```
In [167...
          import pandas as pd
          import numpy as np
          import warnings
          import re
          import nltk
          import seaborn as sns
          import matplotlib.pyplot as plt
          nltk.download('stopwords')
          from nltk.corpus import stopwords
          from nltk.stem import WordNetLemmatizer
          from nltk.tokenize import sent_tokenize, word_tokenize
          import warnings
          warnings.filterwarnings(action = 'ignore')
          import gensim
          from gensim.models import Word2Vec,Phrases
          import keras
          from keras.preprocessing.sequence import pad_sequences
          from keras.models import Sequential
          from keras.layers import Dense, LSTM, Dropout, Bidirectional
          from keras.layers.embeddings import Embedding
          import tensorflow as tf
          from sklearn.model_selection import train_test_split
          from sklearn.metrics import accuracy_score,confusion_matrix
          nltk.download('wordnet')
         [nltk_data] Downloading package stopwords to
         [nltk data]
                         /Users/molly1998/nltk data...
         [nltk data]
                       Package stopwords is already up-to-date!
         [nltk data] Downloading package wordnet to
                         /Users/molly1998/nltk data...
         [nltk data]
         [nltk data]
                       Package wordnet is already up-to-date!
Out[167... True
         trainset=pd.read_csv('/Users/molly1998/Desktop/python//train.csv')
In [118...
          testset=pd.read csv('/Users/molly1998/Desktop/python//test.csv')
          dataset=pd.concat([trainset,testset],ignore index=True)
          dataset comm=np.concatenate([trainset['text'],testset['text']],axis=0)
          dataset tar=np.concatenate([trainset['sentiment'],testset['sentiment']],axis=0)
          len(dataset comm)
          len(dataset tar)
```

Out[118... 50000

```
In [82]: dataset_comm[1]
```

Out[82]: 'The saddest thing about this "tribute" is that almost all the singers (includin g the otherwise incredibly talented Nick Cave) seem to have missed the whole point where Cohen\'s intensity lies: by delivering his lines in an almost tuneless poise, Cohen transmits the full extent of his poetry, his irony, his all-round humanity, laughter and tears in one.<br/>
'> '> To see some of these singer upstarts make convoluted suffering faces, launch their pathetic squeals in the patent effort to scream "I\'m a singer!," is a true pain. It\'s the same feeling many of you probably had listening in to some horrendous operatic versions of simple songs such as Lennon\'s "Imagine." Nothing, simply nothing gets close to the simp licity and directness of the original. If there is a form of art that doesn\'t need embellishments, it\'s Cohen\'s art. Embellishments cast it in the street looking like the tasteless make-up of sex for sale.<br/>
'> The found myself suffering and suffering through pitiful tributes and awful reinterpretations, all of them entirely lacking the original irony of the master

and, if truth be told, several of these singers sounded as if they had been recruited at some asylum talent show. It\'s Cohen doing a tribute to them by letting them sing his material, really, not the other way around: they may have been friends, or his daughter\'s, he could have become very tender-hearted and in the mood for a gift. Too bad it didn\'t stay in the family.<br/>
'><br/>
Fortunately, but only at the very end, Cohen himself performed his majestic "Tower of Song," but even that flower was spoiled by the totally incongruous background of the U2, all of them carrying the expression that bored kids have when they visit their poor grandpa at the nursing home.<br/>
'><br/>
A sad show, really, and sadder if you truly love Cohen as I do.'

```
In [83]: | en_stops = set(stopwords.words('english'))
In [104...
          def clean review(comments: str) -> str:
              #Remove non-letters
              letters_only = re.compile(r'[^A-Za-z\s]').sub(" ", comments)
              #Convert to lower case
              lowercase_letters = letters_only.lower()
              return lowercase_letters
          def lemmatize(tokens: list) -> list:
              #Lemmatize
              tokens = list(map(WordNetLemmatizer().lemmatize, tokens))
              lemmatized tokens = list(map(lambda x: WordNetLemmatizer().lemmatize(x, "v"),
              #Remove stop words
              meaningful words = list(filter(lambda x: not x in en stops, lemmatized token
              return meaningful words
          def preprocess(review: str, total: int, show progress: bool = True) -> list:
              if show progress:
                  global counter
                  counter += 1
                  print('Processing... %6i/%6i'% (counter, total), end='\r')
              review = clean review(review)
              tokens = word tokenize(review)
              lemmas = lemmatize(tokens)
              return lemmas
          counter=0
          all comments=list(map(lambda x: preprocess(x,len(dataset comm)),dataset comm))
         Processing... 50000/ 50000
In [107... all comments[1]
Out[107... ['saddest',
           'thing',
           'tribute',
           'almost',
           'singer',
           'include'
           'otherwise',
           'incredibly',
           'talented',
           'nick',
           'cave',
           'seem',
           'miss',
           'whole',
           'point'
           'cohen',
```

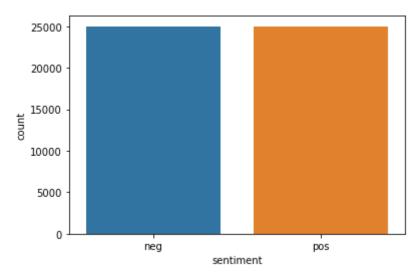
```
'intensity',
'lie',
'deliver',
'line',
'almost'
'tuneless',
'poise',
'cohen',
'transmit',
'full',
'extent',
'poetry',
'irony',
'round',
'humanity',
'laughter',
'tear',
'one',
'br',
'see',
'singer',
'upstart',
'make',
'convolute',
'suffer',
'face',
'launch',
'pathetic',
'squeal',
'patent',
'effort',
'scream',
'singer',
'true',
'pain',
'feel',
'many',
'probably',
'listen',
'horrendous',
'operatic',
'version',
'simple',
'song',
'lennon',
'imagine',
'nothing',
'simply',
'nothing',
'get',
'close',
'simplicity',
'directness',
'original',
'form',
'art',
'need',
'embellishment',
'cohen',
'art',
'embellishment',
'cast',
'street',
'look',
```

```
'like',
'tasteless',
'make',
'sex',
'sale',
'br',
'br',
'cohen',
'tribute',
'find',
'suffer',
'suffer',
'pitiful',
'tribute',
'awful',
'reinterpretation',
'entirely',
'lack',
'original',
'irony',
'master',
'truth',
'tell',
'several',
'singer',
'sound',
'recruit',
'asylum',
'talent',
'show',
'cohen',
'tribute',
'let',
'sing',
'material',
'really',
'way',
'around',
'may',
'friend',
'daughter',
'could',
'become',
'tender',
'hearted',
'mood',
'gift',
'bad',
'stay',
'family',
'br',
'br',
'fortunately',
'end',
'cohen',
'perform',
'majestic',
'tower',
'song',
'even',
'flower',
'wa',
'spoil',
'totally',
'incongruous',
```

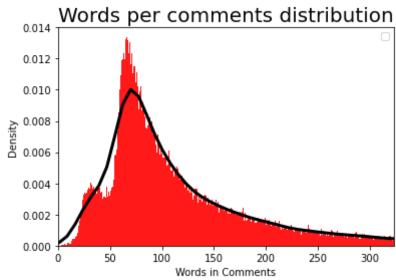
```
'background',
           'u',
           'carry',
           'expression',
           'bore',
           'kid',
           'visit',
           'poor',
           'grandpa',
           'nurse',
           'home',
           'br',
           'br',
           'sad',
           'show',
           'really',
           'sadder',
           'truly',
           'love',
           'cohen']
           dataset tar[dataset tar=="neg"]=0
In [121...
           dataset_tar[dataset_tar=="pos"]=1
           ##check if data is balanced
```

## Out[121... <AxesSubplot:xlabel='sentiment', ylabel='count'>

sns.countplot(x='sentiment', data=dataset)



No handles with labels found to put in legend.



data=data\_senti.copy(deep=True) X\_train=data.loc[:35000,'text']
y\_train=data.loc[:35000,'sentiment'] X\_test=data.loc[35000:,'text']
y\_test=data.loc[35000:,'sentiment'] x\_word2\_train=X\_train.copy() x\_word2\_test=X\_test.copy()

```
In [146... bigrams = Phrases(sentences=all_comments)
    trigrams = Phrases(sentences=bigrams[all_comments])
```

In [147... print(trigrams[bigrams[all\_comments[1]]])

['saddest\_thing', 'tribute', 'almost', 'singer', 'include', 'otherwise', 'incred ibly\_talented', 'nick\_cave', 'seem', 'miss', 'whole', 'point', 'cohen', 'intensi ty', 'lie', 'deliver\_line', 'almost', 'tuneless', 'poise', 'cohen', 'transmit', 'full\_extent', 'poetry', 'irony', 'round', 'humanity', 'laughter\_tear', 'one', 'br', 'see', 'singer', 'upstart', 'make', 'convolute', 'suffer', 'face', 'launch', 'pathetic', 'squeal', 'patent', 'effort', 'scream', 'singer', 'true', 'pain', 'feel', 'many', 'probably', 'listen', 'horrendous', 'operatic', 'versio n', 'simple', 'song', 'lennon', 'imagine', 'nothing', 'simply', 'nothing', 'ge t', 'close', 'simplicity', 'directness', 'original', 'form', 'art', 'need', 'emb ellishment', 'cohen', 'art', 'need', 'emb ellishment', 'cohen', 'art', 'suffer', 'suffer', 'sake', 'sax', 'sale', 'br', 'br', 'cohen', 'tribute', 'find', 'suffer', 'suffer', 'pitiful', 'tribute', 'awful', 'reinterpretation', 'entirely', 'la ck', 'original', 'irony', 'master', 'truth\_tell', 'several', 'singer', 'sound', 'recruit', 'asylum', 'talent', 'show', 'cohen', 'tribute', 'let', 'sing', 'mater ial', 'really', 'way', 'around', 'may', 'friend', 'daughter', 'could', 'become', 'tender', 'hearted', 'mood', 'gift', 'bad', 'stay', 'family', 'br', 'br', 'fortu nately', 'end', 'cohen', 'perform', 'majestic', 'tower', 'song', 'even', 'flowe r', 'wa', 'spoil', 'totally', 'incongruous', 'background', 'u', 'carry', 'expres sion', 'bore', 'kid', 'visit', 'poor', 'grandpa', 'nurse\_home', 'br', 'br', 'sa d', 'show', 'really', 'sadder', 'truly', 'love', 'cohen']

```
('heart', 0.5048468112945557),
          ('chick flick', 0.5030962228775024),
          ('inspirational', 0.5007826089859009),
          ('heartwarming', 0.49950939416885376),
          ('absolutely_love', 0.48973602056503296)]
In [163...
          print(list(w2v model.wv.vocab.keys()).index("love"))
          print(list(w2v_model.wv.vocab.keys()).index("fell_love"))
         157
         6151
In [156...
          %%time
          def vectorize(text, vocabulary):
              keys = list(vocabulary.keys())
              filter unknown = lambda x: vocabulary.get(x, None) is not None
              encode = lambda x: list(map(keys.index, filter(filter_unknown, x)))
              vectorized = list(map(encode, text))
              return vectorized
          padded=pad sequences(vectorize(trigrams[bigrams[all comments]], w2v model.wv.voca
         CPU times: user 3min 33s, sys: 414 ms, total: 3min 34s
         Wall time: 3min 34s
In [200... | X_train, X_test, y_train, y_test = train_test_split(padded,dataset tar,test size
In [201... | X_train=np.asarray(X_train).astype(np.int)
          X_test=np.asarray(X_test).astype(np.int)
          y_test=np.asarray(y_test).astype(np.int)
          y train=np.asarray(y train).astype(np.int)
          from keras.models import Sequential
In [235...
          from keras.layers import Dense, LSTM, Dropout, Bidirectional, BatchNormalization,
          from keras.layers.embeddings import Embedding
          def build model(embedding matrix: np.ndarray, input length: int):
              model = Sequential()
              model.add(Embedding(
                  input dim = embedding matrix.shape[0],
                  output dim = embedding matrix.shape[1],
                  input length = input length,
                  weights = [embedding_matrix],
                  trainable=False))
              #model.add(Bidirectional(LSTM(88, recurrent dropout=0.1)))
              #model.add(Dense(32))
              #model.add(BatchNormalization())
              model.add(Dense(64))
              model.add(Dropout(0.25))
              model.add(Flatten())
              #model.add(Dense(128))
              #model.add(Dropout(0.15))
              model.add(Dense(1, activation='sigmoid'))
              model.summary()
              return model
          model = build model(
              embedding matrix=w2v model.wv.vectors,
              input_length=250)
          model.compile(
              loss="binary crossentropy",
```

```
optimizer='adam',
metrics=['accuracy'])
```

## Model: "sequential 24"

```
Layer (type)
                         Output Shape
                                                Param #
embedding 24 (Embedding)
                         (None, 250, 256)
                                                3269888
                         (None, 250, 64)
dense 51 (Dense)
                                                16448
dropout 44 (Dropout)
                         (None, 250, 64)
flatten 2 (Flatten)
                         (None, 16000)
dense_52 (Dense)
                         (None, 1)
                                                16001
______
Total params: 3,302,337
Trainable params: 32,449
Non-trainable params: 3,269,888
```

```
In [236...
```

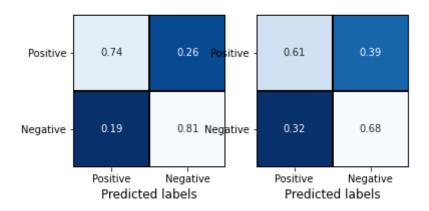
```
mlp_model = model.fit(
    x=X_train,
    y=y_train,
    validation_data=(X_test, y_test),
    batch size=100,
     epochs=20)
```

```
Epoch 1/20
cy: 0.5521 - val loss: 0.6448 - val accuracy: 0.6169
425/425 [============] - 20s 46ms/step - loss: 0.6171 - accura
cy: 0.6572 - val loss: 0.6279 - val accuracy: 0.6497
Epoch 3/20
425/425 [===============] - 19s 44ms/step - loss: 0.5907 - accura
cy: 0.6858 - val loss: 0.6406 - val accuracy: 0.6324
Epoch 4/20
425/425 [============== ] - 19s 44ms/step - loss: 0.5766 - accura
cy: 0.6970 - val loss: 0.6268 - val accuracy: 0.6516
Epoch 5/20
425/425 [============] - 19s 44ms/step - loss: 0.5605 - accura
cy: 0.7103 - val loss: 0.6294 - val accuracy: 0.6540
Epoch 6/20
425/425 [============] - 19s 44ms/step - loss: 0.5511 - accura
cy: 0.7150 - val loss: 0.6547 - val accuracy: 0.6341
425/425 [=============] - 19s 44ms/step - loss: 0.5473 - accura
cy: 0.7175 - val loss: 0.6445 - val accuracy: 0.6477
Epoch 8/20
425/425 [===============] - 19s 44ms/step - loss: 0.5411 - accura
cy: 0.7244 - val_loss: 0.6390 - val_accuracy: 0.6532
Epoch 9/20
425/425 [===============] - 18s 44ms/step - loss: 0.5244 - accura
cy: 0.7366 - val_loss: 0.6510 - val_accuracy: 0.6445
Epoch 10/20
425/425 [===============] - 18s 43ms/step - loss: 0.5247 - accura
cy: 0.7356 - val loss: 0.6511 - val accuracy: 0.6475
Epoch 11/20
425/425 [===============] - 19s 44ms/step - loss: 0.5232 - accura
cy: 0.7376 - val loss: 0.6534 - val accuracy: 0.6447
cy: 0.7480 - val loss: 0.6571 - val accuracy: 0.6492
```

Epoch 13/20

```
425/425 [============== ] - 19s 44ms/step - loss: 0.5142 - accura
        cy: 0.7404 - val loss: 0.6620 - val accuracy: 0.6443
       Epoch 14/20
        cy: 0.7440 - val_loss: 0.6593 - val_accuracy: 0.6493
       Epoch 15/20
        cy: 0.7494 - val_loss: 0.6705 - val_accuracy: 0.6495
       Epoch 16/20
        425/425 [==============] - 19s 44ms/step - loss: 0.4987 - accura
        cy: 0.7492 - val_loss: 0.6792 - val_accuracy: 0.6448
        Epoch 17/20
        425/425 [==============] - 19s 44ms/step - loss: 0.4963 - accura
        cy: 0.7510 - val_loss: 0.6703 - val_accuracy: 0.6492
        425/425 [============== ] - 19s 44ms/step - loss: 0.5032 - accura
        cy: 0.7474 - val_loss: 0.6765 - val_accuracy: 0.6469
       Epoch 19/20
        425/425 [=======================] - 19s 44ms/step - loss: 0.4957 - accura
        cy: 0.7543 - val_loss: 0.6767 - val_accuracy: 0.6516
       Epoch 20/20
        cy: 0.7606 - val loss: 0.6998 - val accuracy: 0.6457
       y train pred = model.predict classes(X train)
In [237...
        y_test_pred = model.predict_classes(X_test)
        def plot_confusion_matrix(y_true, y_pred, ax, class_names, vmax=None,
                               normed=True, title='Confusion matrix'):
            matrix = confusion_matrix(y_true,y_pred)
            if normed:
               matrix = matrix.astype('float') / matrix.sum(axis=1)[:, np.newaxis]
            sns.heatmap(matrix, vmax=vmax, annot=True, square=True, ax=ax,
                     cmap=plt.cm.Blues_r, cbar=False, linecolor='black',
                     linewidths=1, xticklabels=class names)
            ax.set title(title, y=1.20, fontsize=16)
            #ax.set_ylabel('True labels', fontsize=12)
            ax.set_xlabel('Predicted labels', y=1.10, fontsize=12)
            ax.set yticklabels(class names, rotation=0)
        fig, (axis1, axis2) = plt.subplots(nrows=1, ncols=2)
        plot confusion matrix(y test, y test pred, ax=axis2,
                           title='Confusion matrix (test data)',
                           class_names=['Positive', 'Negative'])
        plot confusion matrix(y train, y train pred, ax=axis1,
                           title='Confusion matrix (train data)',
                           class names=['Positive', 'Negative'])
```

## Confusion matrix (train@arta)sion matrix (test data)



In [ ]: