

File Edit View Insert Cell Kernel Widgets Help

Trusted | Python 3 (ipykernel) O

```
In [6]: 1 import os
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from nltk.corpus import stopwords
7 from nltk.stem import WordNetLemmatizer
8 from nltk.tokenize import word_tokenize
9 from wordcloud import WordCloud,STOPWORDS
10 from bs4 import BeautifulSoup
11 import re,string,unicodedata
12
13 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.svm import LinearSVC
16 from sklearn.ensemble import GradientBoostingClassifier
17 from sklearn.naive_bayes import GaussianNB, MultinomialNB
18 from sklearn.model_selection import train_test_split
19 from sklearn.metrics import classification_report, accuracy_score, confusion_matrix, plot_confusion_matrix, plot_roc_curve,
20 from xgboost.sklearn import XGBClassifier
21
22 import tensorflow as tf
23 from tensorflow.keras.preprocessing.text import Tokenizer
24 from tensorflow.keras.preprocessing.sequence import pad_sequences
25 from tensorflow.keras.callbacks import ModelCheckpoint
26 from tensorflow.keras.layers import Dense, Input, Embedding, LSTM, Dropout, Conv1D, MaxPooling1D, GlobalMaxPooling1D, Dropout, Bidirectional
27 from tensorflow.keras.callbacks import EarlyStopping
28 from tensorflow.keras.models import Model
29 from tensorflow.keras.optimizers import Adam
30 from tensorflow.keras.utils import plot_model
31 import transformers
32 import tokenizers
```

importing csv

```
In [7]: 1 df = pd.read_csv(r'C:\Users\LENOVO\Desktop\IMDB Dataset.csv')
2 df.head()
```

```
Out[7]:
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Peter Matier's "Love in the Time of Money" is...	positive

CLEANING

```
In [8]: 1 # Finding Duplicates in the data
2 num_duplicates = df.duplicated().sum()
3 print(f"the number of duplicates in the the data", format(num_duplicates))
4 review = df['review']
5 duplicate_review = df[review.isin(review[review.duplicated()])].sort_values('review')
6 duplicate_review.head()
7
8 # dropping the duplicate values
9 df.drop_duplicates(inplace =True)
```

the number of duplicates in the the data 418

```
In [9]: 1 # Finding Duplicates in the data
2 num_duplicates = df.duplicated().sum()
3 print(f"the number of duplicates in the the data", format(num_duplicates))
4 review = df['review']
5 duplicate_review = df[review.isin(review[review.duplicated()])].sort_values('review')
6 duplicate_review.head()
7
8 # dropping the duplicate values
9 df.drop_duplicates(inplace =True)
```

the number of duplicates in the the data 0

```
In [10]: 1 # Converting the text values of sentiments into numbers
2 from sklearn import preprocessing
3
4 # LabelEncoder object knows how to understand word Labels.
5 label_encoder = preprocessing.LabelEncoder()
6 df['sentiment'] = label_encoder.fit_transform(df['sentiment'])
```

```
In [11]: 1 stop = stopwords.words('english')
2 wl = WordNetLemmatizer()
```

```
In [12]: 1 mapping = {"ain't": "is not", "aren't": "are not", "can't": "cannot",
2           "'cause": "because", "could've": "could have", "couldn't": "could not",
3           "didn't": "did not", "doesn't": "does not", "don't": "do not", "hadn't": "had not",
4           "hasn't": "has not", "haven't": "have not", "he'd": "he would", "he'll": "he will",
5           "he's": "he is", "how'd": "how did", "how'd'y": "how do you", "how'll": "how will",
6           "how's": "how is", "I'd": "I would", "I'd've": "I would have", "I'll": "I will",
7           "I'll've": "I will have", "I'm": "I am", "I've": "I have", "i'd": "i would",
8           "i'd've": "i would have", "i'll": "i will", "i'll've": "i will have",
9           "i'm": "i am", "i've": "i have", "isn't": "is not", "it'd": "i would",
10          "it'd've": "it would have", "it'll": "it will", "it'll've": "i will have",
11          "it's": "it is", "let's": "let us", "ma'am": "madam", "mayn't": "may not",
12          "might've": "might have", "mighn't": "might not", "mighn't've": "might not have",
13          "must've": "must have", "mustn't": "must not", "mustn't've": "must not have",
14          "needn't": "need not", "needn't've": "need not have", "o'clock": "of the clock",
15          "oughtn't": "ought not", "oughtn't've": "ought not have", "shan't": "shall not",
16          "sha'n't": "shall not", "shan't've": "shall not have", "she'd": "she would",
17          "she'd've": "she would have", "she'll": "she will", "she'll've": "she will have",
18          "she's": "she is", "should've": "should have", "shouldn't": "should not", "shouldn't've": "should not have"}
```

