

Configuration Manual

MSc Research Project
Data Analytics

Kartik Sharma
Student ID: 21125813

School of Computing
National College of Ireland

Supervisor: Prof. Christian Horn

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Kartik Sharma
Student ID:	21125813
Programme:	Data Analytics
Year:	2022
Module:	MSc Research Project
Supervisor:	Prof.Christian Horn
Submission Due Date:	15/12/2022
Project Title:	Configuration Manual
Word Count:	355
Page Count:	9

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:	Kartik Sharma
Date:	15th December 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

Kartik Sharma
21125813

1 System Configuration

The project is done on 64-Bit Windows 11 operating system with 16GB ram with AMD Ryzen 9 5900HX processor having a base clock speed of 3.30 GHz.

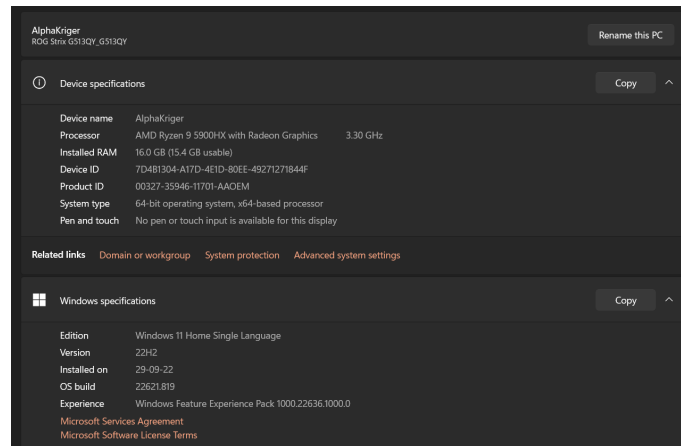


Figure 1: System Configuration

2 Software Requirement

For the project we have used following software :

1. Anaconda 2.3.3
2. Python 3.9.7
3. Jupyter Notebook

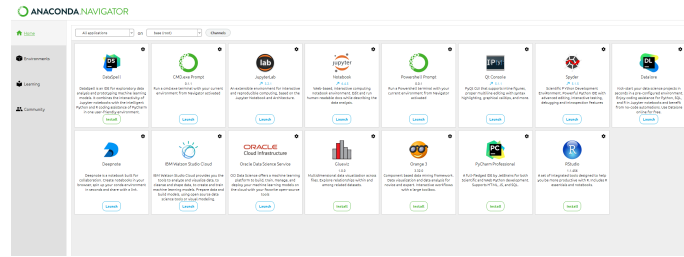


Figure 2: Software Requirement

3 Python Libraries

The project uses following python libraries :

1. TensorFlow
2. Keras
3. numpy
4. matplotlib
5. pandas
6. sklearn
7. itertools
8. nltk
9. transformers
10. scipy

4 Dataset

- The Data is originally created by Edmunds(USA-based car selling business), and for the study, the data is taken from the Kaggle.com where it was publicly published. Original Dataset link: <https://tinyurl.com/bdfjau8y>
- The Data consists of reviews from car owners in the USA from 2002 to 2018 in separate (comma separated value(.csv)) files for each car brand.

5 Data Preprocessing

- All the data is combined into one common dataset to begin the project.
- Data is cleaned and not required columns were dropped from the data frame ,refer to figure 3.

```
In [4]: tempData = [['vehicleTemplateName', 'ReviewTempTitle', 'Review', 'Rating']]

In [5]: reviewDataset = pd.DataFrame()

In [6]: reviewDataset
Out[6]: _

In [7]: for i in range(0, len(item_list)):
tempDataset = pd.read_csv(path+item_list[i], engine='python', error_bad_lines=False)
#print(item_list[i])
tempDataset = tempDataset.drop(['Unnamed: 0', 'Author_Name'], axis=1)
reviewDataset = reviewDataset.append(tempDataset, ignore_index=True)

In [8]: reviewDataset = reviewDataset.dropna()
reviewDataset = reviewDataset.reset_index()
reviewDataset = reviewDataset.drop(['index'], axis=1)

In [9]: reviewDataset.shape
Out[9]: (214808, 5)

In [10]: brands = []
for i in item_list:
i = i.replace('Scrapped_Car_Review_', '')
i = i.replace('Scrapped_Car_Review_', '')
brand_name = i[:4]
brands.append(brand_name)
```

Figure 3: Data Processing step 1

- The text preprocessing and cleaning is done before training the models.

Cleaning the texts

Removing Punctuations , Special Characters and Emojis

```
In [23]: def text_cleaner(review):
review = review.lower()
review = re.sub('[\.\?!\,]', '', review)
review = re.sub('[%s]' % re.escape(string.punctuation), '', review)
review = re.sub('[\u00d7\u2013]', '', review)
review = remove_emojis(review)
return review

def remove_emojis(review):
emoji = re.compile("[
    u'\U0001f600-\U0001f64f' # emoticons
    u'\U0001f300-\U0001f3ff' # symbols & pictographs
    u'\U0001f680-\U0001f6ff' # transport & map symbols
    u'\U0001f1e0-\U0001f1ff' # flags (ios)
    u'\U00002500-\U000025ff' # chinese char
    u'\U00002700-\U000027ff'
    u'\U000024c2-\U000024c3'
    u'\U0001f926-\U0001f927'
    u'\U00010000-\U00010fff'
    u'\u2600-\u2602'
    u'\u2600-\u2605'
    u'\u200d'
    u'\u23cf'
    u'\u23e9'
    u'\u231a'
    u'\ufe0f' # dingbats
    u'\u3030'
    ]+")"
return re.sub(emoji, '', review)

cleaned_review = lambda x : text_cleaner(x)
```

Figure 4: Data Processing step 2

- Figure 5 is the final processed and cleaned dataset that was finally used to train the Neural Networks

Out[59]:

	cleaned_Reviews	Rob_Emotion
0	its been a great delivery vehicle for my cafe...	positive
1	bought this car as a commuter vehicle for a v...	negative
2	this van rocks its the best lots of \rroom i ...	positive
3	great work vehicle drives nice has lots of ro...	positive
4	good solid frame and suspension well equippe...	negative
...
212847	i have owned the tigan for a year and month...	positive
212848	now had months with suv nice suv handles we...	negative
212849	smaller dimensions and my driving experience ...	positive
212850	we have had our tigan for a month prior to ...	positive
212851	the tigan s can be had for a reasonable and...	positive

212852 rows x 2 columns

Figure 5: Processed Dataset for Neural Network

6 Data Analysis

The Sales were compared to the Sentiments of consumers per brand to find the pattern between sales and sentiments.

Finding the Pattern for Ford

```
In [30]: tempData = []
brandDataset = roberta_Senti_Dataset[roberta_Senti_Dataset['brand'] == 'ford']
brandDataset = brandDataset.reset_index()
brandDataset.head()
```

Out[30]:

	index	Rob_Emotion	Review_Date	Brand	Rob_Emo_Quantified
0	8306	positive	2018-05-08	ford	3
1	8307	positive	2017-08-12	ford	3
2	8308	positive	2017-08-15	ford	3
3	8309	positive	2017-05-18	ford	3
4	8300	positive	2018-01-03	ford	3

```
In [31]: from tqdm import tqdm
import numpy

# Getting Average Emotion per year
tempData = []
brandEmotionDataset = pd.DataFrame()

for i in range(0, len(brandDataset.index)):
    tempData.append('temp')
brandEmotionDataset['AverageEmotion'] = tempData
brandEmotionDataset['Year'] = tempData

year=[]
for i in range(0, len(brandDataset.index)):
    year.append(brandDataset['Review_Date'][i].year)
year = list(set(year))
year = numpy.array(year)

for i in range(len(year)):
    emovalue = 0
    count = 0
    averageEmotion = 0
    for j in range(0, len(brandDataset.index)):
        if year[i] == brandDataset['Review_Date'][j].year:
            emovalue = emovalue + brandDataset['Rob_Emo_Quantified'][j]
            count = count+1
    averageEmotion = emovalue/count
    brandEmotionDataset['Year'][i] = year[i]
    brandEmotionDataset['AverageEmotion'][i] = averageEmotion

brandEmotionDataset.drop(brandEmotionDataset[brandEmotionDataset.AverageEmotion == 'temp'].index, inplace=True)
brandEmotionDataset
```

Figure 6: Processing data to generate sentiment graph

Comparing with Ford Sales

```
In [35]: SalesData = pd.read_csv('FordSales.csv')
SalesData = SalesData[['Year', 'Sales']]

for i in range(0, len(SalesData.index)):
    SalesData['Sales'][i] = re.sub(',', '', SalesData['Sales'][i])

SalesData
```

Out[35]:

	Year	Sales
0	2002	2990472
1	2003	2886575
2	2004	2786025
3	2005	2634041
4	2006	2415059
5	2007	2087048
6	2008	1880321
7	2009	1440863
8	2010	1752811
9	2011	2057210
10	2012	2180859
11	2013	2403542
12	2014	2378841
13	2015	2501855
14	2016	2487487
15	2017	2484041
16	2018	2381635
17	2019	2293684

Figure 7: Processing sales data to generate graph

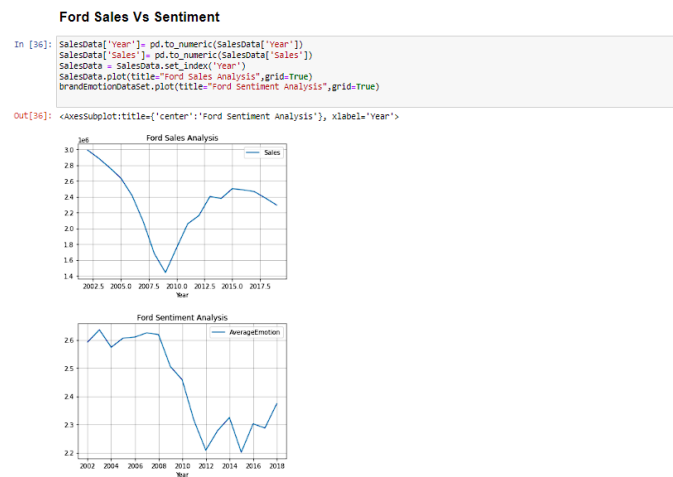


Figure 8: Analysing the Sales and Sentiment Graphs

7 Model Training and Testing

Three approaches were taken to extract the sentiments from the texts, and after evaluation the best approach was combined with neural network to create final model.

1. Logistic Regression

```
Model Training - Model Based On LogisticRegression

In [4]: from sklearn.model_selection import train_test_split

x = dataset.cleaned_Reviews
y = dataset.Happiness_Level_RoundRatingBased
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.2,random_state=444)
print('x_train : ',len(x_train),'\n x_test : ',len(x_test),'\n y_train : ',len(y_train),'\n y_test : ',len(y_test))

x_train : 171846
x_test : 42962
y_train : 171846
y_test : 42962

In [5]: from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression

tfidf = TfidfVectorizer()
lrclassifier = LogisticRegression(solver='saga')

In [6]: from sklearn.pipeline import Pipeline
from sklearn.metrics import confusion_matrix

model_logistic = Pipeline([('vectorizer',tfidf),('classifier',lrclassifier)])
model_lr = model_logistic.fit(x_train,y_train)

predictions = model_lr.predict(x_test)
#confusion_matrix(predictions,y_test)

Model Evaluations

In [7]: modelEvaluator('Logistic Regression',model_lr,predictions,y_test)

Model Evaluations for Logistic Regression
Accuracy : 0.6511568362739165
Precision : 0.707206216235996
Recall : 0.6511568362739165
      precision    recall  f1-score   support

excellent      0.75      0.89      0.81      23700
happy          0.51      0.41      0.45      12405
neutral        0.37      0.31      0.34      3755
sad            0.38      0.26      0.31      2190
very sad       0.43      0.16      0.23       912
```

Figure 9: Logistic Regression Model and Evaluation

2. Vader Sentiment Scoring

VADER Sentiment Scoring - Is Based on words in the sentence and not the overall meaning of sentence

```
In [11]: from nltk.sentiment import SentimentIntensityAnalyzer
from tqdm.notebook import tqdm
sia = SentimentIntensityAnalyzer()

In [12]: sia.polarity_scores('I am somewhat happy ')
Out[12]: {'neg': 0.0, 'neu': 0.37, 'pos': 0.63, 'compound': 0.5279}

In [13]: results = {}
for i, row in tqdm(dataset.iterrows(), total=len(dataset)):
    text = row['cleaned_reviews']
    data_id = row['id']
    results[data_id] = sia.polarity_scores(text)
    results
```

Out[13]:

0%	0/214000 [00:00<, 1it/s]
0: {'neg': 0.0, 'neu': 0.825, 'pos': 0.175, 'compound': 0.743},	
1: {'neg': 0.125, 'neu': 0.776, 'pos': 0.099, 'compound': -0.4793},	
2: {'neg': 0.0, 'neu': 0.787, 'pos': 0.213, 'compound': 0.8519},	
3: {'neg': 0.0, 'neu': 0.87, 'pos': 0.13, 'compound': 0.8689},	
4: {'neg': 0.128, 'neu': 0.757, 'pos': 0.115, 'compound': -0.501},	
5: {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0},	
6: {'neg': 0.004, 'neu': 0.746, 'pos': 0.249, 'compound': 0.9412},	
7: {'neg': 0.0, 'neu': 0.896, 'pos': 0.104, 'compound': 0.8999},	
8: {'neg': 0.0, 'neu': 0.766, 'pos': 0.232, 'compound': 0.9405},	
9: {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0},	
10: {'neg': 0.133, 'neu': 0.499, 'pos': 0.368, 'compound': 0.8749},	
11: {'neg': 0.0, 'neu': 0.908, 'pos': 0.092, 'compound': 0.8834},	
12: {'neg': 0.0, 'neu': 0.985, 'pos': 0.015, 'compound': 0.6136},	
13: {'neg': 0.0, 'neu': 0.886, 'pos': 0.114, 'compound': 0.9785},	
14: {'neg': 0.022, 'neu': 0.891, 'pos': 0.827, 'compound': 0.6187},	
15: {'neg': 0.018, 'neu': 0.721, 'pos': 0.261, 'compound': 0.9758},	
16: {'neg': 0.112, 'neu': 0.821, 'pos': 0.067, 'compound': -0.7618},	
17: {'neg': 0.0, 'neu': 0.991, 'pos': 0.009, 'compound': 0.8325},	

```
In [14]: vader_scoring_reviews = pd.DataFrame(results).T
vader_scoring_reviews = vader_scoring_reviews.reset_index().rename(columns={'index':'id'})
vader_scoring_reviews = vader_scoring_reviews.merge(dataset)
```

Figure 10: Vader Sentiment Scoring

3. RoBERTa Pretrained model

Roberta Pretrained Model (A Deep Learning Model)- Based on the meaning and context of the sentence

```
In [18]: from transformers import AutoTokenizer
from transformers import AutoModelForSequenceClassification
from scipy.special import softmax

In [19]: ## Transfer Learning
pretrainedModel = f"cardiffnlp/twitter-roberta-base-sentiment"
tokenizer = AutoTokenizer.from_pretrained(pretrainedModel)
model = AutoModelForSequenceClassification.from_pretrained(pretrainedModel)

In [20]: def scores_RobertaModel(text):
    encoded_text = tokenizer(text, return_tensors='pt')
    output = model(**encoded_text)
    scores = output[0][0].detach().numpy()
    scores = softmax(scores)
    scores_dict = {
        'Rob_neg': scores[0],
        'Rob_neu': scores[1],
        'Rob_pos': scores[2],
    }
    return scores_dict

In [21]: results = {}
total_errors=0
for i, row in tqdm(dataset.iterrows(), total=len(dataset)):
    try:
        text = row['cleaned_reviews']
        data_id = row['id']
        vader_results = sia.polarity_scores(text)
        vader_result_rename = {}
        for key, value in vader_results.items():
            vader_result_rename[f'vader_{key}'] = value
        roberto_results = scores_RobertaModel(text)
        both = {'vader_results': vader_result_rename, 'roberto_results': roberto_results}
        results[data_id] = both
    except RuntimeError:
        total_errors = total_errors + 1
        print(f"Runtime Error for id : {data_id}")
    except IndexError:
        total_errors = total_errors + 1
        print(f"Index Error for id : {data_id}")

print('Total sentences left out = ', total_errors, '\n Total sentences analyzed successfully : ', len(dataset) - total_errors)
```

Runtime Error for id : 735
Runtime Error for id : 1550
Runtime Error for id : 2732
Runtime Error for id : 3404
Runtime Error for id : 4235
Index Error for id : 4633
Runtime Error for id : 4857

Figure 11: RoBERTa pretrained model

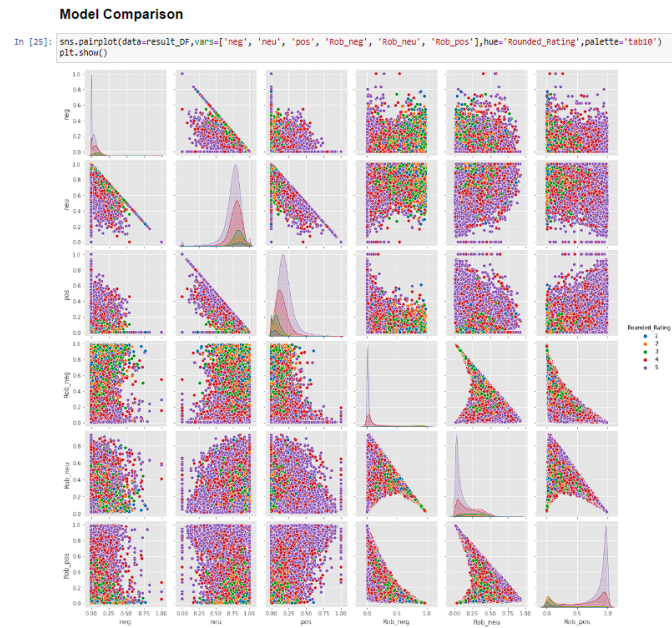


Figure 12: RoBERTa pretrained model Vs Vader Sentiment Scoring

4. Neural Networks

Emotion Detection Using Neural Network and TensorFlow - Model Creation

```
In [58]: import tensorflow as tf
import keras

from nltk.stem.porter import PorterStemmer
from nltk.corpus import stopwords

In [59]: dataset_NeuralNet = result_df[['cleaned_Reviews','Rob_Emotion']]
dataset_NeuralNet

Out[59]:
```

	cleaned_Reviews	Rob_Emotion
0	its been a great delivery vehicle for my cafe...	positive
1	bought this car as a commuter vehicle for a v...	negative
2	this van rocks its the best lots of vroom i ...	positive
3	great work vehicle drives nice has lots of ro...	positive
4	good solid frame and suspension well equippe...	negative
...
212847	i have owned the tiuan for a year and month...	positive
212848	now had months with suv nice suv handles we...	negative
212849	smaller dimensions and my driving experience ...	positive
212850	we have had our tiuan for a month prior to ...	positive
212851	the tiuan a can be had for a reasonable and...	positive
212852 rows x 2 columns		

```
In [60]: dataset_NeuralNet.isna().any(axis=1).sum()
Out[60]: 0

In [61]: ps = PorterStemmer()

def dataPreprocessing(text):
    text = text.split()

    #stemming
    text = [ps.stem(word) for word in text if not word in stopwords.words('english')]
    return " ".join(text)
```

Figure 13: Data Preprocessing for Neural network

```

In [62]: dataset_NeuralNet['cleaned_Reviews'] = dataset_NeuralNet['cleaned_Reviews'].apply(lambda x: dataPreprocessing(x))

In [63]: from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
dataset_NeuralNet['N_Label'] = label_encoder.fit_transform(dataset_NeuralNet['Rob_Emotion'])
dataset_NeuralNet

Out[63]:
   cleaned_Reviews  Rob_Emotion  N_Label
0  great deliver vehicl safe busi good power eco...    positive      2
1  bought car commut vehicl van good disappoint e...    negative      0
2  van rock boat lot room cant lat design nlike bu...    positive      2
3  great work vehicl drive nice lat room east tram...    positive      2
4  good solid frame supports well equip full power...    negative      0
...
212547  own ligan year month bought ligan want smal...    positive      2
212548  month suv nice suv handl well accoder well ama...    negative      0
212549  smaller dream drive experi put ligan class e...    positive      2
212550  ligan month prior trouble free passad overal e...    positive      2
212551  ligan reason get pay averag around superior d...    positive      2
212552 rows x 3 columns

In [64]: from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=5000, ngram_range=(1,3))
dataset_NeuralNet_CV = cv.fit_transform(dataset_NeuralNet['cleaned_Reviews']).toarray()

In [65]: dataset_NeuralNet_CV

Out[65]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 1, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]], dtype=int64)

In [66]: x_train,x_test,y_train,y_test = train_test_split(dataset_NeuralNet_CV,dataset_NeuralNet['N_Label'],test_size=0.36,random_state=42)

In [67]: from tensorflow.keras.layers import Dense
from tensorflow.keras.models import Sequential

In [68]:
model = Sequential()
model.add(Dense(12, input_shape=(x_train.shape[1],), activation='relu'))
model.add(Dense(8, activation='relu'))
model.add(Dense(8, activation='softmax'))
model.compile(loss='sparse_categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
model.fit(x_train,y_train,epochs=10,batch_size=10)

Epoch 1/10
14000/14000 [=====] - 18s 18ms/step - loss: 0.3761 - accuracy: 0.8620
Epoch 2/10
14000/14000 [=====] - 12s 816us/step - loss: 0.3071 - accuracy: 0.8836

```

Figure 14: Layers Creation for Neural network

Roberta Pretrained Model (A Deep Learning Model)- Based on the meaning and context of the sentence

```

In [18]: from transformers import AutoTokenizer
from transformers import AutoModelForSequenceClassification
from scipy.special import softmax

In [19]: ## Transfer Learning
pretrainedModel = "cardiffnlp/twitter-roberta-base-sentiment"
tokenizer = AutoTokenizer.from_pretrained(pretrainedModel)
model_2=AutoModelForSequenceClassification.from_pretrained(pretrainedModel)

In [20]: def scores_RobertaModel(text):
    encoded_text = tokenizer(text,return_tensors='pt')
    output = model_2(**encoded_text)
    scores = output[0][0].detach().numpy()
    scores = softmax(scores)
    scores_dict = {
        'Rob_neg': scores[0],
        'Rob_neu': scores[1],
        'Rob_pos': scores[2],
    }
    return scores_dict

In [21]: results = {}
total_errors=0
for i, row in tqdm(dataset.iterrows(),total = len(dataset)):
    try:
        text = row['cleaned_Reviews']
        data_id = row['id']
        vader_results=sia.polarity_scores(text)
        vader_result_rename = {}
        for key, value in vader_results.items():
            vader_result_rename[f'vader_{key}']=value
        roberta_results = scores_RobertaModel(text)
        both = (vader_results,roberta_results)
        results[data_id]=both
    except RuntimeError:
        total_errors=total_errors+1
        print(f"Runtime Error for id : {data_id}")
    except IndexError:
        total_errors=total_errors+1
        print(f"Index Error for id : {data_id}")

print('Total sentences left out = ',total_errors,"\n Total sentences analyzed successfully : ',len(dataset)-total_errors)

Runtime Error for id : 735
Runtime Error for id : 1558
Runtime Error for id : 2732
Runtime Error for id : 3484
Runtime Error for id : 4235
Index Error for id : 4633
Runtime Error for id : 4800

```

Figure 15: Evaluating the Neural network