



# 1 Packages

In [1]:

```
1  #Standard python Libraries
2  import pandas as pd
3  import seaborn as sns
4  import matplotlib.pyplot as plt
5  import numpy as np
6  import warnings
7  warnings.filterwarnings(action='ignore')
8
9  # Preprocessing tools
10 from sklearn.model_selection import train_test_split, cross_val_predict, cross_val_score
11 from sklearn.preprocessing import MinMaxScaler, StandardScaler, OneHotEncoder
12 scaler = StandardScaler()
13 from sklearn import metrics
14
15 # # Models & Utilities
16 from sklearn.dummy import DummyClassifier
17 from sklearn.linear_model import LogisticRegression, LogisticRegressionCV
18 from sklearn.ensemble import RandomForestClassifier
19 from sklearn.model_selection import train_test_split
20 from sklearn.linear_model import LogisticRegression
21 from sklearn.metrics import classification_report
22 from sklearn.model_selection import cross_val_score
23 from xgboost import XGBClassifier
24 from sklearn.model_selection import GridSearchCV
25 from sklearn.metrics import plot_confusion_matrix
26 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
27
28 # Warnings
29 import warnings
30 warnings.filterwarnings(action='ignore')
31
32 # NLP Libraries
33 import nltk
34 import collections
35 nltk.download('punkt')
36 from sklearn.manifold import TSNE
37 from nltk.tokenize import word_tokenize
38 np.random.seed(0)
39 import re
40 from nltk.corpus import stopwords
41 from nltk.collocations import *
42 from nltk import FreqDist
43 from nltk import word_tokenize
44 from nltk import ngrams
45 import string
46 from sklearn.feature_extraction.text import TfidfVectorizer
47 from sklearn.metrics import accuracy_score
48 from sklearn.datasets import fetch_20newsgroups
49 from sklearn.ensemble import RandomForestClassifier
50 from sklearn.naive_bayes import MultinomialNB
51
52 # Added
53 nltk.download('stopwords')
54
```

executed in 3.06s, finished 12:22:37 2021-06-21

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\Johnny\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\Johnny\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

Out[1]: True

```
In [52]: 1 df = pd.read_csv('data/IMDB Dataset.csv')
          2
          3 df.head()
```

executed in 724ms, finished 13:36:51 2021-06-21

Out[52]:

	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	Petter Mattei's "Love in the Time of Money" is...	positive

```
In [53]: 1 # Taking a look at our columns
          2 print(df.info())
          3
          4 # Checking for NA data
          5 print(df.isna().sum())
```

executed in 44ms, finished 13:36:51 2021-06-21

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  -
0   review      50000 non-null  object
1   sentiment   50000 non-null  object
dtypes: object(2)
memory usage: 781.4+ KB
None
review      0
sentiment   0
dtype: int64
```

```
In [54]: 1 df['sentiment'].value_counts()
```

executed in 21ms, finished 13:36:51 2021-06-21

Out[54]: positive 25000  
negative 25000  
Name: sentiment, dtype: int64

In [55]: 1 df.describe()

executed in 161ms, finished 13:36:53 2021-06-21

Out[55]:

	review	sentiment
count	50000	50000
unique	49582	2
top	Loved today's show!!! It was a variety and not...	positive
freq	5	25000

In [61]: 1 df[df.duplicated(keep=False)]

executed in 265ms, finished 13:39:44 2021-06-21

Out[61]:

	review	sentiment
42	Of all the films I have seen, this one, The Ra...	negative
84	We brought this film as a joke for a friend, a...	negative
140	Before I begin, let me get something off my ch...	negative
219	Ed Wood rides again. The fact that this movie ...	negative
245	I have seen this film at least 100 times and I...	positive
...	...	...
49912	This is an incredible piece of drama and power...	positive
49950	This was a very brief episode that appeared in...	negative
49984	Hello it is I Derrick Cannon and I welcome you...	negative
49986	This movie is a disgrace to the Major League F...	negative
49991	Les Visiteurs, the first movie about the medie...	negative

824 rows × 2 columns

In [63]: 1 df = df.drop\_duplicates()

executed in 253ms, finished 13:40:27 2021-06-21

### ▼ 1.0.1 remove duplicates?

### 1.0.2 change pos to 1 and neg to 0

## ▼ 1.1 Add features

# ▼ 2 Text Preprocessing



## 2.1 Remove

```
In [64]: 1 first_review = df['review'][0][0:500]
          2 first_review
```

executed in 14ms, finished 13:40:33 2021-06-21

Out[64]: "One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be hooked. They are right, as this is exactly what happened with me.<br /><br />The first thing that struck me about Oz was its brutality and unflinching scenes of violence, which set in right from the word GO. Trust me, this is not a show for the faint hearted or timid. This show pulls no punches with regards to drugs, sex or violence. Its is hardcore, in the classic use of the word.<br /><br />It is called OZ"

```
In [65]: 1  ## From tjhe lessons
          2  from nltk import regexp_tokenize
          3  pattern = r"([a-zA-Z]+(?:'[a-z]+)?)"
          4  regexp_tokenize(first_review,pattern)
```

executed in 19ms, finished 13:40:33 2021-06-21

```
Out[65]: ['One',
          'of',
          'the',
          'other',
          'reviewers',
          'has',
          'mentioned',
          'that',
          'after',
          'watching',
          'just',
          'Oz',
          'episode',
          "you'll",
          'be',
          'hooked',
          'They',
          'are',
          'right',
          'as',
          'this',
          'is',
          'exactly',
          'what',
          'happened',
          'with',
          'me',
          'br',
          'br',
          'The',
          'first',
          'thing',
          'that',
          'struck',
          'me',
          'about',
          'Oz',
          'was',
          'its',
          'brutality',
          'and',
          'unflinching',
          'scenes',
          'of',
          'violence',
          'which',
          'set',
          'in',
          'right',
          'from',
```

```
'the',  
'word',  
'GO',  
'Trust',  
'me',  
'this',  
'is',  
'not',  
'a',  
'show',  
'for',  
'the',  
'faint',  
'hearted',  
'or',  
'timid',  
'This',  
'show',  
'pulls',  
'no',  
'punches',  
'with',  
'regards',  
'to',  
'drugs',  
'sex',  
'or',  
'violence',  
'Its',  
'is',  
'hardcore',  
'in',  
'the',  
'classic',  
'use',  
'of',  
'the',  
'word',  
'br',  
'br',  
'It',  
'is',  
'called',  
'OZ']
```

```
In [66]: 1 stop_words_list = stopwords.words('english')  
         2  
         3 stop_words_list.append('br')
```

executed in 14ms, finished 13:40:33 2021-06-21

```
In [67]: 1 tokens = regexp_tokenize(first_review, pattern)
2
3 tokens
4
5 bad_tags = ['br']
6
7 cleaned_tokens = []
8
9 for token in tokens:
10     if token not in bad_tags:
11         cleaned_tokens.append(token)
12
13 cleaned_tokens[0:5]
14
```

executed in 21ms, finished 13:40:33 2021-06-21

Out[67]: ['One', 'of', 'the', 'other', 'reviewers']

```
In [68]: 1 stop_words_list = stopwords.words('english')
2
3 remove_words = ["i've", "i'm", 'br']
4
5 stop_words_list += remove_words
6
7 stop_words_list
```

executed in 16ms, finished 13:40:34 2021-06-21

Out[68]: ['i',  
'me',  
'my',  
'myself',  
'we',  
'our',  
'ours',  
'ourselves',  
'you',  
"you're",  
"you've",  
"you'll",  
"you'd",  
'your',  
'yours',  
'yourself',  
'yourselves',  
'he',  
'him',  
...]

### 3 Tokenizer



```

In [69]: 1 def my_tokenizer(first_review, stop_words=False,
2           stop_words_add=[],
3           remove_words=remove_words, show_full=False):
4
5           pattern = r"([a-zA-Z]+(?:'[a-z]+)?)"
6           # pattern = r"[a-zA-Z.0-9+#+-/*]*[^\s]"
7
8           tokens = regexp_tokenize(first_review, pattern)
9
10          stop_words_list = []
11
12          if stop_words == True:
13              stop_words_list = stopwords.words('english')
14              stop_words_list += stop_words_add
15
16          stop_words_list += remove_words
17          [x for x in stop_words_list]
18          cleaned_tokens = []
19
20          for token in tokens:
21              if token not in stop_words_list:
22                  cleaned_tokens.append(token)
23
24          if show_full == False:
25              return cleaned_tokens
26          else:
27              return " ".join(cleaned_tokens) #, stop_words_list
28
29          my_tokenizer("i'm going to the store", stop_words=True,
30                      stop_words_add=[], show_full=True)
31

```

executed in 27ms, finished 13:40:34 2021-06-21

Out[69]: 'going store'

```

In [70]: 1 df['reviews_t'] = df['review'].apply(lambda text: my_tokenizer(text, stop_wo
2
3         df['reviews_t']
4

```

executed in 56.1s, finished 13:41:31 2021-06-21

```

Out[70]: 0      [One, reviewers, mentioned, watching, Oz, epis...
1      [A, wonderful, little, production, The, filmin...
2      [I, thought, wonderful, way, spend, time, hot,...
3      [Basically, there's, family, little, boy, Jake...
4      [Petter, Mattei's, Love, Time, Money, visually...
        ...
49995  [I, thought, movie, right, good, job, It, crea...
49996  [Bad, plot, bad, dialogue, bad, acting, idioti...
49997  [I, Catholic, taught, parochial, elementary, s...
49998  [I'm, going, disagree, previous, comment, side...
49999  [No, one, expects, Star, Trek, movies, high, a...
Name: reviews_t, Length: 49582, dtype: object

```

## 4 EDA TO DO:

Break out everything by positive and negative FIRST

visualize pos neg bigrams trigrams etc

Are pos negative review shorter or longer

Are longer words more typical of positive or negative reviews

Code to use from:

<https://medium.com/plotly/nlp-visualisations-for-clear-immediate-insights-into-text-data-and-outputs-9ebfab168d5b> (<https://medium.com/plotly/nlp-visualisations-for-clear-immediate-insights-into-text-data-and-outputs-9ebfab168d5b>)

```
fig = px.bar(long_bigram_df_tidy, title='Comparision: ' + ngrams_list[0] + ' | ' + ngrams_list[1],
x='ngram', y='value', color='variable', template='plotly_white',
color_discrete_sequence=px.colors.qualitative.Bold, labels={'variable': 'Company:', 'ngram': 'N-
Gram'}) fig.update_layout(legend_orientation="h") fig.update_layout(legend=dict(x=0.1, y=1.1))
fig.update_yaxes(title="", showticklabels=False) fig.show()
```

### 4.1 Bag of Words

```
In [71]: 1 # Tokens
2
3 tokens = []
4
5 for row in df['reviews_t']:
6     tokens.extend(row)
7
8 print(f'Number of tokens: {len(tokens)}')
```

executed in 331ms, finished 13:41:31 2021-06-21

Number of tokens: 6466732

```
In [72]: 1 # Frequency distribution
2
3 review_freqdist = FreqDist(tokens)
4 most_common = review_freqdist.most_common(50)
```

executed in 6.11s, finished 13:41:37 2021-06-21

```
In [73]: 1 # Number of unique tokens
2
3 len(review_freqdist)
```

executed in 14ms, finished 13:41:37 2021-06-21

Out[73]: 138151

In [74]:

1 `print(most_common)`

executed in 13ms, finished 13:41:37 2021-06-21

```
[('I', 143636), ('The', 89220), ('movie', 85132), ('film', 76206), ('one', 47899), ('like', 38334), ('This', 29312), ('good', 28333), ('time', 24336), ('It', 23967), ('would', 23762), ('really', 22171), ('story', 22133), ('see', 22017), ('even', 21591), ('much', 18586), ('well', 17971), ('get', 17773), ('bad', 17273), ('people', 16838), ('great', 16680), ('made', 15803), ('first', 15573), ('make', 15543), ('also', 15185), ('way', 15130), ('movies', 14819), ('could', 14808), ('But', 14141), ('characters', 14106), ('think', 13936), ('films', 13410), ('And', 13317), ('seen', 13179), ('character', 13119), ('A', 13092), ('watch', 12806), ('plot', 12577), ('many', 12509), ('two', 12451), ('acting', 12382), ('know', 12229), ('never', 12113), ('life', 12046), ('love', 11786), ('It's', 11783), ('In', 11776), ('show', 11710), ('little', 11569), ('best', 11527)]
```

In [75]:

```
1 total = review_freqdist.N()
2 for word in review_freqdist:
3     review_freqdist[word] /= float(total)
4
5 review_freqdist
```

executed in 675ms, finished 13:41:38 2021-06-21

Out[75]: FreqDist({'I': 0.0222115281721896, 'The': 0.013796767826469382, 'movie': 0.013164609264772377, 'film': 0.01178431393167368, 'one': 0.007406987022192972, 'like': 0.005927878254425883, 'This': 0.004532737710484987, 'good': 0.004381347487417138, 'time': 0.0037632609484976338, 'It': 0.0037061996693229285, ...})

In [76]:

```
1 total_word_count = sum(review_freqdist.values())
2 review_top_10 = review_freqdist.most_common(10)
3 print('Word\t\t\tNormalized Frequency')
4 for word in review_top_10:
5     normalized_frequency = word[1] / total_word_count
6     print('{ } \t\t\t {:.4}'.format(word[0], normalized_frequency))
```

executed in 76ms, finished 13:41:38 2021-06-21

Word	Normalized Frequency
I	0.02221
The	0.0138
movie	0.01316
film	0.01178
one	0.007407
like	0.005928
This	0.004533
good	0.004381
time	0.003763
It	0.003706



## 4.2 Word Clouds

```
In [77]: 1 # !pip install wordcloud
2
3 from wordcloud import WordCloud
```

executed in 9ms, finished 13:41:38 2021-06-21

#### ▼ 4.2.1 To Do:

Break out by positive and negative

#### ▼ 4.2.2 Positive

```
In [78]: 1 # Positive tokens
2
3 df_positive = df['reviews_t'].loc[df['sentiment'] == 'positive']
4
5 tokens_positive = []
6
7 for row in df_positive:
8     tokens_positive.extend(row)
9
10 print(f'Number of tokens: {len(tokens_positive)}')
```

executed in 235ms, finished 13:41:38 2021-06-21

Number of tokens: 3278252

