



1 IMDB Review Sentiment Analysis

- Student name: Johnny Dryman
- Student pace: full time
- Scheduled project review date/time: 5/27/2021
- Instructor name: James Irving

1.1 Business Problem

IMDB is the world's most popular and authoritative source for movie, TV and celebrity content, designed to help fans explore the world of movies and shows and decide what to watch.

One of the most popular features on the site is its user reviews. Users are able to give each movie a score between 1-10 along with a written review. The average score is a recognized metric in the industry, and IMDB's top movies based on user reviews is a coveted list for impactful films.

One of the issues inherent with user submitted scores is that the 1-10 rating system might mean different things to different reviewers. For example, one reviewer think that a 'good' movie deserves a 6/10, and another might think a 'good' movie deserves an 8/10. Every user has their own evaluation metrics, and this is often determined internally rather than by a specific standard.

1.2 Natural Language Processing

Using Natural Language Processing (NLP), we can create machine learning models that might help us get closer at the core truth of what people are saying with their written reviews.

Although not as verbose as a scale from 1-10, we found a dataset where 50,000 IMDB reviews were denoted as 'positive' or 'negative.' Using this dataset, we can train a model to classify reviews as positive or negative. In theory, distilling the reviews in this fashion would help us to generate a score based on true sentimentality rather than a sliding scale that has different meanings to different users.



2 Packages

We will first want to install the Python packages we will need to perform data import, exploraty data analysis, machine learning modeling, and natural language processing.

In [54]:

```

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

```

executed in 219ms, finished 15:29:18 2021-06-23

[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!
```

3 Data

The data for this project was sourced from Kaggle, a data science community featuring great datasets for exploration and analysis. The dataset contained 25,000 positive reviews and 25,000 negative reviews, totaling 50,000. The only other data provided were 'positive' and 'negative' classifications. According to the source, the reviews were "highly polar," meaning they were strongly positive or strongly negative.

```
In [2]: 1 # Import .csv file obtained from Kaggle
        2 df = pd.read_csv('data/IMDB Dataset.csv')
        3
        4 # View first 5 rows
        5 df.head()
```

executed in 699ms, finished 14:23:05 2021-06-23

Out[2]:

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

executed in 43ms, finished 14:23:05 2021-06-23

```
<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
```

Fortunately our data has no nulls that we need to worry about.

```
In [4]: 1 # View class balance
        2 df['sentiment'].value_counts()
```

executed in 26ms, finished 14:23:05 2021-06-23

```
Out[4]: positive    25000
        negative    25000
        Name: sentiment, dtype: int64
```

The reviews are also very balanced across positive and negative sentiment.

```
In [5]: 1 # View pandas 'describe' to check for data issues
        2 df.describe()
```

executed in 167ms, finished 14:23:05 2021-06-23

```
Out[5]:
```

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

There are a few hundred duplicates in the data. Removing them will not violate any data science best practices.

```
In [6]: 1 # Drop duplicates from dataframe
        2 df = df.drop_duplicates()
```

executed in 247ms, finished 14:23:06 2021-06-23

▼ 4 Text Preprocessing

One of the most crucial aspects of NLP is preparing the data for machine learning models. To do this, we will need to remove punctuation, symbols, and 'stopwords' - common English language words that while important for communication, are less helpful our models.

```
In [7]: 1 # Taking a glance at the first 500 characters of our first review
        2 first_review = df['review'][0][0:500]
        3 first_review
```

executed in 14ms, finished 14:23:06 2021-06-23

```
Out[7]: "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 unflinchi
ng scenes of violence, which set in right from the word GO. Trust me, this is n
ot a show for the faint hearted or timid. This show pulls no punches with regar
ds to drugs, sex or violence. Its is hardcore, in the classic use of the word.<
br /><br />It is called OZ"
```

Let's take a look at how NLTK, a robust Python library for NLP, will 'tokenize' our first review.

Tokenization is the process of splitting our reviews into a list of strings rather than one large string.

```
In [8]: 1 # Designating pattern for regex tokenizer
        2 pattern = r"([a-zA-Z]+(?:'[a-z]+')?)"
        3
        4 # Tokenize our first review based on the pattern
        5 regexp_tokenize(first_review, pattern)
```

executed in 26ms, finished 14:23:06 2021-06-23

```
Out[8]: ['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']
```

▼ 4.1 Tokenizer

By creating a tokenizer function, we will be able to repeat the tokenization process for without breaking down the entire process. This will also allow us to insert a tokenizer into our TF-IDF and Count Vectorizers later in this notebook. Should we choose, this method will also allow us to evaluate the performance of our models with different tokenizers or a modified version of this one.

