

```

1 import pandas as pd
2 import numpy as np
3 import nltk
4 import sqlite3
5 import os
6 from bs4 import BeautifulSoup
7 from sklearn.model_selection import train_test_split
8 from keras.preprocessing.text import Tokenizer
9 from keras.preprocessing import sequence
10 from keras.models import Sequential
11 from keras.layers import Dense, BatchNormalization, LSTM, Dropout
12 from keras.layers.embeddings import Embedding

```

↳ Using TensorFlow backend.

```

1 from google.colab import drive
2 drive.mount('/content/gdrive')

```

↳ Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=...

Enter your authorization code:

.....

Mounted at /content/gdrive

```

1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import time
5 # https://gist.github.com/greydanus/f6eee59eaf1d90fcb3b534a25362cea4
6 # https://stackoverflow.com/a/14434334
7 # this function is used to update the plots for each epoch and error
8 def plt_dynamic(x, vy, ty, ax, colors=['b']):
9     ax.plot(x, vy, 'b', label="Validation Loss")
10    ax.plot(x, ty, 'r', label="Train Loss")
11    plt.legend()
12    plt.grid()
13    fig.canvas.draw()

```

```

1 con = sqlite3.connect("/content/gdrive/My Drive/Colab Notebooks/database.sqli
2 filtered_data = pd.read_sql_query("""SELECT * FROM Reviews WHERE Score != 3""
3
4 def partition(x):
5     if x < 3:
6         return 0
7     return 1
8
9 #changing reviews with score less than 3 to be positive and vice-versa
10 actualScore = filtered_data['Score']
11 positiveNegative = actualScore.map(partition)
12 filtered_data['Score'] = positiveNegative

```

```

1 #Deduplication of entries
2 final=filtered_data.drop_duplicates(subset={"UserId","ProfileName","Time","Te
3 final.shape

```

↳ (364173, 10)

```

1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]

```

```

1 final = final.sample(n=50000, replace=True)

```

```

1 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours',
2               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
3               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'i',
4               'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',
5               'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
6               'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
7               'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
8               'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off',
9               'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all',
10              'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than',
11              's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shouldn't",
12              've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
13              "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'mustn't',
14              'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wouldn',
15              "wouldn't"])

```

```

1 # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-tags
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\'re", " are", phrase)
12    phrase = re.sub(r"\s", " is", phrase)
13    phrase = re.sub(r"\d", " would", phrase)
14    phrase = re.sub(r"\ll", " will", phrase)
15    phrase = re.sub(r"\t", " not", phrase)
16    phrase = re.sub(r"\ve", " have", phrase)
17    phrase = re.sub(r"\m", " am", phrase)
18    return phrase

```

```

1 from tqdm import tqdm
2 preprocessed_reviews = []
3 # tqdm is for printing the status bar
4 for sentence in tqdm(final['Text'].values):
5     sentence = re.sub(r"http\S+", "", sentence)
6     sentence = BeautifulSoup(sentence, 'lxml').get_text()
7     sentence = decontracted(sentence)
8     sentence = re.sub(r"\S*\d\S*", "", sentence).strip()
9     sentence = re.sub(r"^[A-Za-z]+", "", sentence)
10    # https://gist.github.com/sebleier/554280
11    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
12    preprocessed_reviews.append(sentence.strip())

```

100%|██████████| 50000/50000 [00:16<00:00, 3051.25it/s]

```
1 final['CleanedText'] = preprocessed_reviews
```

```

1 x = final['CleanedText']
2 y = final['Score']
3
4 x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3)

```

```

1 #Source: https://keras.io/preprocessing/text/
2 #Source: https://www.kaggle.com/kredyl0/simple-lstm-for-text-classification
3 #Source: https://keras.io/preprocessing/sequence/
4
5 max_words = 5000
6 max_len = 600
7 tok = Tokenizer(num_words=max_words)
8 tok.fit_on_texts(x_train)