```
In [ ]: 1
```

executed in 186ms, finished 16:56:10 2021-06-18



In [80]: 1 plt.show(wc)

executed in 12ms, finished 13:42:10 2021-06-21

### ▼ 4.2.3 Positive

In [81]:

```
1 # Negative tokens
2
3 df_negative = df['reviews_t'].loc[df['sentiment'] == 'negative']
4
5 tokens_negative = []
6
7 for row in df_negative:
8     tokens_negative.extend(row)
9
10 print(f'Number of tokens: {len(tokens_negative)}')
```

executed in 219ms, finished 13:42:10 2021-06-21

Number of tokens: 3188480



In [ ]:

1

## 4.3 Ngrams

In [83]:

```
1 n_gram = (pd.Series(nltk.ngrams(tokens, 2)).value_counts())[:7]
```

executed in 13.3s, finished 13:42:52 2021-06-21

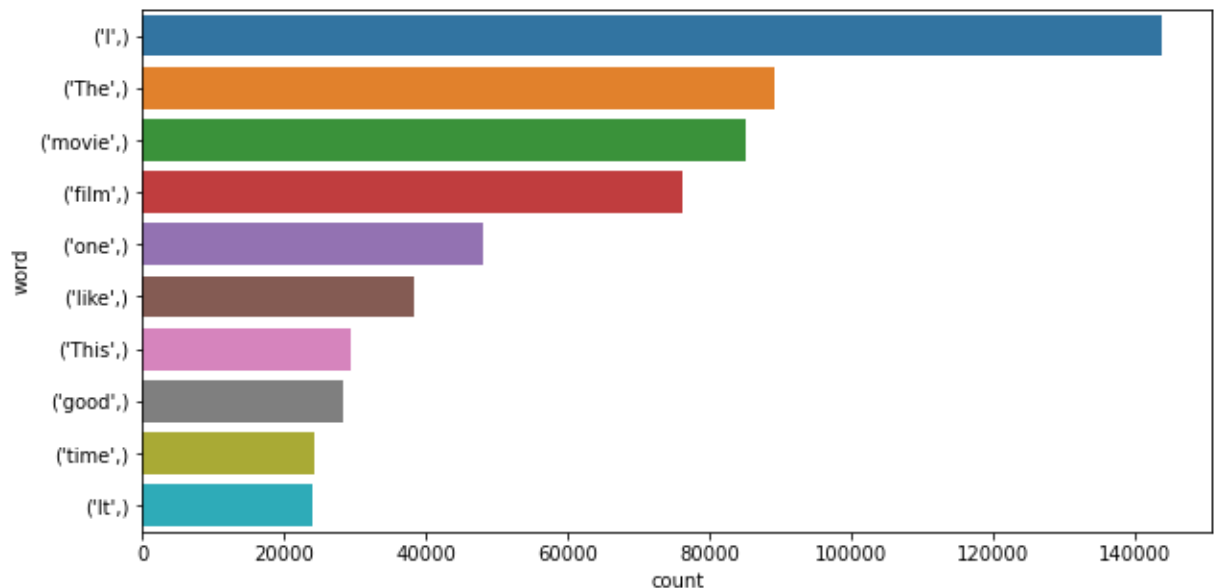
In [84]:

```
1 def plot_ngram(tokens, i):
2     n_gram = (pd.Series(nltk.ngrams(tokens, i)).value_counts())[:10]
3     n_gram_df=pd.DataFrame(n_gram)
4     n_gram_df = n_gram_df.reset_index()
5     n_gram_df = n_gram_df.rename(columns={"index": "word", 0: "count"})
6     print(n_gram_df.head(10))
7     plt.figure(figsize = (10,5))
8     return sns.barplot(x='count',y='word', data=n_gram_df)
9
10 plot_ngram(tokens, 1)
```

executed in 6.19s, finished 13:42:58 2021-06-21

	word	count
0	('I,')	143636
1	('The,')	89220
2	('movie,')	85132
3	('film,')	76206
4	('one,')	47899
5	('like,')	38334
6	('This,')	29312
7	('good,')	28333
8	('time,')	24336
9	('It,')	23967

Out[84]: &lt;AxesSubplot:xlabel='count', ylabel='word'&gt;



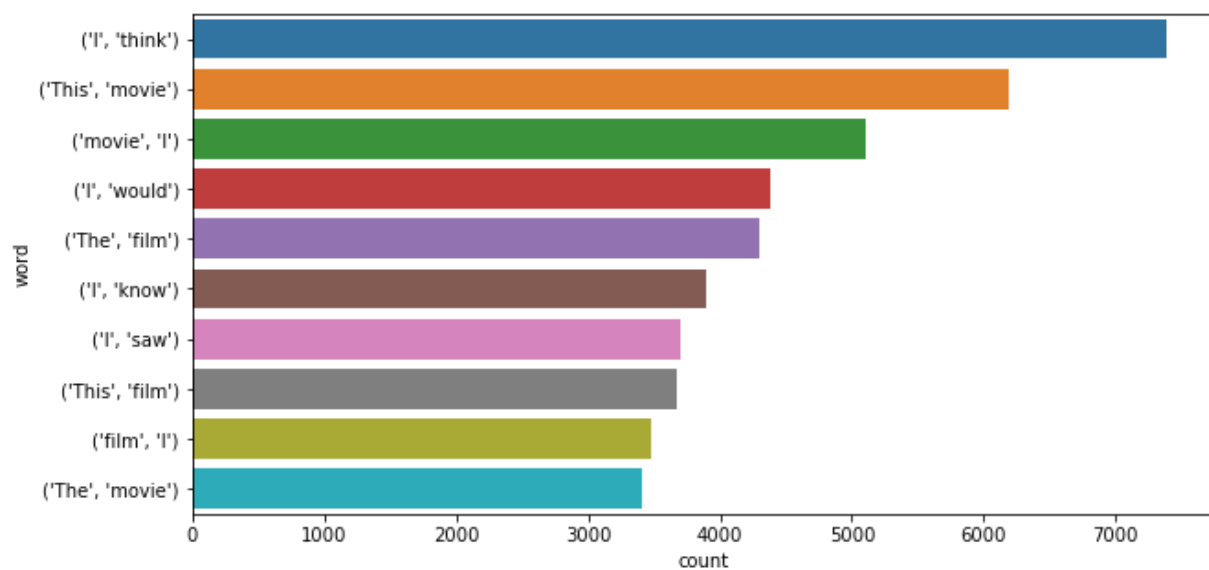


In [85]: 1 plot\_ngram(tokens, 2)

executed in 17.4s, finished 13:43:16 2021-06-21

	word	count
0	(I, think)	7386
1	(This, movie)	6187
2	(movie, I)	5099
3	(I, would)	4387
4	(The, film)	4297
5	(I, know)	3895
6	(I, saw)	3697
7	(This, film)	3670
8	(film, I)	3480
9	(The, movie)	3405

Out[85]: <AxesSubplot:xlabel='count', ylabel='word'>

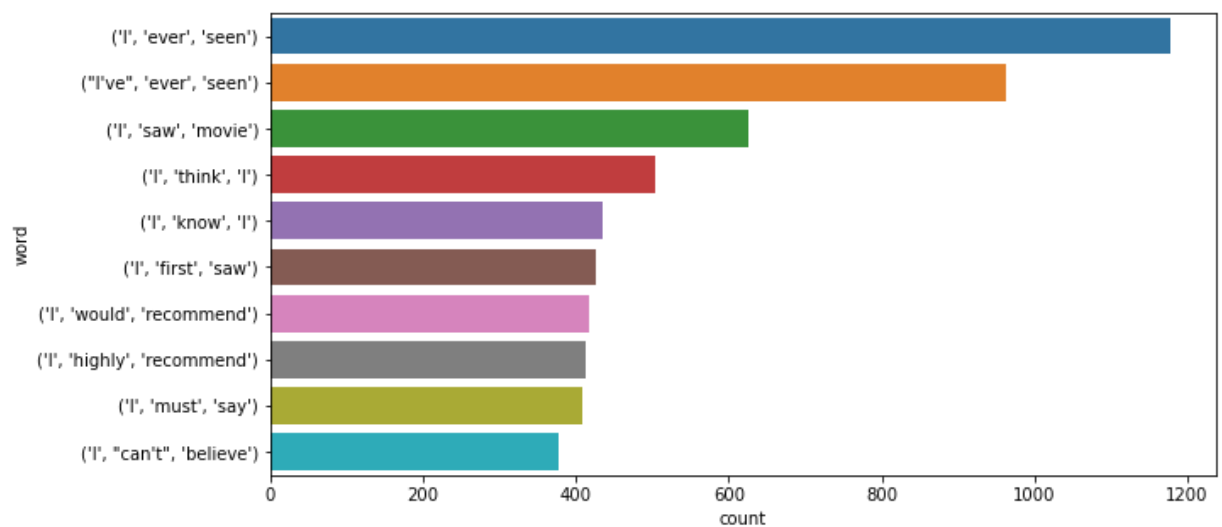


In [86]: 1 plot\_ngram(tokens, 3)

executed in 16.9s, finished 13:43:33 2021-06-21

	word	count
0	(I, ever, seen)	1179
1	(I've, ever, seen)	963
2	(I, saw, movie)	626
3	(I, think, I)	504
4	(I, know, I)	436
5	(I, first, saw)	426
6	(I, would, recommend)	417
7	(I, highly, recommend)	412
8	(I, must, say)	408
9	(I, can't, believe)	377

Out[86]: <AxesSubplot:xlabel='count', ylabel='word'>

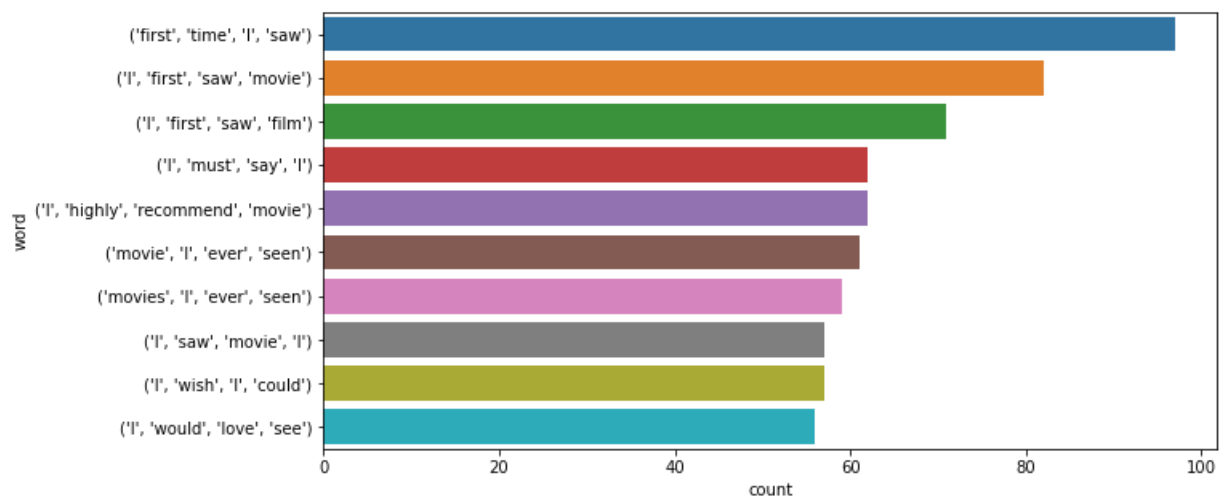


In [87]: 1 plot\_ngram(tokens\_positive, 4)

executed in 9.32s, finished 13:43:42 2021-06-21

	word	count
0	(first, time, I, saw)	97
1	(I, first, saw, movie)	82
2	(I, first, saw, film)	71
3	(I, must, say, I)	62
4	(I, highly, recommend, movie)	62
5	(movie, I, ever, seen)	61
6	(movies, I, ever, seen)	59
7	(I, saw, movie, I)	57
8	(I, wish, I, could)	57
9	(I, would, love, see)	56

Out[87]: <AxesSubplot:xlabel='count', ylabel='word'>



1	tokens
2	

```
Out[88]: ['One',
           'reviewers',
           'mentioned',
           'watching',
           'Oz',
           'episode',
           'hooked',
           'They',
           'right',
           'exactly',
           'happened',
           'The',
           'first',
           'thing',
           'struck',
           'Oz',
           'brutality',
           'unflinching',
           'scenes',
```

1

1

1

1

## ▼ 5.1 TF - IDF

```
In [89]: 1  ## set up text preprocessing pipeline
2
3  tfidf = TfidfVectorizer(strip_accents='unicode',
4                          tokenizer=my_tokenizer,
5                          stop_words=stop_words_list
6                          )
7
8  X = df['review']
9  y = df['sentiment']
10
11  # Do train test split here
12  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,rand
13
14  # Fit both x train and test
15  X_train_tfidf = tfidf.fit_transform(X_train)
16  X_test_tfidf = tfidf.transform(X_test)
17
18
19
```

executed in 17.9s, finished 13:44:00 2021-06-21

```
In [90]: 1  # Logistic Regression with TF-IDF Vectoriser
2  tfidf_log = LogisticRegression(penalty='l2',C=10)
3  tfidf_log.fit(X_train_tfidf, y_train)
```

executed in 5.20s, finished 13:44:05 2021-06-21

Out[90]: LogisticRegression(C=10)

