# Configuration Manual

MSc Research Project
Data Analytics

# Ronan Flannery
Student ID: x19113269

School of Computing
National College of Ireland

Supervisor:    Manaz Kaleel

| Student Name: | Ronan Flannery |
|---|---|
| Student ID: | x19113269 |
| Programme: | Data Analytics |
| Year: | 2020 |
| Module: | MSc Research Project |
| Supervisor: | Manaz Kaleel |
| Submission Due Date: | 17/08/2020 |
| Project Title: | Configuration Manual |
| Word Count: | 955 |
| Page Count: | 8 |

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: | 11th August 2020 |

**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

Ronan Flannery
x19113269

## 1   Introduction

This configuration manual describes the software, environments and settings used in the research project "A Machine Learning Approach to Predicting Gross Domestic Product". This document may be used to replicate the technical work carried out in the research project.

## 2   Hardware Used

This research project was conducted on a Lenovo ideapad 320-15AST with the following configuration:

- **Operating System**: Windows 10

- **Processor**: AMD-A9-9420 RADEON R5, 5 Cores

- **RAM**: 8GB

Google Colab was used for creating and running the models used in this project.

## 3   Environment

Python was used to create the models in this project. The following environment was used in this research project:

- **Google Colaboratory**

Google Colab allows the user to write and run Python code in an online browser notebook. Colab provides a hosted online Jupyter Notebook service. It has advantages for creating and running machine learning models, particularly Artificial Neural Networks. There are less resource constraints than running models on a local machine and there is no requirement to install Python packages locally. Code written in Google Colab is saved to the users Google Drive and can be easily shared if working on a collaborative project. The following is required to use Google Colab:

- Internet Browser, for example Chrome, Firefox or Safari

- Google Account - A google account is required for using Google Colab

Google Colab is accessible here: `https://colab.research.google.com/notebooks/intro.ipynb`
A new notebook for writing code can be accessed through the File - New notebook option. An existing notebook can also be opened or uploaded via the File menu.
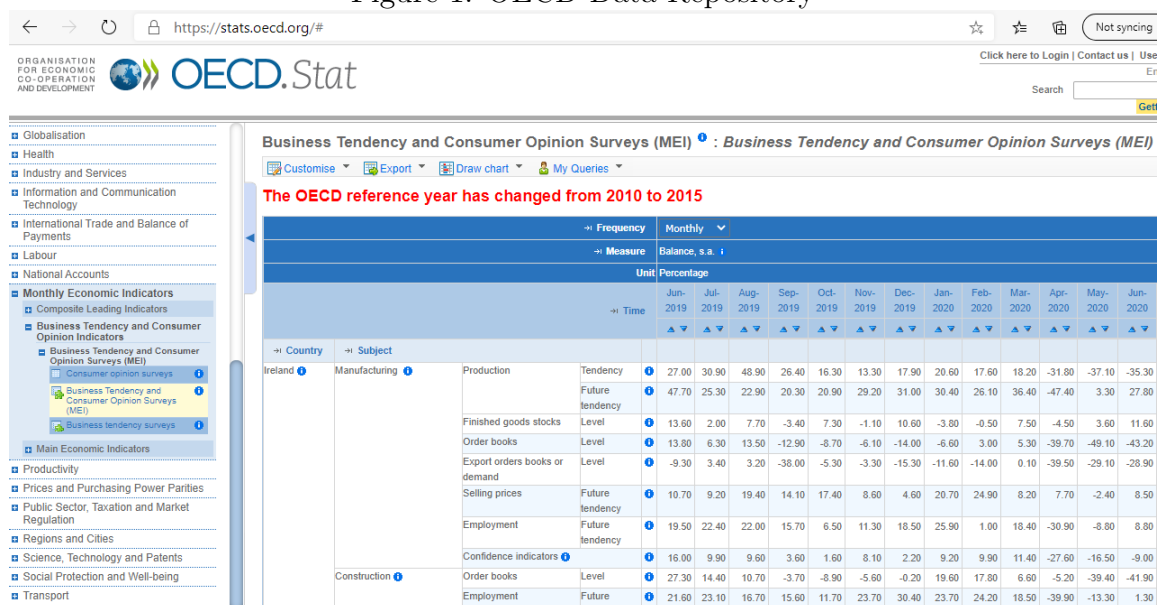
# 4 Implementation

The following section outlines the implementation of the project from a technical perspective.

## 4.1 Dataset

The data used in this research project was gathered from the OECD data repository. The link for this data source is: `https://stats.oecd.org/`

Figure 1 below shows the OECD data repository. The various indicators used in this projects data set can be selected from the menus on the left of the OECD screen. The data returned can be filtered by country, time period etc. The data can be downloaded in various formats, however in this project the data was downloaded in csv format. Data was taken for 10 different countries (Australia, Austria, Belgium, Denmark, Finland, Germany, Ireland, Korea, Sweden and the UK). Individual files for the economic indicators were downloaded, including Current and Constant GDP, Exchange Rates, Consumer Confidence Indicators and Business Confidence Indicators, Long Term and Short Term Interest Rates, Unemployment Rate, Exports, Government Expenditure, Private Consumption, Imports.

Figure 1: OECD Data Repository

In this project, data showing the dates of national elections in various countries was also sourced. This was used to create a variable during the feature engineering portion of the project. The variable indicated the months when an election was held in the country. This data was sourced from the following links: `http://www.parties-and-elections.eu/countries.html`, `http://elections.uwa.edu.au/index.lasso` and `https://en.wikipedia.org/wiki/Elections_in_South_Korea`

## 4.2 Google Colab environment setup

In this project Multilayer Perceptron and Random Forest Models were created. Figure 2 below shows the libraries imported for creating these models in Google Colab.

Figure 2: Libraries used for creating the models

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score

from sklearn.preprocessing import StandardScaler

import pickle
from sklearn.externals import joblib

Random Forest
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import metrics
%matplotlib inline
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestRegressor
```

The data files used in this project were stored on a google drive. In order to connect to the google drive to load/store files an authentication needs to be made with the google drive. Figure 3 below shows the code for authenticating with the Google Drive. A file can be shared from a google drive by marking it as sharable and creating a file id that can be incorporated in the code.

Figure 3: Google drive connection

**Connect to Google Drive where data is stored**

```
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)
```

## 4.3 Data Pre-processing

After downloading the data an exploratory data analysis was carried out and the individual data files were combined into a main data file in Microsoft Excel. Figure 4 below shows the combined data file used in this research project. The main data file used in the project contains data for 10 different countries over a period of 20 years (2000 - 2019).

Data pre-processing was carried out in Python. A number of unnecessary columns were removed from the data set and some other data pre-processing steps were performed.

3

Figure 4: Main data file



Figures 5 and 6 below show an example of the pre-processing code from Google Colab.

Figure 5: Data Pre-processing

```python
#Drop LOCATION, Subject, Measure, Frequency and other columns
del GDPdata['LOCATION']
del GDPdata['Subject']
del GDPdata['Measure']
del GDPdata['Frequency']
del GDPdata['Unit Code']
del GDPdata['Unit']
del GDPdata['PowerCode']
del GDPdata['Flag Codes']
del GDPdata['Flags']
data_top2 = GDPdata.head()
data_top2


# check for NaNs
tempdf1 = tempdf[tempdf.isna().any(axis=1)]
```

Figure 6: Training, test, validation split of the data

```python
# split data into training and testing data
# 15% for testing. 70% training. 15% validation
data_train, data_test = train_test_split(dfwip1,
                                         test_size=0.15,
                                         random_state=1)
```

## 4.4 Models

After the data pre-processing is completed, the the models are created.

### 4.4.1 Multi-Layer Perceptron Model

The Multilayer Perceptron model calls individual modules within the code to perform tasks as follows:

- Options.py - contains parameters for the model

- Dataset.py - this contains the code for loading the dataset

- Model.py - this contains code for the MLP model

- plot.py - this shows the final plots created for analysing the output of the model

Figure 7 below shows a sample of the code from the implementation of the MLP model. Figure 8 shows some of the hyperparameters configured for the MLP model.

Figure 7: Multi-Layer Perceptron Implementation

```
# MLP Model

# Importing the required Keras libraries and packages for the MLP Model
import keras
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import Dropout
from keras.layers import LSTM
batch = 1
from sklearn.neural_network import MLPRegressor

from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

my_module = drive.CreateFile({'id':'1oA5m72fjBwQ74TCX4dRoy955Jgd5O-LO'})
my_module.GetContentFile('Options.py')

import Options

model = MLPRegressor(hidden_layer_sizes=(16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16), early_stopping=True, validation_fraction=0.17647, max_iter
```

Figure 8: MLP hyperparameter Options

```
params = {
 'hidden_layer_sizes': (16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16)
 'early_stopping': True,
 'validation_fraction': 0.17647,
 'max_iter': 5000,
 'n_iter_no_change': 10,
 'verbose': True
}
```

### 4.4.2 Random Forest Model

Figure 9 below shows an example of the Random Forest model implementation, including the parameters of the Random Forest Regressor and some of the evaluation metrics.

Figure 9: Random Forest Implementation

```
rf1 = RandomForestRegressor(random_state=50, n_estimators=128)

rf1.fit(X, Y)

RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                      max_depth=None, max_features='auto', max_leaf_nodes=None,
                      max_samples=None, min_impurity_decrease=0.0,
                      min_impurity_split=None, min_samples_leaf=1,
                      min_samples_split=2, min_weight_fraction_leaf=0.0,
                      n_estimators=128, n_jobs=None, oob_score=False,
                      random_state=50, verbose=0, warm_start=False)


rf1_pred = rf1.predict(Xtest)


print('Random Forest Model Performance:')
print('MAE:', metrics.mean_absolute_error(Ytest, rf1_pred))
print('RMSE:', np.sqrt(metrics.mean_squared_error(Ytest, rf1_pred)))
print('R2_Score: ', metrics.r2_score(Ytest, rf1_pred))

Random Forest Model Performance:
MAE: 430.6834835907653
RMSE: 978.0751804655268
R2_Score:  0.9810473584483614
```

# 5  Output and Evaluation

This project uses the Python 'pickle' module for saving the MLP model in order to allow it to be called in a separate instance if required. The output of the MLP is also redirected to a log file, which can be used for further evaluation of the training iterations, if required.

Figure 10: Saving the model with pickle

```
filename = 'MLPgdp_model.sav'
scalerfile = 'MLPgdp_scale.sav'


pickle.dump(Model.model, open(filename, 'wb'))
pickle.dump(scaler, open(scalerfile, 'wb'))


# using JobLib
joblib.dump(Model.model, 'gdp.model')
joblib.dump(scaler, 'gdp.scaler')
```

Figure 11: Calling the Model and Displaying result

```
loaded_model = pickle.load(open(filename, 'rb'))
loaded_scaler = pickle.load(open(scalerfile, 'rb'))
result = loaded_model.score(Xtest, Ytest)
print(result)
```

```
0.9416148238964995
```

Figure 12: Redirecting the model output to a log file

```
#redirect output of model to log file
import sys
orig_stdout = sys.stdout
f = open('LogModelOutput.txt', 'w')
sys.stdout = f
print(Model.model.fit(X, Y))
sys.stdout = orig_stdout
f.close()
```

Figure 13: Output log showing the iterations of training

```
Iteration 436, loss = 1436349.88648591
Validation score: 0.876606
Iteration 437, loss = 1436322.31156765
Validation score: 0.882990
Iteration 438, loss = 1427336.52073444
Validation score: 0.883278
Iteration 439, loss = 1410908.82039794
Validation score: 0.880669
Iteration 440, loss = 1447898.66052808
```

## 5.1 Plots

This research project uses the Python 'matplotlib' library for creating graphs from the output of the models. Various graphs were created from the output, for example graphs showing the predictions of GDP vs the actual GDP figures. Figure 14 below shows an example of one of the output graphs.

Figure 14: Actual vs Predicted GDP MLP model