```

```

9 sequences = tok.texts_to_sequences(x_train)
10 train_sequences_matrix = sequence.pad_sequences(sequences,maxlen=max_len)

1 test_sequences = tok.texts_to_sequences(x_test)
2 test_sequences_matrix = sequence.pad_sequences(test_sequences,maxlen=max_len)

1 embedding_vecor_length = 32
2 model = Sequential()
3 model.add(Embedding(max_words+1, embedding_vecor_length, input_length=max_len
4 model.add(LSTM(100, dropout = 0.3, recurrent_dropout=0.3, return_sequences=Tr
5 model.add(BatchNormalization())
6 model.add(LSTM(100, dropout = 0.4, recurrent_dropout=0.4))
7 model.add(Dense(1, activation='sigmoid'))
8 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accurac
9 print(model.summary())

```

⏏ WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/tensorflow/Instructions for updating:
Colocations handled automatically by placer.
WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backerInstructions for updating:
Please use `rate` instead of `keep_prob`. Rate should be set to `rate = 1`.

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 600, 32)	160032
lstm_1 (LSTM)	(None, 600, 100)	53200
batch_normalization_1 (Batch Normalization)	(None, 600, 100)	400
lstm_2 (LSTM)	(None, 100)	80400
dense_1 (Dense)	(None, 1)	101
Total params: 294,133		
Trainable params: 293,933		
Non-trainable params: 200		
None		

```
1 history = model.fit(train_sequences_matrix, y_train, epochs=10, batch_size=10
```

⏏

```
WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/tensorflow/
Instructions for updating:
Use tf.cast instead.
Train on 35000 samples, validate on 15000 samples
Epoch 1/10
35000/35000 [=====] - 771s 22ms/step - loss: 0.292
Epoch 2/10
35000/35000 [=====] - 771s 22ms/step - loss: 0.190
Epoch 3/10
```

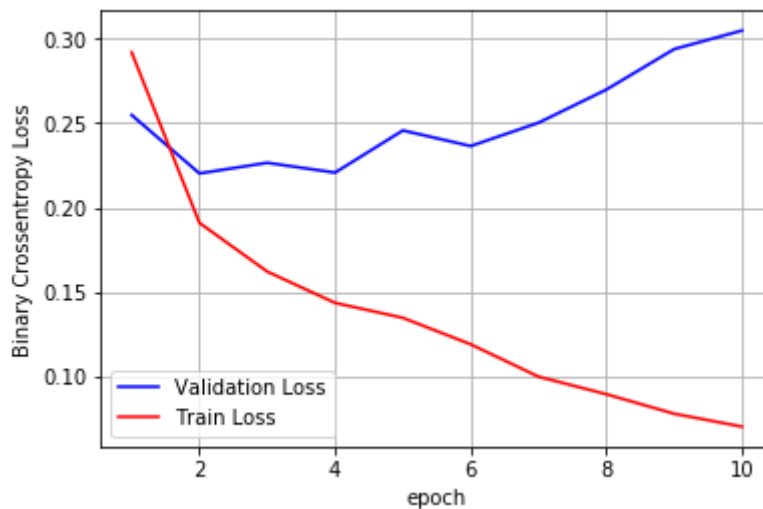
```
1 scores = model.evaluate(test_sequences_matrix, y_test, verbose=0)
2 model2_acc = scores[1]
3 print("Accuracy: %.2f%%" % (scores[1]*100))
```

Accuracy: 91.69%

```
35000/35000 [=====] - 773s 22ms/step - loss: 0.115
```

```
1 print('Test score:', scores[0])
2 print('Test accuracy:', scores[1])
3
4 fig,ax = plt.subplots(1,1)
5 ax.set_xlabel('epoch') ; ax.set_ylabel('Binary Crossentropy Loss')
6
7 # list of epoch numbers
8 x = list(range(1,11))
9
10 vy = history.history['val_loss']
11 ty = history.history['loss']
12
13
14 plt_dynamic(x, vy, ty, ax)
```

Test score: 0.3049351974586646
Test accuracy: 0.9169333333651225



```
1 embedding_vector_length = 32
2 model = Sequential()
3 model.add(Embedding(max_words+1, embedding_vector_length, input_length=max_len
4 model.add(LSTM(100, dropout = 0.3, recurrent_dropout=0.3, return_sequences=True)
5 model.add(BatchNormalization())
6 model.add(LSTM(100, dropout = 0.5, recurrent_dropout=0.5))
7 model.add(Dense(1, activation='sigmoid'))
8 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
9 print(model.summary())
```



Layer (type)	Output Shape	Param #
embedding_2 (Embedding)	(None, 600, 32)	160032
lstm_3 (LSTM)	(None, 600, 100)	53200
batch_normalization_2 (Batch Normalization)	(None, 600, 100)	400
lstm_4 (LSTM)	(None, 100)	80400
dense_2 (Dense)	(None, 1)	101
Total params: 294,133		
Trainable params: 293,933		

```
1 history = model.fit(train_sequences_matrix, y_train, epochs=10, batch_size=10
```

☞ Train on 35000 samples, validate on 15000 samples

```
Epoch 1/10
35000/35000 [=====] - 778s 22ms/step - loss: 0.291
Epoch 2/10
35000/35000 [=====] - 778s 22ms/step - loss: 0.189
Epoch 3/10
35000/35000 [=====] - 780s 22ms/step - loss: 0.166
Epoch 4/10
35000/35000 [=====] - 778s 22ms/step - loss: 0.147
Epoch 5/10
35000/35000 [=====] - 779s 22ms/step - loss: 0.145
Epoch 6/10
35000/35000 [=====] - 772s 22ms/step - loss: 0.131
Epoch 7/10
35000/35000 [=====] - 777s 22ms/step - loss: 0.114
Epoch 8/10
35000/35000 [=====] - 776s 22ms/step - loss: 0.106
Epoch 9/10
35000/35000 [=====] - 775s 22ms/step - loss: 0.099
Epoch 10/10
35000/35000 [=====] - 777s 22ms/step - loss: 0.099
```

```
1 scores = model.evaluate(test_sequences_matrix, y_test, verbose=0)
2 model3_acc = scores[1]
3 print("Accuracy: %.2f%%" % (scores[1]*100))
```

☞ Accuracy: 91.53%

```
1 print('Test score:', scores[0])
2 print('Test accuracy:', scores[1])
3
4 fig,ax = plt.subplots(1,1)
5 ax.set_xlabel('epoch') ; ax.set_ylabel('Categorical Crossentropy Loss')
6
7 # list of epoch numbers
8 x = list(range(1,11))
9
10
11 vy = history.history['val_loss']
12 ty = history.history['loss']
13
14 plt_dynamic(x, vy, ty, ax)
```

☞

Test score: 0.3190554770519336

Test accuracy: 0.9152666666348775



```

1 embedding_vecor_length = 32
2 model = Sequential()
3 model.add(Embedding(max_words+1, embedding_vecor_length, input_length=max_len
4 model.add(LSTM(100, dropout = 0.7, recurrent_dropout=0.7, return_sequences=Tr
5 model.add(BatchNormalization())
6 model.add(LSTM(100, dropout = 0.5, recurrent_dropout=0.5, return_sequences=Tr
7 model.add(BatchNormalization())
8 model.add(LSTM(100, dropout = 0.3, recurrent_dropout=0.3))
9 model.add(Dense(1, activation='sigmoid'))
10 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accurac
11 print(model.summary())

```



Layer (type)	Output Shape	Param #
=====		
embedding_3 (Embedding)	(None, 600, 32)	160032
lstm_5 (LSTM)	(None, 600, 100)	53200
batch_normalization_3 (Batch Normalization)	(None, 600, 100)	400
lstm_6 (LSTM)	(None, 600, 100)	80400
batch_normalization_4 (Batch Normalization)	(None, 600, 100)	400
lstm_7 (LSTM)	(None, 100)	80400
dense_3 (Dense)	(None, 1)	101
=====		
Total params: 374,933		
Trainable params: 374,533		
Non-trainable params: 400		
None		

```
1 history = model.fit(train_sequences_matrix, y_train, epochs=10, batch_size=10
```



Train on 35000 samples, validate on 15000 samples

```
Epoch 1/10
35000/35000 [=====] - 1169s 33ms/step - loss: 0.38
Epoch 2/10
35000/35000 [=====] - 1167s 33ms/step - loss: 0.26
Epoch 3/10
35000/35000 [=====] - 1164s 33ms/step - loss: 0.23
Epoch 4/10
35000/35000 [=====] - 1169s 33ms/step - loss: 0.23
Epoch 5/10
35000/35000 [=====] - 1165s 33ms/step - loss: 0.19
Epoch 6/10
35000/35000 [=====] - 1172s 33ms/step - loss: 0.20
Epoch 7/10
```

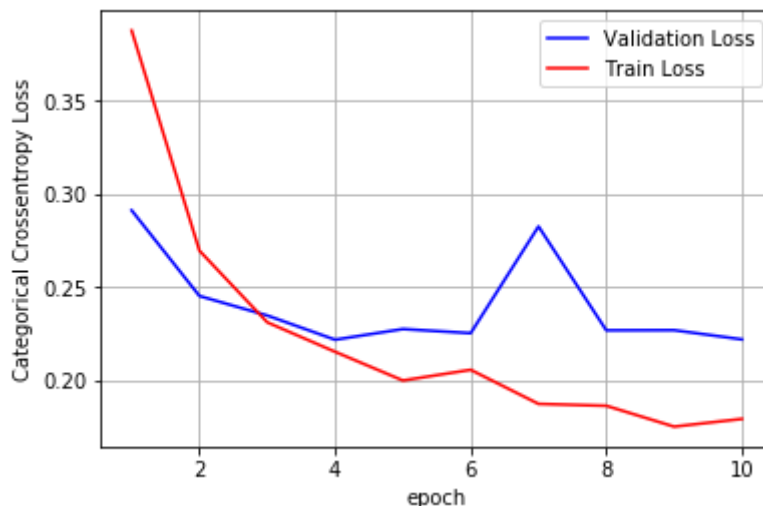
```
1 scores = model.evaluate(test_sequences_matrix, y_test, verbose=0)
2 model4_acc = scores[1]
3 print("Accuracy: %.2f%%" % (scores[1]*100))
```

Accuracy: 91.63%

```
35000/35000 [=====] - 1170s 33ms/step - loss: 0.17
```

```
1 print('Test score:', scores[0])
2 print('Test accuracy:', scores[1])
3
4 fig, ax = plt.subplots(1, 1)
5 ax.set_xlabel('epoch') ; ax.set_ylabel('Categorical Crossentropy Loss')
6
7 # list of epoch numbers
8 x = list(range(1, 11))
9
10
11 vy = history.history['val_loss']
12 ty = history.history['loss']
13
14 plt_dynamic(x, vy, ty, ax)
```

Test score: 0.22179906602303187
Test accuracy: 0.9163333333015442



```
1 from prettytable import PrettyTable
2
3 model_metric = PrettyTable()
4
5 model_metric = PrettyTable(["Model", "# LSTM Layers", 'Test Accuracy'])
6
7 model_metric.add_row(['LSTM', '2 layers', '%.4f%%' % model2_acc])
8 model_metric.add_row(['LSTM', '2 layer', '%.4f%%' % model3_acc])
9 model_metric.add_row(['LSTM', '3 layers', '%.4f%%' % model4_acc])
```

```
10 print(model_metric.get_string(start=0, end=8))
```

↳

Model	# LSTM Layers	Test Accuracy
LSTM	2 layers	0.9169%
LSTM	2 layer	0.9153%
LSTM	3 layers	0.9163%