# Conclusion Report: COVID-19 Time Series Analysis

## 1. Introduction

The COVID-19 pandemic, a global health crisis that emerged in late 2019, has significantly impacted countries worldwide. In the United States, the virus has posed formidable challenges, leading to widespread infections and fatalities. As of the latest available data, the USA has witnessed a substantial number of confirmed cases and deaths due to COVID-19. This report focuses on developing an accurate time series forecasting model to predict the spread of COVID-19, specifically in the USA.

### 1.1 Historical and Statistical Overview
The COVID-19 outbreak in the USA began in early 2020, with a surge in confirmed cases and deaths. Over time, various interventions and public health measures were implemented to mitigate the impact. As of March 9, 2023, the dataset spans from January 22, 2020, capturing the evolving dynamics of the pandemic. Descriptive statistics reveal a mean of approximately 47 million confirmed cases and 624,562 deaths, emphasizing the severity and magnitude of the crisis.

## 2. Problem Statement

The goal is to develop a robust time series forecasting model to predict the future trajectory of COVID-19 confirmed cases and deaths within the USA. Accurate predictions are crucial for informed decision n-making, resource allocation, and implementing effective public health measures.

## 3. Approach

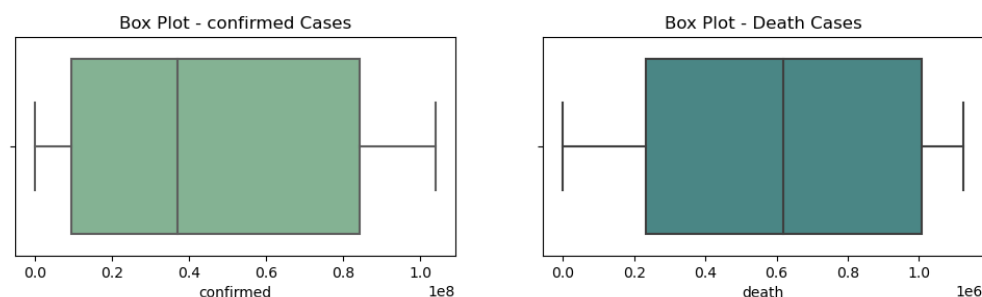### 3.1 Data Collection and Cleaning

3.1.1 Data Sources:
The primary data source is the COVID-19 dataset specific to the USA, obtained from Johns Hopkins University of Medicine, Corona Virus resource center. The dataset includes information on confirmed cases and deaths across states (source). https://coronavirus.jhu.edu/about/how-to-use-our-data

3.1.2 Cleaning Process:
Two distinct datasets for confirmed cases and deaths were obtained and concatenated into a unified dataframe named 'df.'
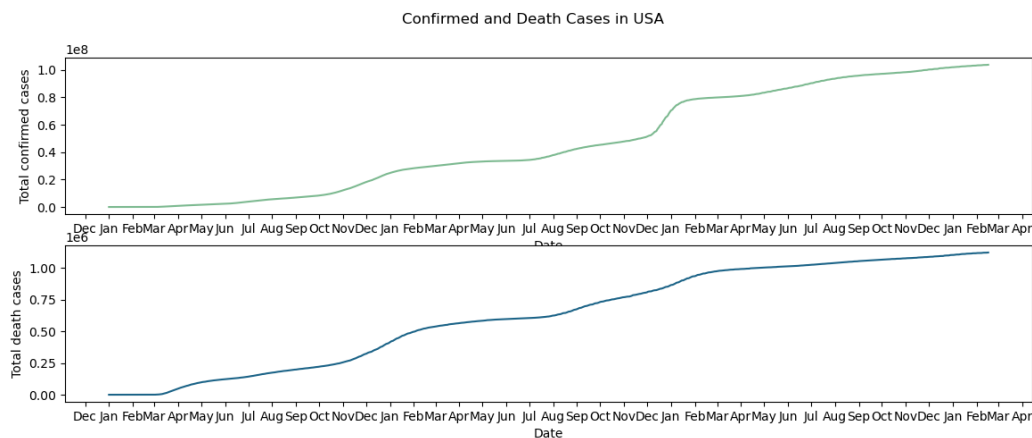Redundant columns were removed, and the data was transposed for clarity.

