

# STA 137 Final Project

Zhuorui He, Shiyu Wu, Sulei Wang, Zhan Shi

2024-03-12

## Contents

<b>Introduction</b>	<b>2</b>
<b>Background</b>	<b>2</b>
<b>Exploratory Data Analysis</b>	<b>2</b>
Data Description . . . . .	2
Missing Value Plots . . . . .	3
Visualization . . . . .	3
<b>Inferential Analysis</b>	<b>3</b>
ADF and KPSS Tests for Stationary . . . . .	3
Decomposition . . . . .	4
First-order Difference . . . . .	4
Box-Cox Transformation . . . . .	4
Log Transformation . . . . .	4
ARIMA Model . . . . .	4
Forecast . . . . .	5
<b>Conclusion</b>	<b>5</b>
<b>Reference</b>	<b>7</b>
<b>Code Appendix</b>	<b>7</b>

```
knitr::opts_chunk$set(  
  include = FALSE  
)
```

## Introduction

This study examines the annual exports data of the Central African Republic, represented as a percentage of its total GDP from 1960 to 2017. Tackling real-world datasets often involves navigating their inherent complexity, as they seldom present themselves in a straightforward manner. The methodology employed here includes applying decomposition and transformations to facilitate more effective analysis, identifying the optimal ARIMA parameters, and performing residual diagnostics to ensure model reliability. The ultimate goal is to develop an ARIMA model capable of forecasting future exports with a reasonable degree of precision.

## Background

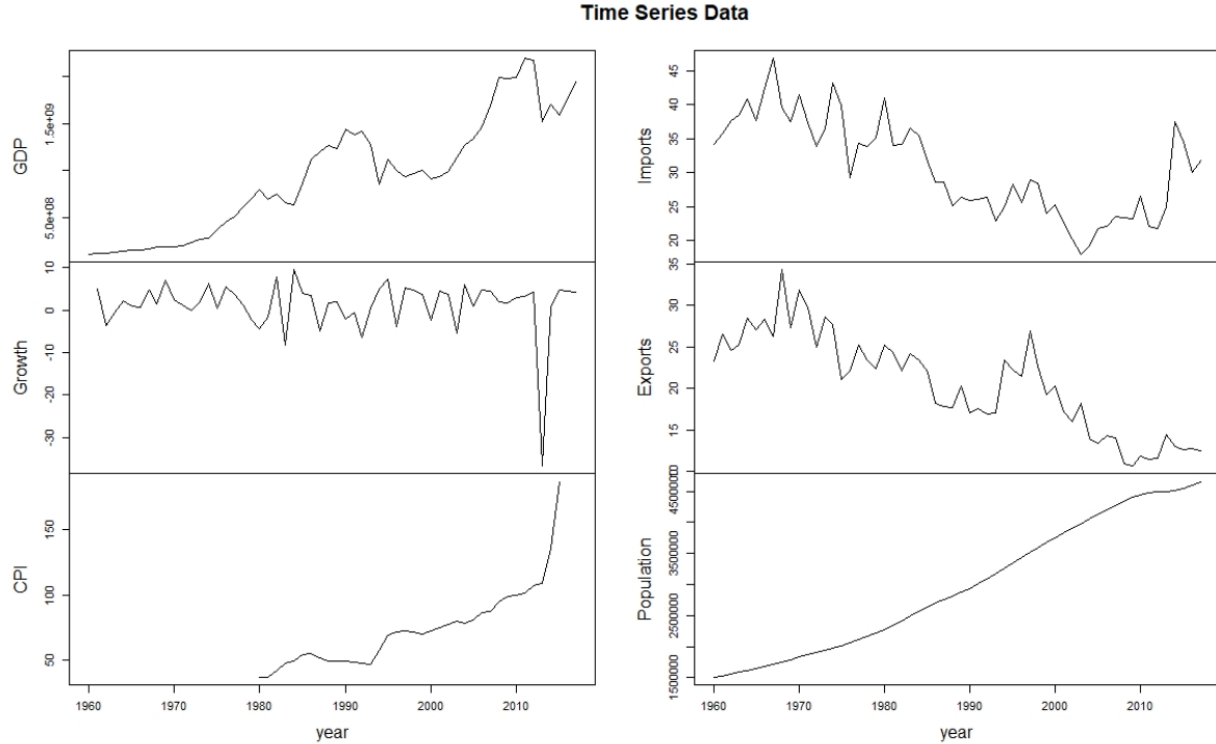
The Central African Republic (CAR), a landlocked nation at Africa's heart, ranked 183rd globally in exports in 2022 per the Observatory of Economic Complexity (OEC). Its economy, rich in natural resources, focuses on agriculture, services, and exporting minerals, oil, timber, and agricultural products, with gold and diamonds among its top exports. Despite agriculture being pivotal, its contribution to exports is lesser. From 1960 to 2017, export volumes rose, but their GDP percentage has been declining since 1968. Major importers include the UAE, Italy, Pakistan, China, and France.

