

# Assignment\_14\_LSTM\_IMDB

June 18, 2019

```
In [2]: # Credits: https://machinelearningmastery.com/sequence-classification-lstm-recurrent-n
# LSTM for sequence classification in the IMDB dataset
import numpy
from keras.datasets import imdb
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import LSTM,
from keras.layers import Dropout
from keras.layers.embeddings import Embedding
from keras.preprocessing import sequence
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
import os
import sqlite3

from sklearn.metrics import roc_auc_score
import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
# fix random seed for reproducibility
numpy.random.seed(7)

!apt-get install -y -qq software-properties-common python-software-properties module-i
!add-apt-repository -y ppa:alessandro-strada/ppa 2>&1 > /dev/null
!apt-get update -qq 2>&1 > /dev/null
!apt-get -y install -qq google-drive-ocamlfuse fuse
from google.colab import auth
auth.authenticate_user()
from oauth2client.client import GoogleCredentials
creds = GoogleCredentials.get_application_default()
import getpass
!google-drive-ocamlfuse -headless -id={creds.client_id} -secret={creds.client_secret} .
```



```
Out[4]:
```

	Id	...	Text
0	1	...	I have bought several of the Vitality canned d...
1	2	...	Product arrived labeled as Jumbo Salted Peanut...
2	3	...	This is a confection that has been around a fe...

[3 rows x 10 columns]

```
In [0]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

```
In [6]: print(display.shape)
display.head()
```

```
(80668, 7)
```

```
Out[6]:
```

	UserId	...	COUNT(*)
0	#oc-R115TNMSPFT9I7	...	2
1	#oc-R11D9D7SHXIJB9	...	3
2	#oc-R11DNU2NBKQ23Z	...	2
3	#oc-R1105J5ZVQE25C	...	3
4	#oc-R12KPBODL2B5ZD	...	2

[5 rows x 7 columns]

```
In [7]: display['COUNT(*)'].sum()
```

```
Out[7]: 393063
```

```
In [0]: #Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False)
```

```
In [9]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first')
final.shape
```

```
Out[9]: (348262, 10)
```

```
In [10]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[10]: 69.6524
```

```
In [0]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [12]: #Before starting the next phase of preprocessing lets see the number of entries left  
print(final.shape)
```

```
#How many positive and negative reviews are present in our dataset?  
final['Score'].value_counts()
```

```
(348260, 10)
```

```
Out[12]: 1    293516  
        0    54744  
        Name: Score, dtype: int64
```

```
In [0]: # https://stackoverflow.com/a/47091490/4084039  
import re
```

```
def decontracted(phrase):  
    # specific  
    phrase = re.sub(r"won't", "will not", phrase)  
    phrase = re.sub(r"can't", "can not", phrase)  
  
    # general  
    phrase = re.sub(r"n't", " not", phrase)  
    phrase = re.sub(r"\ 're", " are", phrase)  
    phrase = re.sub(r"\ 's", " is", phrase)  
    phrase = re.sub(r"\ 'd", " would", phrase)  
    phrase = re.sub(r"\ 'll", " will", phrase)  
    phrase = re.sub(r"\ 't", " not", phrase)  
    phrase = re.sub(r"\ 've", " have", phrase)  
    phrase = re.sub(r"\ 'm", " am", phrase)  
    return phrase
```

```
In [18]: # printing some random reviews  
sent_0 = final['Text'].values[0]  
print(sent_0)  
print("="*50)  
  
sent_1000 = final['Text'].values[1000]  
print(sent_1000)  
print("="*50)  
  
sent_1500 = final['Text'].values[1500]  
print(sent_1500)  
print("="*50)  
  
sent_4900 = final['Text'].values[4900]  
print(sent_4900)  
print("="*50)
```

This book was purchased as a birthday gift for a 4 year old boy. He squealed with delight and l  
=====

I've purchased both the Espressione Espresso (classic) and the 100% Arabica. My vote is defin  
=====

This is a great product. It is very healthy for all of our dogs, and it is the first food that  
=====

I find everything I need at Amazon so I always look there first. Chocolate tennis balls for a t  
=====