```

19 "shouldn't've": "should not have", "so've": "so have", "so's": "so as", "this's": "this is",
20 "that'd": "that would", "that'd've": "that would have", "that's": "that is",
21 "there'd": "there would", "there'd've": "there would have", "there's": "there is",
22 "here's": "here is", "they'd": "they would", "they'd've": "they would have",
23 "they'll": "they will", "they'll've": "they will have", "they're": "they are",
24 "they've": "they have", "to've": "to have", "wasn't": "was not", "we'd": "we would",
25 "we'd've": "we would have", "we'll": "we will", "we'll've": "we will have",
26 "we're": "we are", "we've": "we have", "weren't": "were not",
27 "what'll": "what will", "what'll've": "what will have", "what're": "what are",
28 "what's": "what is", "what've": "what have", "when's": "when is", "when've": "when have",
29 "where'd": "where did", "where's": "where is", "where've": "where have", "who'll": "who will",
30 "who'll've": "who will have", "who's": "who is", "who've": "who have", "why's": "why is",
31 "why've": "why have", "will've": "will have", "won't": "will not", "won't've": "will not have",
32 "would've": "would have", "wouldn't": "would not", "wouldn't've": "would not have",
33 "y'all": "you all", "y'all'd": "you all would", "y'all'd've": "you all would have",
34 "y'all're": "you all are", "y'all've": "you all have", "you'd": "you would",
35 "you'd've": "you would have", "you'll": "you will", "you'll've": "you will have",
36 "you're": "you are", "you've": "you have" }

```

```

In [13]: 1 clean_data
2 rt BeautifulSoup
3 ing,unicodedata
4 t(text,lemmatize = True):
5 autifulSoup(text, "html.parser") #remove html tags
6 up.get_text()
7 ''.join([mapping[t] if t in mapping else t for t in text.split(" ")]) #expanding chatwords and contracts clearing contractions
8 an= re.compile("("
9         "\u20001F600-\u20001F64F" # emoticons
10        "\u20001F300-\u20001FFF" # symbols & pictographs
11        "\u20001F680-\u20001F6FF" # transport & map symbols
12        "\u20001F1E0-\u20001FF" # flags (ios)
13        "\u200002702-\u200002780"
14        "\u2000024C2-\u20001F251"
15    "]+", flags=re.UNICODE)
16 oji_clean.sub(r'',text)
17 .sub('.(?\s)', ' ',text) #add space after full stop
18 .sub('httpS?', ' ',text) #remove urls
19 .join([word.lower() for word in text if word not in string.punctuation]) #remove punctuation
20 re.split('W+', text) #create tokens
21 size:
22 = " ".join([wl.lemmatize(word) for word in text.split() if word not in stop and word.isalpha()]) #lemmatize
23 |
24 = " ".join([word for word in text.split() if word not in stop and word.isalpha()])
25 xt
26

```

```
In [14]: 1 data_copy = df.copy()
```

```
In [16]: 1 # df['review']=df['review'].apply(clean_text,lemmatize = True)
```

```
In [17]: 1 df.head()
```

```
Out[17]:
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	1
1	A wonderful little production. The...	1
2	I thought this was a wonderful way to spend ti...	1
3	Basically there's a family where a little boy ...	0
4	Petter Mattel's "Love in the Time of Money" is...	1

text processing

```

In [18]: 1 def convert_lowercase(column):
2     column = column.str.lower()
3     return column
4
5 df['review'] = convert_lowercase(df['review'])

```

```

In [19]: 1 def remove_html_tags(text):
2     re_html = re.compile('<.*?>')
3     return re_html.sub(r'', text)
4
5 df['review'] = df['review'].apply(remove_html_tags)

```

```

In [20]: 1 def remove_url(text):
2     re_url = re.compile('https://\S+|www.\S+')
3     return re_url.sub('', text)
4
5 df['review'] = df['review'].apply(remove_url)

```

```
In [21]: 1 df.head()
```

```
Out[21]:
```

	review	sentiment
0	one of the other reviewers has mentioned that ...	1
1	a wonderful little production. the filming tec...	1
2	i thought this was a wonderful way to spend ti...	1
3	basically there's a family where a little boy ...	0
4	petter mattel's "love in the time of money" is...	1

```
In [22]: 1 df.to_csv('imdb_rough_practicemetrics.csv')
```