```

In [11]: 1 # Custom 'remove words' list based on experimentation
2 # Removes a few contractions as well as the HTML tag 'br'
3 remove_words = ["i've", "i'm", 'br']
4
5 # Create stop words list and append remove_words list
6 stop_words_list = stopwords.words('english')
7 stop_words_list += remove_words
8
9 # Define tokenizer function
10 def my_tokenizer(review, stop_words=False,
11                 stop_words_add=[],
12                 remove_words=remove_words, show_full=False):
13
14     # Determine pattern for regexp_tokenize
15     pattern = r"([a-zA-Z]+(?:'[a-z]+)?)"
16
17     # Convert review into tokens based on our pattern
18     tokens = regexp_tokenize(review, pattern)
19
20     # Instantiating empty stopwords list
21     stop_words_list = []
22
23     # Option to insert stopwords for testing
24     # Stop words will be listed separately for our vectorizer
25     if stop_words == True:
26         stop_words_list = stopwords.words('english')
27         stop_words_list += stop_words_add
28
29     # Remove additional words from our tokens based on remove_words list
30     stop_words_list += remove_words
31     [x.lower() for x in stop_words_list]
32     cleaned_tokens = []
33
34     # Add word to token list if word not in stop words list
35     for token in tokens:
36         if token.lower() not in stop_words_list:
37             cleaned_tokens.append(token.lower())
38
39     # Return tokens with option to preview
40     if show_full == False:
41         return cleaned_tokens
42     else:
43         return " ".join(cleaned_tokens) + "#, stop_words_list
44
45 # Print before and after tokenization
46 print("First review before tokenization:")
47 print('')
48 print(first_review)
49 print('')
50 print("First review after tokenization:")
51 print('')
52 print(my_tokenizer(first_review, stop_words=True,
53                   stop_words_add=[], show_full=True))

```

executed in 25ms, finished 14:24:57 2021-06-23

First review before tokenization:

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.

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.

It is called OZ

First review after tokenization:

one reviewers mentioned watching oz episode hooked right exactly happened first thing struck oz brutality unflinching scenes violence set right word go trust s how faint hearted timid show pulls punches regards drugs sex violence hardcore classic use word called oz

Tokenization has cut our word count significantly and has removed elements that will not be helpful for our machine learning models.

We will leave our original review data intact while also creating a separate column with tokenized reviews. This will make it easier to perform our exploratory data analysis.

```
In [12]: 1 # Create new column with tokenized reviews
2 df['reviews_t'] = df['review'].apply(lambda text: my_tokenizer(text, stop_wo
3
4 # Preview first 5
5 df['reviews_t'].head())
```

executed in 53.9s, finished 14:25:51 2021-06-23

```
Out[12]: 0 [one, reviewers, mentioned, watching, oz, epis...
1 [wonderful, little, production, filming, techn...
2 [thought, wonderful, way, spend, time, hot, su...
3 [basically, there's, family, little, boy, jake...
4 [petter, mattei's, love, time, money, visually...
Name: reviews_t, dtype: object
```

5 Exploratory Data Analysis

Now that we have our tokenized reviews, we can dig into the reviews to see what we can learn.

First, we'll split our reviews between positive and negative classifiers. There will be instances where we'll want to see how they're different and also instances where we'll want to examine the entire corpus.

```
In [13]: 1 # Create new dataframes separating positive and negative reviews
2 df_pos = df['reviews_t'].loc[df['sentiment'] == 'positive']
3 df_neg = df['reviews_t'].loc[df['sentiment'] == 'negative']
4
5 # Instantiating empty positive / negative / total token lists
6 tokens = []
7 tokens_pos = []
8 tokens_neg = []
9
10 # Populating token lists from p/n/t dataframes
11 for row in df['reviews_t']:
12     tokens.extend(row)
13 for row in df_pos:
14     tokens_pos.extend(row)
15 for row in df_neg:
16     tokens_neg.extend(row)
17
18 # Print number of p/n tokens
19 print(f'Total corpus tokens: {len(tokens)}')
20 print(f'Number of positive tokens: {len(tokens_pos)}')
21 print(f'Number of negative tokens: {len(tokens_neg)}')
```

executed in 508ms, finished 14:25:51 2021-06-23

Total corpus tokens: 5903613
Number of positive tokens: 3005174
Number of negative tokens: 2898439

5.1 Frequency Distribution

Frequency distributions will show us how often tokens appear in our reviews. The results can be informative especially when comparing between classifiers.