```
In [91]: 1 def evaluate_model(model, X_train=X_train, X_test=X_test, y_train=y_train,
2             y_test=y_test, cmap='Greens', normalize='true',
3             classes=None,figsize=(10,4)):
4
5     # Print model accuracy
6     print(f'Training Accuracy: {model.score(X_train,y_train):.2%}')
7     print(f'Test Accuracy: {model.score(X_test,y_test):.2%}')
8     print('')
9
10    # Print classification report
11    y_test_predict = model.predict(X_test)
12    print(metrics.classification_report(y_test, y_test_predict,
13                                       target_names=classes))
14
15    # Plot confusion matrix
16    fig,ax = plt.subplots(ncols=2,figsize=figsize)
17    metrics.plot_confusion_matrix(model, X_test,y_test,cmap=cmap,
18                                 normalize=normalize,display_labels=classes
19                                 ax=ax[0])
20
21    #Plot ROC curves
22    with sns.axes_style("darkgrid"):
23        curve = metrics.plot_roc_curve(model,X_train,y_train,ax=ax[1])
24        curve2 = metrics.plot_roc_curve(model,X_test,y_test,ax=ax[1])
25        curve.ax_.grid()
26        curve.ax_.plot([0,1],[0,1],ls=':')
27        fig.tight_layout()
28        plt.show()
```

executed in 27ms, finished 13:44:05 2021-06-21

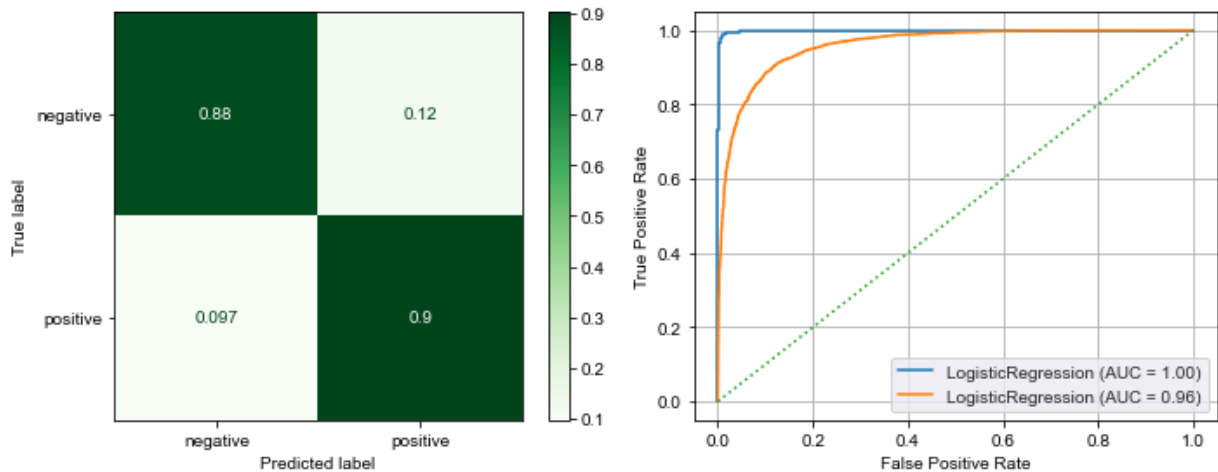
```
In [92]: 1 evaluate_model(tfidf_log, X_train=X_train_tfidf, X_test=X_test_tfidf)
```

executed in 1.40s, finished 13:44:06 2021-06-21

Training Accuracy: 98.81%

Test Accuracy: 89.27%

	precision	recall	f1-score	support
negative	0.90	0.88	0.89	7422
positive	0.88	0.90	0.89	7453
accuracy			0.89	14875
macro avg	0.89	0.89	0.89	14875
weighted avg	0.89	0.89	0.89	14875



In [93]:

```

1 feature_names = tfidf.get_feature_names()
2 feature_names
3
4 df_coef = pd.DataFrame()
5
6 df_coef['features'] = feature_names
7 df_coef['coefficients'] = tfidf_log.coef_.flatten()
8
9 df_coef.sort_values(by='coefficients', ascending=False)

```

executed in 235ms, finished 13:44:07 2021-06-21

Out[93]:

	features	coefficients
<b>27700</b>	excellent	11.641220
<b>34870</b>	great	11.243282
<b>61781</b>	perfect	10.045060
<b>38009</b>	hilarious	10.004170
<b>2309</b>	amazing	9.765448
...	...	...
<b>82963</b>	terrible	-10.316837
<b>9456</b>	boring	-11.348241
<b>5218</b>	awful	-13.200745
<b>90478</b>	waste	-15.084933
<b>92413</b>	worst	-18.461313

93909 rows × 2 columns



## 5.2 Model Iteration



```

In [94]: 1 run = False
2
3 # Initiate new model and perform grid search
4 tfidf_log_hp = LogisticRegression(random_state=8)
5
6
7 if run == True:
8
9 # Define lists of parameters to compare
10     params = {'C':[0.01,0.1,1,10,100],
11               'penalty':['l1','l2','elastic_net'],
12               'solver':['liblinear', 'newton-cg', 'lbfgs', 'sag','saga']}
13
14
15 else:
16     params = {'C':[100],
17               'penalty':['l2'],
18               'solver':['newton-cg']}
19
20
21 # 'C': 10, 'penalty': 'l2', 'solver': 'newton-cg'
22
23 # Run the grid search with a focus on accuracy
24 log_grid_search = GridSearchCV(tfidf_log_hp,params,scoring='accuracy',
25                                verbose=100,
26                                n_jobs=-1)
27
28 # Fit grid search to training data and display best parameters
29 log_grid_search.fit(X_train_tfidf, y_train)
30
31 # Print best parameters
32 log_grid_search.best_params_

```

executed in 23.9s, finished 13:44:31 2021-06-21

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 1 tasks      | elapsed: 13.7s
[Parallel(n_jobs=-1)]: Done 2 out of 5 | elapsed: 14.2s remaining: 21.3
s
[Parallel(n_jobs=-1)]: Done 3 out of 5 | elapsed: 14.7s remaining: 9.8
s
[Parallel(n_jobs=-1)]: Done 5 out of 5 | elapsed: 18.3s remaining: 0.0
s
[Parallel(n_jobs=-1)]: Done 5 out of 5 | elapsed: 18.3s finished

```

Out[94]: {'C': 100, 'penalty': 'l2', 'solver': 'newton-cg'}

```

In [95]: 1 # 'C': 10, 'penalty': 'l2', 'solver': 'newton-cg'

```

executed in 12ms, finished 13:44:31 2021-06-21

```

In [96]: 1 log_grid_search.best_score_

```

executed in 13ms, finished 13:44:31 2021-06-21

Out[96]: 0.8879478226386113

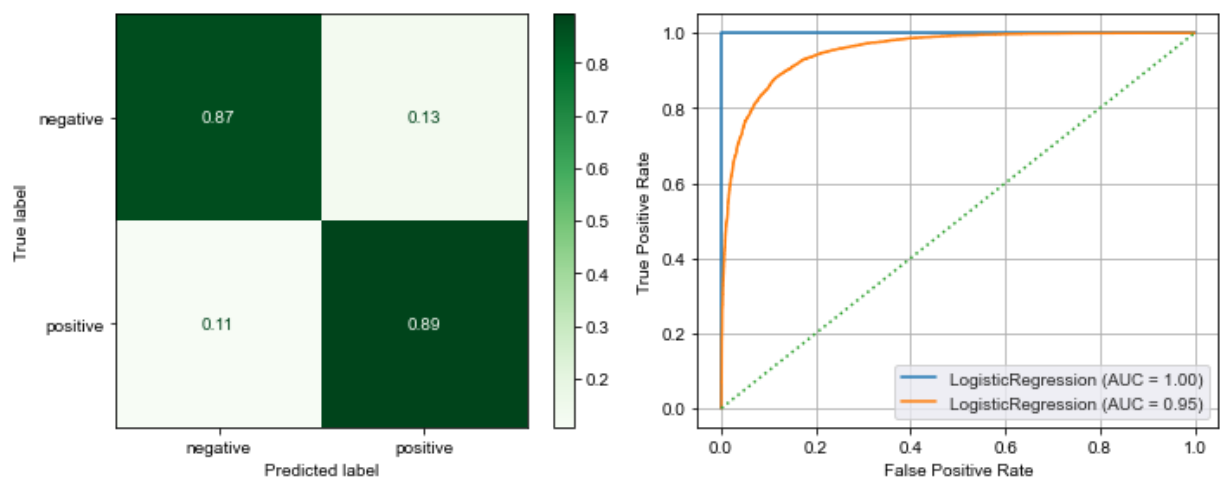
In [97]: 1 evaluate\_model(log\_grid\_search.best\_estimator\_, X\_train=X\_train\_tfidf, X\_test=X\_test\_tfidf)

executed in 1.43s, finished 13:44:32 2021-06-21

Training Accuracy: 100.00%

Test Accuracy: 88.16%

	precision	recall	f1-score	support
negative	0.89	0.87	0.88	7422
positive	0.87	0.89	0.88	7453
accuracy			0.88	14875
macro avg	0.88	0.88	0.88	14875
weighted avg	0.88	0.88	0.88	14875