TEST AND TRAIN

```

In [23]: 1 #####
2 from sklearn.model_selection import train_test_split
3 train_text, test_text, train_labels, test_labels = train_test_split(df["review"],df["sentiment"], test_size=0.33)

```

```

In [24]: 1 #####
2 from tensorflow.keras.preprocessing.text import Tokenizer
3 from tensorflow.keras.preprocessing.sequence import pad_sequences
4
5 #Creates and fits a TensorFlow Tokenizer
6 tokenizer = Tokenizer(num_words=10_000, oov_token='<OOV>')
7 tokenizer.fit_on_texts(train_text)

```

```

8
9 #Creates sequences of numeric representations of words
10 training_sequences = tokenizer.texts_to_sequences(train_text)
11 #pads sequences so they all have the same length
12 training_sequences = pad_sequences(training_sequences, maxlen=20)
13
14 #Process test data in the same way for later evaluation
15 testing_sequences = tokenizer.texts_to_sequences(test_text)
16 testing_sequences = pad_sequences(testing_sequences, maxlen=20)
17
18
19 #prints a sample of the new sequences
20 print('PROCESSED TEXT DATA')
21 print('=====')
22 for i in range(5):
23     print(training_sequences[i], '\n')

PROCESSED TEXT DATA
=====
[ 4 1391 11 19 7 27 5 2 88 312 199 122 106 3
 199 106 4 168 5 98]
[6230 5 1 48 1 7 636 11 19 7 27 5 57 342
 867 874 104 5 29 55]
[ 4 406 1316 15 10 4650 29 344 16 4 337 726 839 62
 11 15 17 12 305 396]
[ 19 17 29 2 4739 3 1044 5 32 1899 1 527 4 2928
 1 21 4 1 1 828]
[ 20 44 1559 388 37 12 12 94 2 394 658 21 39 4
 3293 2061 18 246 14 70]

```

In [25]: 1 train_labels.reset_index(inplace=True, drop=True)

callback

```

In [85]: 1 # CUSTOM CALLBACK:
2 from tensorflow.keras.callbacks import Callback
3
4 class CustomCallback(Callback):
5     def on_epoch_end(self, epoch, logs={}):
6         if logs.get('accuracy') > 0.93:
7             print("Accuracy over 95%... Stopping training")
8             self.model.stop_training = True
9
10 my_callback = CustomCallback()

```

```

In [86]: 1 # PRE-DEFINED CALLBACK:
2 from tensorflow.keras.callbacks import LearningRateScheduler
3
4 #creates a function that updates the learning rate based on the epoch number
5 def scheduler(epoch, lr):
6     if epoch < 2:
7         return 0.01
8     else:
9         return lr * 0.99
10
11 lr_scheduler = LearningRateScheduler(scheduler)

```

model architecture:

```

In [87]: 1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import Embedding, LSTM, Dense, Bidirectional, Dropout
3 from tensorflow.keras.optimizers import Adam
4
5 #input dimensions is equal to number of words tokenized (defined above)
6 input_dim = 10_000
7 # input length will be the length of our padded sequences
8 input_length = 20
9
10
11 #defines a text classifier model
12 model = Sequential([
13     Embedding(input_dim=input_dim, output_dim=64, input_length=input_length),
14     Bidirectional(LSTM(150)),
15     Dropout(0.4),
16     Dense(512, activation='relu'),
17     Dropout(0.5),
18     Dense(1, activation='sigmoid')
19 ])
20
21 model.compile(
22     loss='binary_crossentropy',
23     optimizer=Adam(weight_decay=0.08),
24     metrics=['accuracy']
25 )
26
27 model.summary()

```

Model: "sequential_3"

Layer (type)	Output Shape	Param #
<hr/>		
embedding_3 (Embedding)	(None, 20, 64)	640000
<hr/>		
bidirectional_3 (Bidirectional)	(None, 300)	258000
<hr/>		
dropout_6 (Dropout)	(None, 300)	0
<hr/>		
dense_6 (Dense)	(None, 512)	154112
<hr/>		
dropout_7 (Dropout)	(None, 512)	0
<hr/>		
dense_7 (Dense)	(None, 1)	513
<hr/>		
Total params: 1,052,625		
Trainable params: 1,052,625		
Non-trainable params: 0		

```
In [88]: 1 from sklearn import preprocessing
2
3 # LabelEncoder object knows how to understand word labels.
4 label_encoder = preprocessing.LabelEncoder()
5 df['sentiment']=label_encoder.fit_transform(df['sentiment'])
```

