Nefisa hassen

D213 task 2

**Part I: Research Question.**

**A1. Summarize one research question that you will answer using neural network models and NLP techniques**

Can we analyze customer reviews of a product and predict whether a review is positive or negative based on its text using a neural network?

**A2. Define the objectives or goals of the data analysis.**

The objective of this analysis is to develop a neural network model capable of predicting whether a review is positive or negative based on its text.

**A3. Identify a type of neural network capable of performing a text classification task that can be trained to produce useful predictions on text sequences on the selected data set.**

To tackle this task, I'll use a Bidirectional LSTM network. This type of neural network works well with text because it reads the sequence of words in both directions, helping it understand the context and sentiment more effectively. By training the model on labeled data, such as positive and negative reviews, it can pick up patterns in the text that help it predict the sentiment behind each review.

**Part II: Data Preparation**

B. Summarize the data cleaning process by doing the following:

1. Perform exploratory data analysis on the chosen data set, and include an explanation of *each* of the following elements:
   • presence of unusual characters (e.g., emojis, non-English characters)
   • vocabulary size
   • proposed word embedding length
   • statistical justification for the chosen maximum sequence length

```
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
import pandas as pd
import numpy as np
import re
import seaborn as sns
import csv
```

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, LSTM, Dense
from tensorflow.keras.optimizers import Adam

from google.colab import files
uploaded = files.upload()

df1 = pd.read_csv('amazon_cells_labelled.txt', sep='\t', header=None)

df2 = pd.read_csv('imdb_labelled.txt', sep='\t', header=None)

df3 = pd.read_csv('yelp_labelled.txt', sep='\t', header=None)

df = pd.concat([df1, df2, df3])

df.columns = ['review', 'sentiment']

display(df)
```

| | review | sentiment |
|---|---|---|
| 0 | So there is no way for me to plug it in here i... | 0 |
| 1 | Good case, Excellent value. | 1 |
| 2 | Great for the jawbone. | 1 |
| 3 | Tied to charger for conversations lasting more... | 0 |
| 4 | The mic is great. | 1 |
| ... | ... | ... |
| 995 | I think food should have flavor and texture an... | 0 |
| 996 | Appetite instantly gone. | 0 |
| 997 | Overall I was not impressed and would not go b... | 0 |

| | | |
|---|---|---|
| 99 8 | The whole experience was underwhelming, and I ... | 0 |
| 99 9 | Then, as if I hadn't wasted enough of my life ... | 0 |

2748 rows × 2 columns

Df.shape
(2748, 2)

df.head()

```
#count of positive  and negative  reviews
df.sentiment.value_counts()
```

|  | count |
|---|---|
| sentiment | |
| 1 | 1386 |
| 0 | 1362 |

dtype: int64

```
import matplotlib.pyplot as plt
sns.countplot(x='sentiment', data=df)
plt.title('Sentiment Distribution')
plt.xlabel('Sentiment')
plt.ylabel('Count')
```

Sentiment Distribution

df.isna().sum()

0

review        0

sentiment     0

dtype: int64

df.dropna(inplace=True)

df.isnull().sum()

0

review        0

sentiment     0

dtype: int64

**Part II: Data Preparation**

**B1. Perform exploratory data analysis on the chosen data set, and include an explanation of each of the following elements:**

- **presence of unusual characters (e.g., emojis, non-English characters)**

- **vocabulary size**
- **proposed word embedding length**
- **Statistical justification for the chosen maximum sequence length**

---

[ ]

```python
#identify and display all the unique characters present in the 'review' column

list_of_chars = set(''.join(df['review'].astype(str)))

print(list(list_of_chars))
```

['3', 'F', '"', 'ê', '(', 'b', '-', 'd', '%', '4', 'm', '\x85', 'R', 'O', '*', ']', 'K', '0', '?', '+', 'B', '\n', 'I', '9', "''", 'å', 'c', 'X', 'u', '\t', 'é', 'r', 'G', 'A', 'i', 'Q', 'J', 'z', 't', 'x', 'k', ' ', '2', 'L', 'Y', '5', 'C', 'j', '/', 'M', '1', '#', 'N', 'V', '7', 'w', '\x96', ',', 'l', '.', 'e', 'Z', 'g', 'U', 'p', 'D', ')', 'v', 'E', '6', '\x97', 'y', '!', ';', 'a', '&', 'P', 'q', 'o', '$', 'W', 'T', 'f', 'n', ':', 'S', 'H', 's', '[', '8', 'h']

```python
#remove empty rows

df = df[df['review'].str.strip().astype(bool)]

#  no emojis in the in data review dataset
emojis = [re.findall(r'[^\w\s,]', review) for review in df['review']]
print(emojis)


no non english words in the data
non_english = [re.findall(r'[^\x00-\x7F]+', review) for review in df['review']]
print(non_english)


# Initialize stopwords
stop_words = set(stopwords.words('english'))


# Preprocessing function
def preprocess_text(description):
    # Remove punctuation but keep spaces
    description = re.sub("[^a-zA-Z\s]", "", description)
    # Convert to lowercase
    description = description.lower()
    # Perform tokenization
```

```python
    description = nltk.word_tokenize(description)
    # Perform lemmatization
    lemma = nltk.WordNetLemmatizer()
    description = [lemma.lemmatize(word) for word in description]
    # Remove stopwords
    description = [word for word in description if word not in stop_words]
    return ' '.join(description)


# Apply preprocessing to all reviews
description_list = [preprocess_text(desc) for desc in df['review']]

# Tokenization and vectorization
vocab_size = 10000
tokenizer = Tokenizer(num_words=vocab_size, oov_token="<OOV>")
tokenizer.fit_on_texts(description_list)

# Convert texts to sequences
sequences = tokenizer.texts_to_sequences(description_list)

print("Sample tokenized sequence:", sequences[0])
print("Word index:", tokenizer.word_index)
```
Sample tokenized sequence: [47, 256, 107, 536, 27, 1989]

Word index: {'<OOV>': 1, 'wa': 2, 'good': 3, 'movie': 4, 'great': 5, 'film': 6, 'phone': 7, 'one': 8, 'time': 9, 'like': 10, 'food': 11, 'place': 12, 'work': 13, 'service': 14, 'really': 15, 'bad': 16, 'ha': 17, 'well': 18, 'dont': 19, 'would': 20, 'best': 21, 'even': 22, 'ever': 23, 'also': 24, 'back': 25, 'get': 26, 'go': 27, 'quality': 28, 'love': 29, 'make': 30, 'ive': 31, 'made': 32, 'character': 33, 'product': 34, 'im': 35, 'headset': 36, 'could': 37, 'nice': 38, 'thing': 39, 'better': 40, 'excellent': 41, 'sound': 42, 'never': 43, 'recommend': 44, 'much': 45, 'use': 46, 'way': 47, 'battery': 48, 'think': 49, 'first': 50,]

**Only showing 50 values here but i have the full values on notebook**

**B2. Describe the goals of the tokenization process, including any code generated and packages that are used to normalize text during the tokenization process.**

The goal of Tokenization is to break raw text into smaller pieces, like words or subwords, making it easier to work with in machine learning. It standardizes the text by splitting it into

tokens and handles out-of-vocabulary (OOV) words by assigning them a special token, so the model can deal with new or unseen words without issues. The text is then converted into sequences of integers, where each number represents a word's position in the vocabulary.

```python
# Vocabulary size
vocab_size = len(tokenizer.word_index) + 1
print("Vocabulary size:", vocab_size)
Vocabulary size: 4757


#determine the min max and  length of reviews
review_length = []
for char_length in df['review']:
  review_length.append(len(char_length))
  #print(review_length)


max_review_length = max(review_length)
min_review_length = min(review_length)
avg_review_length = sum(review_length)/len(review_length)
print("Max review length:",max_review_length)
print("Min review length:",min_review_length)
print("Avg review length:",avg_review_length)
Max review length: 7944
Min review length: 7
Avg review length: 71.52838427947599


#split the data into train and test
import numpy as np
from sklearn.model_selection import train_test_split


X = np.array(description_list)
y = df.sentiment.values


# Split the data into training and testing sets
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=15,
stratify=y)
y_train = pd.Series(y_train)
y_test = pd.Series(y_test)

print("X_train size:",X_train.shape)
print("X_tests size:",X_test.shape)
print("y_train size:",y_train.shape)
print("y_test size:",y_test.shape)
```

X_train size: (2198,)
X_tests size: (550,)
y_train size: (2198,)
y_test size: (550,)

```
# Define max_length, padding_type, and trunc_type
max_length = int(np.percentile(review_length, 95))
padding_type = 'post'
trunc_type = 'post'
```

**B3. Explain the padding process used to standardize the length of sequences. Include the following in your explanation:**
- **if the padding occurs before or after the text sequence**
- **a screenshot of a single padded sequence**

Padding helps ensure all sequences in the dataset have the same length, making them suitable for input into neural network models like LSTMs, which require consistent input shapes. For this task, padding is used to extend shorter sequences to the desired length without changing their meaning, so the model can process the data effectively. Below is the padded result.

