

Sentiment Analysis on Covid-19 Vaccination Reviews Using BERT and Comparative Study with LSTM, Vader, and Text blob Models - Configuration Manual

MSc Research Project
Data Analytics

Sourav Ramalingam
Student ID: x20199911

School of Computing
National College of Ireland

Supervisor: Mr. Taimur Hafeez

National College of Ireland
MSc Project Submission Sheet
School of Computing



Student Name: MR. Sourav Ramalingam
Student ID: x20199911
Programme: Data Analytics **Year:** 2022
Module: MSc Research Project
Lecturer: Mr. Taimur Hafeez
Submission Due Date: 15/12/2022.....
Project Title: Sentiment Analysis on Covid-19 Vaccination Reviews Using BERT and Comparative Study with LSTM, Vader, and Text blob Models.
Word Count:1142..... **Page Count:**10.....

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: Sourav Ramalingam.....
Date: 15/12/2022.....

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

Sourav Ramalingam
Student ID: x20199911

1. Introduction

In the manual configuration document, the system specification that has been used for the application development is documented. Documenting the configuration details would help to set up the code in any other system for evaluation, making enhancements and so on. On top of the system configuration, the document also contains the packages that are used in the chosen programming language. These details are very important as multiple versions are available for the same package and the combination of the right ones needs to be used.

The document also highlights the steps like pre-processing the content, applying the model and so on.

2. System Configuration

In this section, the hardware, software, and programming configuration that is used for model development are documented.

2.1 Hardware & Base Specification:

Operating System	Windows 11 Home Edition
Installed RAM	8.00 GB
Processor	Intel Core i5 CPU @2.50 GHz
System Type	64-Bit Operating System
Programming Language	Python Programming
Package Management	PIP
Development Environment	PyCharm Free Community Edition

2.2 Software Specification:

The application is developed using the Python programming language and the base python package version used is 3.10.x.

The application also uses the below components:

- HTML for the web user interface development.
- Bootstrap CSS for designing the web page more easily.
- jQuery version 3.4.1 for the client-side scripting operations.
- Roboto Slab Font from online CDN.
- Python Programming for the Core Model Development.
- Python Flask module for developing the WEB API URL access to interact from the client-side HTML pages.

2.3 Package Specification:

The below-mentioned packages are used in Python programming for application development. The packages are downloaded from the package management system called 'PIP'.

Package Name	Version	Package Description
NLTK	3.7	Package used for natural language processing.
Word2number	1.1	Package used to convert the numbers in words to numeric.
Pandas	1.5.1	Package provides a fast and flexible way of working with the data.
TensorFlow	2.10.0	Package for open-source machine learning framework.
Gensim	4.2.0	Package used for document indexing and similarity retrieval from the huge text.
Keras	2.10.0	An application programming interface (API) developed by Google for implementing neural networks.
Numpy	1.23.4	Package used for performing powerful operations on the array object with multiple dimensions.
Scikit-learn	1.1.3	Package that is used to perform machine learning tasks which are built on top of SciPy.
VaderSentiment	3.3.2	Package to consume the VADER (Valence Aware Dictionary and sEntiment Reasoner) which is a lexicon and rule-based sentiment analysis tool.
TextBlob	0.17.1	Library that processes the textual data and find the sentiment.
Torch	1.13.0	Package that provides features like tensor computation and building deep neural networks.
Transformers	4.23.1	Package that has pretrained models to perform tasks on text, vision, and audio data.
Flask	2.2.2	It is a lightweight WSGI web application framework, used for API development.

3. Data Source

The vaccination-related tweets extracted from the Twitter feed are used for the model development. The source has N Number of columns, anyways for finding the sentiment score, the user review text and the rating would be taken into consideration. The rest of the fields are dropped while processing the same in the python coding.

The dataset has total records of 1,25,906.

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	id	user_name	user_locat	user_desc	user_creat	user_follo	user_frien	user_favo	user_verifi	date	text	hashtags	source
2	1.34E+18	Rachel Rol	La Crescer	Aggregato	#####	405	1692	3247	FALSE	#####	Same folks	['PfizerBio	Twitter f
3	1.34E+18	Albert Fon	San Franci	Marketing	#####	834	666	178	FALSE	#####	While the world has k		Twitter V
4	1.34E+18	eliðŸ±±ðŸ!	Your Bed	heil, hydra	#####	10	88	155	FALSE	#####	#coronavi	['coronavi	Twitter f
5	1.34E+18	Charles Ad	Vancouver	Hosting "C	#####	49165	3933	21853	TRUE	#####	Facts are immutable,		Twitter V
6	1.34E+18	Citizen News	Channel	Citizen Ne	#####	152	580	1473	FALSE	#####	Explain to	['whereare	Twitter f
7	1.34E+18	Dee	Birmingham	Gastroent	#####	105	108	106	FALSE	#####	Does anyone have an		Twitter f
8	1.34E+18	Gunther Fe	Austria, Uk	End North	#####	2731	5001	69344	FALSE	#####	it is a bit se	['vaccinati	Twitter V
9	1.34E+18	Dr.Krutika	Kuppalli	ID, Global	#####	21924	593	7815	TRUE	#####	There	['BidenHar	Twitter f
10	1.34E+18	Erin Despas		Designing&	#####	887	1515	9639	FALSE	#####	Covid	['CovidVac	Twitter V
11	1.34E+18	Ch.Amjad .	Islamabad	#ProudPa	#####	671	2368	20469	FALSE	#####	#CovidVa	['CovidVac	Twitter V
12	1.34E+18	Tamer Yaz	Turkey-Isra	Im	#####	1302	78	339	FALSE	#####	while	['PfizerBio	Twitter V
13	1.34E+18	VoiceM		campaigne	#####	2	25	20	FALSE	#####	@cnnbrk #	['COVID19	Twitter V
14	1.34E+18	WION	India	#WION: W	#####	292510	91	7531	TRUE	#####	The agency also relea		TweetDe
15	1.34E+18	Dr.Krutika	Kuppalli	ID, Global	#####	21924	593	7815	TRUE	#####	For all the	['PfizerBio	Twitter f
16	1.34E+18	Opoyi		High-qualit	#####	10332	49	16	FALSE	#####	"Expect 145 sites acrc		TweetDe
17	1.34E+18	City A.M.	London, Er	London's k	#####	66224	603	771	TRUE	#####	Trump	['vaccine']	Twitter f
18	1.34E+18	STOPCOM	Global	'Trust' is n	#####	406	176	479	FALSE	#####	UPDATED	['YellowFe	Twitter V
19	1.34E+18	ILKHA	TÃ¼rkiye	Official Tw	#####	4056	6	3	TRUE	#####	Coronavirt	['Iran', 'coi	TweetDe
20	1.34E+18	Braderz73	Bristol, UK	One of	#####	6430	6292	45007	FALSE	#####	.@Pfizer w	['CovidVac	Twitter f
21	1.34E+18	Alex Vie	Los Angele	Marine vet	#####	125	442	5401	FALSE	#####	The trump	['COVIDIO	Twitter f

4. Code Setup – Step-by-Step

- Download the source code and place it in a local path under an empty folder.
- Find the /dataset folder and place the downloaded dataset file from the provided link if the folder is empty.
- Ensure Python 3.10.X is installed.
- Ensure PyCharm Community Version IDE or any similar one is installed.
- Open the source code with IDE and it would prompt to create a Virtual Environment, create the same.
- To activate the virtual environment, the below code needs to be executed in the Terminal:
 - o Cd env/Scripts activate
- Install the required packages with the below comment:
 - o Pip3 install -r requirements.txt
- Execute the below command to initiate the web app and interact with the UI by providing sample inputs and viewing the results:
 - o Python -m flask run

5. Source Code Modules

The below code snippet is to perform the pre-processing on the tweet text before it is taken for processing by the model.

```
preprocess.py x
1 import re
2 import string
3 import nltk
4 import unicodedata
5 from word2number import w2n
6 from nltk.tokenize import word_tokenize
7 from nltk.corpus import stopwords, wordnet
8 from nltk.stem import WordNetLemmatizer
9
10
11 # Learn From https://www.w3schools.com/python/python\_regex.asp
12 class TweetTextCleaner(object):
13
14     # Remove 'RT' or 'rt' in text
15     def remove_retweets(self, text):
16         cleaned_text = re.sub(r'\bRT\b', '', text)
17         return cleaned_text
18
19     # Remove URLs in text
20     def remove_urls(self, text):
21         cleaned_text = re.sub('(?:http?:\:\/\/|http?:\/|https?:\:\/\/|https?:\/|https?:\:\/\/www)\S+', '', text)
22         return cleaned_text
23
24     # Remove username with '@' in text
25     def remove_mentions(self, text):
26         cleaned_text = re.sub('@[\^s]+', '', text)
27         return cleaned_text
28
29     # Remove hashtags '#' in text
30     def remove_hashtags(self, text):
31         cleaned_text = re.sub('#[\^s]+', '', text)
32         return cleaned_text
```

The below code snippet is to calculate the sentiment score using the Vader Sentiment Package.

```
VaderSentimentModel.py x
1 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
2
3
4 def VaderSentiment_Predict(input_sentence):
5     sid_obj = SentimentIntensityAnalyzer()
6     score = sid_obj.polarity_scores(input_sentence)
7     if float(score['compound']) >= 0.05:
8         result = 'Positive'
9     elif float(score['compound']) <= - 0.05:
10        result = 'Negative'
11    else:
12        result = 'Neutral'
13
14    print(f"Predicted Result from Vader Sentiment Package: {result}")
15    return result
16
```

The below code snippet is to calculate the sentiment score using the TextBlob Sentiment Package.

```
TextBlobModel.py x
1   from textblob import TextBlob
2
3
4   def TextBlob_Predict(input_sentence):
5       data = TextBlob(input_sentence)
6       result = ''
7       if data.sentiment.polarity < 0:
8           result = 'Negative'
9       elif data.sentiment.polarity == 0:
10          result = 'Neutral'
11       elif data.sentiment.polarity > 0:
12          result = 'Positive'
13
14          print(f"Predicted Result from TextBlog Package: {result}")
15          return result
16
```

The below code snippet is the LSTM Sentiment Analysis Model that performs both the model training and prediction on based of the value of the mode variable: 'test/train'.

```
SentimentAnalysis_LSTM.py x
123  model1.add(layers.Embedding(max_words, 20))
124  model1.add(layers.LSTM(15, dropout=0.5))
125  model1.add(layers.Dense(3, activation='softmax'))
126
127  model1.compile(optimizer='rmsprop', loss='categorical_crossentropy', metrics=['accuracy'])
128  checkpoint1 = ModelCheckpoint("model_weight/LSTM_Model.hdf5", monitor='val_accuracy', verbose=1,
129                               save_best_only=True, mode='auto',
130                               period=1, save_weights_only=False)
131  model1.fit(X_train, y_train, epochs=10, validation_data=(X_test, y_test), callbacks=[checkpoint1])
132
133  best_model = keras.models.load_model("model_weight/LSTM_Model.hdf5")
134  model_trained = True
135  lstmpredict = model1.predict(X_test)
136  print(lstmpredict)
137
138  test_loss, test_acc = best_model.evaluate(X_test, y_test, verbose=2)
139  print('Model accuracy: ', test_acc)
140  print(classification_report(np.argmax(y_test, axis=1), np.argmax(lstmpredict, axis=1)))
141  return test_acc
142
143  elif mode == 'test':
144      if not model_trained:
145          train_predict('train')
146          print('testing inputs')
147          sentiment = ['Neutral', 'Negative', 'Positive']
148          sequence = tokenizer.texts_to_sequences([input])
149          test = pad_sequences(sequence, maxlen=max_len)
150          result = sentiment[np.around(best_model.predict(test), decimals=0).argmax(axis=1)[0]]
151          print(f"Predicted Result from Long Short Term Memory (LSTM) Sentiment Analysis Module: {result}")
152          return result
```

The below code snippet is the BERT Model training code:

```
BERT_ModelTraining.py x
144     best_accuracy = 0
145
146     for epoch in tqdm(range(EPOCHS)):
147         print(f"Epoch {epoch + 1}/{EPOCHS}")
148         print("-" * 10)
149
150         train_acc, train_loss = train_epoch(
151             model, train_data_loader, loss_fn, optimizer, device, len(df_train)
152         )
153
154         print(f"Epoch: {epoch}, Train loss: {train_loss}, accuracy: {train_acc}")
155
156         val_acc, val_loss = eval_model(
157             model, val_data_loader, loss_fn, device, len(df_val)
158         )
159
160         print(f"Epoch: {epoch}, Val loss: {val_loss}, accuracy: {val_acc}")
161
162         history["train_acc"].append(train_acc)
163         history["train_loss"].append(train_loss)
164         history["val_acc"].append(val_acc)
165         history["val_loss"].append(val_loss)
166
167         if val_acc > best_accuracy:
168             torch.save(model.state_dict(), "model_weight/BERT_Model.hdf5")
169             best_accuracy = val_acc
170
171     print("Training completed")
172
```

The below code snippet is the BERT Model Prediction code:

```
BERT_ModelPredict.py x
1  import keras
2  import numpy as np
3  import pickle
4  from keras_preprocessing.sequence import pad_sequences
5
6
7  def predict(input_sentence=''):
8      max_len = 200
9      with open('model_weight/tokenizer.pickle', 'rb') as handle:
10         tokenizer = pickle.load(handle)
11         data = [input_sentence]
12         tokenizer.fit_on_texts(data)
13         best_model = keras.models.load_model("model_weight/BERT_Model.hdf5")
14         sentiment = ['Neutral', 'Negative', 'Positive']
15         sequence = tokenizer.texts_to_sequences(data)
16         test = pad_sequences(sequence, maxlen=max_len)
17         result = sentiment[np.argmax(best_model.predict(test), decimals=0).argmax(axis=1)[0]]
18         print(f"Predicted Result from BERT Sentiment Analysis Module: {result}")
19         return result
20
```


The below code snippet is for exposing the Python Flask API:

```
app.py x
4 from LSTM import SentimentAnalysis_LSTM
5
6 app = Flask(__name__)
7 app.secret_key = "secret key"
8
9
10 @app.route('/')
11 def index():
12     return render_template('index.html')
13
14
15 @app.route('/modelexec', methods=['GET', 'POST'])
16 def model_execution():
17     if request.method == 'GET':
18         input_Text = request.args.get("input_Text")
19         print(input_Text)
20         result = {
21             "textblob": TextBlobModel.TextBlob_Predict(input_Text),
22             "vader": VaderSentimentModel.VaderSentiment_Predict(input_Text),
23             "bert": BERT_ModelPredict.predict(input_Text),
24             "lstm": SentimentAnalysis_LSTM.train_predict('test', input_Text)
25         }
26         return result
27
28
29 if __name__ == '__main__':
30     app.run(debug=True)
```

Related Links:

- Python Download Link: <https://www.python.org/downloads/release/python-3100/>
- Pip: <https://pypi.org/project/pip/>
- PyCharm IDE Download: <https://www.jetbrains.com/pycharm/download/>
- Covid-19 Vaccination Dataset: <https://www.kaggle.com/datasets/gpreda/all-covid19-vaccines-tweets>
- NLTK - <https://pypi.org/project/nltk/>
- Word2number - <https://pypi.org/project/word2number/>
- Pandas - <https://pypi.org/project/pandas/>
- TensorFlow - <https://www.tensorflow.org/install/pip>
- Gensim - <https://pypi.org/project/gensim/>
- Keras - <https://pypi.org/project/keras/>
- NumPy - <https://pypi.org/project/numpy/>

- Scikit-learn - <https://pypi.org/project/scikit-learn/>
- Vader Sentiment - <https://pypi.org/project/vaderSentiment/>
- Text Blob - <https://pypi.org/project/textblob/>
- Torch - <https://pypi.org/project/torch/>
- Transformers - <https://pypi.org/project/transformers/>
- Flask - <https://pypi.org/project/Flask/>