

NLP Assignment 3 (Sentiment Analysis)
 Name : Krishna Kant Verma
 Roll No : 2211cs19
 Name : Gourob Chatterjee
 Roll No : 2211cs08

Importing All Required Libraries

```
1 import spacy
2 import re
3 import sys
4 import os
5 import pandas as pd
6 import numpy as np
7 from matplotlib import pyplot as plt
8 from keras.layers import LSTM,Dense,Flatten
9 from keras.utils import to_categorical
10 from sklearn.model_selection import train_test_split
11 from tensorflow import keras
12 from progressbar import progressbar
13 import pandas as pd
14 import numpy as np
```

▼ Zip File Extraction

Extracting the Positive Zip Folder and Negative Zip Folder to POS/NEG SubFolder

```
1 import shutil
2 pos_zip = "/content/pos.zip"
3 neg_zip = "/content/neg.zip"
4 shutil.unpack_archive(pos_zip, "/content/POS", "zip")
5 shutil.unpack_archive(neg_zip, "/content/NEG", "zip")
```

Import OS Library and Read File

```
1 import os
2 def read_file(file_location):
3     with open(file_location, 'r', encoding='utf-8') as f:
4         return f.readlines()
```

Reading Text Data and Processing it Make it Operable

```
1 def read_text(folder_locations, max_files_toread):
2     textData = []
3     textFileLocations = []
4     if not isinstance(folder_locations, str):
5         for location in folder_locations:
6             textFileLocations.append([os.path.join(location, file_name) for file_name in os.listdir(location)][0:max_files_toread])
7     else:
8         print("Folder locations should be in list or tuple format.\nExample: [Folder loc1, Folder loc2, Folder loc3]")
9         return
10
11     classNames = [0 for _ in folder_locations]
12     currLen = 0
13
14     for location_ind, text_file_location in enumerate(textFileLocations):
15         for text_file in text_file_location:
16             try:
17                 textData.append(read_file(text_file))
18             except Exception as e:
19                 print(e, text_file)
20         classNames[location_ind] = len(textData) - currLen
21         currLen = len(textData)
22
```

```

23 classLabels = []
24 for index, no_of_files_inclass in enumerate(classNames):
25     classLabels += [index for number in range(no_of_files_inclass)]
26     return textData, classLabels
27
28 textData, classLabels = read_text(["/content/NEG/neg", "/content/POS/pos"], 5000)

```

Printing textData 1 output (with labels)

```
1 textData[1],classLabels[1]
```

```

(["Who did the research for this film? It's set in Baghdad in 2004, however all the Soldiers are wearing ACUs and have
all Universal Camouflage Pattern gear. No one was wearing that stuff in 04. <br /><br />I just saw this film while
deployed overseas and I can say that the overwhelming feeling from the audience was WTF? This movie made no sense, had
characters come and go with no explanation, and people doing ridiculous things that would NEVER happen in real life. I
realize that it's a movie, but it's obviously trying to portray something realistic. It fails miserably, but it's
trying. <br /><br />It's like someone came up with a bunch of random ideas, chewed them up and swallowed, then vomited
out a film. I would not recommend this film to anyone. I'm still not sure why I sat through the whole thing. GI Joe
was one that really made you think compared to this. STAY AWAY!"],
0)

```

```
1 textData[0]
```

```

["From the creators of Shrek\x85\x85\x85\x85.. OK, that grabbed my attention.<br /><br />Well the creators of Shrek
also made Madagascar. Madagascar was half as good as Shrek.<br /><br />And now Flushed Away is half as good as
Madagascar.<br /><br />That means Flushed Away isn't good. The animation and all that special effects were extremely
good but the movie wasn't.<br /><br />The story of this movie was only meant for kids. It's seriously not possible for
adults to actually love this flick.<br /><br />But there were many jokes meant for adults. I bet kids dint understand
the jokes.<br /><br />Despite that I dint like this flick.<br /><br />I am completely disappointed. 4/10"]

```

Eliminating all unnecessary HTML and other hyper Texts

```

1 regex = re.compile(r'<[^>]+>')
2 def removeHyperTags(string):
3     return regex.sub(' ', string)