```
#Apply Padding to training data
sequences_train = tokenizer.texts_to_sequences(X_train)
padded_train = pad_sequences(sequences_train, maxlen= max_length, padding
=padding_type,truncating=trunc_type)
array([[  2, 219, 282, ...,   0,   0,   0],
       [ 43, 224, 636, ...,   0,   0,   0],
       [394, 515, 150, ...,   0,   0,   0],
```

```
      ...,
    [ 600, 1536,   32, ...,    0,   0,   0],
    [2413,  984,  622, ...,    0,   0,   0],
    [  24,    5,   53, ...,    0,   0,   0]], dtype=int32)
```

#apply padding to test data
sequences_test = tokenizer.texts_to_sequences(X_test)
padded_test = pad_sequences(sequences_test, maxlen= max_length, padding
=padding_type,truncating=trunc_type)
```
array([[  88,   53,    4, ...,    0,   0,   0],
    [1486, 2988,  246, ...,    0,   0,   0],
    [1121,  601, 2806, ...,    0,   0,   0],

      ...,
    [ 273, 3812, 1496, ...,    0,   0,   0],
    [1732,   53,    3, ...,    0,   0,   0],
    [  62,  480,   23, ...,    0,   0,   0]], dtype=int32)
```
from sys import maxsize
#Display the padded sequence
np.set_printoptions(threshold=maxsize)
# Show first 50 tokens
print(padded_train[1][:50])
```
[  43  224  636 4122   43   23 1127 4123    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    0    0    0    0    0    0
    0    0    0    0    0    0]
```
#convert padded data to numpy array to be used in model
padded_train = np.array(padded_train)
padded_test = np.array(padded_test)
train_labels = np.array(y_train)
test_labels = np.array(y_test)

**B4.  Identify how many categories of sentiment will be used and an activation function for the final dense layer of the network.**

For this model, I use sigmoid activation in the output layer because it's perfect for binary classification tasks, like distinguishing between Positive and Negative reviews, by producing outputs between 0 and 1. The binary_crossentropy loss function is chosen because it works well for binary classification, measuring the difference between predicted probabilities and the actual class labels. Adam is selected as the optimizer since it's efficient and adapts the learning rate, making it suitable for text classification tasks. I  set num_epochs = 25 to allow the model enough time to learn without overfitting, while monitoring performance to adjust if needed.

**B5.  Explain the steps used to prepare the data for analysis, including the size of the training, validation, and test set split (based on the industry average).**

To get the data ready for analysis, I first cleaned it up by removing any rows with missing values to ensure everything was complete and reliable. Then, I split the data into 80% for training and 20% for testing to make sure the model gets enough data to learn while still having some set aside to check how well it performs on new, unseen data.

B6.a copy of the prepared data set

```python
#export the data to CSV file
pd.DataFrame(padded_train).to_csv('padded_train.csv', index=False)
pd.DataFrame(padded_test).to_csv('padded_test.csv', index=False)
pd.DataFrame(train_labels).to_csv('train_labels.csv', index=False)
pd.DataFrame(test_labels).to_csv('test_labels.csv', index=False)
```

**Part III: Network Architecture**

**C1.  Provide the output of the model summary of the function from TensorFlow.**

```python
#Building The Neural Network Model
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, GlobalAveragePooling1D, Dense, Dropout
from tensorflow.keras.layers import LSTM, Bidirectional
from tensorflow.keras.regularizers import l2
vocab_size= 4757
```

```python
embedding_dim = 32
max_length = 20
trunc_type='post'
padding_type='post'
oov_tok = "<OOV>"
activation = 'sigmoid'
loss = 'binary_crossentropy'
optimizer = Adam(learning_rate=0.0005)
num_epochs = 25
callback=tf.keras.callbacks.EarlyStopping(monitor='val_loss', patience=3)


# Build the model

model = Sequential([
    Embedding(vocab_size, embedding_dim, input_length=max_length),
    Bidirectional(LSTM(32, dropout=0.2, recurrent_dropout=0.2)),
    Dropout(0.5),
    Dense(100, activation='relu', kernel_regularizer=l2(0.01)),
    Dropout(0.5),
    Dense(50, activation='relu'),
    Dense(1, activation='sigmoid')
])
 # Compile the model
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
# Build the model by providing input shape
model.build(input_shape=(None,max_length))
model.summary()
```

**Model: "sequential"**

| Layer (type) | Output Shape | Param # |

| Layer | Output Shape | Param # |
|---|---|---|
| embedding_1 (Embedding) | (None, 20, 32) | 152,224 |
| bidirectional_1 (Bidirectional) | (None, 64) | 16,640 |
| dropout (Dropout) | (None, 64) | 0 |
| dense (Dense) | (None, 100) | 6,500 |
| dropout_1 (Dropout) | (None, 100) | 0 |
| dense_1 (Dense) | (None, 50) | 5,050 |
| dense_2 (Dense) | (None, 1) | 51 |

**Total params:** 180,465 (704.94 KB)

**Trainable params:** 180,465 (704.94 KB)

**Non-trainable params:** 0 (0.00 B)

```
early_stopping_monitor = tf.keras.callbacks.EarlyStopping(
    monitor='val_loss', patience=5, mode='min', restore_best_weights=True
)
```

history = model.fit(padded_train, train_labels, epochs=num_epochs, validation_data=(padded_test, test_labels), callbacks=[early_stopping_monitor])

```
poch 1/25
69/69 ———————————————————— 14s 122ms/step - accuracy: 0.4966
- loss: 1.3059 - val_accuracy: 0.5036 - val_loss: 0.9335
Epoch 2/25
69/69 ———————————————————— 10s 129ms/step - accuracy: 0.5466
- loss: 0.8718 - val_accuracy: 0.7200 - val_loss: 0.7272
Epoch 3/25
69/69 ———————————————————— 9s 118ms/step - accuracy: 0.8094
- loss: 0.5938 - val_accuracy: 0.7873 - val_loss: 0.4922
Epoch 4/25
69/69 ———————————————————— 8s 113ms/step - accuracy: 0.9337
- loss: 0.2550 - val_accuracy: 0.7818 - val_loss: 0.5047
Epoch 5/25
69/69 ———————————————————— 9s 125ms/step - accuracy: 0.9513
- loss: 0.1705 - val_accuracy: 0.7891 - val_loss: 0.5683
Epoch 6/25
69/69 ———————————————————— 7s 103ms/step - accuracy: 0.9705
- loss: 0.1115 - val_accuracy: 0.7964 - val_loss: 0.6117
Epoch 7/25
69/69 ———————————————————— 11s 107ms/step - accuracy: 0.9750
- loss: 0.0848 - val_accuracy: 0.7709 - val_loss: 0.6468
Epoch 8/25
69/69 ———————————————————— 12s 127ms/step - accuracy: 0.9829
- loss: 0.0813 - val_accuracy: 0.7873 - val_loss: 0.6796
```

**C2. Discuss the number of layers, the type of layers, and the total number of parameters.**

- The model begins with an Embedding layer to transform words into vector representations. It is followed by a Bidirectional LSTM layer that captures contextual information from both directions in a sequence. Dropout layers are incorporated to help prevent overfitting, and two Dense layers with ReLU activation are used to refine the extracted features. The final Dense layer utilizes a sigmoid activation to produce probabilities for binary classification. In total, the model has 180,465 parameters, all of which are trainable, with no non-trainable parameters.

**C3 Justify the choice of hyperparameters, including the following elements:**

**• activation functions number of nodes per layer,loss function,optimizer,stopping criteria,evaluation metric**

- The model uses ReLU activation for efficient learning and sigmoid for binary classification. The LSTM has 32 nodes, and the Dense layers (100 and 50 nodes) simplify feature extraction. Binary cross-entropy is the chosen loss function for accurate probability-based error calculation, and Adam optimizer ensures steady and efficient weight updates. Early stopping avoids overfitting by stopping training when validation loss stops improving, and accuracy is used as the evaluation metric to measure performance.

**Part IV: Model Evaluation**

**D1. Discuss the impact of using stopping criteria to include defining the number of epochs**

- Early stopping helps prevent overfitting by stopping the training once the model's performance on the validation data stops improving. This is useful because it saves time and resources by avoiding unnecessary computation and the risk of the model overfitting to the training data. Setting a limit on the number of epochs, like 25 in this case, defines the maximum training duration. Early stopping ensures that the model can stop earlier if it reaches its best performance before hitting that limit, helping the model generalize better without wasting time.

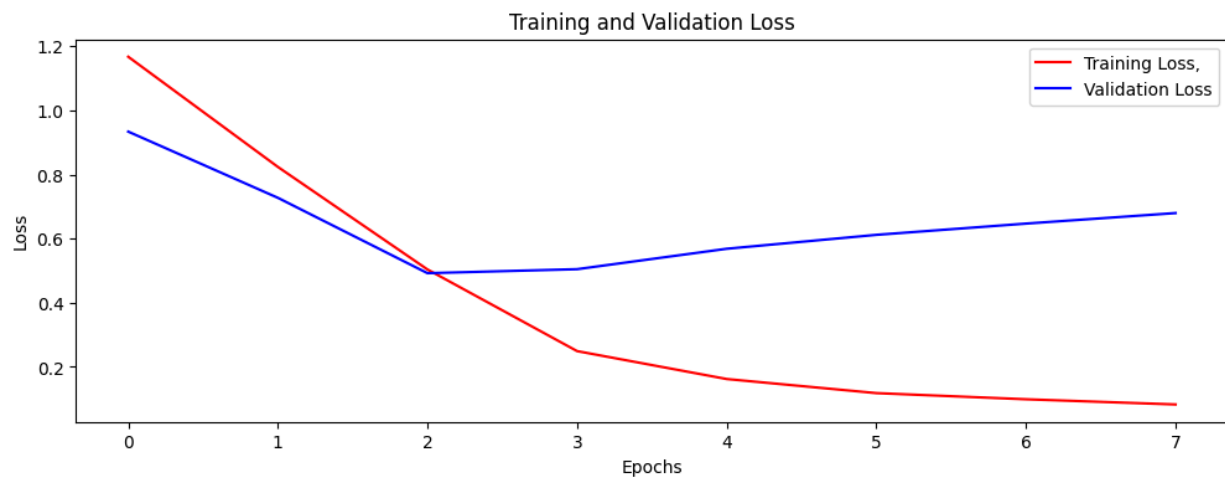**D2. Assess the fitness of the model and any actions taken to address overfitting.**

- The model's fitness is assessed by its performance on both the training and validation data. To prevent overfitting, early stopping is used to halt training when the validation loss stops improving. Dropout layers and L2 regularization are also implemented to reduce reliance on specific neurons and penalize large weights, encouraging generalization. These strategies help the model perform well on unseen data, avoiding overfitting while maintaining efficiency.

**D3. Provide visualizations of the model's training process, including a line graph of the loss and chosen evaluation metric**

```
import matplotlib.pyplot as plt
train_loss = history.history['loss']
val_loss = history.history['val_loss']
```

```python
train_acc = history.history['accuracy']
val_acc = history.history['val_accuracy']

plt.figure(figsize=(12, 4))
plt.plot(train_loss, label='Training Loss,', color = 'red')
plt.plot(val_loss, label='Validation Loss', color = 'blue')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.title('Training and Validation Loss')
plt.show()
```
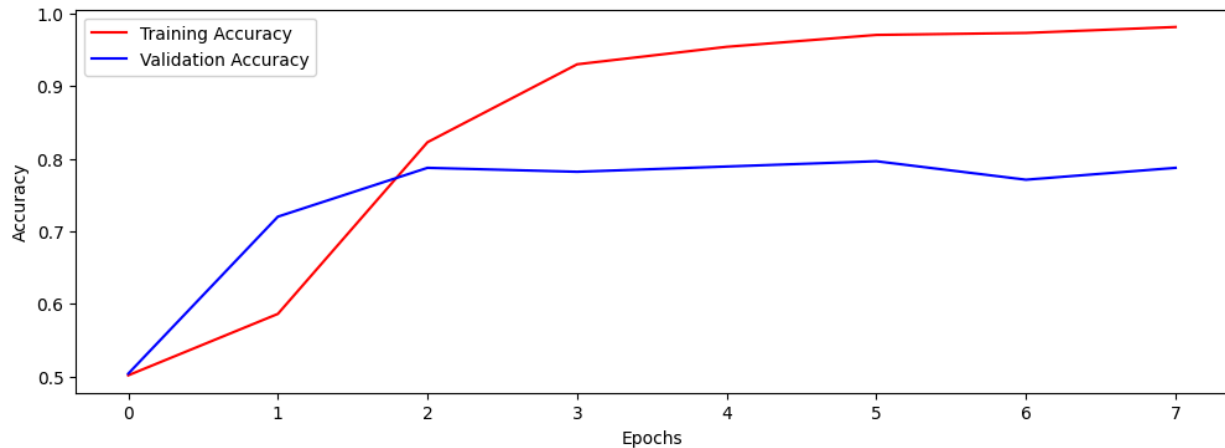


Training and Validation Loss The plot above shows that the training loss decreases steadily, indicating that the model is effectively learning and fitting the training data. Additionally, the decrease in validation loss suggests that the model generalizes well on unseen data.

