

Configuration Manual

MSc Research Project
Data Analytics

Sharon George
Student ID: 21245240

School of Computing
National College of Ireland

Supervisor: Mayank Jain

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Sharon George
Student ID:	21245240
Programme:	Data Analytics
Year:	2023
Module:	MSc Research Project
Supervisor:	Mayank Jain
Submission Due Date:	14/12/2023
Project Title:	Configuration Manual
Word Count:	602
Page Count:	6

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:	Sharon George
Date:	31st January 2024

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
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.	<input type="checkbox"/>

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

Sharon George
21245240

1 Introduction

This config manual has all the information needed to use the project. It tells you the least and best things you need to do in the thesis. It connects the project and the thesis, explaining the software and hardware you need and showing the important parts with code bits. It gives steps for getting data, running the project and showing the main results. All of this is focused on an AI-based Product Recommendation System that looks at YouTube comments.

2 Hardware Configuration

①	Device specifications
Device name	DESKTOP-M9R8QHR
Processor	12th Gen Intel(R) Core(TM) i5-1235U 1.30 GHz
Installed RAM	16.0 GB (15.7 GB usable)
Device ID	71272EFB-F0E1-4248-802F-FEC17B07614B
Product ID	00342-42621-98970-AAOEM
System type	64-bit operating system, x64-based processor
Pen and touch	Pen and touch support with 10 touch points
Related links	Domain or workgroup System protection Advanced system s
⚙	Windows specifications
Edition	Windows 11 Home Single Language
Version	22H2
Installed on	01-02-2023
OS build	22621.2715
Experience	Windows Feature Experience Pack 1000.22677.1000.0
	Microsoft Services Agreement
	Microsoft Software License Terms

Figure 1: Hardware specifications

Figure 1 above outlines the required software specifications for implementing the project on the local machine.

3 Software Configuration

The following software must be installed for the project to run successfully:

1. Anaconda Navigator for Windows

2. Jupyter Notebook
3. Python
4. YouTube

4 Data Selection

The dataset for the project is available from an open data source called Kaggle. The URL is <https://www.kaggle.com/datasets/kazanova/sentiment140>. It has around 1.6 lakh tweets from a twitter api.

5 Code files of the project

- 1_GPT2_Feature_Extraction.ipynb - The aim of this code is to collection the data, data cleaning, data preprocessing and create a new csv file with the extracted features. The below Figure 2 displays all the libraries used by this file. Foundation (year of the documentation)

```
import warnings
warnings.filterwarnings("ignore")

import pandas as pd
pd.set_option("display.max_columns", None)
import matplotlib.pyplot as plt
%matplotlib inline
from tqdm import tqdm
from lib_file import lib_path
import re
import torch
from transformers import GPT2Model, GPT2Tokenizer
import re
import contractions
from langdetect import detect
import wordcloud
```

Figure 2: Imported Libraries for 1_GPT2_Feature_Extraction.ipynb file

- 2_ModelTraining.ipynb - A CLSTM model is created on the extracted features and saved. Also three common evaluation metrics used for classification tasks are accuracy score, classification report, and confusion matrix. The below Figure 3 displays all the libraries used by this file.

```

import warnings
warnings.filterwarnings("ignore")

import pandas as pd
pd.set_option("display.max_columns", None)
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
from sklearn.utils import resample
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from lib_file import lib_path

```

Figure 3: Imported Libraries for 2_ModelTraining.ipynb file

- 3_Inference.ipynb - The inference code begins with the user providing a YouTube link. Once the user provides a YouTube link, the system proceeds to collect all comments associated with the video. Following this, the collected data undergoes cleaning and preparation before being fed into the model, which ultimately produces the result. The below Figure 4 displays all the libraries used by this file.

```

import warnings
warnings.filterwarnings("ignore")

import pandas as pd
pd.set_option("display.max_columns", None)
import numpy as np
from tensorflow.keras.models import load_model
import re
from bs4 import BeautifulSoup
from tqdm import tqdm
import torch
from transformers import GPT2Model, GPT2Tokenizer
from transformers import pipeline
from urllib.parse import urlparse
from lib_file import lib_path
import googleapiclient
from langdetect import detect
from IPython import display

