# Configuration Manual

MSc Research Project
Data Analytics

# RAMANAN LOGANATHAN

Student ID: 21202702

School of Computing
National College of Ireland

Supervisor:     FURQAN RUSTAM

| Student Name: | RAMANAN LOGANATHAN |
|---|---|
| Student ID: | 21202702 |
| Programme: | Data Analytics |
| Year: | 2023 |
| Module: | MSc Research Project |
| Supervisor: | FURQAN RUSTAM |
| Submission Due Date: | 14/08/2023 |
| Project Title: | Configuration Manual |
| Word Count: | 1077 |
| Page Count: | 5 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| Signature: | |
|---|---|
| Date: | 18th September 2023 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

RAMANAN LOGANATHAN

21202702

# 1 Introduction

The configuration manual outlines the sequential guide for executing the time series models associated with the research project. The project involves the implementation and evaluation of four different time series models, namely RNN (Recurrent Neural Network), ARIMA (AutoRegressive Integrated Moving Average), SARIMA (Seasonal AutoRegressive Integrated Moving Average), and LSTM (Long Short-Term Memory).

# 2 System Configuration

## 2.1 Software Requirements

The research work was executed utilising Google Colab, an open-source framework for AI/ML projects in the Google ecosystem. This setting is powered by a Python module. Installing each of these packages is necessary before the project can be built.

## 2.2 Hardware specifications

- System Name: suryas_zenbook

- Processor: 11th Gen Intel(R) Core(TM) i7-1165G7 @ 2.80GHz 2.80 GHz

- Installed RAM: 16.0 GB (15.7 GB usable)

- Storage Size: 1TB SSD (109,951,162,7776 bytes)

- OS type: 64-bit operating system, x64-based processor

# 3 Installation and Environment Setup

- **Python**

This project made use of a Python package. Since the majority of Deep Learning and Machine Learning Projects are supported by its numerous built-in libraries. With a variety of plots, it makes developing and analysing models easier. Since python's latest version is already installed in google colab, no need to install it again.

- **Google Colab**

Figure 1: Latest version in Google colab

Google offers a Jupyter notebook environment on the cloud called Colaboratory, also known as Google Colab. As It provides a practical platform for research, machine learning, data analysis, and coding tasks and Colab serves as a popular option for students, academics, and professionals because it avoids the need for complicated settings by allowing Python code to execute easily in a browser.

One of Colab's key benefits is its interface with Google Drive, which enables simple storing, sharing, and notebook collaboration. Additionally, it offers open access to GPUs and TPUs, facilitating the quicker completion of resource-intensive operations. Because Colab notebooks are interactive, data can be visualised and analysed in real-time, providing immediate insights. This can accessed from [1]as shown in the Figure 2 .
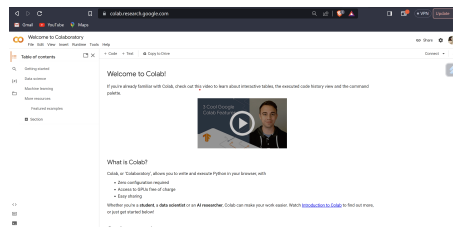


Figure 2: Google Colab's user interface

# 4 Data Collection

The quarterly data from the Ministry of Statistics and Implementation Programme [2] had the quarterly GDP data of different sectors. This dataset had Gross Domestic Product (GDP) figures for various industries across different quarters for a span of 15 years.

# 5 Implementation

## 5.1 Importing Libraries

The sections below demonstrate how to implement the project using Python. Run them as instructed, step by step. Preprocessing the provided data is the first step before we start the implementation. The libraries required for startup are shown in the picture below 3. ,

---

[1]https://colab.research.google.com/
[2]https://mospi.gov.in/data

```
[ ]  import os
     import pandas as pd
     import numpy as np
     import glob
     import seaborn as sns
     import matplotlib.pyplot as plt
     from sklearn.ensemble import RandomForestRegressor
     from sklearn.linear_model import Lasso
     from statsmodels.tsa.holtwinters import ExponentialSmoothing
     from tensorflow.keras.models import Sequential
     from tensorflow.keras.layers import Dense, LSTM
     from sklearn.preprocessing import MinMaxScaler
     from keras.models import Sequential
     from keras.layers import LSTM, Dense
     from sklearn.model_selection import train_test_split
     from statsmodels.tsa.arima.model import ARIMA
     from sklearn.metrics import mean_squared_error, mean_absolute_error
     from statsmodels.tsa.statespace.sarimax import SARIMAX
```

Figure 3: Importing Required Libraries

## 5.2 Data Preprocessing

The industry column specifies the particular industry or sector to which the GDP values pertain. This data is further splitted into 9 different CSV files as it represents 9 different GDP data.To do this operation the code shown in figure 4 need to be executed.

```
# Read the dataset
df = pd.read_csv(input_path)

# Iterate over each row in the dataframe
for index, row in df.iterrows():
    industry = row['INDUSTRY']
    industry_data = row[2:]  # Extract the data for the industry, excluding the first two columns

    # Drop all columns named 'Annual'
    industry_data = industry_data.drop(industry_data.filter(like='Annual').index, axis=0)

    # Create a separate DataFrame for the industry
    industry_df = pd.DataFrame(columns=['Time', 'GDP'])

    # Iterate over each quarter in the row
    for i, quarter in enumerate(industry_data):
        quarter_number = i % 4 + 1  # Calculate the quarter number
        year = 1996 + i // 4  # Calculate the year

        # Filter out NaN values
        if not pd.isna(quarter):
            date = pd.Timestamp(year=year, month=(quarter_number - 1) * 3 + 3, day=30)

            # Check if the value is a float or an integer
            if isinstance(quarter, float):
                value = int(quarter)
            else:
                value = str(quarter).replace(',', '')  # Remove any commas in the value

            # Append the values to the industry DataFrame
            industry_df = industry_df.append({'Time': date.strftime('%Y/%m/%d'), 'GDP': value}, ignore_index=True)

    # Save the industry DataFrame to a separate file
    file_name = f"{index+1}. {industry}.csv"
    file_path = os.path.join(output_path, file_name)
    industry_df.to_csv(file_path, index=False)
```

Figure 4: Data Preprocessing

## 5.3 Splitting of Data

The data has been splitted into 80-20 ratio for train and test respectively for ARIMA model as shown in the figure 5 and for RNN, LSTM and SARIMA the split was 70-30 as shown in the below figure 6 7 8 respectively.

3

```
# Calculate the number of records for the training set
num_records = len(df)
num_train_records = int(0.8 * num_records)
```

Figure 5: Train and Test Split for ARIMA

```
num_records = len(df)
num_train_records = int(0.7 * num_records)
train = df[:num_train_records]
test = df[num_train_records:]
```

Figure 6: Train and Test Split for RNN

## 5.4 Implementing ARIMA Model

The ARIMA model was implemented by executing the code in the below picture 9. All the necessary steps for building the model has been included in the code. This code gives a result with actual GDP, predicted GDP, difference between actual and predicted GDP values and the difference percentage.

## 5.5 Implementing SARIMA Model

The SARIMA model was implemented by executing the code in the below picture 10. All the necessary steps for building the model has been included in the code. This code gives a result with actual GDP, predicted GDP, difference between actual and predicted GDP values and the difference percentage.

## 5.6 Implementing LSTM Model

The LSTM model was implemented by executing the code in the below picture 11. All the necessary steps for building the model has been included in the code. This code gives a result with actual GDP, predicted GDP, difference between actual and predicted GDP values and the difference percentage.

## 5.7 Implementing RNN Model

The RNN model was implemented by executing the code in the below picture 12. All the necessary steps for building the model has been included in the code. This code gives a result with actual GDP, predicted GDP, difference between actual and predicted GDP values and the difference percentage.

```
num_records = len(df)
num_train_records = int(0.7 * num_records)
train = df[:num_train_records]
test = df[num_train_records:]
```

Figure 7: Train and Test Split for LSTM

```
# Calculate the number of records for the training set
num_records = len(df)
num_train_records = int(0.7 * num_records)
```

Figure 8: Train and Test Split for SARIMA



Figure 9: ARIMA Model



Figure 10: SARIMA Model



Figure 11: LSTM Model



Figure 12: RNN Model