```

calling removeHyperText Function

```
1 textData = [removeHyperTags(text[0]) for text in textData]
```

Text After Processing all the removal of hypertexts

```
1 textData[0]
```

```

'From the creators of Shrek\x85\x85\x85\x85.. OK, that grabbed my attention.
Well the creators of Shrek also made Madagascar. Madagascar was half as good
as Shrek. And now Flushed Away is half as good as Madagascar. That means F
lushed Away isn't good. The animation and all that special effects were extr
emely good but the movie wasn't. The story of this movie was only meant for

```

▼ Tokenizing Texts

Function to Tokenize text

```

1 def tokenize(texts):
2     for text_ind,text in enumerate(texts):
3         texts[text_ind]=text.lower()
4     nlp = spacy.load('en_core_web_sm')
5     tokenizedTexts = []
6     for ind,text in enumerate(texts):
7         print(ind,len(texts))
8         doc = nlp(text)
9         tokens = [token.text for token in doc]
10        tokenizedTexts.append(tokens)
11    return tokenizedTexts

```

Calling Tokenize Text Function over First Dataset

```
1 tokenizedTexts = tokenize(textData)
```

```
9942 10000
9943 10000
9944 10000
9945 10000
9946 10000
9947 10000
9948 10000
9949 10000
9950 10000
9951 10000
9952 10000
9953 10000
9954 10000
9955 10000
9956 10000
9957 10000
9958 10000
9959 10000
9960 10000
9961 10000
9962 10000
9963 10000
9964 10000
9965 10000
9966 10000
9967 10000
9968 10000
9969 10000
9970 10000
9971 10000
9972 10000
9973 10000
9974 10000
9975 10000
9976 10000
9977 10000
9978 10000
9979 10000
9980 10000
9981 10000
9982 10000
9983 10000
9984 10000
9985 10000
9986 10000
9987 10000
9988 10000
9989 10000
9990 10000
9991 10000
9992 10000
9993 10000
9994 10000
9995 10000
9996 10000
9997 10000
9998 10000
9999 10000
```

▼ Finding Out Most Frequent Tokens In DataSet

Function to find Out most frequent tokens

```
1 def maxFrequentTokens( textData, frequency=5 ):
2     counter={}
3     for text in textData:
4         for token in text:
5             if token in counter:
6                 counter[token]+=1
7             else:
8                 counter[token]=1
9     FrequentTokens=[]
10    for i in counter.items():
```

```

11     if i[1]>frequency:
12         FrequentTokens.append(i[0])
13     return FrequentTokens
14
15 FrequentTokens = maxFrequentTokens(tokenizedTexts, frequency=300)
16
17 f" Most frequent Tokens length  {len(FrequentTokens)}"

' Most frequent Tokens length  814'

```

```

1 tokenizedTexts[1][:10]

['who', 'did', 'the', 'research', 'for', 'this', 'film', '?', 'it', "s"]

```

▼ Finding Out Padding Sequences

```

1 def padSeq(sequences):
2     averageLen=0
3     for text in sequences:
4         averageLen+=len(text)
5     averageLen=int(averageLen/len(tokenizedTexts))
6     for text_ind,text in enumerate(sequences):
7         if len(text)>=averageLen:
8             sequences[text_ind]=text[:averageLen]
9         else:
10             for i in range(averageLen-len(text)):
11                 sequences[text_ind].append('PAD')
12

```

```

1 padSeq(tokenizedTexts)

```

▼ Encoding Words to One-Hot-Encoding

```

1 def oneHotEncoding(tokenizedTexts):
2     token_to_number={x:to_categorical(ind,num_classes=len(FrequentTokens)+1,dtype='uint8') for ind,x in enumerate(FrequentTokens)}
3     token_to_number['unk']=np.array([0 for x in range(len(FrequentTokens))]+[1])
4     for text in tokenizedTexts:
5         for token_ind,token in enumerate(text):
6             if token in token_to_number:
7                 text[token_ind]=token_to_number[token]
8             else:
9                 text[token_ind]=token_to_number['unk']
10 oneHotEncoding(tokenizedTexts)

```

```

1 X_array = np.array(tokenizedTexts,dtype='uint8')

