

# NeuralNetworkSentimentAnalysis

August 8, 2023

## 1 Lab 7: Implement a Neural Network for Sentiment Analysis

```
[1]: import pandas as pd
import numpy as np
import os
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
import tensorflow.keras as keras
import time
```

In this lab assignment, you will implement a neural network that performs sentiment analysis for a binary classification problem. You will:

1. Load the book review data set.
2. Create training and test datasets.
3. Transform the training and test text data using a TF-IDF vectorizer.
4. Construct a neural network
5. Train the neural network.
6. Compare the model's performance on the training data vs test data.
7. Improve its generalization performance.

For this lab, you may use the following demos as reference: Transforming Text into Features for Sentiment Analysis and Implementing a Neural Network in Keras.

**Note: some of the code cells in this notebook may take a while to run**

### 1.1 Part 1: Load the Data Set

We will work with the book review data set that contains book reviews taken from Amazon.com reviews.

Task: In the code cell below, use the same method you have been using to load the data using `pd.read_csv()` and save it to DataFrame `df`.

You will be working with the file named "bookReviews.csv" that is located in a folder named "data".

```
[2]: # YOUR CODE HERE
df=pd.read_csv("data/bookReviews.csv")
```

```
[3]: df.head()
```

```
[3]:
```

	Review	Positive Review
0	This was perhaps the best of Johannes Steinhof...	True
1	This very fascinating book is a story written ...	True
2	The four tales in this collection are beautifu...	True
3	The book contained more profanity than I expec...	False
4	We have now entered a second time of deep conc...	True

```
[4]: df.shape
```

```
[4]: (1973, 2)
```

## 1.2 Part 2: Create Training and Test Data Sets

### 1.2.1 Create Labeled Examples

Task: Create labeled examples from DataFrame df. In the code cell below carry out the following steps:

- Get the Positive\_Review column from DataFrame df and assign it to the variable y. This will be our label.
- Get the Review column from DataFrame df and assign it to the variable X. This will be our feature.

```
[5]: # YOUR CODE HERE
y=df['Positive Review']
X=df['Review']
```

```
[6]: X.head()
```

```
[6]: 0    This was perhaps the best of Johannes Steinhof...
1    This very fascinating book is a story written ...
2    The four tales in this collection are beautifu...
3    The book contained more profanity than I expec...
4    We have now entered a second time of deep conc...
Name: Review, dtype: object
```

```
[7]: X.shape
```

```
[7]: (1973,)
```

### 1.2.2 Split Labeled Examples into Training and Test Sets

Task: In the code cell below, create training and test sets out of the labeled examples.

1. Use scikit-learn's `train_test_split()` function to create the data sets.
2. Specify:

- A test set that is 20 percent of the size of the data set.
- A seed value of '1234'.

```
[8]: # YOUR CODE HERE
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=20,random_state=1234)
```

```
[9]: X_train.head()
```

```
[9]: 1216    "Our attention is the most precious resou...
661    I read Desperation and enjoyed it immensely mu...
183    I love Mr. Koontz and have long read his novel...
1428    I got this book when it first came out and rea...
1130    I have read Baldacci's first four novels and h...
Name: Review, dtype: object
```

### 1.3 Part 3: Implement TF-IDF Vectorizer to Transform Text

In the code cell below, you will transform the features into numerical vectors using `TfidfVectorizer`.

Task: Follow the steps to complete the code in the cell below:

1. Create a `TfidfVectorizer` object and save it to the variable `tfidf_vectorizer`.
2. Call `tfidf_vectorizer.fit()` to fit the vectorizer to the training data `X_train`.
3. Call the `tfidf_vectorizer.transform()` method to use the fitted vectorizer to transform the training data `X_train`. Save the result to `X_train_tfidf`.
4. Call the `tfidf_vectorizer.transform()` method to use the fitted vectorizer to transform the test data `X_test`. Save the result to `X_test_tfidf`.

```
[10]: # 1. Create a TfidfVectorizer object
# YOUR CODE HERE
tfidf_vectorizer=TfidfVectorizer()

# 2. Fit the vectorizer to X_train
# YOUR CODE HERE
tfidf_vectorizer.fit(X_train)

# 3. Using the fitted vectorizer, transform the training data
# YOUR CODE HERE
X_train_tfidf=tfidf_vectorizer.transform(X_train)

# 4. Using the fitted vectorizer, transform the test data
# YOUR CODE HERE
X_test_tfidf=tfidf_vectorizer.transform(X_test)
```

When constructing our neural network, we will have to specify the `input_shape`, meaning the dimensionality of the input layer. This corresponds to the dimension of each of the training examples, which in our case is our vocabulary size. Run the code cell below to see the vocabulary size.

```
[11]: vocabulary_size = len(tfidf_vectorizer.vocabulary_)

print(vocabulary_size)
```

21189

## 1.4 Part 4: Construct a Neural Network

### 1.4.1 Step 1. Define Model Structure

Next we will create our neural network structure. We will create an input layer, three hidden layers and an output layer:

- Input layer: The input layer will have the input shape corresponding to the vocabulary size.
- Hidden layers: We will create three hidden layers of widths (number of nodes) 64, 32, and 16. They will utilize the ReLu activation function.
- Output layer: The output layer will have a width of 1. The output layer will utilize the sigmoid activation function. Since we are working with binary classification, we will be using the sigmoid activation function to map the output to a probability between 0.0 and 1.0. We can later set a threshold and assume that the prediction is class 1 if the probability is larger than or equal to our threshold, or class 0 if it is lower than our threshold.

To construct the neural network model using Keras, we will do the following: \* We will use the Keras Sequential class to group a stack of layers. This will be our neural network model object. \* We will use the Dense class to create each layer. \* We will add each layer to the neural network model object.

Task: Follow these steps to complete the code in the cell below:

1. Create the neural network model object.
  - Use `keras.Sequential()` to create a model object, and assign the result to the variable `nn_model`.
2. Create the input layer:
  - Call `keras.layers.Dense()` with the argument `input_shape=(vocabulary_size,)` to specify the dimension of the input.
  - Assign the results to the variable `input_layer`.
  - Use `nn_model.add(input_layer)` to add the layer `input_layer` to the neural network model object.
3. Create the first hidden layer:
  - Call `keras.layers.Dense()` with the arguments `units=64` and `activation='relu'`.
  - Assign the results to the variable `hidden_layer_1`.
  - Use `nn_model.add(hidden_layer_1)` to add the layer `hidden_layer_1` to the neural network model object.
4. Create the second hidden layer using the same approach that you used to create the first hidden layer, specifying 32 units and the `relu` activation function.
  - Assign the results to the variable `hidden_layer_2`.

- Add the layer to the neural network model object.
5. Create the third hidden layer using the same approach that you used to create the first two hidden layers, specifying 16 units and the relu activation function.
    - Assign the results to the variable `hidden_layer_3`.
    - Add the layer to the neural network model object.
  6. Create the output layer using the same approach that you used to create the hidden layers, specifying 1 unit and the sigmoid activation function.
    - Assign the results to the variable `output_layer`.
    - Add the layer to the neural network model object.

```
[30]: from tensorflow.keras import layers
# 1. Create model object
# YOUR CODE HERE
nn_model=keras.Sequential()

# 2. Create the input layer and add it to the model object:

# Create input layer:
input_layer = layers.
→Dense(units=64,activation='relu',input_shape=(vocabulary_size,))

# Add input_layer to the model object:
# YOUR CODE HERE
nn_model.add(input_layer)

# 3. Create the first hidden layer and add it to the model object:

# Create input layer:
hidden_layer_1 = layers.Dense(units=64,activation='relu')

# Add hidden_layer_1 to the model object:
# YOUR CODE HERE
nn_model.add(hidden_layer_1)

# 4. Create the second layer and add it to the model object:

# Create input layer:
hidden_layer_2 = layers.Dense(units=32,activation='relu')

# Add hidden_layer_2 to the model object:
# YOUR CODE HERE
nn_model.add(hidden_layer_2)

# 5. Create the third layer and add it to the model object:
```

```

# Create input layer:
hidden_layer_3 = layers.Dense(units=16,activation='relu')

# Add hidden_layer_3 to the model object:
# YOUR CODE HERE
nn_model.add(hidden_layer_3)

# 6. Create the output layer and add it to the model object:

# Create input layer:
output_layer = layers.Dense(units=1,activation='sigmoid')

# Add output_layer to the model object:
# YOUR CODE HERE
nn_model.add(output_layer)

# Print summary of neural network model structure
nn_model.summary()