In [ ]: 1

In [ ]: 1

## ▼ 5.3 Count Vectorizer

```
In [98]: 1  ## set up text preprocessing pipeline
2
3  from sklearn.feature_extraction.text import CountVectorizer
4
5  cv = CountVectorizer(strip_accents='unicode',
6                      tokenizer=my_tokenizer,
7
8                      stop_words=stop_words_list
9                      )
10
11  # X = df['review']
12  # y = df['sentiment']
13
14  # X
15
16  # # Do train test split here
17  # X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, ra
18
19  # Fit both x train and test
20  X_train_cv = cv.fit_transform(X_train)
21  X_test_cv = cv.transform(X_test)
22
23
24
```

executed in 17.5s, finished 13:44:50 2021-06-21

```
In [99]: 1  # X_train_cv
```

executed in 14ms, finished 13:44:50 2021-06-21

```
In [100]: 1  # Logistic Regression with Count Vectoriser
2  cv_log = LogisticRegression(penalty='l2',C=10)
3  cv_log.fit(X_train_cv, y_train)
```

executed in 5.64s, finished 13:44:55 2021-06-21

Out[100]: LogisticRegression(C=10)

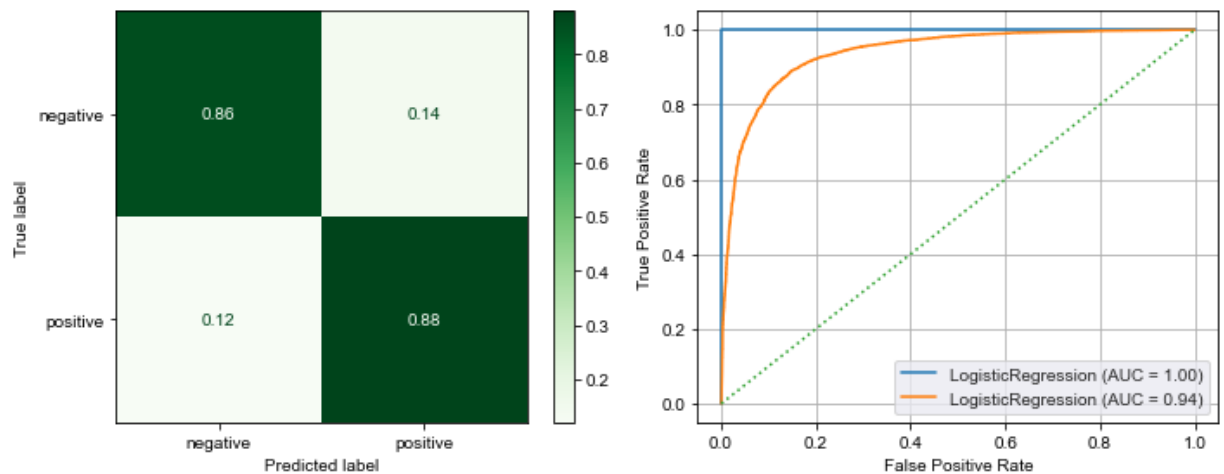
```
In [101]: 1 evaluate_model(cv_log, X_train=X_train_cv, X_test=X_test_cv)
```

executed in 1.61s, finished 13:44:57 2021-06-21

Training Accuracy: 99.99%

Test Accuracy: 86.98%

	precision	recall	f1-score	support
negative	0.88	0.86	0.87	7422
positive	0.86	0.88	0.87	7453
accuracy			0.87	14875
macro avg	0.87	0.87	0.87	14875
weighted avg	0.87	0.87	0.87	14875



## ▼ 5.4 XGBoost

```
In [ ]: 1 model_xgb_tf = XGBClassifier() #Default XGBoost Model with TF_IDF
2 model_xgb_tf.fit(X_train_tfidf, y_train)
3
4 evaluate_model(model_xgb_tf, X_train=X_train_tfidf, X_test=X_test_tfidf)
```

executed in 1m 15.7s, finished 11:00:04 2021-06-21

```
In [ ]: 1
```

```
In [ ]: 1 model_xgb_cv = XGBClassifier() #Default XGBoost Model with TF_IDF
2 model_xgb_cv.fit(X_train_tfidf, y_train)
3
4 evaluate_model(model_xgb_cv, X_train=X_train_cv, X_test=X_test_cv)
```

executed in 1m 12.9s, finished 11:01:17 2021-06-21

```
In [ ]: 1 # Create compare model function
2
3 def model_compare(base_model, grid_search_model):
4
5     # Calculate accuracies
6     base_score = base_model.score(X_test, y_test)
7     grid_score = grid_search_model.score(X_test, y_test)
8
9     #Print accuracies
10    print("--- Base Model ---")
11    model_accuracy(base_model)
12    print('')
13    print("--- Grid Search Model ---")
14    model_accuracy(grid_search_model)
15    print('')
16
17    # If/else function to display best model and score improvement
18    if base_score < grid_score:
19        print(f'Our grid search model outperformed our base model by {(grid_
20    else:
21        print(f'Our base model outperformed our grid search model by {(base_
22
23    # model_compare(model_log, log_grid_search.best_estimator_)
```

executed in 12ms, finished 11:01:17 2021-06-21

```

In [ ]: 1 # If run = True, code will perform full grid search
        2 # If run = False, code will use previously calculated best parameters
        3 run = False
        4
        5 # Instantiate new model for hyperparameter tuning
        6 model_xgb_hp = XGBClassifier(random_state=8)
        7
        8 # Define grid search parameters
        9 if run == True:
        10     param_grid = {
        11         'learning_rate': [0.0001, 0.001, 0.01, 0.1],
        12         'max_depth': [3, 5, 7, 9],
        13         'min_child_weight': [1, 2],
        14         'subsample': [0.5, 0.7, 1],
        15         'n_estimators': [10, 100, 1000]}
        16 else:
        17     param_grid = {
        18         'learning_rate': [0.01],
        19         'max_depth': [5],
        20         'min_child_weight': [2],
        21         'subsample': [0.5],
        22         'n_estimators': [100]}
        23
        24 # Create grid search and train
        25 xgb_grid_search = GridSearchCV(model_xgb_hp, param_grid, scoring='accuracy',
        26                                cv=5, n_jobs=-1, verbose=100)
        27 xgb_grid_search.fit(X_train_tfidf, y_train)
        28
        29 # # Print metrics
        30 # model_compare(model_xgb_tf, xgb_grid_search.best_estimator_)
        31 # print("")
        32 # print(f"Cross Validated Score: {xgb_grid_search.best_score_ :.2%}")
        33 # print("")
        34 # print(f"Optimal Parameters: {xgb_grid_search.best_params_}")