```python
#plot accuracy
plt.figure(figsize=(12, 4))
plt.plot(train_acc, label='Training Accuracy', color = 'red')
plt.plot(val_acc, label='Validation Accuracy', color = 'blue')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend()
```

The training accuracy steadily improves throughout the epochs, reaching close to 99% by the end, which shows that the model is learning the training data effectively. On the other hand, the validation accuracy levels off at around 78% after the third epoch, indicating that the model's performance on unseen data stops improving beyond this point.

**4. Discuss the predictive accuracy of the trained network using the chosen evaluation metric from part D3.**

The trained network achieved a test accuracy of approximately 79.09%, indicating reliable performance in classifying sentiment on unseen data. While the loss suggests room for improvement, the predictions align well with actual labels, as demonstrated by correctly identifying the sentiment in sample reviews.

#verify model accuracy on test data

```
#verify model accuracy on test data
score = model.evaluate(padded_test, test_labels, verbose=0)
print('Test loss:', score[0])
print('Test accuracy:', score[1])
Test loss: 0.49218860268592834
Test accuracy: 0.7872727513313293


#perform model prediction
predictions = model.predict(padded_test)
i = 45
print("Review:", X_test[i], "\n")
print("Predicted:", "Negative" if predictions[i][0] < 0.5 else "Positive",
"reviews")
```

```
print("Actual:", "Negative" if test_labels[i] == 0 else "Positive",
"reviews")
Review: thing really worth watching wa scenery house beautiful

Predicted: Positive reviews
Actual: Positive reviews
```

## Part V: Summary and Recommendations

**E. Provide the code you used to save the trained network within the neural network.**

```
model.save('my_model.keras')
```

**F. Discuss the functionality of your neural network, including the impact of the network architecture.**

- The neural network architecture includes an Embedding layer that converts words into vector representations, followed by a Bidirectional LSTM layer that captures contextual information from both directions in the sequence. Dropout layers are incorporated to mitigate overfitting, while Dense layers with ReLU activation help refine the learned features. The final layer uses a sigmoid activation to output probabilities for binary classification. This architecture effectively captures the sequential nature of text, while techniques like dropout and L2 regularization help improve generalization. Overall, the design strikes a balance between performance and efficiency, making it well-suited for sentiment analysis tasks

**G. Recommend a course of action based on your results.**

- The goal of this analysis was to develop a neural network model capable of predicting whether a review is positive or negative. Based on the results, I recommend that the company implement this model for predicting sentiment in their data. By leveraging this model, the company can gain valuable insights into customer feedback, enabling them to make more informed decisions and enhance their products or services. Additionally, the model can be fine-tuned and scaled to address different business needs, ensuring its effectiveness in providing accurate sentiment predictions.

## Part VI: Reporting

H.  Show your neural network in an industry-relevant interactive development environment (e.g., a Jupyter Notebook). Include a PDF or HTML document of your executed notebook presentation.

I.  Denote specific web sources you used to acquire segments of third-party code that was used to support the application.


J.  Acknowledge sources, using in-text citations and references, for content that is quoted, paraphrased, or summarized.

GeeksforGeeks. (n.d.). *What is sentiment analysis?* Retrieved from https://www.geeksforgeeks.org/what-is-sentiment-analysis/
Idrees, H. (n.d.). *RNN vs. LSTM vs. GRU: A comprehensive guide to sequential data modeling*. Medium. Retrieved from https://medium.com/@hassaanidrees7/rnn-vs-lstm-vs-gru-a-comprehensive-guide-to-sequential-data-modeling-03aab16647bb
Towards Data Science. (n.d.). *Simple guide to hyperparameter tuning in neural networks*. Retrieved from https://towardsdatascience.com/simple-guide-to-hyperparameter-tuning-in-neural-networks-3fe03dad8594

K.  Demonstrate professional communication in the content and presentation of your submission.