```

Model: "sequential\_2"

Layer (type)	Output Shape	Param #
dense_10 (Dense)	(None, 64)	1356160
dense_11 (Dense)	(None, 64)	4160
dropout_6 (Dropout)	(None, 64)	0
dense_12 (Dense)	(None, 32)	2080
dropout_7 (Dropout)	(None, 32)	0
dense_13 (Dense)	(None, 16)	528
dropout_8 (Dropout)	(None, 16)	0
dense_14 (Dense)	(None, 1)	17

Total params: 1,362,945  
 Trainable params: 1,362,945  
 Non-trainable params: 0

[ ]:

### 1.4.2 Step 2. Define the Optimization Function

Task: In the code cell below, create a stochastic gradient descent optimizer using `keras.optimizers.SGD()`. Specify a learning rate of 0.1 using the `learning_rate` parameter. Assign the result to the variable `sgd_optimizer`.

```
[31]: # YOUR CODE HERE
sgd_optimizer=keras.optimizers.SGD(learning_rate=0.1)
```

### 1.4.3 Step 3. Define the Loss Function

Task: In the code cell below, create a binary cross entropy loss function using `keras.losses.BinaryCrossentropy()`. Use the parameter `from_logits=False`. Assign the result to the variable `loss_fn`.

```
[32]: # YOUR CODE HERE
loss_fn=keras.losses.BinaryCrossentropy(from_logits=False)
nn_model.compile(optimizer=sgd_optimizer,loss=loss_fn,metrics=['accuracy'])
```

### 1.4.4 Step 4. Compile the Model

Task: In the code cell below, package the network architecture with the optimizer and the loss function using the `compile()` method.

You will specify the optimizer, loss function and accuracy evaluation metric. Call the `nn_model.compile()` method with the following arguments: \* Use the optimizer parameter and assign it your optimizer variable: `optimizer=sgd_optimizer` \* Use the loss parameter and assign it your loss function variable: `loss=loss_fn` \* Use the metrics parameter and assign it the accuracy evaluation metric: `metrics=['accuracy']`

```
[33]: # YOUR CODE HERE
nn_model.compile(optimizer=sgd_optimizer,loss=loss_fn,metrics=['accuracy'])
```

## 1.5 Part 5. Fit the Model on the Training Data

We will define our own callback class to output information from our model while it is training. Make sure you execute the code cell below so that it can be used in subsequent cells.

```
[34]: class ProgBarLoggerNEpochs(keras.callbacks.Callback):

    def __init__(self, num_epochs: int, every_n: int = 50):
        self.num_epochs = num_epochs
        self.every_n = every_n

    def on_epoch_end(self, epoch, logs=None):
        if (epoch + 1) % self.every_n == 0:
            s = 'Epoch [{}/ {}]' .format(epoch + 1, self.num_epochs)
            logs_s = ['{}: {:.4f}' .format(k.capitalize(), v)
                      for k, v in logs.items()]
            s_list = [s] + logs_s
            print(', '.join(s_list))
```

Task: In the code cell below, fit the neural network model to the vectorized training data.

1. Call `nn_model.fit()` with the training data `X_train_tfidf` and `y_train` as arguments. Note that `X_train_tfidf` is currently of type sparse matrix. The Keras `fit()` method requires that input data be of specific types. One type that is allowed is a NumPy array. Convert `X_train_tfidf` to a NumPy array using the `toarray()` method.
2. In addition, specify the following parameters:
  - Use the `epochs` parameter and assign it the variable to `epochs=num_epochs`
  - Use the `verbose` parameter and assign it the value of 0: `verbose=0`
  - Use the `callbacks` parameter and assign it a list containing our logger function: `callbacks=[ProgBarLoggerNEpochs(num_epochs_M, every_n=50)]`
  - We will use a portion of our training data to serve as validation data. Use the `validation_split` parameter and assign it the value 0.2
3. Save the results to the variable `history`.