However, CAR faces challenges like political instability, violence, and inadequate resource management, hindering sustainable growth and global market competitiveness. Reliance on international aid is critical due to slow growth and a rising population, with poverty affecting 65.7% of the populace by 2021.

## Exploratory Data Analysis

### Data Description

This dataset, sourced from the World Bank, encompasses detailed yearly data on the Gross Domestic Product (GDP), imports, and exports as a percentage of total GDP, population figures, and the GDP growth of the Central African Republic from 1960 to 2017. The focal point of our report is the Exports data. In this dataset, exports are measured as a percentage of GDP, indicating the total value of the country's exports of goods and services in relation to the size of its Gross Domestic Product. The figure below provides a compelling overview of the country's economic trends over the past decades.



## Missing Value Plots

Upon examining the dataset for missing values, it was found that the Exports data is free from any gaps, affirming the dataset's completeness and its suitability for conducting time series model analyses focused on Exports.

## Visualization

The preliminary analysis commenced with an examination of the time series data for Exports. A subtle downward trend observed in exports, as illustrated in the figure below, suggests the presence of a non-constant mean, indicating that the time series is non-stationary. Further investigation was conducted on the dataset's Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF). These analyses are crucial for understanding the direct influence of past data points on future values and for determining the appropriate order of the time series model.

## Inferential Analysis

### ADF and KPSS Tests for Stationary

To numerically evaluate the stationarity of the dataset, we applied two distinct methods: the Augmented Dickey-Fuller (ADF) test and the Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test, both set at a significance level of 0.05. The ADF test deems a time series as stationary if the p-values are less than 0.05, whereas the KPSS test considers a time series stationary if the p-values exceed 0.05. Results from both tests led to the rejection of the hypothesis that the time series is stationary, confirming that the dataset is indeed

non-stationary. All our models are predicated on the foundation of stationary data, which imparts numerous characteristics and properties crucial for prediction. Consequently, transforming the data into a time-independent format is essential to harness these advantages effectively.

## Decomposition

To address the non-stationarity of the dataset, four different methods were employed to detrend the data, aiming to isolate a suitable residual component. The initial approach involved Kernel smoothing, a technique utilized to eliminate the overarching decreasing trend observed over time. This process was intended to distill the data down to its cyclical trend component, as depicted in figures below. The goal was to refine the data in such a way that the remaining variability could be attributed primarily to cyclical fluctuations, thereby facilitating a more focused analysis of the time series.

To delve deeper into the model exploration, calculations of the Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) were conducted. ACF and PACF are instrumental in identifying the lags with significant impact on the current value. They play a crucial role in delineating the pattern of the time series, thereby aiding in the selection of an appropriate model. The ACF revealed a very strong correlation at a lag of 1, while the PACF displayed two pronounced correlations at lags 1 and 2. These findings suggest the potential suitability of AR(2) and MA(1) models for the de-trended data.

## First-order Difference

Another transformation attempt was made using the first-order difference, which similarly indicated the appropriateness of an AR(2), MA(1) model for the dataset. Despite this, the first-order difference transformation failed to pass the ADF test, indicating that it did not achieve stationarity. Given this outcome, the decision was made to abandon the first-order difference approach.