```

executed in 3m 37s, finished 11:04:54 2021-06-21

```

In [ ]: 1 print("")
        2 print(f"Cross Validated Score: {xgb_grid_search.best_score_ :.2%}")
        3 print("")
        4 print(f"Optimal Parameters: {xgb_grid_search.best_params_}")

```

executed in 13ms, finished 11:04:54 2021-06-21

```

In [ ]: 1 # model_compare(model_xgb_tf, xgb_grid_search.best_estimator_)
        2 # print("")
        3 # print(f"Cross Validated Score: {xgb_grid_search.best_score_ :.2%}")
        4 # print("")
        5 # print(f"Optimal Parameters: {xgb_grid_search.best_params_}")

```

executed in 15ms, finished 11:06:26 2021-06-21

```

In [ ]: 1 X_train_tfidf

```

executed in 10ms, finished 11:06:26 2021-06-21

In [ ]:

1



## 6 Word analysis

### 6.1 Frequency Distribution

### 6.2 Normalized word frequency

### 6.3 Bigrams

### 6.4 Mutual information scores

#### 6.4.1 TF-IDF vectorization

Visualize Vector



#### 6.4.2 bag of words

In [ ]:

1

In [ ]:

1



## 6.5 Models

In [ ]:

1

In [ ]:

1

In [ ]:

1



## 7 Scratch



### 7.1 Random Forest

```
In [ ]: 1 # Initiate a random forest model
        2 model_rf = RandomForestClassifier(random_state=8)
        3 model_rf.fit(X_train_tfidf, y_train)
        4
        5 # model_accuracy(model_rf)
```

```
In [ ]: 1 evaluate_model(model_rf, X_train=X_train_tfidf, X_test=X_test_tfidf)
```

```
In [ ]: 1 # params = {
        2 #     'n_estimators': [10, 50, 100, 150],
        3 #     'max_depth': [10, 20, 50, None],
        4 #     'min_samples_split': [5, 30]
        5 # }
        6
        7 # clf_rf = RandomForestClassifier()
        8 # gs_rf = GridSearchCV(clf_rf, param_grid=params, scoring='f1_macro', cv=3,
        9 # gs_rf.fit(X_train_tfidf, y_train)
```

```
In [ ]: 1 # gs_rf.best_score_
```

```
In [ ]: 1 # gs_rf.best_params_
```

```
In [ ]: 1
```

## ▼ 7.2 LSTM

```
In [ ]: 1 embedding_vecor_length = 32
        2 callback = EarlyStopping(monitor='val_loss', patience=2)
        3
        4 model = Sequential()
        5 model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim, input_length=
        6 model.add(LSTM(100, dropout=0.2, recurrent_dropout=0.2))
        7 model.add(Dense(1, activation='sigmoid'))
        8
        9 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
        10
        11 print(model.summary())
```

```
In [ ]: 1 model.fit(X_train1, y_train, epochs=10, batch_size=256, verbose = 1, validation_data=(X_test1, y_test))
```

```
In [ ]: 1 accuracy_score(y_test, model.predict_classes(X_test1))
```

```
In [ ]: 1 plot_history(history)
```

```
In [ ]: 1
```



```
In [ ]: 1 import pandas as pd
2 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
3
4 sentence_1="This is a good job.I will not miss it for anything"
5 sentence_2="This is not good at all"
6
7
8
9 #without smooth IDF
10 print("Without Smoothing:")
11 #define tf-idf
12 tf_idf_vec = TfidfVectorizer(use_idf=True,
13                               smooth_idf=False,
14                               ngram_range=(1,1),stop_words='english') # to use only
15 #transform
16 tf_idf_data = tf_idf_vec.fit_transform([sentence_1,sentence_2])
17
18 #create dataframe
19 tf_idf_dataframe=pd.DataFrame(tf_idf_data.toarray(),columns=tf_idf_vec.get_f
20 print(tf_idf_dataframe)
21 print("\n")
22
23 #with smooth
24 tf_idf_vec_smooth = TfidfVectorizer(use_idf=True,
25                                     smooth_idf=True,
26                                     ngram_range=(1,1),stop_words='english')
27
28
29 tf_idf_data_smooth = tf_idf_vec_smooth.fit_transform([sentence_1,sentence_2]
30
31 print("With Smoothing:")
32 tf_idf_dataframe_smooth=pd.DataFrame(tf_idf_data_smooth.toarray(),columns=tf
33 print(tf_idf_dataframe_smooth)
34
```

```

In [ ]: 1 def plot_coefficients(classifier, feature_names, top_features=20):
2
3     # Access the coefficients from classifier
4     coef = classifier.coef_
5
6     # Access the classes
7     classes = classifier.classes_
8
9     # Iterate the loop for number of classes
10    for i in range(len(classes)):
11
12        i = 0
13        print(classes[i])
14
15        # Access the row containing the coefficients for this class
16        class_coef = coef[i]
17
18
19        # Below this, I have just replaced 'i' in your code with 'class_coef'
20        # Pass this to get top and bottom features
21        top_positive_coefficients = np.argsort(class_coef)[-top_features:]
22        top_negative_coefficients = np.argsort(class_coef)[:top_features]
23
24        # Concatenate the above two
25        top_coefficients = np.hstack([top_negative_coefficients,
26                                     top_positive_coefficients])
27
28        # create plot
29        plt.figure(figsize=(10, 3))
30
31        colors = ["red" if c < 0 else "blue" for c in class_coef[top_coefficients]]
32        plt.bar(np.arange(2 * top_features), class_coef[top_coefficients], c=colors)
33        feature_names = np.array(feature_names)
34
35        # Here I corrected the start to 0 (Your code has 1, which shifted the ticks)
36        plt.xticks(np.arange(0, 1 + 2 * top_features),
37                   feature_names[top_coefficients], rotation=60, ha="right")
38        plt.show()
39    plot_coefficients(tfidf_log, tfidf.get_feature_names(), top_features=20)

```

In [ ]: 1

In [ ]: 1

In [ ]: 1

### ▼ 7.2.1 25 most common words in positive and negative reviews

In [ ]:

```
1  
2 # ax = neg_freq_df.set_index('Word').sort_values('Frequency').plot(kind='bar'  
3 # ax.set(title="25 Most Common Words in Negative Tweets")
```

executed in 11ms, finished 10:54:37 2021-06-21

## ▼ 7.2.2 Word clouds pos / neg

## ▼ 7.2.3 visualize TSNE

## 7.2.4 mean word length postive or negative?. bucket and histogram?

## 7.2.5 Same with word unique percentage

Sentiment across word count

Word count histogram

In [ ]:

```
1 df
```

executed in 56ms, finished 10:54:37 2021-06-21

```

In [ ]: 1 run = False
        2
        3 if run == True:
        4
        5     ## Indirect features
        6     eng_stopwords = set(stopwords.words("english"))
        7
        8     df['count_sent']=df["review"].apply(lambda x: len(re.findall("\n",str(x)
        9     #Word count in each comment:
        10    df['count_word']=df["review"].apply(lambda x: len(str(x).split()))
        11    #Unique word count
        12    df['count_unique_word']=df["review"].apply(lambda x: len(set(str(x).spli
        13    #Letter count
        14    df['count_letters']=df["review"].apply(lambda x: len(str(x)))
        15    #punctuation count
        16    df["count_punctuations"] =df["review"].apply(lambda x: len([c for c in s
        17    #upper case words count
        18    df["count_words_upper"] = df["review"].apply(lambda x: len([w for w in s
        19    #title case words count
        20    df["count_words_title"] = df["review"].apply(lambda x: len([w for w in s
        21    #Number of stopwords
        22    df["count_stopwords"] = df["review"].apply(lambda x: len([w for w in str
        23    #Average length of the words
        24    df["mean_word_len"] = df["review"].apply(lambda x: np.mean([len(w) for w
        25    #Word count percent in each comment:
        26    df['word_unique_percent']=df['count_unique_word']*100/df['count_word']
        27    #Punct percent in each comment:
        28    df['punct_percent']=df['count_punctuations']*100/df['count_word']
        29    #derived features
        30    #Word count percent in each comment:
        31    df['word_unique_percent']=df['count_unique_word']*100/df['count_word']
        32    #derived features
        33    #Punct percent in each comment:
        34    df['punct_percent']=df['count_punctuations']*100/df['count_word']
        35
        36    df.to_json('data/df_json_store.json')
        37
        38 else:
        39     df = pd.read_json('data/df_json_store.json')
        40
        41 df.head(2)

```

executed in 1.58s, finished 10:54:37 2021-06-21



## 7.3 Neural Networks

```
In [ ]: 1 from keras.preprocessing.text import Tokenizer
2 from keras.preprocessing.sequence import pad_sequences
3 from keras.callbacks import ModelCheckpoint
4 from keras.layers import LSTM,Dropout
5 from keras.layers.embeddings import Embedding
6 from keras.preprocessing import sequence
7 from keras.layers import LSTM, Conv1D, MaxPooling1D, Dropout
8 from keras.callbacks import EarlyStopping
9 from keras.models import Sequential
10 from keras.layers import Dense, Dropout
11 from sklearn.metrics import classification_report,confusion_matrix,accuracy_
12 from keras.wrappers.scikit_learn import KerasClassifier
13 import string
14 from keras.preprocessing import text, sequence
15 from keras import layers, models, optimizers
```

executed in 6.98s, finished 11:06:35 2021-06-21

```
In [ ]: 1 def plot_history(history):
2     acc = history.history['accuracy']
3     val_acc = history.history['val_accuracy']
4     loss = history.history['loss']
5     val_loss = history.history['val_loss']
6     x = range(1, len(acc) + 1)
7
8     plt.figure(figsize=(12, 5))
9     plt.subplot(1, 2, 1)
10    plt.plot(x, acc, 'b', label='Training acc')
11    plt.plot(x, val_acc, 'r', label='Validation acc')
12    plt.title('Training and validation accuracy')
13    plt.legend()
14    plt.subplot(1, 2, 2)
15    plt.plot(x, loss, 'b', label='Training loss')
16    plt.plot(x, val_loss, 'r', label='Validation loss')
17    plt.title('Training and validation loss')
18    plt.legend()
```

executed in 26ms, finished 11:06:35 2021-06-21

```
In [ ]: 1 df['sentiment'] = df['sentiment'].replace('positive',1)
2 df['sentiment'] = df['sentiment'].replace('negative',0)
```

executed in 58ms, finished 11:06:35 2021-06-21

```
In [ ]: 1 X_train, X_test, y_train, y_test = train_test_split(df['reviews_t'], df['sen
2
3 X_train, X_valid, y_train, y_valid = train_test_split(X_train, y_train,test_
```

executed in 43ms, finished 11:06:35 2021-06-21

```
In [ ]: 1 [x.shape for x in [X_train,X_valid,X_test]]
```

executed in 13ms, finished 11:06:35 2021-06-21

```
In [ ]: 1 tokenizer = Tokenizer(num_words=5000)
2 tokenizer.fit_on_texts(df['reviews_t'])
3
4 X_train1 = tokenizer.texts_to_sequences(X_train)
5 X_valid1 = tokenizer.texts_to_sequences(X_valid)
6 X_test1 = tokenizer.texts_to_sequences(X_test)
7
8 vocab_size = len(tokenizer.word_index) + 1 # Adding 1 because of reserved 0
9
10 print(X_train[2])
11 print(X_train1[2])
```

executed in 12.0s, finished 11:06:47 2021-06-21

```
In [ ]: 1 seq_lens = [len(s) for s in X_train1]
2 print("average length: %0.1f" % np.mean(seq_lens))
3 print("max length: %d" % max(seq_lens))
```

executed in 25ms, finished 11:06:47 2021-06-21

```
In [ ]: 1 maxlen = 150
2
3 X_train1 = pad_sequences(X_train1, padding='post', maxlen=maxlen)
4 X_valid1 = pad_sequences(X_valid1, padding='post', maxlen=maxlen)
5 X_test1 = pad_sequences(X_test1, padding='post', maxlen=maxlen)
6
7 print(X_train1[2, :])
```

executed in 784ms, finished 11:06:48 2021-06-21

```
In [ ]: 1 vocab_size
```

executed in 12ms, finished 11:06:48 2021-06-21

```
In [ ]: 1 embedding_dim = 50
2 callback = EarlyStopping(monitor='val_loss', patience=2)
3
4 model = Sequential()
5 model.add(layers.Embedding(input_dim=vocab_size, output_dim=embedding_dim, i
6 model.add(layers.Flatten()))
7 model.add(layers.Dense(10, activation='relu'))
8 model.add(layers.Dense(1, activation='sigmoid'))
9
10 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy
11
12 model.summary()
```

executed in 250ms, finished 11:06:48 2021-06-21

```
In [ ]: 1 y_train
```

executed in 14ms, finished 11:06:48 2021-06-21

```
In [ ]: 1 history = model.fit(X_train1, y_train, epochs=10, verbose=True, validation_data
```

executed in 26.5s, finished 11:07:14 2021-06-21

```
In [ ]: 1 accuracy_score(y_test, model.predict_classes(X_test1))
```

executed in 547ms, finished 11:07:15 2021-06-21

```
In [ ]: 1 plot_history(history)
```

executed in 1.07s, finished 11:07:16 2021-06-21

## ▼ 7.4 CNN Model

```
In [ ]: 1 embedding_vecor_length = 32
2 callback = EarlyStopping(monitor='val_loss', patience=2)
3
4 model = Sequential()
5 model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim, input_le
6 model.add(Conv1D(filters=32, kernel_size=3, padding='same', activation='relu
7 model.add(MaxPooling1D(pool_size=2))
8 model.add(LSTM(100))
9 model.add(Dense(1, activation='sigmoid'))
10
11 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accura
12
13 print(model.summary())
```

executed in 617ms, finished 11:07:17 2021-06-21

```
In [ ]: 1 model.fit(X_train1, y_train, epochs=10, batch_size=256, verbose = 1, validatio
```

executed in 4m 9s, finished 11:11:26 2021-06-21

```
In [ ]: 1 accuracy_score(y_test, model.predict_classes(X_test1))
```

executed in 5.61s, finished 11:11:31 2021-06-21

```
In [ ]: 1 plot_history(history)
```

executed in 584ms, finished 11:11:32 2021-06-21

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

In [ ]:

1

In [ ]:

1