Note: This may take a while to run.

```
[35]: t0 = time.time() # start time

X_train_np=X_train_tfidf.toarray()

num_epochs = 50 #epochs

history=nn_model.fit(X_train_tfidf.toarray(),y_train,
    ↳epochs=num_epochs,verbose=0,callbacks=[ProgBarLoggerNEpochs(num_epochs,every_n=50)],validat
    ↳2)

t1 = time.time() # stop time

print('Elapsed time: %.2fs' % (t1-t0))
```

```
Epoch [50/ 50], Loss: 0.1644, Accuracy: 0.9616, Val_loss: 0.5126, Val_accuracy:
0.8107
```

```
Elapsed time: 27.55s
```

### 1.5.1 Visualize the Model's Performance Over Time

The code above outputs both the training loss and accuracy and the validation loss and accuracy. Let us visualize the model's performance over time:

```
[36]: # Plot training and validation loss
plt.plot(range(1, num_epochs + 1), history.history['loss'], label='Training
    ↳Loss')
plt.plot(range(1, num_epochs + 1), history.history['val_loss'],
    ↳label='Validation Loss')
```



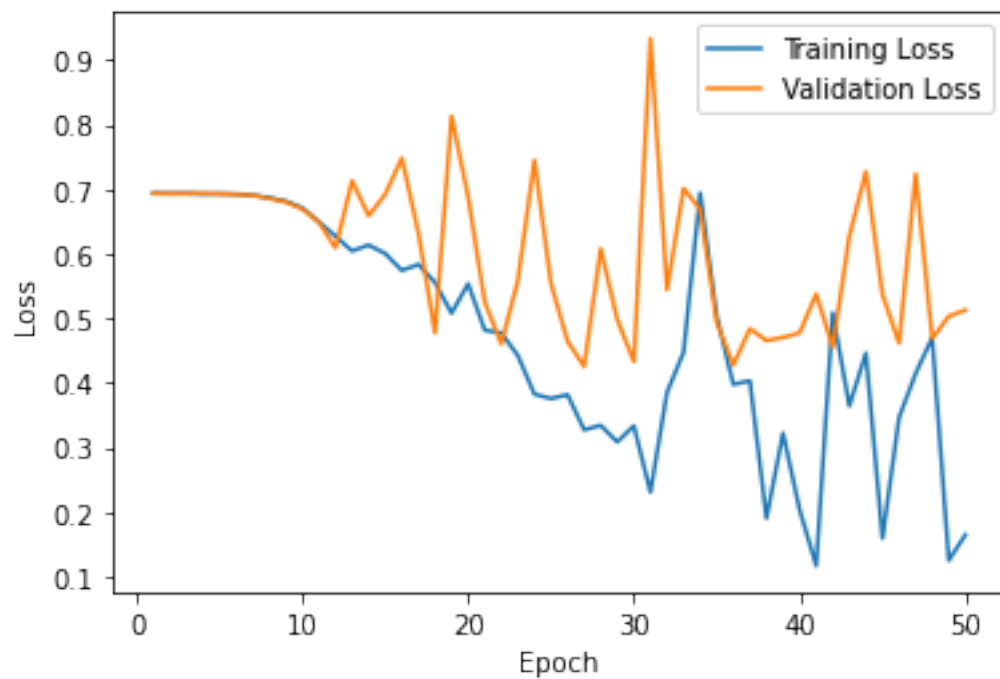
```

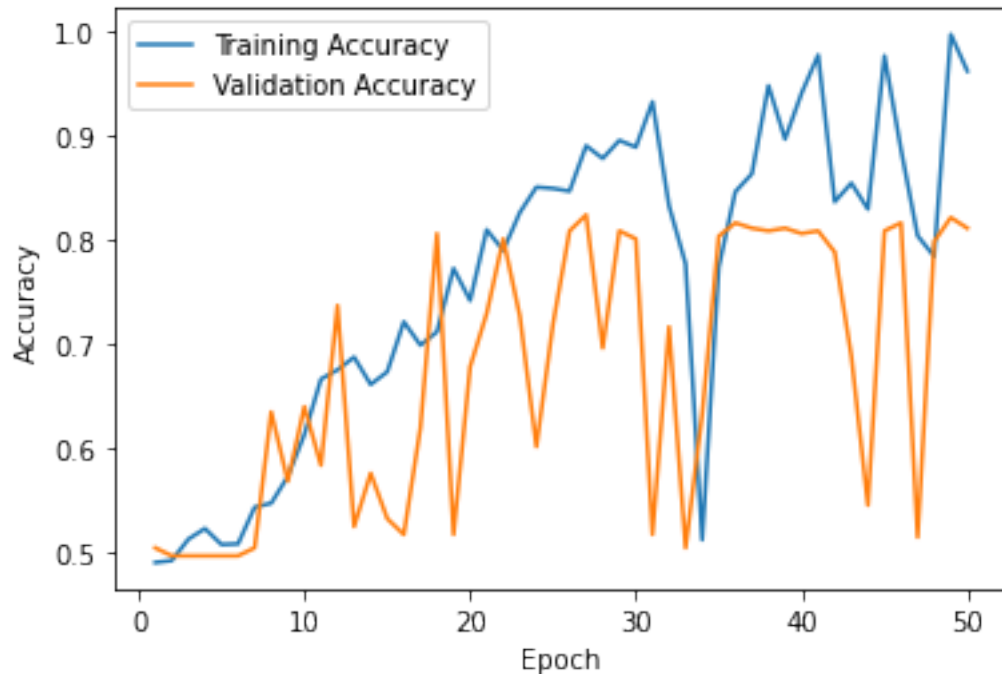
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.show()

# Plot training and validation accuracy
plt.plot(range(1, num_epochs + 1), history.history['accuracy'], label='Training_
→Accuracy')
plt.plot(range(1, num_epochs + 1), history.history['val_accuracy'],
→label='Validation Accuracy')

plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend()
plt.show()