## Box-Cox Transformation

The application of both Box-Cox Transformation and Log Transformation was explored to address the non-stationary nature of the time series. However, subsequent analysis revealed that both transformations resulted in time series that remained non-stationary, as evidenced in figure below. Following the application of differencing to these transformed datasets, the ACF and PACF analyses of the differences suggested characteristics akin to a random series, resembling white noise. Consequently, due to these outcomes, both the Box-Cox and Log Transformations were deemed unsuitable for further consideration and were thus discarded from the analysis.

## Log Transformation

After thorough comparison and careful consideration, the decision was made to select the detrending method as the preferred approach for preparing our data for model fitting in the subsequent phase of analysis.

## ARIMA Model

For the ARIMA(q,d,p) model

$$X'_t = \sum_{i=1}^p \alpha_p X'_{t-p} + \sum_{i=1}^q \theta_q \omega_{t-q}$$

where  $X'_t = (1-B)^d X_t$  The model exhibits an order in the form of (p, d, q) where: P = The order of the Autoregressive Model d = the order of differencing Q = The order of the Moving Average Following the detrending

of the Exports dataset, the next step involved establishing the ARIMA model. To accurately determine the optimal parameters for the ARIMA(p,d,q) model, an automatic selection process was employed. This method suggested ARIMA(2,0,1) as the most fitting model for our data.

## Model Comparison

To discern the most effective model, Mean Squared Error (MSE) and Akaike Information Criterion (AIC) were employed as the primary metrics. MSE helps estimate the average of squared errors during model training, while AIC assesses both the prediction error and the extent of information loss during model fitting.

$$MSE = \frac{\sum (Residuals)^2}{n}$$

$$AIC = n \times \ln\left(\frac{\sum (Residuals)^2}{n}\right) + 2k$$

where n is the number of observations in the model, and k is the number of parameters.

The SARIMA model was utilized for evaluation, showing that the ACF of residuals mostly fell within the expected confidence range, with the Q-Q plot revealing a good alignment of most points along the fit line, albeit with some evidence of heavy tails. Comparisons were made between ARIMA(2,1,1) and ARIMA(3,0,1), against the backdrop of ARIMA(2,0,1). While ARIMA(2,1,1) exhibited a lower MSE of 0.0128 compared to ARIMA(3,0,1)'s MSE of 0.0735 and ARIMA(2,0,1)'s MSE of 0.00430, it was disqualified due to its p-values falling below the significance line for lags less than 5, indicating autocorrelation and dependency within white noise, rendering it suboptimal.

Despite ARIMA(3,0,1) presenting a slightly better AIC of 2.544 against ARIMA(2,0,1)'s AIC of 2.613, the minimal difference and the preference for simplicity guided the choice towards ARIMA(2,0,1) for prediction purposes, with the specific equation provided in the following segment.

## Forecast

After comparing models and conducting diagnostics, ARIMA(2,0,1) with detrended data was identified as the optimal choice. The model formula is

$$(1 - \phi_1 B - \phi_2 B^2)x_t = (1 + B)w_t$$

Utilizing this model, we forecasted the white noise in exports for the next six years with a significance level  $\alpha = 0.05$ .

The forecast depicted focuses solely on the random white noise component of exports over the next six years, not the actual future exports. To accurately forecast total exports, a linear fit and prediction of the trend using a separate linear model would be necessary. However, this paper concentrates primarily on analyzing the random noise within the time series, not on predicting the trend.

## Conclusion

In this project, detrending the data was chosen as the necessary step to adhere to the stationary assumption and reduce errors, leading to the establishment of an ARIMA(2,0,1) model. However, the forecast process encounters a limitation due to its inability to predict future trends, offering only a 95% confidence interval for predictions that exclude any trend component. This constitutes a significant constraint of the study.

Looking ahead, to address this limitation and enhance the project, an initial step could involve identifying suitable functions or distributions to estimate the trend, which would allow for more comprehensive confidence intervals that incorporate the trend. Additionally, the introduction of more data and the exploration of alternative modeling approaches, such as the Long Short-Term Memory (LSTM) model, could offer improved outcomes that include trend analysis, thus broadening the scope and applicability of the findings.

## Reference

(n.d.). Exports of goods and services (% of GDP) - Central African Republic. THE WORLD BANK.  
<https://data.worldbank.org/indicator/NE.EXP.GNFS.ZS?locations=CF>

## Code Appendix