The dataset spans from January 22, 2020, to March 9, 2023, with no missing values. Descriptive statistics confirm the absence of outliers, ensuring dataset cleanliness for subsequent analyses or modeling.

**3.2 Exploratory Data Analysis (EDA)**

3.2.1 Confirmed and Death Cases Visualization

Per State and in the USA: Visualizations were created to depict the distribution of confirmed and death cases across different states and the overall USA. This aids in identifying geographical patterns and variations in the impact of the pandemic.
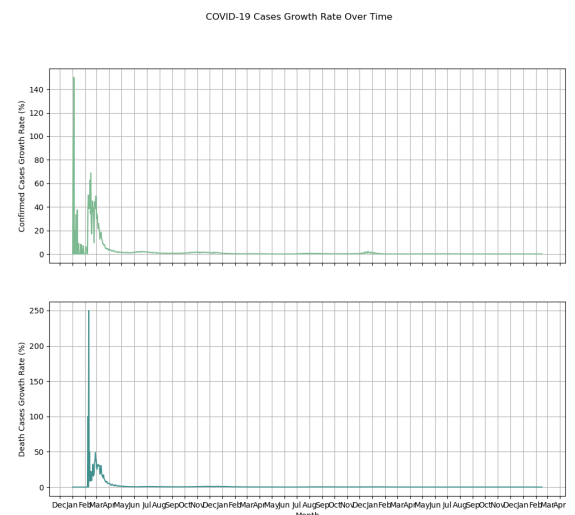


3.2.2 Mortality Rate Calculation

Mortality Rate in USA: The mortality rate, calculated as the ratio of deaths to confirmed cases, was determined to be approximately 1.33. This metric provides a quantitative measure of the severity of the disease's impact and helps in assessing the effectiveness of healthcare systems.

3.2.3 Growth Rate Analysis

Confirmed and Death Cases Growth Rate: The growth rate in time series data signifies the relative change in a value over time, often expressed as a percentage increase or decrease between two consecutive data points. By analyzing the growth rates of confirmed COVID-19 cases and related deaths, temporal patterns were identified.

- Initial Surge and Decline: Significant surges in positive cases and deaths were observed from January to February 2020, marking the onset of the pandemic. A subsequent decline in deaths in mid-February could be attributed to interventions like public health measures and increased testing.
- Effectiveness of Quarantine Measures: A decrease in the peaks of positive cases in mid-January suggested the effectiveness of quarantine measures. However, a slight increase in deaths in mid-February indicated a potential lag effect.
- Consistent Positive Growth: From April onward, the metrics showed a consistent positive growth, reflecting a sustained increase in confirmed cases. This period may be influenced by factors such as emerging variants, changes in testing strategies, or adaptations in public health responses.

# 4. Feature Engineering and Modeling

## 4.1 Data Stationarity
Checked stationarity with KPSS method and Dickey-Fuller test, indicating non-stationarity.
Applied Log transformation to achieve stationarity for both confirmed and death cases.

## 4.2 Hyperparameter Tuning for ARIMA Model
Explored the best parameters for the ARIMA model (pdq: 1,1,1) for confirmed and death cases.
Training data cover 2020-01-22 to 2022-11-01, while the test data covers the subsequent period.

## 4.3 ARIMA Model Summary Statistics and Conclusion
Confirmed Cases Model (ARIMA (1, 1, 1)):

- ARIMA (1, 1, 1) structure with AR and MA coefficients of 0.9860 and -0.7311, respectively.
- Goodness of Fit: Lower AIC, BIC, and HQIC values suggest a better fit.
- Residuals Analysis: No significant autocorrelation but departure from normality. Heteroskedasticity present.
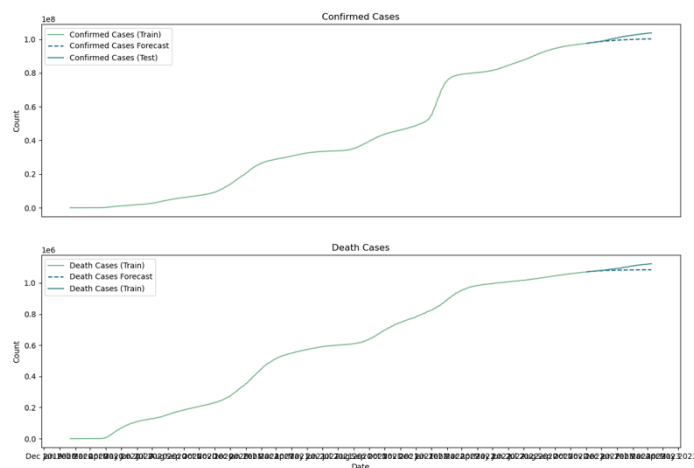
Death Cases Model (ARIMA (1, 1, 1)):

- ARIMA (1, 1, 1) structure with AR and MA coefficients of 0.9745 and -0.7578, respectively.
- Goodness of Fit: Similar to confirmed cases, lower AIC, BIC, and HQIC values suggest a good fit.
- Residuals Analysis: No significant autocorrelation but departure from normality. Heteroskedasticity present.

## 4.4 Conclusion
Both models fit the data well based on AIC, BIC, and HQIC values. Residuals show no significant autocorrelation but exhibit departure from normality. Heteroskedasticity is present, indicating variability in residuals.

## 4.5 Forecasting
The forecasting period spans from 2022-11-01 to 2023-03-09 covering a total of 5 months.

## 5. Personal Conclusion

In hindsight, analyzing the COVID-19 time series data up to March 9, 2023, offers a retrospective view of a challenging period that has now transitioned into history. As of 2024, the pandemic is a chapter of resilience and adaptation.

*Key Reflections:*
**Early Surge and Response:** The initial surge in cases in early 2020 was met with swift responses, leading to a decline in deaths. This highlighted the adaptability of healthcare systems.
**Ongoing Dynamics:** From April 2020 onward, a sustained positive growth in cases reflected the evolving nature of the virus. Factors like emerging variants and changes in testing influenced the pandemic's trajectory.
**Modeling and Forecasting:** ARIMA models provided insights into future trajectories, emphasizing the dynamic impact of the pandemic. The variability in residuals underscores the complexity of the situation.

*Looking Forward:*
As we move beyond the pandemic, lessons learned pave the way for continuous monitoring, research, and adaptation. Future efforts should explore long-term consequences, assess vaccination effectiveness, and contribute to a more resilient global health infrastructure.

In conclusion, the insights gained contribute to a broader understanding of global health challenges, shaping strategies for a resilient post-pandemic future.

Regarding the forecasting, at the beginning, the models worked excellently in predicting the trajectory of the pandemic. However, over time, deviations started to emerge, reflecting the evolving nature of the virus, including the impact of emerging variants and changes in testing. It became evident that continuous adaptation of the forecasting models was necessary to maintain accuracy. Recurrent training every month proved to be an effective strategy in resolving these deviations and ensuring the models remained insightful in providing future trajectories. This highlights the dynamic nature of the pandemic and the importance of ongoing adjustments in modeling approaches to capture the complexity of the situation.

## 6. Ideas for Further Research

- Longitudinal Data Analysis: Explore the integration of longitudinal data to understand the dynamic nature of COVID-19 progression over time.
- Incorporating External Factors: Investigate the inclusion of external factors, such as vaccination rates and public health interventions, to enhance predictive capabilities.
- Ensemble Modeling: Consider the utilization of ensemble models for a more comprehensive understanding of the pandemic's complexity.
- 

This report serves as a foundational step in developing predictive models for COVID-19, contributing to data-driven decision-making and public health efforts. Continuous research and collaboration are imperative to address the evolving challenges posed by the pandemic.