```
In [19]: # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

This book was purchased as a birthday gift for a 4 year old boy. He squealed with delight and l

```
In [20]: # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all
from bs4 import BeautifulSoup

soup = BeautifulSoup(sent_0, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1000, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1500, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_4900, 'lxml')
text = soup.get_text()
print(text)
```

This book was purchased as a birthday gift for a 4 year old boy. He squealed with delight and l  
=====

I've purchased both the Espressione Espresso (classic) and the 100% Arabica. My vote is defin  
=====

This is a great product. It is very healthy for all of our dogs, and it is the first food that  
=====

I find everything I need at Amazon so I always look there first. Chocolate tennis balls for a

```
In [0]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have reumoved in the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves',
'you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'i',
'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that's",
'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had',
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through',
'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over',
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any',
'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too',
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'no',
've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
"hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
'mustn't', 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'v',
'won', "won't", 'wouldn', "wouldn't"])
```

```
In [21]: # Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

100%|| 348260/348260 [02:52<00:00, 2016.77it/s]

```
In [22]: preprocessed_reviews[1500]
```

```
Out[22]: 'great product healthy dogs first food love eat helped older dog lose weight year old
```

Splitting the data :

```
In [0]: from sklearn.model_selection import train_test_split
final['Text'] = preprocessed_reviews
```

```

finalp = final[final.Score == 1].sample(25000,random_state =2)
finaln = final[final.Score == 0].sample(25000,random_state =2)
finalx = pd.concat([finalp,finaln],ignore_index=True)
finalx = finalx.sort_values('Time')
y = finalx.Score.values
X = finalx.Text.values
X_train, X_test , y_train, y_test = train_test_split(X,y,test_size=0.2)
X_tr, X_cv , y_tr, y_cv = train_test_split(X_train,y_train,test_size=0.2)

```

```

In [0]: from keras.preprocessing.text import Tokenizer
tokens = Tokenizer(num_words=5000)
tokens.fit_on_texts(X_tr)

```

```

X_tr = tokens.texts_to_sequences(X_tr)
X_cv = tokens.texts_to_sequences(X_cv)
X_test = tokens.texts_to_sequences(X_test)

```

```

In [66]: print(X_tr[1])
         print(type(X_tr[1]))
         print(len(X_tr[1]))

```

```

[17, 15, 167, 65, 246, 164, 135, 167, 31, 471, 95, 598, 2285, 309, 213, 441, 3, 221, 67]
<class 'list'>

```

19

```

In [67]: # truncate and/or pad input sequences
max_review_length = 600
X_tr = sequence.pad_sequences(X_tr, maxlen=max_review_length)
X_cv = sequence.pad_sequences(X_cv, maxlen=max_review_length)
X_test = sequence.pad_sequences(X_test, maxlen=max_review_length)

```

```

print(X_tr.shape)
print(X_tr[1])
print(X_cv.shape)
print(X_cv[1])
print(X_test.shape)
print(X_test[1])

```

(32000, 600)

```

[ 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  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  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  0  0  0  0  0  0
  0  0  0  0  0  0  0  0  0  0  0  0  0  0

```

167    31  
(8000, 600)







```

nb_epoch = 10
history = model.fit(X_tr, y_tr, nb_epoch=nb_epoch, batch_size=64, validation_data=(X_val, y_val))
# Final evaluation of the model
scores = model.evaluate(X_test, y_test, verbose=0)
print("Accuracy: %.2f%%" % (scores[1]*100))
fig, ax = plt.subplots(1, 1)
ax.set_xlabel('epoch')
ax.set_ylabel('Categorical Crossentropy Loss')
x = list(range(1, nb_epoch+1))
vy = history.history['val_loss']
ty = history.history['loss']
plt_dynamic(x, vy, ty, ax)
model()

```

Layer (type)	Output Shape	Param #
embedding_16 (Embedding)	(None, 600, 32)	160032
lstm_23 (LSTM)	(None, 100)	53200
dense_15 (Dense)	(None, 1)	101

```

Total params: 213,333
Trainable params: 213,333
Non-trainable params: 0

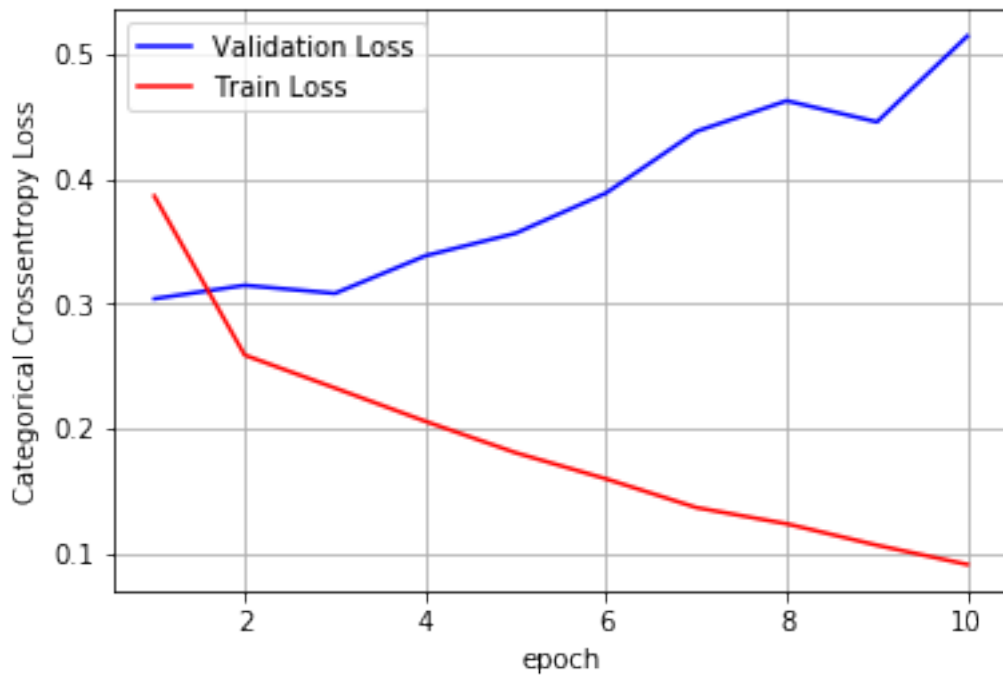
```

```

None
Train on 32000 samples, validate on 8000 samples
Epoch 1/10
32000/32000 [=====] - 419s 13ms/step - loss: 0.3861 - acc: 0.8264 - val_loss: 0.2588 - val_acc: 0.8973
Epoch 2/10
32000/32000 [=====] - 420s 13ms/step - loss: 0.2588 - acc: 0.8973 - val_loss: 0.2324 - val_acc: 0.9090
Epoch 3/10
32000/32000 [=====] - 418s 13ms/step - loss: 0.2324 - acc: 0.9090 - val_loss: 0.2058 - val_acc: 0.9208
Epoch 4/10
32000/32000 [=====] - 417s 13ms/step - loss: 0.2058 - acc: 0.9208 - val_loss: 0.1806 - val_acc: 0.9324
Epoch 5/10
32000/32000 [=====] - 419s 13ms/step - loss: 0.1806 - acc: 0.9324 - val_loss: 0.1598 - val_acc: 0.9395
Epoch 6/10
32000/32000 [=====] - 419s 13ms/step - loss: 0.1598 - acc: 0.9395 - val_loss: 0.1367 - val_acc: 0.9488
Epoch 7/10
32000/32000 [=====] - 419s 13ms/step - loss: 0.1367 - acc: 0.9488 - val_loss: 0.1238 - val_acc: 0.9555
Epoch 8/10
32000/32000 [=====] - 417s 13ms/step - loss: 0.1238 - acc: 0.9555 - val_loss: 0.1066 - val_acc: 0.9602
Epoch 9/10
32000/32000 [=====] - 417s 13ms/step - loss: 0.1066 - acc: 0.9602 - val_loss: 0.1066 - val_acc: 0.9602
Epoch 10/10
32000/32000 [=====] - 417s 13ms/step - loss: 0.1066 - acc: 0.9602 - val_loss: 0.1066 - val_acc: 0.9602

```

32000/32000 [=====] - 417s 13ms/step - loss: 0.0912 - acc: 0.9677 - v  
 Accuracy: 85.83%



Model with 2 LSTM Layers and dropouts:

```
In [73]: def model():
    model = Sequential()
    model.add(Embedding(top_words+1, embedding_vector_length, input_length=max_review_length))
    model.add(LSTM(100,return_sequences=True))
    model.add(Dropout(0.5))
    model.add(LSTM(50))
    model.add(Dense(1, activation='sigmoid'))
    model.add(Dropout(0.2))
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
    print(model.summary())
    #Refer: https://datascience.stackexchange.com/questions/10615/number-of-parameters-
    nb_epoch = 10
    history = model.fit(X_tr, y_tr, nb_epoch=nb_epoch, batch_size=64, validation_data=(X_test, y_test))
    # Final evaluation of the model
    scores = model.evaluate(X_test, y_test, verbose=0)
    print("Accuracy: %.2f%%" % (scores[1]*100))
    fig,ax = plt.subplots(1,1)
    ax.set_xlabel('epoch')
    ax.set_ylabel('Categorical Crossentropy Loss')
    x = list(range(1,nb_epoch+1))
```

```

        vy = history.history['val_loss']
        ty = history.history['loss']
        plt_dynamic(x, vy, ty, ax)
    model()

```

Layer (type)	Output Shape	Param #
embedding_17 (Embedding)	(None, 600, 32)	160032
lstm_24 (LSTM)	(None, 600, 100)	53200
dropout_10 (Dropout)	(None, 600, 100)	0
lstm_25 (LSTM)	(None, 50)	30200
dense_16 (Dense)	(None, 1)	51
dropout_11 (Dropout)	(None, 1)	0

```

Total params: 243,483
Trainable params: 243,483
Non-trainable params: 0

```

```
None
```

```
Train on 32000 samples, validate on 8000 samples
```

```
Epoch 1/10
```

```
32000/32000 [=====] - 760s 24ms/step - loss: 1.9702 - acc: 0.7525 - va
```

```
Epoch 2/10
```

```
32000/32000 [=====] - 766s 24ms/step - loss: 1.8459 - acc: 0.8067 - va
```

```
Epoch 3/10
```

```
32000/32000 [=====] - 767s 24ms/step - loss: 1.8191 - acc: 0.8213 - va
```

```
Epoch 4/10
```

```
32000/32000 [=====] - 761s 24ms/step - loss: 1.8311 - acc: 0.8253 - va
```

```
Epoch 5/10
```

```
32000/32000 [=====] - 760s 24ms/step - loss: 1.7983 - acc: 0.8317 - va
```

```
Epoch 6/10
```

```
32000/32000 [=====] - 761s 24ms/step - loss: 1.7668 - acc: 0.8395 - va
```

```
Epoch 7/10
```

```
32000/32000 [=====] - 760s 24ms/step - loss: 1.7360 - acc: 0.8458 - va
```

```
Epoch 8/10
```

```
32000/32000 [=====] - 761s 24ms/step - loss: 1.7572 - acc: 0.8512 - va
```

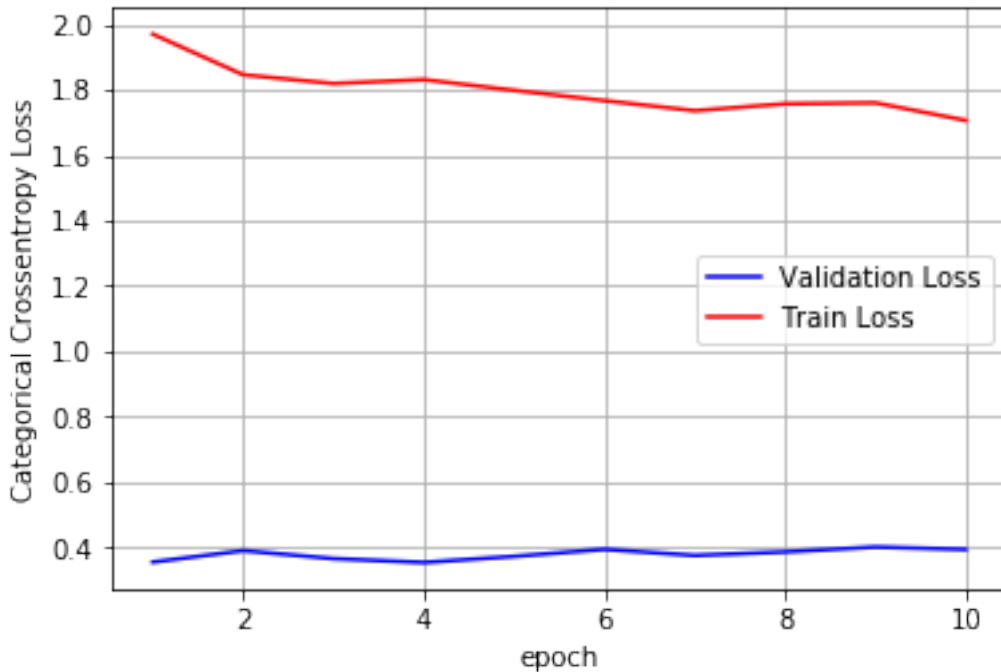
```
Epoch 9/10
```

```
32000/32000 [=====] - 760s 24ms/step - loss: 1.7601 - acc: 0.8546 - va
```

```
Epoch 10/10
```

```
32000/32000 [=====] - 761s 24ms/step - loss: 1.7065 - acc: 0.8627 - va
```

```
Accuracy: 87.30%
```



Model with 3 LSTM layers and dropouts(Using ReLu here):

```
In [84]: def model():
    model = Sequential()
    model.add(Embedding(top_words+1, embedding_vector_length, input_length=max_review_length))
    model.add(LSTM(100,return_sequences=True))
    model.add(Dropout(0.5))
    model.add(LSTM(40,return_sequences=True))
    model.add(Dropout(0.5))
    model.add(LSTM(20))
    model.add(Dense(1,activation='relu',kernel_initializer='he_normal'))
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
    print(model.summary())
    #Refer: https://datascience.stackexchange.com/questions/10615/number-of-parameters-in-a-lstm-model
    nb_epoch = 10
    history = model.fit(X_tr, y_tr, nb_epoch=nb_epoch, batch_size=64,validation_data=(X_test, y_test))
    # Final evaluation of the model
    scores = model.evaluate(X_test, y_test, verbose=0)
    print("Accuracy: %.2f%%" % (scores[1]*100))
    fig,ax = plt.subplots(1,1)
    ax.set_xlabel('epoch')
    ax.set_ylabel('Categorical Crossentropy Loss')
    x = list(range(1,nb_epoch+1))
    vy = history.history['val_loss']
    ty = history.history['loss']
```

```
plt_dynamic(x, vy, ty, ax)
model()
```

Layer (type)	Output Shape	Param #
embedding_28 (Embedding)	(None, 600, 32)	160032
lstm_51 (LSTM)	(None, 600, 100)	53200
dropout_34 (Dropout)	(None, 600, 100)	0
lstm_52 (LSTM)	(None, 600, 40)	22560
dropout_35 (Dropout)	(None, 600, 40)	0
lstm_53 (LSTM)	(None, 20)	4880
dense_24 (Dense)	(None, 1)	21

```
Total params: 240,693
Trainable params: 240,693
Non-trainable params: 0
```

```
None
```

```
Train on 32000 samples, validate on 8000 samples
```

```
Epoch 1/10
```

```
32000/32000 [=====] - 921s 29ms/step - loss: 0.5708 - acc: 0.7575 - va
```

```
Epoch 2/10
```

```
32000/32000 [=====] - 913s 29ms/step - loss: 0.5093 - acc: 0.7662 - va
```

```
Epoch 3/10
```

```
32000/32000 [=====] - 912s 28ms/step - loss: 0.4287 - acc: 0.8385 - va
```

```
Epoch 4/10
```

```
32000/32000 [=====] - 910s 28ms/step - loss: 0.3734 - acc: 0.8599 - va
```

```
Epoch 5/10
```

```
32000/32000 [=====] - 909s 28ms/step - loss: 0.3504 - acc: 0.8782 - va
```

```
Epoch 6/10
```

```
32000/32000 [=====] - 910s 28ms/step - loss: 0.3836 - acc: 0.8439 - va
```

```
Epoch 7/10
```

```
32000/32000 [=====] - 909s 28ms/step - loss: 0.3303 - acc: 0.8679 - va
```

```
Epoch 8/10
```

```
32000/32000 [=====] - 911s 28ms/step - loss: 0.3206 - acc: 0.8698 - va
```

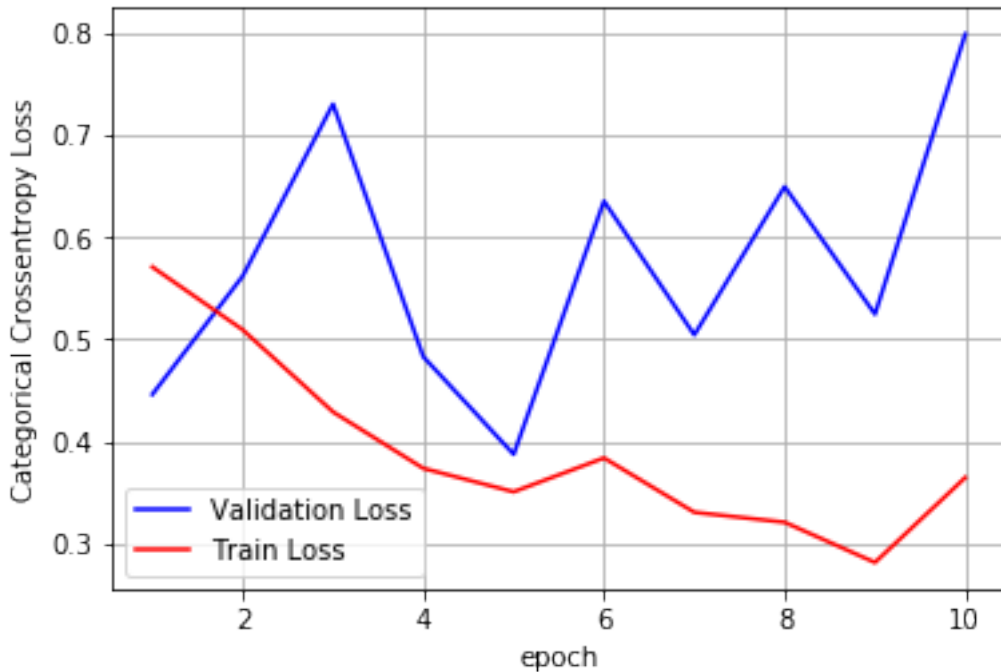
```
Epoch 9/10
```

```
32000/32000 [=====] - 912s 29ms/step - loss: 0.2811 - acc: 0.8729 - va
```

```
Epoch 10/10
```

```
32000/32000 [=====] - 916s 29ms/step - loss: 0.3644 - acc: 0.8334 - va
```

```
Accuracy: 84.19%
```



```
In [88]: from prettytable import PrettyTable
x = PrettyTable()
x.title = "LSTM"
x.field_names = [ "Architecture", "Overfitting", "Accuracy" ]
x.add_row([ "Embedding->LSTM->Sigmoid", "less", "85.83%" ])
x.add_row([ "Embedding->LSTM->Dropout->LSTM->Sigmoid->Dropout", "More", "87.3%" ])
x.add_row([ "Embedding->LSTM->Dropout->LSTM->Dropout->LSTM->ReLU", "less", "84.19%" ])
print(x)
```

Architecture	Overfitting	Accuracy
Embedding->LSTM->Sigmoid	less	85.83%
Embedding->LSTM->Dropout->LSTM->Sigmoid->Dropout	More	87.3%
Embedding->LSTM->Dropout->LSTM->Dropout->LSTM->ReLU	less	84.19%

CONCLUSIONS: Tried 3 different achitectures: single LSTM layer, 2 LSTM layers, and 3 LSTM layers.

Observed that dropout decreases overfitting, but not much of a difference in these models.

Though the second model gave more accuracy, it has comparatively more overfitting.

ReLU worked better than sigmoid as observed from last model.

CNN worked better than LSTM