```

Figure 4: Imported Libraries for 3_Inference.ipynb file

6 Code snippets

6.1 1_GPT2_Feature_Extraction file

The below code snippet in Figure 5 depicts the text_cleaning function which is used for preprocessing textual data.

```
def remove_emojicons(text):
    emojicons_to_remove = [":)", ":(", ":D", ":s", ":)", ":/", ":O:", ":P", ":O", "&)", "\u2013", ">O", ":3", ">:(\"", "s"]
    for emojicon in emojicons_to_remove:
        text = text.replace(emojicon, '')
    return text

def text_cleaning(text):
    text = remove_emojicons(text)
    emoji_pattern = re.compile("[
        \u0001F600-\u0001F64F" # Emojis
        \u0001F300-\u0001F5FF" # Symbols & pictographs
        \u0001F680-\u0001F6FF" # Transport & map symbols
        \u0001F700-\u0001F7FF" # Alphabetic presentation forms
        \u0001F780-\u0001F7FF" # Geometric shapes
        \u0001F800-\u0001F8FF" # Miscellaneous symbols
        \u0001F900-\u0001F9FF" # Supplemental symbols & pictographs
        \u0001FA00-\u0001FA6F" # Extended-A
        \u0001FA70-\u0001FAFF" # Extended-B
        \u0001F004-\u0001F0CF" # Mahjong tiles
        \u0001F170-\u0001F251" # Enclosed characters
        \u00020000-\u0002073F" # Chinese, Japanese, and Korean characters
        \u000E0000-\u000E007F" # Tags
    ]+", flags=re.UNICODE)
    text = re.sub(emoji_pattern, '', text)
    text = re.sub("http[s+|www[s+|https[s+", '', text, flags=re.MULTILINE)
    text = re.sub(r'(@|w+|w+)', '', text)
    text = contractions.fix(text)
    text = re.sub(r'^[a-zA-Z0-9]+', '', text)
    try:
        lang = detect(text)
        if lang == 'en':
            return text.lower()
        else:
            return 'empty'
    except:
        return 'empty'

sentences = []
for text in tqdm(df['text'].values):
    sentences.append(text_cleaning(text))

100%|██████████ 40000/40000 [07:05<00:00, 93.951t/s]
```

Figure 5: Text cleaning function

This code shown in Figure 6 essentially processes a list of cleaned text data using a pre-trained GPT-2 model, extracting features from the hidden states of the model for each text.

```
# Load a pre-trained GPT-2 model
gpt2_model = GPT2Model.from_pretrained('gpt2')
tokenizer = GPT2Tokenizer.from_pretrained('gpt2')
```

Figure 6: Applying GPT-2 to the cleaned text

The code snippet below in Figure 7 is the result after applying the GPT-2 Model

[illegible]

Figure 7: Extracted features

6.2 2_ModelTraining file

The code snippet below is the model architecture to build the CLSTM Model

Algorithm : Convolutional Long Short-Term Memory

```

]: ▶ from tensorflow.keras import Sequential
    from tensorflow.keras.layers import Conv1D,MaxPool1D,Flatten,Dense,Dropout,LSTM,BatchNormalization

]: ▶ clstm_model = Sequential()
    clstm_model.add(Conv1D(filters=32, kernel_size=3, activation='relu', padding='same', input_shape=(x_t
    clstm_model.add(Conv1D(filters=32, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(Conv1D(filters=64, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(Conv1D(filters=64, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(Conv1D(filters=128, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(Conv1D(filters=128, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(Conv1D(filters=256, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(Conv1D(filters=256, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(Conv1D(filters=512, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(Conv1D(filters=512, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(Conv1D(filters=1024, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(Conv1D(filters=1024, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(Conv1D(filters=2048, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(Conv1D(filters=2048, kernel_size=3, activation='relu', padding='same'))
    clstm_model.add(MaxPool1D())
    clstm_model.add(BatchNormalization())
    clstm_model.add(LSTM(units=256, return_sequences=True))
    clstm_model.add(LSTM(units=256, return_sequences=True))
    clstm_model.add(Flatten())
    clstm_model.add(Dense(512, activation='relu'))
    clstm_model.add(Dropout(0.4))
    clstm_model.add(Dense(2, activation='sigmoid'))
    clstm_model.compile(loss='binary_crossentropy', optimizer="adam", metrics=['accuracy'])

```

Figure 8: Model architecture - CLSTM

The Figure 9 below shows the Accuracy Score and Classification report of the model.

Accuracy Score

```

▶ clstm_model_accuracy=accuracy_score(y_true=y_true,y_pred=pred)
print("ConvolutionallongShortTermMemory model's Validation accuracy is {
ConvolutionallongShortTermMemory model's Validation accuracy is 89.05%

```

Classification Report

```

▶ print(classification_report(y_true=y_true,y_pred=pred, target_names=clas

```

	precision	recall	f1-score	support
negative	0.84	0.97	0.90	2000
positive	0.96	0.81	0.88	2000
accuracy			0.89	4000
macro avg	0.90	0.89	0.89	4000
weighted avg	0.90	0.89	0.89	4000

Figure 9: Accuracy Score and Classification report

6.3 3 Inference.ipynb file

The Figure 10 below depicts the code snippet to enter the youTube link and verify it

```
video_url = input("Paste youtube's appropriate link...\n")

Paste youtube's appropriate link...
https://youtu.be/t-URmQbR5FY?si=YVG98lQErypub8v7

video_id = get_video_id_from_url(video_url)
if video_id is not None:
    print(f"The video ID, {video_id}, has been detected.")
else:
    print("The video ID could not be retrieved from the URL. Please try a different URL.")

The video ID, t-URmQbR5FY, has been detected.
```

Figure 10: Entering link and verifying youtube link

The below Figure 11 show show to pick the youTube link that would be added to the above code snippet in Figure 10

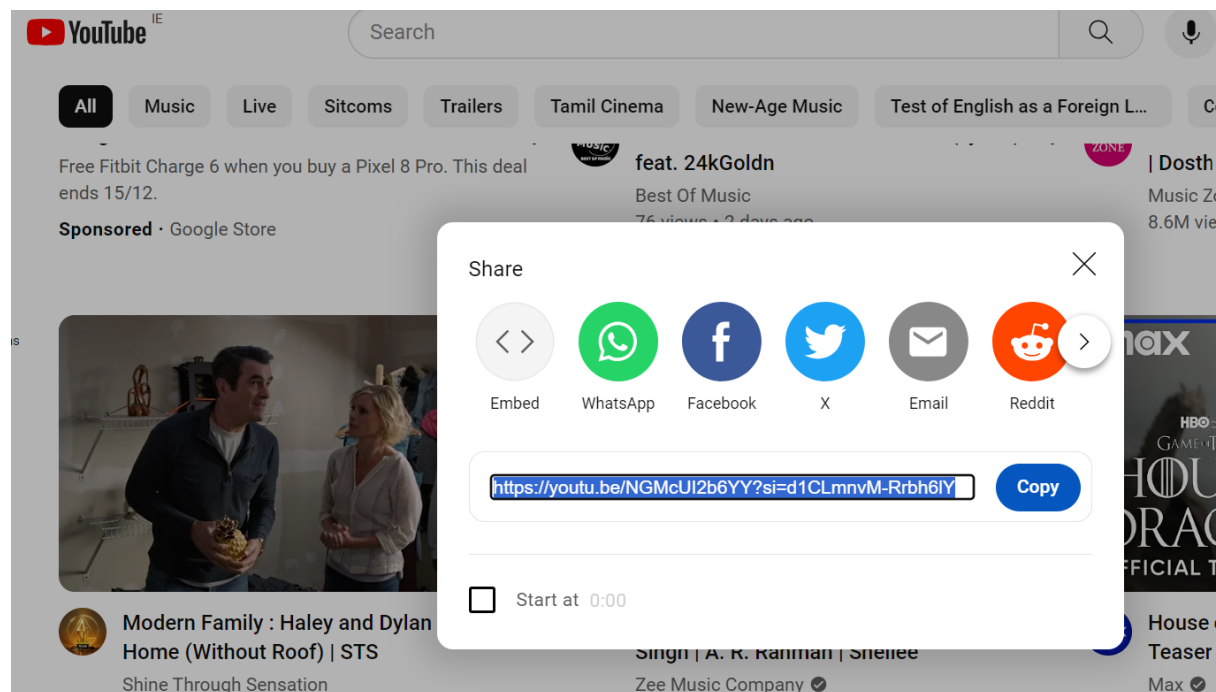


Figure 11: Pick a link from YouTube

References

Foundation, P. S. (year of the documentation). Python documentation. Accessed on: 31st January 2024.

URL: <https://docs.python.org/3/library/index.html>