```

```

1 dataSet_1 = X_array

```

▼ Second Dataset (SemEval Tweet Dataset)

```

1 data_2 = pd.read_csv("/content/2013semeval_train.csv")

```

```

1 data_2.tweet

```

```

0      Gas by my house hit $3.39!!!! I\u2019m going t...
1      Theo Walcott is still shit\u002c watch Rafa an...
2      its not that I\u2019m a GSP fan\u002c i just h...
3      Iranian general says Israel\u2019s Iron Dome c...
4      Tehran\u002c Mon Amour: Obama Tried to Establi...
...
9679   RT @MNFootNg It's monday and Monday Night Foot...
9680   All I know is the road for that Lomardi start ...
9681   "All Blue and White fam, we r meeting at Golde...
9682   @DariusButler28 Have a great game against Tam...
9683   "I'm pisseedddd that I missed Kid Cudi's show...
Name: tweet, Length: 9684, dtype: object

```

▼ unicode Cleaning

```
1 uniEscRegx = re.compile(r'\\u([0-9a-fA-F]{4})')
2 def convert_escape_sequence(match):
3     return chr(int(match.group(1), 16))
```

```
1 result = uniEscRegx.sub(convert_escape_sequence, data_2.tweet[1])
2 def changeString(string):
3     return uniEscRegx.sub(convert_escape_sequence,string)
```

```
1 print(result)
```

Theo Walcott is still shit, watch Rafa and Johnny deal with him on Saturday.

```
1 data_2.tweet = data_2.tweet.apply(changeString)
```

▼ Tokenizing Second Tweet DataSet

```
1 tokenizedTexts = tokenize( data_2.tweet )
```

```
9626 9684
9627 9684
9628 9684
9629 9684
9630 9684
9631 9684
9632 9684
9633 9684
9634 9684
9635 9684
9636 9684
9637 9684
9638 9684
9639 9684
9640 9684
9641 9684
9642 9684
9643 9684
9644 9684
9645 9684
9646 9684
9647 9684
9648 9684
9649 9684
9650 9684
9651 9684
9652 9684
9653 9684
9654 9684
9655 9684
9656 9684
9657 9684
9658 9684
9659 9684
9660 9684
9661 9684
9662 9684
9663 9684
9664 9684
9665 9684
9666 9684
9667 9684
9668 9684
9669 9684
9670 9684
9671 9684
9672 9684
9673 9684
9674 9684
9675 9684
9676 9684
9677 9684
9678 9684
9679 9684
9680 9684
9681 9684
9682 9684
9683 9684
```

▼ Finding Out Most Freuent Tokens

```
1 FrequentTokens = maxFrequentTokens(tokenizedTexts,frequency=200)
2 f" Most frequent Tokens length {len(FrequentTokens)}"

' Most frequent Tokens length 138'
```

▼ Conversion to One Hot Encoding after Padding Sequences

```
1 padSeq(tokenizedTexts)
2 oneHotEncoding(tokenizedTexts)
```

```
1 X_array = np.array(tokenizedTexts,dtype='uint8')
```

```
1 dataSet_2 = X_array
```

```
1 dataSet_1.shape,dataSet_2.shape

((10000, 273, 815), (9684, 23, 139))
```

```
# converting first dataset labels to numpy array
```

```
1 dataSet_1_labels=np.array(classLabels)
```

```
# converting second dataset labels to numpy array
```

```
1 uniqLabels={x:ind for ind,x in enumerate(data_2.label.unique())}
```

```
1 data_2.label = data_2.label.apply(lambda x:uniqLabels[x])
```

```
1 dataSet_2_labels = data_2.label.values
```

▼ Finalizing Our DataSet (Shaping Our Dataset)

```
1 dataSet_1.shape
2 dataSet_1_labels.shape
3 dataSet_2.shape
4 dataSet_2_labels.shape

(9684,)
```

▼ Creating Model Using RNN

```
Creating Model For First Dataset
```

```
1 model_1_RNN=keras.Sequential([
2     LSTM(256,input_shape=dataSet_1.shape[1:]),
3     Dense(2,activation='sigmoid')
4 ])
```

```
1 model_1_RNN.compile(loss='sparse_categorical_crossentropy',optimizer='adam',metrics=['accuracy'])
```

```
# Creating Model For Second Dataset
```

```
1 model_2_RNN=keras.Sequential([
2     LSTM(256,input_shape=dataSet_2.shape[1:]),
```

```
3 Dense(3,activation='sigmoid')
4 ])
```

```
1 model_2_RNN.compile(loss='sparse_categorical_crossentropy',optimizer='adam',metrics=['accuracy'])
```

```
1 from sklearn.model_selection import train_test_split
2 print(len(dataSet_1_labels))
3
4 xTrain_1, xTest_1, yTrain_1, yTest_1 = train_test_split(dataSet_1, dataSet_1_labels, test_size=0.25, random_state=42)
5 xTrain_2, xTest_2, yTrain_2, yTest_2 = train_test_split(dataSet_2, dataSet_2_labels, test_size=0.25, random_state=42)

10000
```

▼ Excution of Model 1

```
1 execModel_1 = model_1_RNN.fit(xTrain_1,yTrain_1,epochs=10)
```

```
Epoch 1/10
235/235 [=====] - 15s 27ms/step - loss: 0.6940 - accuracy: 0.5060
Epoch 2/10
235/235 [=====] - 6s 27ms/step - loss: 0.6943 - accuracy: 0.5444
Epoch 3/10
235/235 [=====] - 6s 27ms/step - loss: 0.7051 - accuracy: 0.4980
Epoch 4/10
235/235 [=====] - 6s 27ms/step - loss: 0.6981 - accuracy: 0.5063
Epoch 5/10
235/235 [=====] - 6s 27ms/step - loss: 0.6948 - accuracy: 0.5151
Epoch 6/10
235/235 [=====] - 6s 27ms/step - loss: 0.6924 - accuracy: 0.5169
Epoch 7/10
235/235 [=====] - 7s 28ms/step - loss: 0.6891 - accuracy: 0.5331
Epoch 8/10
235/235 [=====] - 6s 27ms/step - loss: 0.6880 - accuracy: 0.5333
Epoch 9/10
235/235 [=====] - 7s 29ms/step - loss: 0.6819 - accuracy: 0.5409
Epoch 10/10
235/235 [=====] - 6s 28ms/step - loss: 0.6768 - accuracy: 0.5481
```

▼ Excution of Model 2

```
1 execModel_2 = model_2_RNN.fit(xTrain_2,yTrain_2,epochs=10)
```

```
Epoch 1/10
227/227 [=====] - 7s 5ms/step - loss: 0.9602 - accuracy: 0.5372
Epoch 2/10
227/227 [=====] - 1s 4ms/step - loss: 0.8888 - accuracy: 0.5918
Epoch 3/10
227/227 [=====] - 1s 4ms/step - loss: 0.8743 - accuracy: 0.5974
Epoch 4/10
227/227 [=====] - 1s 4ms/step - loss: 0.8643 - accuracy: 0.6042
Epoch 5/10
227/227 [=====] - 1s 5ms/step - loss: 0.8502 - accuracy: 0.6142
Epoch 6/10
227/227 [=====] - 1s 4ms/step - loss: 0.8355 - accuracy: 0.6165
Epoch 7/10
227/227 [=====] - 1s 4ms/step - loss: 0.8307 - accuracy: 0.6188
Epoch 8/10
227/227 [=====] - 1s 4ms/step - loss: 0.8131 - accuracy: 0.6292
Epoch 9/10
227/227 [=====] - 1s 4ms/step - loss: 0.8084 - accuracy: 0.6329
Epoch 10/10
227/227 [=====] - 1s 5ms/step - loss: 0.7920 - accuracy: 0.6468
```

▼ Calculating Precision Accuracy and Recall for First Dataset

```
1 from sklearn.metrics import precision_recall_fscore_support
2 y_pred_1=model_1_RNN.predict(xTest_1)
3 scores=precision_recall_fscore_support(np.argmax(y_pred_1,axis=1),yTest_1,average='macro')
4 print(f"""\n\nRNN Model_1 For DataSet1
5     Precision = {scores[0]}
6     Recall    = {scores[1]}
7     f1_score  = {scores[2]}
8 """)
```

