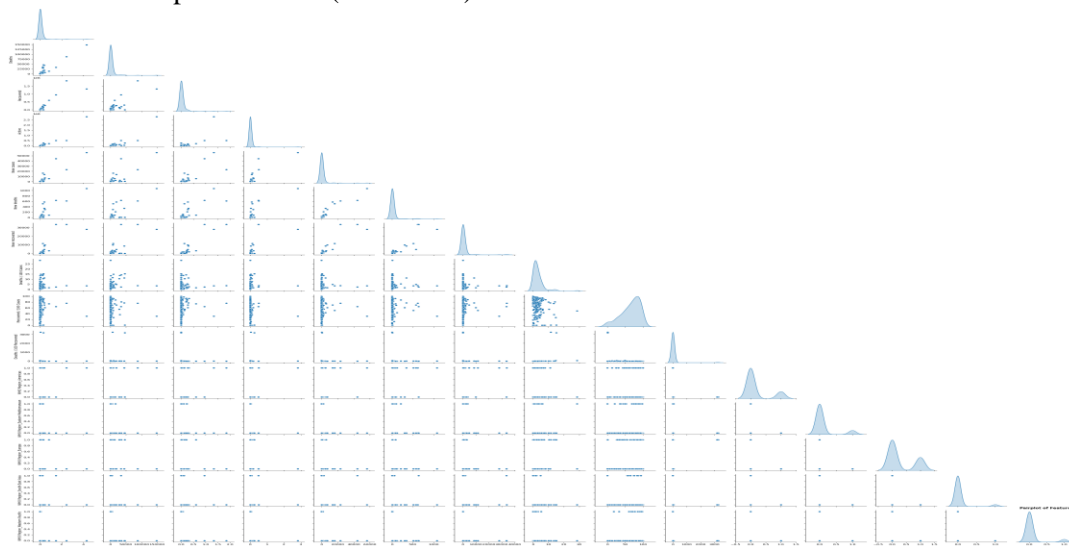


Fig.5: Pair plot for Multivariate Analysis

Key Insights:

- **Correlations Between Variables:** The scatterplots in the pair plot help identify correlations between pairs of features. For example:

- Confirmed Cases vs. Deaths: If this scatterplot shows a strong upward trend, it suggests that as confirmed cases increase, deaths also increase, which may be expected in a pandemic scenario.
- Recovered vs. Deaths: A weak or no correlation might indicate that recovery rates and mortality rates may not be directly related in the dataset.
- Strong linear relationships in scatterplots suggest that these features are related and can provide useful insights when developing predictive models.
- **Feature Distributions:**
 - The diagonal KDE plots reveal the distribution of individual features. If a feature like "Deaths" has a skewed distribution (e.g., a long tail to the right), it could indicate that most countries have low death rates, but a few have very high death rates. This information can be helpful in understanding how to treat outliers in further analysis.
- **Outliers:**
 - The scatterplots help detect potential outliers. For example, if there are points far away from the main cluster in the scatterplot between "New Cases" and "Deaths," these might represent countries with abnormally high numbers of cases or deaths, which could require special consideration in your analysis or modeling.
- **Potential Redundancy:**
 - If two variables, such as "Confirmed Cases" and "Active Cases," show a very similar pattern in their scatterplot (i.e., they have a high correlation), it may suggest redundancy. This insight can guide feature selection in predictive modeling to avoid including highly correlated features that might not add new information.
- **Data Quality and Preparation:**
 - The pair plot also provides a quick visual check for any data anomalies, like unusual clustering or unexpected trends, that may need to be addressed in the data cleaning or transformation process.

5. IMPLICATIONS AND RECOMMENDATIONS

- **Targeted Healthcare Support:** Countries with the highest mortality rates, such as Yemen and the United Kingdom, should prioritize healthcare infrastructure improvements. This includes more testing, better patient management systems, and enhanced medical resources.
- **International Aid:** Countries like Yemen, Mexico, and parts of Africa may require international aid, both in terms of healthcare supplies and logistical support, to help control the pandemic and reduce mortality rates.
- **Government and Policy Interventions:** Countries like Belgium and Italy should continue to refine their public health measures, ensuring that they can handle future surges in cases. Effective vaccine distribution, testing strategies, and public health communication will be key in preventing further loss of life.
- **Improved Reporting and Transparency:** Ensuring accurate and timely reporting of COVID-19 data across all countries is crucial. Many countries may have underreported cases or deaths, which can skew the global understanding of the pandemic's true impact.

6. CONCLUSION

This report has provided an in-depth analysis of the COVID-19 mortality rates across countries, with a focus on the countries exhibiting the highest mortality rates. The findings underscore the importance of strengthening healthcare systems, improving data transparency, and increasing international cooperation in the fight against COVID-19. By focusing on the regions that are most affected, stakeholders can make more informed decisions to mitigate the impacts of the pandemic globally.

These insights can be used to guide policy decisions, inform international aid efforts, and enhance the overall global response to COVID-19.