

Configuration Manual

MSc Research Project Data Analytics

Avis Massey Student ID: x21199752

School of Computing National College of Ireland

Supervisor: Aaloka Anant

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Avis Massey
Student ID:	x21199752
Programme:	Data Analytics
Year:	2023
Module:	MSc Research Project
Supervisor:	Aaloka Anant
Submission Due Date:	14/12/2023
Project Title:	Configuration Manual
Word Count:	404
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:	
Date:	14th December 2023

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

Avis Massey x21199752

1 Introduction

This configuration manual describes the hardware and the software requirements which is essential to build the counterfeit review identifier using Neural Network - LSTM & CNN, Classifier - SVM and Transfer learning models BERT, RoBERTa, DistilBERT and ALBERT.

2 Hardware and Software Requirement

To train the LSTM, CNN and SVM models along with the transfer learning models such as BERT, RoBERTa, DistilBERT, and ALBERT, Google Collaboratory cloud is used which is a cloud machine since we have a vast dataset to train. For which you can find below the specification of the host device.

Device specifications

Device name	DESKTOP-QNDC301					
Processor	Intel(R) Core(TM) i5-6200U CPU @ 2.30GHz 2.4 GHz					
Installed RAM	8.00 GB					
Device ID	AA00322E-F555-4682-ABA9-55F5DC04D1B6					
Product ID	00331-10000-00001-AA009					
System type	64-bit operating system, x64-based processor					
Pen and touch	No pen or touch input is available for this display					

Figure 1: Device Specifications

The following table 1 provides information regarding the programming languages, computaional unit, and libraries employed throughout the project.

Specification	Value
IDE	Google Colab
Computation	GPU
Number of GPU	1
Programming	Python
language	
Modeling library	SimpleTransformer, HuggingFace Transformer, Sklearn, Pandas, Numpy,
	Matplotlib, Seaborn, Wandb, Keras, PyTorch, tqdm, SciPy

Table 1: System Specifications for the Project

3 Dataset

The dataset which we have implied in the project is in the CSV format. Originally these reviews are from amazon and the fake are AI generated using GPT. There are namely 5 column, Category, Label, Rating, Reviews and Fake_Review_flag. The Reviews and Fake_Review_flag are the most important columns. Salminen et al. (2022)

A	0	- C	
Category	Rating	Label	Reviews
Home_and_Kitchen_5	5	CG	Love this! Well made, sturdy, and very comfortable. I love it!Very pretty
Home_and_Kitchen_5	5	CG	love it, a great upgrade from the original. I've had mine for a couple of years
Home_and_Kitchen_5	5	CG	This pillow saved my back. I love the look and feel of this pillow.
Home_and_Kitchen_5	1	CG	Missing information on how to use it, but it is a great product for the price! I
Home_and_Kitchen_5	5	CG	Very nice set. Good quality. We have had the set for two months now and have not been
Home_and_Kitchen_5	3	CG	I WANTED DIFFERENT FLAVORS BUT THEY ARE NOT.
Home_and_Kitchen_5	5	CG	They are the perfect touch for me and the only thing I wish they had a little more space.
Home_and_Kitchen_5	3	CG	These done fit well and look great. I love the smoothness of the edges and the extra
Home_and_Kitchen_5	5	CG	Great big numbers & easy to read, the only thing I didn't like is the size of the
Home_and_Kitchen_5	5	CG	My son loves this comforter and it is very well made. We also have a baby
Home and Kitchen 5	5	CG	As advertised. 5th one I've had. The only problem is that it's not really a

Figure 2: Dataset overview

4 Implementation of the project

In the Google Colab we chage the runtime to make use of the GPU provided by the google which is the T4 GPU showing exceptional computation capabilities.

Pyth	ion 3		*					
Hardware a	accelerat	tor ၇						
0	CPU	О Т4	GPU	\bigcirc	A100 GPL	C) V10	0 GPU
0	TPU							
Want acc	ess to pr	emium GPI	Us? Pur	chase a	additional	compute	units	

Figure 3: change run time to GPU

4.1 Installation of Libraries

The 4 shows install ion of required libraries such as numpy pandas sentencepiece scipy torch transformers.

```
[1] !pip install numpy pandas sentencepiece scipy torch transformers
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packag
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packag
Collecting contengeniese
```

Figure 4: Dataset overview

Once installed the imported libraries are shown below in 5

```
[ ]
from keras.layers import Conv1D, Dense, Embedding, GlobalMaxPooling1D, LSTM
from keras.models import Sequential
from keras.preprocessing.sequence import pad_sequences
from keras.preprocessing.text import Tokenizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import accuracy_score, classification_report
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from torch.utils.data import DataLoader, RandomSampler, SequentialSampler, TensorDataset
from tqdm import tqdm
from transformers import AlbertTokenizer, AlbertForSequenceClassification, BertTokenizer,
from scipy.stats import zscore
```

Figure 5: Import Libraries

4.2 Importing the dataset

We import the data into the dataframe df using the pandas library.



Figure 6: Store the dataset in dataframe.

4.3 Tokenize for LSTM and CNN

We tokenize the dataset for LSTM and CNN models.

```
[ ] # Tokenizer for LSTM and CNN
max_features = 10000 # Number of words to consider as features
tokenizer = Tokenizer(num_words=max_features)
tokenizer.fit_on_texts(X_train)
X_train_seq = tokenizer.texts_to_sequences(X_train)
X_test_seq = tokenizer.texts_to_sequences(X_test)
# Padding sequences for LSTM and CNN
maxlen = 100 # Cuts off reviews after 100 words
X_train_pad = pad_sequences(X_train_seq, maxlen=maxlen)
X_test_pad = pad_sequences(X_test_seq, maxlen=maxlen)
# TF-IDF for SVM
vectorizer = TfidfVectorizer(max_features=max_features)
X_test_tfidf = vectorizer.transform(X_train)
X_test_tfidf = vectorizer.transform(X_test)
```

Figure 7: Tokenization for the models

4.4 Building the LSTM model and CNN model

The LSTM and CNN models are built and executed, for which the outputs are given below.



Figure 8: Building the LSTM Model

4.5 Predicting the LSTM and CNN models

The CNN models outperform the LSTM model after evaluation.

4.6 SVM building and Evaluation

SVM is built and executed, the evaluation states that it can be an alternative.

Figure 9: Building the CNN model

```
[ ] # Predicting with LSTM model
y_pred_lstm = lstm_model.predict(X_test_pad)
y_pred_lstm_classes = (y_pred_lstm > 0.5).astype("int32")
# Predicting with CNN model
y_pred_cnn = cnn_model.predict(X_test_pad)
y_pred_cnn_classes = (y_pred_cnn > 0.5).astype("int32")
# Evaluation for LSTM model
print("LSTM Model Evaluation")
print(classification_report(y_test, y_pred_lstm_classes))
# Evaluation for CNN model
print("CNN Model Evaluation")
print(classification_report(y_test, y_pred_lstm_classes))
253/253 [=======] - 8s 31ms/step
253/253 [======] - 1s 2ms/step
LSTM Model Evaluation
```

Figure 10: Predicting the models

5 Transfer Models

The transformer model hyperparameter are shown in 15 which is used to fine tune the model.

6 Results

The BERT function model is executed for which the reference is given below.

The evaluation of the model is given in 16

Similarly, we build all the other transformer models.

References

Salminen, J., Kandpal, C., Kamel, A., Jung, S. G. and Jansen, B. (2022). Creating and detecting fake reviews of online products, *Journal of Retailing and Consumer Services* 64.

∋	253/253 [=== 253/253 [===	31ms/step 2ms/step			
	LSTM MODEL E	precision	recall	f1-score	support
	0	0.95	0.94	0.94	4118
	1	0.94	0.95	0.94	3969
	accuracy			0.94	8087
	macro avg	0.94	0.94	0.94	8087
	weighted avg	0.94	0.94	0.94	8087
	CNN Model Ev				
		precision	recall	f1-score	support
	0	0.94	0.93	0.93	4118
	1	0.93	0.94	0.93	3969
	accuracv			0.93	8087
	macro avg	0.93	0.93	0.93	8087
	weighted avg	0.93	0.93	0.93	8087

Figure 11: LSTM and CNN evaluation

```
[ ] # Building and training the SVM model
svm_model = SVC(kernel="linear")
svm_model.fit(X_train_tfidf, y_train)
# Predicting and evaluating on the test set
y_pred_svm = svm_model.predict(X_test_tfidf)
accuracy_svm = accuracy_score(y_test, y_pred_svm)
print(f"SVM Model Accuracy: {accuracy_svm}")
```

Figure 12: SVM model building

[]	<pre># Evaluation for SVM model print("SVM Model Evaluation") print(classification_report(y_test, y_pred_svm))</pre>						
	SVM Model Evalua	ation recision	recall	f1-score	support		
	0 1	0.91 0.90	0.90 0.91	0.91 0.90	4118 3969		
	accuracy macro avg weighted avg	0.91 0.91	0.91 0.91	0.91 0.91 0.91	8087 8087 8087		

Figure 13: SVM Model Evaluation


```
[ ] # Hyperparameters
    bert_batch_size = 32
    roberta_batch_size = 32
    distilbert_batch_size = 32
    albert_batch_size = 32
    bert_epochs = 5
    roberta_epochs = 5
    distilbert_epochs = 5
    albert_epochs = 5
    bert_max_len = 64
    roberta_max_len = 64
    distilbert_max_len = 64
    albert_max_len = 64
```

Figure 14: Hyperparameters for the transfer models

```
[ ] # Function for BERT Tokenization
    def bert_encode(data, tokenizer, max_len):
        input_ids = []
        attention_masks = []
        for i in tqdm(range(len(data))):
            encoded = tokenizer.encode_plus(
                data.iloc[i],
                add_special_tokens=True,
                max_length=max_len,
                pad_to_max_length=True,
                return_attention_mask=True,
                return_tensors="pt",
            input_ids.append(encoded["input_ids"])
            attention_masks.append(encoded["attention_mask"])
        input_ids = torch.cat(input_ids, dim=0)
        attention_masks = torch.cat(attention_masks, dim=0)
           tunn innut ide attantion marks
```

Figure 15: BERT Model Function

	print(classif	ication_repo	rt(bert_f	lat_true_la	abels, bert_	flat_prediction		
1	0% /usr/local/li warnings.wa 100% 100% Some weights You should pr /usr/local/li	0/32345 b/python3.10, rn(32345/323 8087/8087 of BertForSec obably TRAIN b/python3.10,	ncation was sformers/tok 54.19it/s] .63it/s] n were not i wn-stream ta sformers/opt	not explicitly enization_utils nitialized from sk to be able t imization.py:41				
	warnings.warn(precision recall f1-score support							
	0 1	0.99 0.95	0.95 0.99	0.97 0.97	4118 3969			
	accuracy macro avg weighted avg	0.97 0.97	0.97 0.97	0.97 0.97 0.97	8087 8087 8087			