79/79 [=====] - 1s 13ms/step

```
RNN Model_1 For DataSet1
Precision = 0.497400794870219
Recall    = 0.4802216538789429
f1_score  = 0.3542091504669417
```

▼ Calculating Precision Accuracy and Recall for Second Dataset

```
1 from sklearn.metrics import precision_recall_fscore_support
2 y_pred_2=model_2_RNN.predict(xTest_2)
3 scores=precision_recall_fscore_support(np.argmax(y_pred_2,axis=1),yTest_2,average='macro')
4 print(f"""\n\nRNN Model_2 For DataSet_2
5     Precision = {scores[0]}
6     Recall    = {scores[1]}
7     f1_score  = {scores[2]}
8 """)
```

76/76 [=====] - 1s 3ms/step

```
RNN Model_2 For DataSet_2
Precision = 0.4864269280775426
Recall    = 0.545716590829335
f1_score  = 0.48656399547127904
```

▼ Logs after Each Epocs

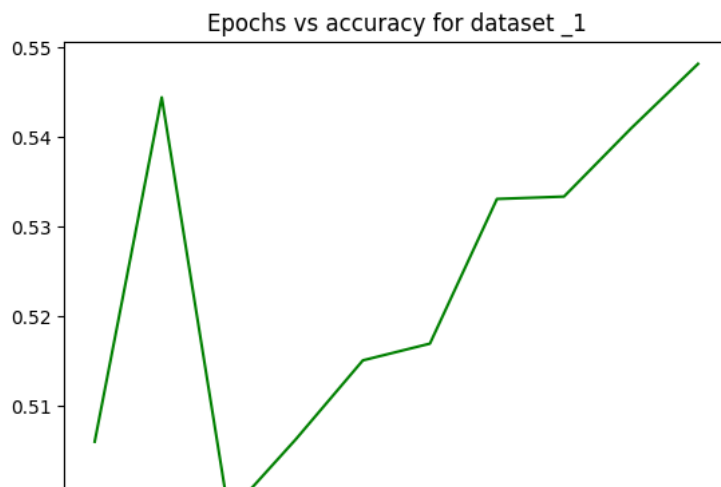
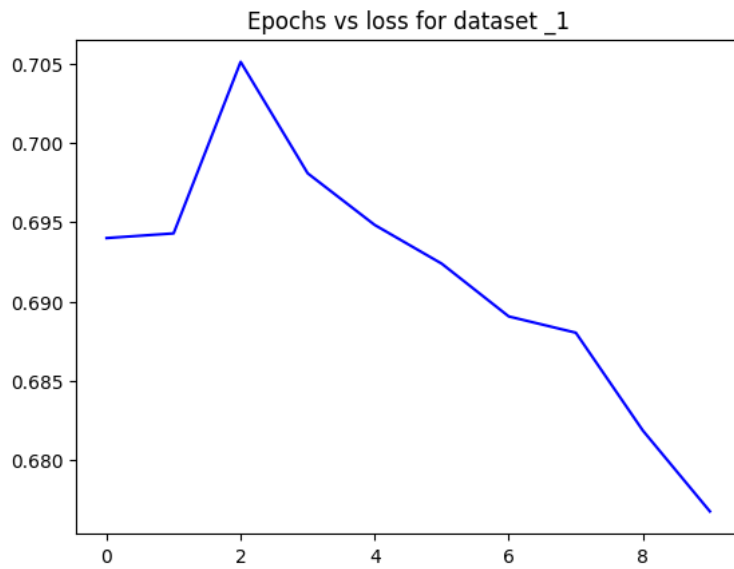
```
1 execModel_1.history

{'loss': [0.6939970254898071,
0.6942936182022095,
0.7051113247871399,
0.6980825662612915,
0.6948240995407104,
0.6923822164535522,
0.6890540719032288,
0.6880289316177368,
0.6818565726280212,
0.6767584681510925],
'accuracy': [0.5059999823570251,
0.5443999767303467,
0.49799999594688416,
0.5062666535377502,
0.5150666832923889,
0.5169333219528198,
0.5330666899681091,
0.5333333611488342,
0.5409333109855652,
0.5481333136558533]}
```

▼ Plots

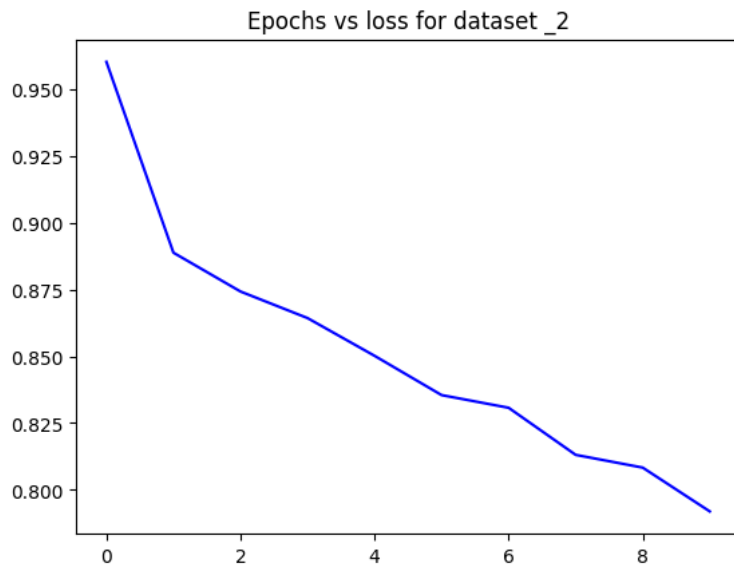
plotting epocs Vs Loss For Model 1

```
1 plt.plot([x for x in range(10)],execModel_1.history['loss'],c='blue')
2 plt.title("Epochs vs loss for dataset _1")
3 plt.xlabel="Epochs"
4 plt.ylabel="Loss"
5 plt.show()
6 plt.plot([x for x in range(10)],execModel_1.history['accuracy'],c='g')
7 plt.title("Epochs vs accuracy for dataset _1")
8 plt.xlabel="Epochs"
9 plt.ylabel="Loss"
10 plt.show()
```

Plotting Epocs Vs Loss for Model 2

```
1 plt.plot([x for x in range(10)],execModel_2.history['loss'],c='blue')
2 plt.title("Epochs vs loss for dataset _2")
3 plt.xlabel="Epochs"
4 plt.ylabel="Loss"
5 plt.show()
6 plt.plot([x for x in range(10)],execModel_2.history['accuracy'],c='green')
7 plt.title("Epochs vs accuracy for dataset _2")
8 plt.xlabel="Epochs"
9 plt.ylabel="Loss"
10 plt.show()
```



▼ Feed forward Neural Networks Architecture

0.044 |

✓

|

My proposed feed-forward neural network (FFNN) architecture consists of two hidden layers connected by non-linear activation functions. The first hidden layer, which has 256 neurons, is coupled to the input layer. The second hidden layer, which includes 128 neurons, is connected to the first hidden layer. The output layer, whose size is determined by the number of classes in the problem, is then connected to the second hidden layer.

Here is a visual representation of the architecture:

" Input Layer (Size depends on input data) --> Hidden Layer 1 (Size: 256) with Non-linearity --> Hidden Layer 2 (Size: 128) with Non-linearity --> Output Layer (Size depends on the number of classes) "

There are several prominent options for non-linearity in this architecture, including the Rectified Linear Unit (ReLU), the hyperbolic tangent (tanh), and the Gaussian Error Linear Unit (GELU).

The binary cross-entropy loss function is frequently used for binary classification issues. The categorical cross-entropy loss function is frequently used for multi-class classification issues.

Feed-Forward There are various uses for neural networks:

They are able to simulate intricate non-linear input-output interactions.
 They are relatively easy to comprehend and use.
 They can be trained well on sizable datasets using methods like backpropagation and stochastic gradient descent.
 They can be used to solve many different types of issues, like as classification, regression, and prediction.

However, because FFNNs do not account for the temporal correlations between inputs, they are only partially capable of handling sequential or time-series data. Additionally, if the dataset is too little or the model is too complicated, they are more likely to overfit.

Thanking You So Much