```





## 1.6 Part 6. Improve the Model and Evaluate the Performance

We just evaluated our model's performance on the training and validation data. Let's now evaluate its performance on our test data and compare the results.

Keras makes the process of evaluating our model very easy. Recall that when we compiled the model we specified the metric we wanted to use to evaluate the model: accuracy. The Keras method `evaluate()` will return the loss and accuracy score of our model on our test data.

**Task:** In the code cell below, call `nn_model.evaluate()` with `X_test_tfidf` and `y_test` as arguments. You must convert `X_test_tfidf` to a NumPy array using the `toarray()` method.

**Note:** The `evaluate()` method returns a list containing two values. The first value is the loss and the second value is the accuracy score.

```
[38]: X_test_np=X_test_tfidf.toarray()
      loss, accuracy = nn_model.evaluate(X_test_np,y_test)

      print('Loss: ', str(loss) , 'Accuracy: ', str(accuracy))
```

```
1/1 [=====] - 0s 28ms/step - loss: 1.0788 - accuracy:
0.6000
Loss: 1.0787767171859741 Accuracy: 0.6000000238418579
```

### 1.6.1 Prevent Overfitting and Improve Model's Performance

Neural networks can be prone to overfitting. Notice that the training accuracy is 100% but the test accuracy is around 82%. This indicates that our model is overfitting; it will not perform as well to

new, previously unseen data as it did during training. We want to have an accurate idea of how well our model will generalize. Our goal is to have our training and testing accuracy scores be as close as possible.

While there are different techniques that can be used to prevent overfitting, for the purpose of this exercise, we will focus on two methods:

1. Changing the number of epochs. Too many epochs can lead to overfitting of the training dataset, whereas too few epochs may result in underfitting.
2. Adding dropout regularization. During training, the nodes of a particular layer may always become influenced only by the output of a particular node in the previous layer, causing overfitting. Dropout regularization is a technique that randomly drops a number of nodes in a neural network during training as a way to adding randomization and prevent nodes from becoming dependent on one another. Adding dropout regularization can reduce overfitting and also improve the performance of the model.

Task:

1. Tweak the variable `num_epochs` above and restart and rerun all of the cells above. Evaluate the performance of the model on the training data and the test data.
2. Add Keras Dropout layers after one or all hidden layers. Add the following line of code after you add a hidden layer to your model object: `nn_model.add(keras.layers.Dropout(.25))`. The parameter `.25` is the fraction of the nodes to drop. You can experiment with this value as well. Restart and rerun all of the cells above. Evaluate the performance of the model on the training data and the test data.

Analysis: In the cell below, specify the different approaches you used to reduce overfitting and summarize which configuration led to the best generalization performance.

Did changing the number of epochs prevent overfitting? Which value of `num_epochs` yielded the closest training and testing accuracy score? Recall that too few epochs can lead to underfitting (both poor training and test performance). Which value of `num_epochs` resulted in the best accuracy score when evaluating the test data?

Did adding dropout layers prevent overfitting? How so? Did it also improve the accuracy score when evaluating the test data? How many dropout layers did you add and which fraction of nodes did you drop?

Record your findings in the cell below.

Changing the number of epochs helped to reduce overfitting. Values returned for 10 epochs were Loss: 0.6878030896186829 Accuracy: 0.6000000238418579, for 20: Loss: 0.6714824438095093 Accuracy: 0.6000000238418579, for 30: Loss: 0.8673494458198547 Accuracy: 0.6499999761581421, and for 50: Loss: 2.0127642154693604 Accuracy: 0.550000011920929. This shows that out of these options, 30 epochs performed the best as lower amounts showed indications that underfitting was occurring, while with 50 epochs the data was overfit. However, the performance of the model is not ideal. Adding dropout layers assisted in handling this data slightly. I added this line to one hidden layer at first, resulting in Loss: 2.5163350105285645 Accuracy: 0.550000011920929, adding it to another 2.043381929397583 Accuracy: 0.6499999761581421, which was only a slight change. However, adding it to all three resulted in Loss: 1.5367114543914795 Accuracy: 0.6499999761581421, which resulted in a better loss value for 50 epochs.

```

[27]: #Code with nn_model.add(keras.layers.Dropout(.25)) lines added
from tensorflow.keras import layers
# 1. Create model object
# YOUR CODE HERE
nn_model=keras.Sequential()

# 2. Create the input layer and add it to the model object:

# Create input layer:
input_layer = layers.
    ↳Dense(units=64,activation='relu',input_shape=(vocabulary_size,))

# Add input_layer to the model object:
# YOUR CODE HERE
nn_model.add(input_layer)

# 3. Create the first hidden layer and add it to the model object:

# Create input layer:
hidden_layer_1 = layers.Dense(units=64,activation='relu')

# Add hidden_layer_1 to the model object:
# YOUR CODE HERE
nn_model.add(hidden_layer_1)
nn_model.add(keras.layers.Dropout(.25))

# 4. Create the second layer and add it to the model object:

# Create input layer:
hidden_layer_2 = layers.Dense(units=32,activation='relu')

# Add hidden_layer_2 to the model object:
# YOUR CODE HERE
nn_model.add(hidden_layer_2)
nn_model.add(keras.layers.Dropout(.25))

# 5. Create the third layer and add it to the model object:

# Create input layer:
hidden_layer_3 = layers.Dense(units=16,activation='relu')

# Add hidden_layer_3 to the model object:
# YOUR CODE HERE

```

```

nn_model.add(hidden_layer_3)
nn_model.add(keras.layers.Dropout(.25))

# 6. Create the output layer and add it to the model object:

# Create input layer:
output_layer = layers.Dense(units=1,activation='sigmoid')

# Add output_layer to the model object:
# YOUR CODE HERE
nn_model.add(output_layer)

# Print summary of neural network model structure
nn_model.summary()

```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_5 (Dense)	(None, 64)	1356160
dense_6 (Dense)	(None, 64)	4160
dropout_3 (Dropout)	(None, 64)	0
dense_7 (Dense)	(None, 32)	2080
dropout_4 (Dropout)	(None, 32)	0
dense_8 (Dense)	(None, 16)	528
dropout_5 (Dropout)	(None, 16)	0
dense_9 (Dense)	(None, 1)	17

Total params: 1,362,945  
 Trainable params: 1,362,945  
 Non-trainable params: 0

## 1.6.2 Make Predictions on the Test Set

Now that we have our best performing model that can generalize to new, previously unseen data, let us make predictions using our test data.