train the model

```
In [89]: 1 history = model.fit(
2     np.array(training_sequences), #must convert to numpy array before sending to model
3     np.array(train_labels),      #must convert to numpy array before sending to model
4     epochs=100,
5     batch_size=128,
6     callbacks=[my_callback, lr_scheduler],verbose=1)
Epoch 1/100
260/260 [=====] - 17s 58ms/step - loss: 0.5283 - accuracy: 0.7336 - lr: 0.0100
Epoch 2/100
260/260 [=====] - 15s 58ms/step - loss: 0.4026 - accuracy: 0.8172 - lr: 0.0100
Epoch 3/100
260/260 [=====] - 15s 58ms/step - loss: 0.3127 - accuracy: 0.8663 - lr: 0.0099
Epoch 4/100
260/260 [=====] - 15s 58ms/step - loss: 0.2394 - accuracy: 0.9024 - lr: 0.0098
Epoch 5/100
260/260 [=====] - 15s 58ms/step - loss: 0.1907 - accuracy: 0.9227 - lr: 0.0097
Epoch 6/100
259/260 [=====] - ETA: 0s - loss: 0.1674 - accuracy: 0.9336Accuracy over 95%... Stopping training
260/260 [=====] - 15s 58ms/step - loss: 0.1674 - accuracy: 0.9336 - lr: 0.0096
```

```
In [90]: 1 # precision for training data prediction
2 from sklearn.metrics import precision_score
3 predicted_train_labels = model.predict(np.array(training_sequences))
4 precision_train = precision_score(train_labels, predicted_train_labels.round())
5 print('Training data precision ',precision_train)

1039/1039 [=====] - 5s 5ms/step
Training data precision = 0.9648892338926375
```

```
In [91]: 1 # recall for training data prediction
2 from sklearn.metrics import recall_score
3 predicted_train_labels = model.predict(np.array(training_sequences))
4 recall_train = recall_score(train_labels, predicted_train_labels.round())
5 print('Training recall data ',recall_train)

1039/1039 [=====] - 5s 5ms/step
Training recall data = 0.9724963890226288
```

```
In [92]: 1 # f1 score for training data prediction
2 from sklearn.metrics import f1_score
3 predicted_train_labels = model.predict(np.array(training_sequences))
4 f1_score_train = f1_score(train_labels, predicted_train_labels.round())
5 print('Training data f1_score ',f1_score_train)

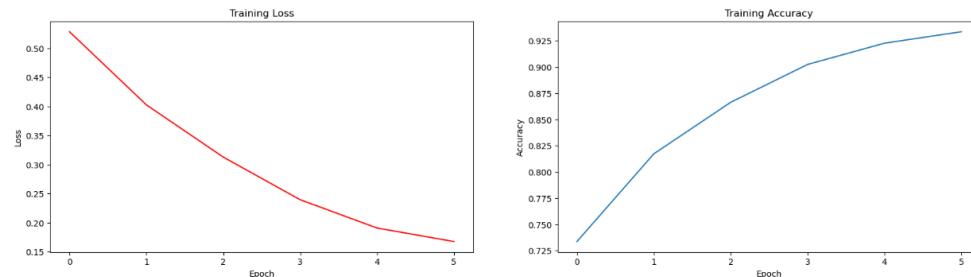
1039/1039 [=====] - 5s 5ms/step
Training data f1_score = 0.968677876689746
```

```
In [ ]: 1
```

plotting the model loss and accuracy:

```
In [93]: 1 #Plots history of model training
2 plt.rcParams["figure.figsize"] = (20,5)
3 fig, axs = plt.subplots(1, 2)
4
5 axs[0].plot(history.history['loss'], color='red')
6 axs[0].set_xlabel('Epoch')
7 axs[0].set_ylabel('Loss')
8 axs[0].set_title('Training Loss')
9
10 axs[1].plot(history.history['accuracy'])
11 axs[1].set_xlabel('epoch')
12 axs[1].set_ylabel('Accuracy')
13 axs[1].set_title('Training Accuracy')
```

Out[93]: Text(0.5, 1.0, 'Training Accuracy')



let's test:

```
In [94]: 1 NEW REVIEW =\
2 """
3 he committed a crime
4 """
```

```
In [95]: 1 #Process the new review the same way the test text was processed
2 new_review_sequence = tokenizer.texts_to_sequences([NEW REVIEW])
3 new_review_sequence = pad_sequences(new_review_sequence, maxlen=20)
4
5 #sends new review to be predicted by the model
6 new_review_prediction = round(model.predict(np.array(new_review_sequence))[0][0])
7 sentiment = "NEGATIVE" if new_review_prediction == 0 else "POSITIVE"
8
9 #displays what the model thinks the sentiment of the review was
```

```

10 print("MOVIE REVIEW:", NEW REVIEW)
11 print("MODEL PREDICTED SENTIMENT:", sentiment)

1/1 [=====] - 0s 17ms/step
MOVIE REVIEW:
he committed a crime

MODEL PREDICTED SENTIMENT: NEGATIVE

eval model:¶

In [98]: 1 loss, accuracy = model.evaluate(testing_sequences, test_labels)
          512/512 [=====] - 3s 5ms/step - loss: 0.7837 - accuracy: 0.7427

In [99]: 1 # recall for training data prediction
          2 from sklearn.metrics import recall_score
          3 predicted_test_labels = model.predict(np.array(testing_sequences))
          4 recall_test = recall_score(test_labels, predicted_test_labels.round())
          5 print('Testing recall data =',recall_test)

          512/512 [=====] - 3s 5ms/step
Testing recall data = 0.7624576681180455

In [100]: 1 f1_score_test = f1_score(test_labels, predicted_test_labels.round())
          2 print('Testing data f1_score =',f1_score_test)

          Testing data f1_score = 0.7496283964563886

In [101]: 1 precision_score_test = precision_score(test_labels, predicted_test_labels.round())
          2 print('Testing data precision_score =', precision_score_test)

          Testing data precision_score = 0.7372237165243831

In [112]: 1 y_pred = model.predict(testing_sequences)
          2

          512/512 [=====] - 2s 5ms/step

In [113]: 1 print(y_pred)
[[9.2250073e-01]
[9.9997121e-01]
[3.7033834e-02]
...
[5.5549435e-06]
[8.0901891e-08]
[9.2241624e-03]]]

In [114]: 1 from sklearn.metrics import confusion_matrix
          2
          3 # Get the true labels for the test set
          4 y_true = test_labels
          5
          6

In [116]: 1 # Generate predicted probabilities for the test set
          2 y_pred_prob = model.predict(testing_sequences)
          3
          4 # Convert probabilities to binary values
          5 y_pred = (y_pred_prob > 0.5).astype(int)

          512/512 [=====] - 2s 5ms/step

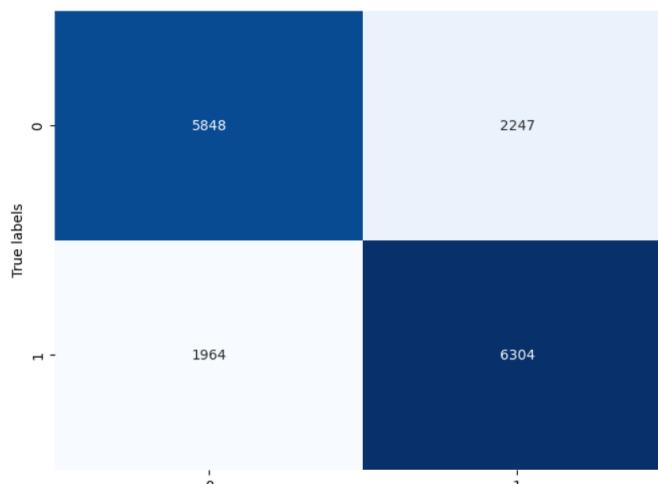
In [117]: 1
          2 # Compute the confusion matrix
          3 cm = confusion_matrix(y_true, y_pred)
          4
          5 # Print the confusion matrix
          6 print(cm)

[[5848 2247]
[1964 6304]]

In [120]: 1 model.save('IMDB_LSTM(hulu).h5')

In [121]: 1 import matplotlib.pyplot as plt
          2 import seaborn as sns
          3 # Plot the confusion matrix
          4 plt.figure(figsize=(8, 6))
          5 sns.heatmap(cm, annot=True, cmap='Blues', fmt='g', cbar=False)
          6 plt.xlabel('Predicted labels')
          7 plt.ylabel('True labels')
          8 plt.show()

```



Predicted labels

```
In [122]: 1 # Extract true positives, true negatives, false positives, false negatives
2 tn, fp, fn, tp = cm.ravel()
3
4 # Print the counts
5 print("True positives:", tp)
6 print("True negatives:", tn)
7 print("False positives:", fp)
8 print("False negatives:", fn)
```

```
True positives: 6304
True negatives: 5848
False positives: 2247
False negatives: 1964
```

```
In [ ]: 1
```