

Configuration Manual

MSc Research Project Programme Name

Shubham Garg Student ID: x19205295

School of Computing National College of Ireland

Supervisor: Martin Alain

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Shubham Garg
Student ID:	x19205295
Programme:	Programme Name
Year:	2021
Module:	MSc Research Project
Supervisor:	Martin Alain
Submission Due Date:	16/12/2021
Project Title:	Configuration Manual
Word Count:	604
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.

<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:	Shubham Garg
Date:	31st January 2022

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

Shubham Garg x19205295

1 Introduction

This paper is the manual of configuration which enlists the all the requirements is needed to run the whole code of the project on the local machine environment. this manual also includes the screenshot of the hardware requirements as well as the software requirements. This file also includes the pictures of code of the Exploratory Data Analysis, Data loading , Data pre-processing, Data cleaning, Model applied which are Long Short Term Memory.

2 Environment

In this section of this module we have explained the requirements of hardware as well as the software and their version to re run the code on the local environment.

2.1 Hardware requirements

To perform research code we need hardware of minimum 11th generation of intel core i5 of version 1135G7 with the processor speed of 2.40GHz., 16GB RAM of DDR4 version with the speed of 3200 MHz, with the operating system of windows 11 of version 64 bit and the hard drive of 512GB SSD.Figure 3

í	Device specifica	tions	
	Device name	DESKTOP-JPE2J7A	
	Processor	11th Gen Intel(R) Core(TM) i5-1135G7 @ 2.40GHz 1.38 GHz	
	Installed RAM	24.0 GB (23.7 GB usable)	
	Device ID	32D59860-29A1-4599-BEF1-5E5F3DF3C5DD	
	Product ID	00327-35922-44027-AAOEM	
	System type 64-bit operating system, x64-based processor		
	Pen and touch	No pen or touch input is available for this display	
	Windows specifi	ications	
	Edition	Windows 11 Home Single Language	
	Version	21H2	
	Installed on	28-10-2021	
	OS build	22000.318	
	Experience	Windows Feature Experience Pack 1000.22000.318.0	
	Mining the Court		

Figure 1: Hardware Requirements

2.2 Software Requirements

- Anaconda Navigator for Windows (Version 1.9.12)
- Jupyter Notebook (Version 6.0.3)
- Python (Version 3.9)

3 Data Collection

The data set for the research project is taken from Kaggle which is an open directory. Link for dataset is https://www.kaggle.com/eliasdabbas/flights-serps-and-landing-pages? select=flights_serp_scrape.csv

https://www.kaggle.com/eliasdabbas/search-engine-results-flights-tickets-keywords? select=flights_tickets_serp2019-08-15.csv both are the same file but the second link has the data set which has been taken twice in every month from december 2018 to 1st of april 2020 what makes it the data set of 24 different csv with same script of code to scraping. it has 4000 rows of first 10 ranked pages for the same query. The rank has been changed every time.

4 Data Exploration

4.1 Importing The Libraries

There are many libraries which is required for processing the the code which has predefined functions shown below

```
from mpl_toolkits.mplot3d import Axes3D
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt # plotting
import numpy as np # linear algebra
import os # accessing directory structure
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
```

Figure 2: Libraries

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from keras.preprocessing.text import Tokenizer
from keras.preprocessing.sequence import pad_sequences
from keras.models import Sequential
from keras.layers import Dense, Embedding, LSTM, SpatialDropout1D
from sklearn.model_selection import train_test_split
from keras.utils.np_utils import to_categorical
from keras.callbacks import EarlyStopping
from keras.layers import Dropout
import re
import nltk
from nltk.corpus import stopwords
from nltk import word_tokenize
nltk.download('stopwords')
STOPWORDS = set(stopwords.words('english'))
from bs4 import BeautifulSoup
import plotly.graph_objs as go
#import plotly.plotly as py
import cufflinks
from IPython.core.interactiveshell import InteractiveShell
import plotly.figure_factory as ff
InteractiveShell.ast_node_interactivity = 'all'
from plotly.offline import iplot
cufflinks.go_offline()
cufflinks.set_config_file(world_readable=True, theme='pearl')
```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.

Figure 3: Libraries

4.2 Exploratory Data Analysis

4.2.1 Data description

	rank	totalResults	count	startIndex	searchTime	formattedSearchTime
count	4000.00000	4.000000e+03	4000.0	4000.0	4000.000000	4000.000000
mean	5.50000	5.538642e+07	10.0	1.0	0.321426	0.321350
std	2.87264	8.004667e+07	0.0	0.0	0.061791	0.062045
min	1.00000	9.940000e+04	10.0	1.0	0.162895	0.160000
25%	3.00000	9.217500e+06	10.0	1.0	0.278504	0.280000
50%	5.50000	2.300000e+07	10.0	1.0	0.314640	0.310000
75%	8.00000	7.097500e+07	10.0	1.0	0.354942	0.352500
max	10.00000	5.420000e+08	10.0	1.0	0.533649	0.530000

Figure 4: DATA DESCRIPTION

4.2.2 checking Null values

This code will check the null or Nan values in the dataset.



<Figure size 720x432 with 0 Axes><matplotlib.axes_subplots.AxesSubplot at 0x7f54f5c08790>Text(0.5, 1.0, 'Missing values'
Missing values?



Figure 5: Null values

4.2.3 Percentage of traffic by rank

This is percentage of traffic on rank based regardless the websites.

<pre>traffic = {1: 0.33,</pre>	
# Display the percentage of traffic with colors	
colors = []	
for 1 in range(10):	
$x_{c} = round(0, 7-0, 05^{*}i, 2)$	
$c = (x_c, x_c, \theta, 5)$	
colors, append(c)	
colors = colors + ["red"]	
x = [str(t) for t in range(1,11)] + [">10"]	
<pre>v = [traffic[kev]*100 for key in traffic.kevs()]</pre>	
v = [100 - sum(v)]	
plt.figure(figsize=(8,6))	
plt.bar(x,y, color = colors)	
plt.title("Estimation of percentage of traffic by ranking\nin Google search result	ts", fontsize = 16)
plt.xticks(x)	
plt.xlabel("# rank")	
plt.ylabel("% of traffic")	
plt.show()	

Figure 6: Percentage of traffic by rank

4.2.4 Percentage of traffic by websites

This code will give the percentage of traffic on the websites regardless the rank.



Figure 7: Percentage of traffic by websites

4.2.5 keyword concentration in titles

This code will get the percentage of the titles over the rank same code will be used for snippet column to know the concentration of key words in snippet.

```
# Calculate the % of keywords/search terms in the titles
df["%search_term_in_title"] = df["searchTerms"].apply(lambda x: len(x.split(" "))) / df["title"].apply(lambda x: len(x.split(" ")))
df["%search_term_in_title"] = 100 * df["%search_term_in_title"] # Convert the result in %
# Display the result
proc_searchterm_rank = pd.pivot_table(df, values = "%search_term_in_title", index = "rank", aggfunc = "mean").sort_index(ascending = False)
proc_searchterm_rank.plot.barh(figsize = (8,5), color = (0.32, 0.32, 0.5))
plt.legend("")
plt.xlabel("Exact keyword concentration in %", fontsize = 12)
plt.ylabel("# rank", fontsize = 14)
plt.title("Average keyword concentration in titles\nby rank", fontsize = 16)
```

Figure 8: keyword concentration in titles

4.2.6 concatenating all data

This function concatenate the all the data sets and formed one data set and make the new data frame.

```
# Import and merge all the csv files
lst_df_path = []
for dirname, _, filenames in os.walk('New_folder'):
    for filename in filenames:
        lst_df_path.append(os.path.join(dirname, filename))
# List the dates of the csv files
# Be aware that some months are missing
sorted([d[-14:-4] for d in lst_df_path ])
```

Figure 9: concatenating all data

4.2.7 Most visitors on websites over the time

this code helps to know the most visited websites as the ranks changes over the time so what will be the effect of it on the traffic of websites.



Figure 10: Most visitors on websites over the time

5 Pre Processing

5.0.1 Removing the impurities and stop words

this code help to remove all the symbols , numeric values, stop words etc. that menns it will remove everything except the text words.

df= df.reset_index(drop=True)	
REPLACE_BY_SPACE_RE = re.compile('[/(){}\[\]\[@,;]')	
BAD_SYMBOLS_RE = re.compile('[^0-9a-z #+_]')	
STOPWORDS = set(stopwords.words('english'))	
def clean text(text):	
text: a string	
return: modified initial string	
text = text.lower() # lowercase text	
tayt = REPLACE BY SPACE RE sub(' ' tayt) # penjace REPLACE BY SPACE RE symb	ols by space in text
#substitute the matched string	n REPLACE BY SPACE RE with space
text - RAD SYMBOLS RE sub('' text) # nemove symbols which are in RAD SYMBOL	S RE from text
text = bab_stibles_tersub(, text) = text symbols when the not show a single fill	AD SYMPOLE BE with pathing
# substitute the matched string in a	AD_STHBOLS_KE WICH HOCHINg.
text = text.replace(x,)	
text = text.replace(0,)	
text = text.replace(1,)	
text = text.replace(2 ,)	
text = text.replace('3', ')	
<pre>text = text.replace('4','')</pre>	
text = text.replace('5','')	
text = text.replace('6','')	
text = text.replace('7','')	
<pre>text = text.replace('8','')</pre>	
text = text.replace('9','')	
<pre># text = re.sub(r'\W+', '', text)</pre>	
text = ' '.join(word for word in text.split() if word not in STOPWORDS) # re	move stopwors from text
return text	
df['title'] = df['title'].apply(clean text)	

Figure 11: Removing the impurities and stop words

5.1 tokenizing the words

5.1.1 tokenizer function

```
MAX_NB_WORDS = 50000
MAX_SEQUENCE_LENGTH = 250
EMBEDDING_DIM = 100
tokenizer = Tokenizer(num_words=MAX_NB_WORDS, filters='!"#$%&()*+,-./;;<=>?@[\]^_`{|}~', lower=True)
tokenizer.fit_on_texts(fu['combined'].values)
word_index_title = tokenizer.word_index
```

Figure 12: tokenizer function

5.1.2 tokenizing the words in the data frame

```
Y = pd.get_dummies(rank_df).values
print('Shape of label tensor:', Y.shape)
X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size = 0.10, random_state = 42)
print(X_train.shape,Y_train.shape)
print(X_test.shape,Y_test.shape)
```

Figure 13: tokenizing the words in the data frame

5.2 splitting into test and train

This code split the data into the ratio of 90 and 10 in train and test data respectively

```
Y = pd.get_dummies(rank_df).values
print('Shape of label tensor:', Y.shape)
```

```
X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size = 0.10, random_state = 42)
print(X_train.shape,Y_train.shape)
print(X_test.shape,Y_test.shape)
```

Figure 14: splitting into test and train

6 Modeling

this is the model which is applied on the data set for prediction the target variable "RANK" this model is then changed by adding and deleting the dense layer and by changing the weights.

```
model = Sequential()
model.add(Embedding(MAX_NB_WORDS, EMBEDDING_DIM, input_length=X.shape[1]))
model.add(SpatialDropout1D(0.2))
model.add(LSTM(100, dropout=0.2, recurrent_dropout=0.2))
model.add(Dense(10, activation='softmax'))
model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
print(model.summary())
```



Here the model has been fitted and will provide the result with the epochs history.

epochs = 50
batch_size = 25
history = model.fit(X train, Y train, epochs=epochs, batch size=batch size,validation split=0.1,callbacks=[EarlyStopping(monitor='val loss', patience=6, min delta=0.0001)])

Figure 16: fitting the model