In the cell below, we will make a prediction on our test set using the `predict()` method, receive

a probability between 0.0 and 1.0, and then apply a threshold to obtain the the predicted class for each example. We will use a threshold of 0.5.

For the first 10 examples, we will output their probabilities and the corresponding classes. Examine the output to see how this works.

```
[20]: probability_predictions = nn_model.predict(X_test_tfidf.toarray())

print("Predictions for the first 10 examples:")
print("Probability\t\t\t\tClass")
for i in range(0,10):
    if probability_predictions[i] >= .5:
        class_pred = "Good Review"
    else:
        class_pred = "Bad Review"
    print(str(probability_predictions[i]) + "\t\t\t\t" + str(class_pred))
```

```
Predictions for the first 10 examples:
Probability          Class
[0.99993044]         Good Review
[0.99242747]         Good Review
[0.11647522]         Bad Review
[0.09843245]         Bad Review
[0.98083043]         Good Review
[0.72994316]         Good Review
[0.07349026]         Bad Review
[0.05808163]         Bad Review
[0.9998639]          Good Review
[0.9998069]          Good Review
```

Let's check two book reviews and see if our model properly predicted whether the reviews are good or bad reviews.

```
[21]: print('Review #1:\n') #I tested 16 instead of 56 as I got an error saying 56_
      →was out of bounds
print(X_test.to_numpy()[16])

goodReview = True if probability_predictions[16] >= .5 else False

print('\nPrediction: Is this a good review? {}'.format(goodReview))

print('Actual: Is this a good review? {}'.format(y_test.to_numpy()[16]))
```

Review #1:

I forced myself to finish this book, though it was touch and go in several places, just so I could feel able to review it fairly. It seems I liked it even less than reviewer Ms. Trieste "CF", below, but I am in general agreement with her points. There are just so many things wrong with this book that even the grating "famous names" don't really stand out for me. Let's see, where to

begin: the naive political ranting, the unspeakable dialogue, the corny love interest, the wholly implausible technology, the absurd coincidences, the distractingly disjointed structure, the proliferation of minor characters, the talky explanations of the detective's thought processes, the total incompetence and corruption of the police - I can't go on; it's hackneyed and poorly written, and that's all there is to it.

Prediction: Is this a good review? False

Actual: Is this a good review? False

```
[22]: print('Review #2:\n') #I tested 19 instead of 24 as I got an error saying 24␣
      ↪was out of bounds

print(X_test.to_numpy()[19])

goodReview = True if probability_predictions[19] >= .5 else False

print('\nPrediction: Is this a good review? {}\n'.format(goodReview))

print('Actual: Is this a good review? {}\n'.format(y_test.to_numpy()[19]))
```

Review #2:

We went to a book signing to get this book. Truthfully, I wasn't sure if the book was going to live up to the hype. T.O. graciously signed books and all proceeds went to charity. I was surprised this book was so well written. My four year old wants me to read it over and over again. He doesn't like for me to read to him but he loves this book! It is the perfect book for children. Its not too short but not too long. I have read this book at least 10 times in the past 3 days. I would recommend this book to anyone with small children. It is great book for little boys because of the football theme. However, the lesson is something everyone can appreciate. Great Book

Prediction: Is this a good review? True

Actual: Is this a good review? False

### 1.6.3 Deep Dive (Ungraded):

Experiment with the vectorizer and neural network implementation above and compare your results every time you train the network. Pay attention to the time it takes to train the network, and the resulting loss and accuracy on both the training and test data.

Below are some ideas for things you can try:

- Adjust the learning rate.
- Add more hidden layers and/or experiment with different values for the `unit` parameter in the hidden layers to change the number of nodes in the hidden layers.
- Fit your vectorizer using different document frequency values and a different n-gram ranges. When creating a `TfidfVectorizer` object, use the parameter `min_df` to specify the minimum 'document frequency' and use `gram_range=(1,2)` to change the default n-gram range of (1,1).

[ ]: