

Configuration Manual

MSc Research Project Data Analytics

Harsh Chudasama Student ID: X18187340

School of Computing National College of Ireland

Supervisor: Dr. Muhammad Iqbal

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Harsh Chudasama
Student ID:	X18187340
Programme:	Data Analytics
Year:	2020
Module:	MSc Research Project
Supervisor:	Dr. Muhammad Iqbal
Submission Due Date:	17/08/2020
Project Title:	Configuration Manual
Word Count:	650
Page Count:	11

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.

<u>ALL</u> 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:	17th 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
or□

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

Harsh Chudasama X18187340

1 Introduction

This configuration manual specifies the hardware and software requirements and the code snippets explaining the implementation of the below research project in detail:

"Forecasting the Novel Coronavirus(COVID-19) using Time Series Model"

2 System Configurations

2.1 Hardware

- Processor: Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz 1.60 GHz
- RAM: 8 GB
- System Type: Windows 10 (64 bit)
- GPU: Intel(R) UHD Graphics 620 , 2 GB
- Storage: 256 GB SSD

2.2 Software

• Microsoft Excel 2019

I am currently using the latest version of MS Excel(2019) Version for the data handling for Comma Separated Values(CSV) and Plotting the graphs. Additionally, the NumXL Plugin¹ is installed for ARIMA and ARMA Forecasting.

• Anaconda Distribution-Jupyter Notebook

Jupyter Notebook² is an interactive cell-based Integrated development environment(IDE). In this software, Python or R machine learning codes can be executed in a cell and output will be displayed. For this project, Exploratory data analysis, Data preprocessing, Data Transformation, Visualization, and Model Implementation is done using Jupyter Notebook (**Python Version 3.8.3**)

¹NumXL Tool Plugin : https://www.numxl.com/products/numxl

²Anaconda Distribution https://www.anaconda.com/distribution/

3 Project Development

The research project was executed successfully and all the four models(LSTM, ARMA, ARIMA, and Prophet) were deployed. Python language was used for scripting in the Jupyter Notebook IDE. The Project Development lifecycle consists of following steps Data preparation, Data Transformation, and Time Series Modelling.

3.1 Data preparation

For this project, Data(University; 2020) is downloaded and loaded in the Jupyter Notebook using Numpy libraries, and the Data manipulation task is done using the Pandas library.

3.1.1 Symptom Dataset

The main objective of this project was to extract the key symptoms for early prediction of COVID-19. The dataset was loaded successfully using Numpy CSV reader function and it consisted of 32 columns with 14,216 entries (Refer to Figure 1).

1	#Importing the dataset											
2	<pre>open_line = pd.read_csv('COVID19_open_line_list.csv')</pre>											
3	<pre>open_line = open_line.iloc[:, :-12]</pre>											
4	# Check for missing information, datatypes and shape											
5	open	_line.info()										
6	#Out	put										
7	<cla< th=""><th>ss 'pandas.core.frame.Data</th><th>Frame'></th><th></th></cla<>	ss 'pandas.core.frame.Data	Frame'>									
8	RangeIndex: 14126 entries, 0 to 14125											
9	Data	columns (total 33 columns):									
10	#	Column	Non-Null Count	Dtype								
11												
12	0	ID	13173 non-null	float64								
13	1	age	1349 non-null	object								
14	2	sex	1264 non-null	object								
15	3	city	10194 non-null	object								
16	4	province	12906 non-null	object								
17	5	country	13148 non-null	object								
18	6	wuhan(0)_not_wuhan(1)	13170 non-null	float64								
19	7	latitude	13147 non-null	float64								
20	8	longitude	13147 non-null	float64								
21	9	geo_resolution	13147 non-null	object								
22	10	date_onset_symptoms	746 non-null	object								
23	11	date_admission_hospital	730 non-null	object								
24	12	date_confirmation	13089 non-null	object								
25	13	symptoms	493 non-null	object								
26	14	lives_in_Wuhan	565 non-null	object								
27	15	travel_history_dates	503 non-null	object								
28	16	travel_history_location	758 non-null	object								
29	17	reported_market_exposure	35 non-null	object								
30	18	additional_information	2412 non-null	object								
31	19	chronic_disease_binary	<pre>18 non-null</pre>	float64								
32	20	chronic_disease	<pre>13 non-null</pre>	object								
33	21	source	12950 non-null	object								
34	22	sequence_available	1 non-null	object								
35	23	outcome	184 non-null	object								
36	24	date_death_or_discharge	93 non-null	object								
37	25	notes_for_discussion	187 non-null	object								
38	26	location	1024 non-null	object								
39	27	admin3	1159 non-null	object								
40	28	admin2	9068 non-null	object								
41	29	admin1	12877 non-null	object								
42	30	country_new	13079 non-null	object								
43	31	admin_id	13103 non-null	object								
44	L 32	data_moderator_initials	17 non-null	object								
45	dtyp	es: float64(5), object(28)										

Figure 1: Symptom Dataset

3.1.2 COVID-19 Dataset

This is the main dataset which consists of the Number of confirmed cases, deaths, and recovered cases that are mapped to their respective states (Refer to Figure 2).



Figure 2: COVID-19 Dataset

3.2 Data Transformation

3.2.1 Symptom Dataset

Dataset was loaded successfully with 32 features consisting of 14,216 rows. But, not all features were required for the analysis, and hence a sneak peek of data was done with the first four rows of each column with its data types and respective values counter. Below are a few observations (Refer to Figure 3):

- Feature 'sex' contains some ambiguities and requires unifying the values.
- Feature 'Wuhan(0)_not_wuhan(1)' indicates that all cases originate from outside of Wuhan the i.e. epicenter of the disease.
- Features ('age', 'additional_information', 'reported_market_exposure') are currently not planning for the project scope.
- The following features are removed as they are inconsistent and provide ambiguous information:



Figure 3: Symptom Dataset Analsis

age, chronic_disease_binary, chronic_disease, sequence_available, outcome, date_death_or_discharge, notes_for_discussion, location, admin3, admin2, admin1, country_new, admin_id, data_moderator_initials, lives_in_Wuhan, travel_history_dates, travel_history_location

• As the feature 'sex 'consisted of ambiguity it was normalized and the graph was plotted to look at the date(Refer to Figure 4).



Figure 4: Normalized Feature 'Sex'

• Based on the feature 'Symptom', Text analysis was performed to extract the Top 10 Key Symptoms by defining the custom function(Refer to Figure 5).

```
In [11]: # A function to extract symptoms
         def find_symptoms(word):
             word_split = word.replace('()',',').split(',')
word_split = [word.strip().rstrip(',') for word in word_split]
             key_symptoms.extend(word_split)
In [12]: # creating a dataframe of major symptoms (Top 10)
         key_symptoms = []
         symptoms_df['symptoms'].dropna().apply(find_symptoms)
         1
         print("Top 10 Major Sympyoms identified are :")
         major_symptoms[:10]
         Top 10 Major Sympyoms identified are :
Out[12]: fever
                         290
         cough
                         158
         sore throat
                          27
         pneumonitis
                          19
                          17
         fatigue
         chills
                          16
         pneumonia
                          16
         .
headache
                          13
         runny nose
                          13
         malaise
                          12
         dtype: int64
```

Figure 5: Top 10 Key Symptoms Feature Extraction

• The Outcome of the key Symptom was summarized using a Word cloud pack-age(Refer to Figure 6).



Figure 6: Word Cloud of Key Symptoms

3.2.2 COVID-19 Dataset

- Some discrepancies were observed in the state names and hence they were corrected and stored in new data frame 'merged' (Refer to Figure 7).
- Additional feature 'merged['dcratio']' was derived using the following formula and applied to each row of the data frame.

1.2 Covid-19 Exploratory Data Analysis

In [15]:	#current date today = df[df.date == '2020-08-04']								
In [16]:	<pre># Data Transformation gdf = gpd.read_file('Indian_States.shp') #renaming state names gdf['st_nm'].replace({"Andaman & Nicobar Island": "Andaman and Nicobar Islands",</pre>								
In [17]:	merged.info()								
	<pre><cli><cli><cli><cli><cli><cli><cli><cli< th=""></cli<></cli></cli></cli></cli></cli></cli></cli></pre>								

Figure 7: COVID-19 Dataset Transformation

merged['dcratio'] = merged['deaths'] / merged['confirmed'] * 100

• Seaborn library was used for plotting the graph of Fatality Rate Per State that required geopandas library for plotting the layout of India map (Refer to Figure 8).

3.3 Time Series Modelling

3.3.1 LSTM Model

- LSTM model was implemented which belongs to the following Recurrent Neural Network(RNN) class (Ayyoubzadeh et al.; 2020).
- Model Consisted of following layers deep neural network layer followed by dense and drop out layer(Refer to Figure 9).

3.3.2 Prophet Model

- The simplest and straight forward model to implement.
- Simple provide the input time series and give trend-setting information.(Refer to Figure 10)(Taylor and Letham; 2018).

3.3.3 ARIMA & ARMA Model

- These models were implemented using the MS Excel NumXL Plugin(Zhang; 2003).
- Initially, data were log-transformed and provided as input and the stationary test was performed.



Figure 8: Fatality Rate Per State

- Stationary test results proved that model shows ARCH Effect i.e. lags are not constant rather exponential (Refer to Figure 11) (Kelvin et al.; 2020).
- In conclusion, The model is aligned dependent on beginning segments and coefficients, and the result is anticipated with the ideal strides to conjecture(Refer to Figure 12).

References

- Ayyoubzadeh, S. M., Ayyoubzadeh, S. M., Zahedi, H., Ahmadi, M. and R Niakan Kalhori, S. (2020). Predicting covid-19 incidence through analysis of google trends data in iran: Data mining and deep learning pilot study, *JMIR Public Health Surveill* 6(2): e18828. URL: https://doi.org/10.2196/18828
- Kelvin, D., Rubino, S. and Kelvin, N. (2020). Similarity in case fatality rates (cfr) of covid-19/sars-cov-2 in italy and china, J Infect Dev Ctries 14(2): 125–128. URL: https://jidc.org/index.php/journal/article/view/32146445
- Taylor, S. J. and Letham, B. (2018). Forecasting at Scale, The American Statistician 72(1): 37–45.
 URL: https://ideas.repec.org/a/taf/amstat/v72y2018i1p37-45.html
- University, J. H. (2020). Covid-19 dataset. URL: https://github.com/CSSEGISandData/COVID-19

#Importing the required from sklearn.preprocess	Package ing import MinMaxScaler			
<pre>scaler = MinMaxScaler()</pre>				
#Scaling the date				
<pre>scaler.fit(train data)</pre>				
scaled train data = sca	ler.transform(train data)	1		
<pre>scaled_test_data = scale</pre>	er.transform(test_data)			
#Creating Dataframe of	Scaled Data as Model Innu	ı t		
data = pd.DataFrame(col	umns = ['ds'.'v'])			
data['ds'] = train data	index			
data['v'] = scaled train	n data			
#Creating Model	-			
#from keras.preprocessi	ng.sequence import Timese	riesGenerator		
from keras.models impor	t Sequential			
from keras.layers impor	t Dense			
from keras.layers import	t LSTM			
from keras.layers impor	t Dropout			
import tensorflow as tf				
n input = 12 # Number o	F Forecasting Days			
n_features= 1 #Input Co	lumns			
#generating hatches of	temporal data			
generator = TimeseriesG	enerator(scaled train dat	a. #Scaled Input Trai	n Data	
	scaled train dat	a. #Taraet Data		
	length=n input,	#Lenath as per input	size=12	
	<pre>batch_size=1) #B</pre>	Batches generated		
#Starting of Model in s	equential manner			
<pre>lstm model = Sequential</pre>	0			
#First LSTM Laver of si	ze (12.1) with 500 hidden	Lavers		
lstm model.add(LSTM(500	, activation='tanh', inpu	nt shape=(n input, n f	eatures)))	
#Drop out Laver	,, ,,			
lstm model.add(Dropout(0.10))			
#Adding Dense Laver				
lstm model.add(Dense(1))			
#Compiling the model wi	th Optimiszer as adam			
<pre>lstm_model.compile(opti)</pre>	<pre>nizer='adam', loss='mse')</pre>	1		
#Displaying Summary of	the model			
lstm_model.summary()				
#Generating the model w	ith 50 Epochs			
istm_model.fit_generato	(generator,epochs=50)			
Model: "sequential_18"				
Layer (type)	Output Shape	Param #		
	·····			
Istm_19 (LSTM)	(None, 500)	1004000		
dropout_10 (Dropout)	(None, 500)	0		
dense_13 (Dense)	(None, 1)	501		
_ ,				
Total params: 1,004.501				
Trainable params: 1.004	,501			
Non-trainable params: 0	,			

Figure 9: Code for LSTM Model

```
In [130]: #importing fbprophet
from fbprophet import Prophet
#model
m = Prophet()
#fitting the model
m.fit(df3)
#forecasting Future dates
future = m.make_future_dataframe(periods= 12)
future.tail(12)
#rename the column
fb_res.columns = ['ds', 'FBProphet']
fb_res['FBProphet'] = fb_res['FBProphet'].astype(int)
result = pd.concat([df2,fb_res],axis=1)
del result['deaths']
del result['ds']
result.FBProphet = result.FBProphet.replace(np.nan, 0)
out = result.tail(12)
```

Out[130]:

	date	confirmed	FBProphet
176	2020-07-24	1287945	1148408.0
177	2020-07-25	1336861	1171513.0
178	2020-07-26	1385522	1195112.0
179	2020-07-27	1435453	1218749.0
180	2020-07-28	1483156	1242036.0
181	2020-07-29	1531669	1265485.0
182	2020-07-30	1583792	1288950.0
183	2020-07-31	1638870	1307538.0
184	2020-08-01	1695988	1330644.0
185	2020-08-02	1750723	1354242.0
186	2020-08-03	1803695	1377879.0
187	2020-08-04	1855745	1401167.0

Figure 10: Code for FB Prophet Model



Figure 11: Stationary Test Check for ARIMA & ARMA Model

ARMA(1,1)				Goodness-of-fit				Re	siduals (standardized) Analysis						
	Param	Value	-	LLF A	AIC CH	неск		_	AVG	STDEV	Skew	Kurtosis	Noise?	Normal?	ARCH?
	μ	8.65	-	-246.64	501.27	1.			0.07	0.13	8.05	84.15	FALSE	FALSE	FALSE
	Φ1	1.00					Target	t	0.00	1.00	0.00	0.00			
	θ,	0.09					SIG?	?	TRUE	TRUE	TRUE	TRUE			
	σ	1.01													
ARMA(1,3)				Goodness-of-fit				Re	siduals (standardized) Analysis						
	Param	Value	-	LLF A	AIC CI	НЕСК		_	AVG	STDEV	Skew	Kurtosis	Noise?	Normal?	ARCH?
	μ	8.65	-	-254.50	517.00	1.		_	0.06	0.13	7.57	80.21	FALSE	FALSE	FALSE
	Φ1	1.00					Target	t	0.00	1.00	0.00	0.00			
	θ1	0.28					SIG?	?	TRUE	TRUE	TRUE	TRUE			
Step	Mean	STD	UL	ш											
1	14.03207963	1.012822498	16.01717525	12.04698401											
2	14.03153587	1.496123649	16.96388434	11.0991874											
3	14.03099217	1.857610735	17.67184231	10.39014203											
4	14.03044852	2.15935279	18.26270222	9.798194824											
5	14.02990493	2.423767008	18.78040097	9.279408889											
6	14.02936139	2.661999963	19.24678545	8.81193734											
7	14.02881791	2.880554877	19.67460173	8.383034097											
8	14.02827448	3.083618711	20.0720561	7.984492868											
9	14.02773111	3.274075094	20.44480038	7.610661846											
10	14.02718779	3.454010211	20.79692341	7.257452179											
11	14.02664453	3.624991153	21.13149664	6.921792429											
12	14.02610132	3.788230554	21.45089678	6.601305873											

Figure 12: ARMA Model Predictions

Zhang, G. (2003). Time series forecasting using a hybrid arima and neural network model, Neurocomputing 50: 159 - 175.
URL: http://www.sciencedirect.com/science/article/pii/S0925231201007020