```
In [14]: 1 # Instantiating p/n FreqDists
2 corpus_freqdist = FreqDist(tokens)
3 pos_freqdist = FreqDist(tokens_pos)
4 neg_freqdist = FreqDist(tokens_neg)
```

executed in 10.0s, finished 14:26:01 2021-06-23

In [15]:

```

1 # View top 20 most frequent terms in p/n
2 print("Top 20 most frequent terms in positive reviews:")
3 print(pos_freqdist.most_common(20))
4 print('')
5 print("Top 20 most frequent terms in negative reviews:")
6 print(neg_freqdist.most_common(20))

```

executed in 76ms, finished 14:26:01 2021-06-23

Top 20 most frequent terms in positive reviews:

```

[('film', 40890), ('movie', 37300), ('one', 26920), ('like', 17651), ('good', 14965), ('great', 12888), ('story', 12835), ('time', 12693), ('well', 12682), ('see', 12212), ('also', 10761), ('really', 10707), ('would', 10425), ('even', 9574), ('first', 9186), ('much', 9176), ('love', 8649), ('people', 8519), ('best', 8478), ('get', 8234)]

```

Top 20 most frequent terms in negative reviews:

```

[('movie', 49173), ('film', 36318), ('one', 25776), ('like', 22192), ('even', 15095), ('good', 14576), ('bad', 14563), ('would', 13611), ('really', 12216), ('time', 12197), ('see', 10589), ('story', 10032), ('get', 9990), ('much', 9977), ('make', 9263), ('people', 9185), ('could', 8958), ('made', 8707), ('well', 8424), ('first', 8246)]

```

Interestingly, the most frequently used language between positive and negative reviews is mostly shared. Even the word 'good' appears in our negative reviews 14,576 times. There are some words that aren't quite as useful such as 'like,' 'would,' 'even,' etc. but these are still useful in computing ngrams, so we will leave them in our dataset.

It's worth taking a brief look at the normalized frequency as well.

In [16]:

```

1 # Define and calculate total word count for token values
2 total_word_count = sum(corpus_freqdist.values())
3
4 # Create a variable for just the top 10 most common words
5 review_top_10 = corpus_freqdist.most_common(10)
6
7 # Print top 10 words with their highest frequencies
8 print('Word\t\t\tNormalized Frequency')
9 for word in review_top_10:
10     normalized_frequency = word[1] / total_word_count
11     print('{ } \t\t\t {:.4}'.format(word[0], normalized_frequency))

```

executed in 58ms, finished 14:26:01 2021-06-23

Word	Normalized Frequency
movie	0.01465
film	0.01308
one	0.008926
like	0.006749
good	0.005004
time	0.004216
even	0.004179
would	0.004071
really	0.003883
story	0.003873

This isn't particularly helpful, but we can see that 'movie' and 'film' both appear a little more than once every 100 words, which isn't too surprising.

5.2 Word Clouds

While not the most informative in a statistical sense, word clouds are engaging visuals that allow for a much more digestible interpretation of the most frequently used words in any given corpus. It's worth taking a look to see if anything strikes us.

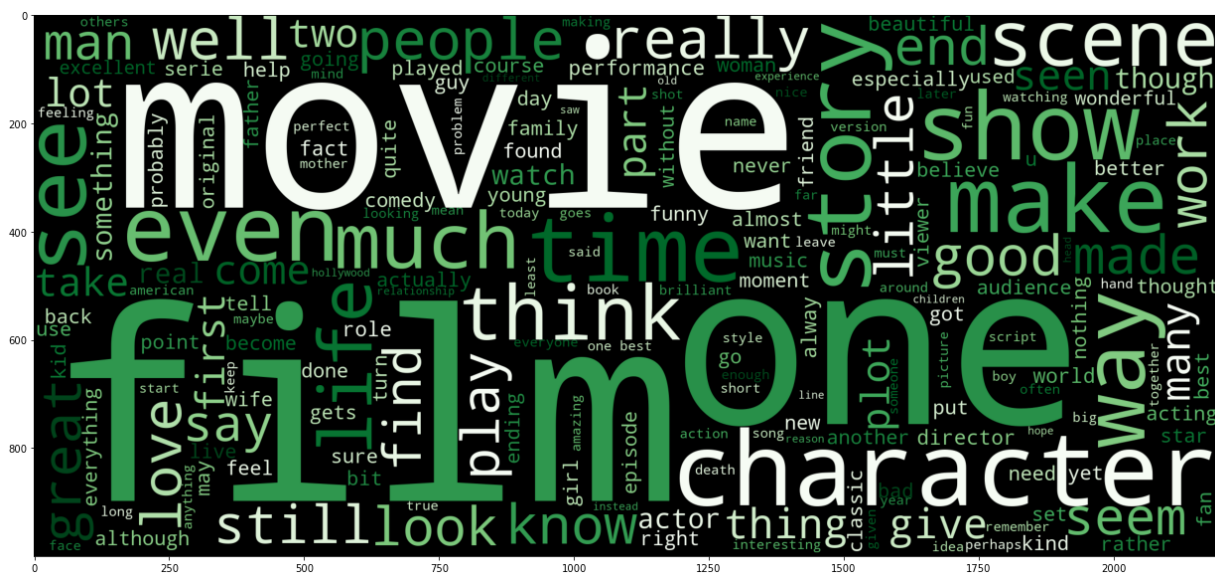
In [17]:

```

1  # Define function for generating word clouds
2  def draw_wordcloud(tokens, colormap):
3
4      # Instantiate plot
5      plt.figure(figsize = (22,10))
6
7      # Instantiate wordcloud
8      wc = WordCloud(max_words = 200,
9                      width = 2200, height = 1000,
10                     colormap=colormap).generate(" ".join(tokens_pos))
11
12     #Show wordcloud
13     plt.imshow(wc)
14
15     # Plot word cloud with positive tokens
16     draw_wordcloud(tokens_pos, "Greens")

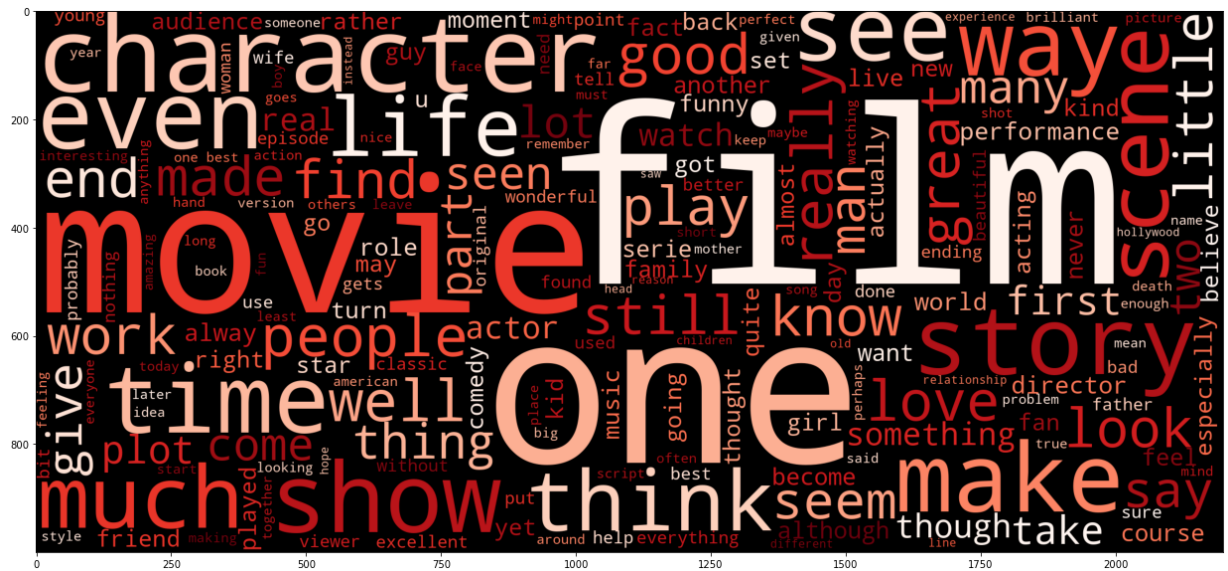
```

executed in 36.2s, finished 14:26:38 2021-06-23



```
In [18]: 1 # Plot word cloud with negative tokens
          2 draw_wordcloud(tokens_neg, "Reds")
```

executed in 35.9s, finished 14:27:13 2021-06-23



Aside from the color palettes, the results match what we found in our frequency distributions.

▼ 5.3 Ngrams

Ngrams are another useful tool for exploratory analysis. Instead of singular tokens, ngrams demonstrate the frequency of specific phrases. Let's take at ngrams with one, two, three, and four words with our positive and negative dataframes.

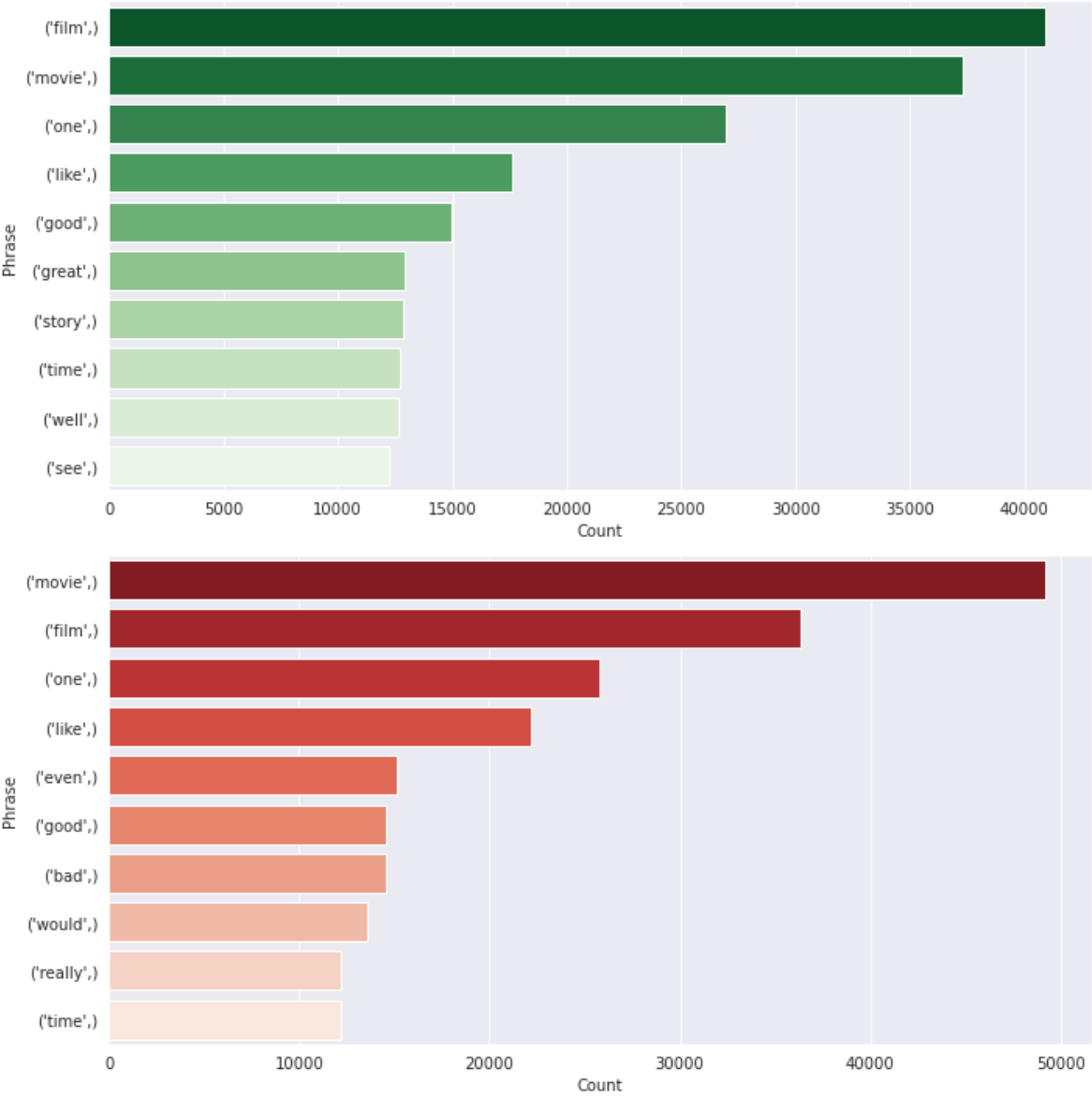
In [20]:

```

1  # Define function for plotting horizontal bar charts based on ngrams
2  def plot_ngram(i, tokens_pos=tokens_pos, tokens_neg=tokens_neg):
3
4      # Setting up ngrams depending on our specified value for 'i'
5      n_gram_pos = (pd.Series(nltk.ngrams(tokens_pos, i)).value_counts())[:10]
6      n_gram_neg = (pd.Series(nltk.ngrams(tokens_neg, i)).value_counts())[:10]
7
8      # Creating p/n dataframes
9      n_gram_df_pos = pd.DataFrame(n_gram_pos)
10     n_gram_df_neg = pd.DataFrame(n_gram_neg)
11
12     # Resetting index for labeling on plots
13     n_gram_df_pos = n_gram_df_pos.reset_index()
14     n_gram_df_neg = n_gram_df_neg.reset_index()
15
16     # Renaming plots
17     n_gram_df_pos = n_gram_df_pos.rename(columns={'index': 'Phrase', 0: 'Cou
18     n_gram_df_neg = n_gram_df_neg.rename(columns={'index': 'Phrase', 0: 'Cou
19
20     # Setting seaborn grid style to 'darkgrid'
21     with sns.axes_style('darkgrid'):
22
23         # Setting up two figures to stack on top of each other
24         fig = plt.figure(figsize = (10,10))
25         ax1 = fig.add_subplot(211)
26         ax2 = fig.add_subplot(212)
27
28         # Assigning each barplot to positive and negative ngram sets
29         sns.barplot(ax=ax1, x='Count', y='Phrase', data=n_gram_df_pos, palett
30         sns.barplot(ax=ax2, x='Count', y='Phrase', data=n_gram_df_neg, palett
31
32     # Return plot with tight layout
33     plt.tight_layout()
34
35     # Plot one word ngrams
36     plot_ngram(1)

```

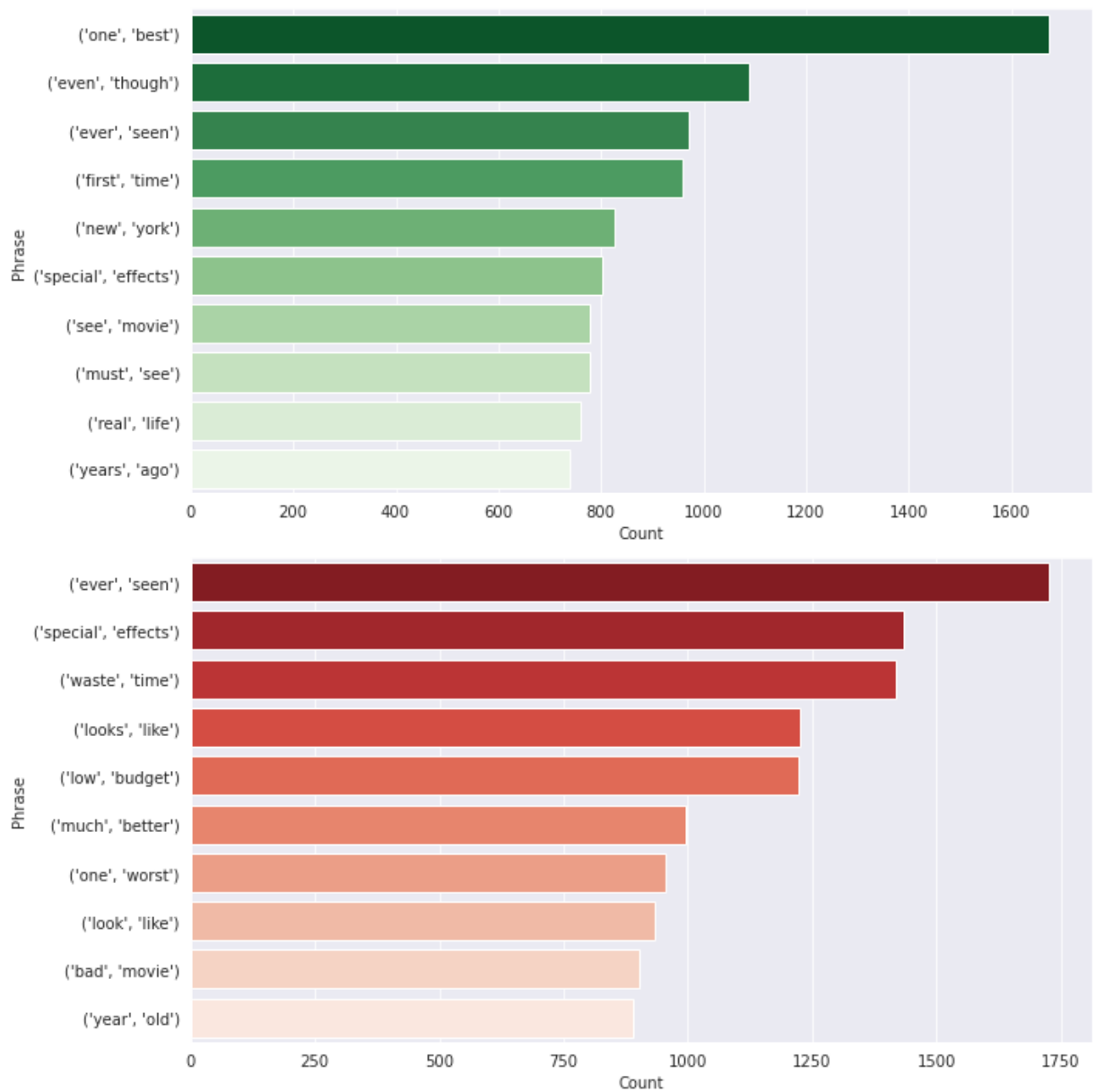
executed in 5.44s, finished 14:27:29 2021-06-23



In [21]:

```
1 # Plot two word ngrams
2 plot_ngram(2)
```

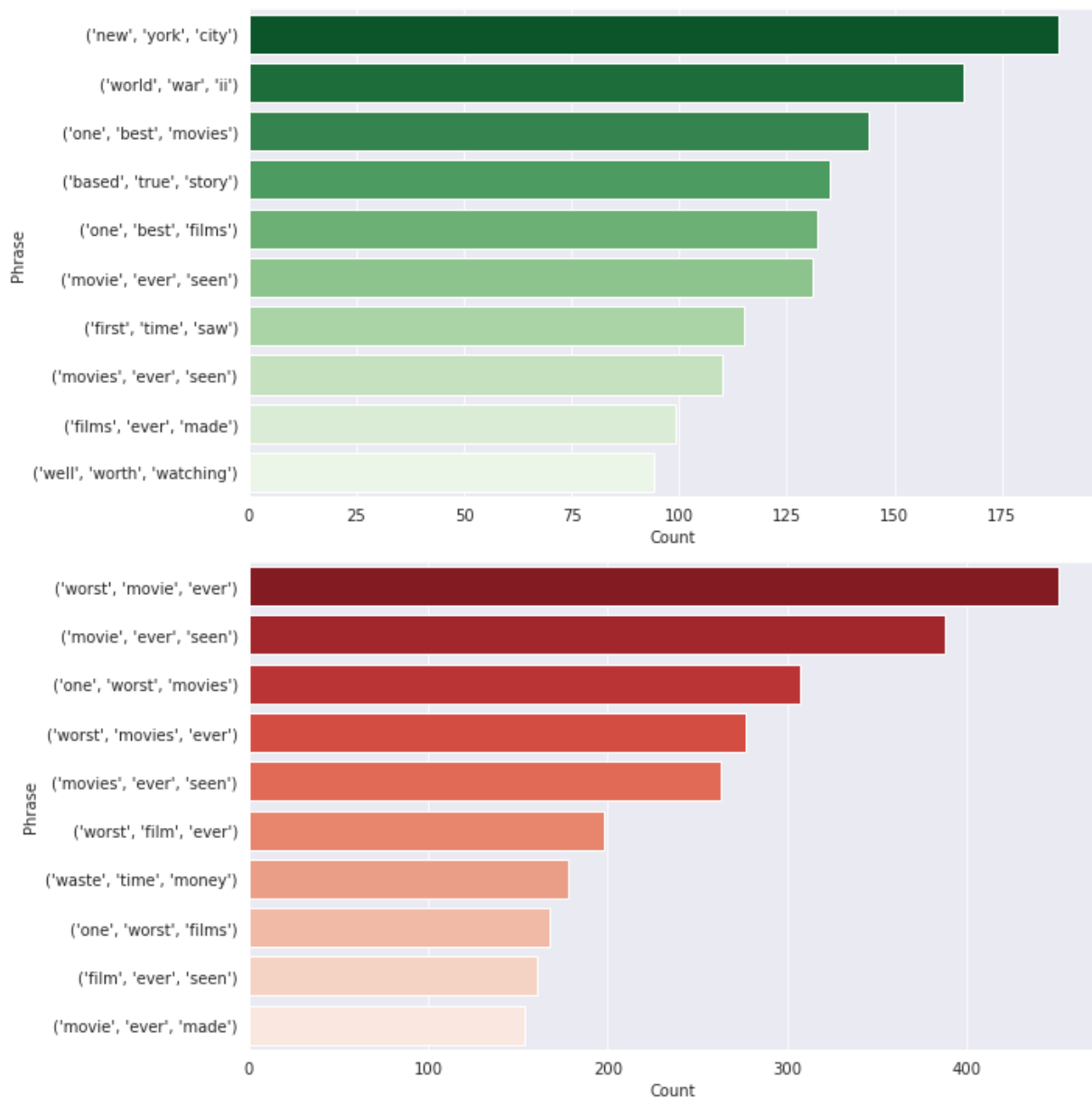
executed in 12.9s, finished 14:27:42 2021-06-23



In [22]:

```
1 # Plot three word ngrams
2 plot_ngram(3)
```

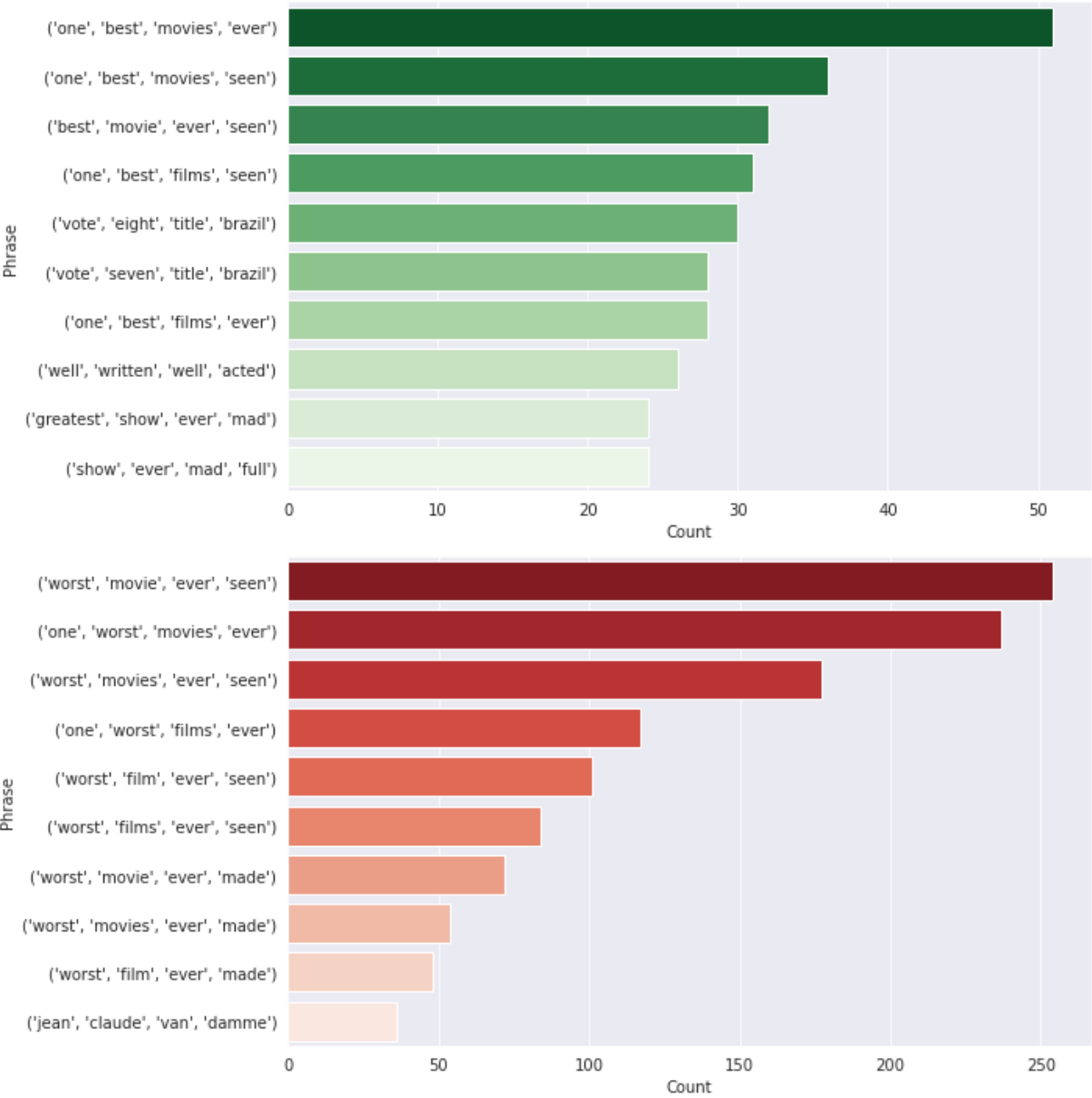
executed in 13.8s, finished 14:27:56 2021-06-23



In [23]:

```
1 ## Plot four word ngrams
2 plot_ngram(4)
```

executed in 13.9s, finished 14:28:10 2021-06-23



```
1 Finally, the differences in
  positive and negative tokens are
  becoming more evident.
2
```

Finally, the differences in positive and negative tokens are becoming more evident.

```

3 Ngrams(3) shows us more
  interesting information than
  Ngrams(1) and Ngrams(2). New York
  City appears as the top most
  frequent among trigrams, as does
  World War II. The rest are
  superlative - "one best films,"
  "movie ever seen," "films ever
  made."
4
5 The negative trigrams are more
  superlative and don't include
  things we wouldn't normally expect
  like New York City and World War
  II. This might lead us to
  consider whether or not movies
  based in New York City are more
  well reviewed, and if World War II
  movies also receive a similar
  boost. Perhaps fans of movies
  portraying NYC and WWII are more
  likely to write positive reviews.
6
7 Adorably, Jean Claude Van Damme
  appears frequently in the negative
  ngram(4) plot.

```

Ngrams(3) shows us more interesting information than Ngrams(1) and Ngrams(2). New York City appears as the top most frequent among trigrams, as does World War II. The rest are superlative - "one best films," "movie ever seen," "films ever made."

The negative trigrams are more superlative and don't include things we wouldn't normally expect like New York City and World War II. This might lead us to consider whether or not movies based in New York City are more well reviewed, and if World War II movies also receive a similar boost. Perhaps fans of movies portraying NYC and WWII are more likely to write positive reviews.

Adorably, Jean Claude Van Damme appears frequently in the negative ngram(4) plot.

▼ 5.4 Mutual Information Scores

Pairwise mutual information scores are also worth a quick look. This will give us a measure of association between two tokens.

```
In [24]: 1 # Instantiating NLTK's Bigram Association Measures
2 bigram_measures = nltk.collocations.BigramAssocMeasures()
3
4 # Instantiating Bigram Collocation Finder
5 tokens_pmi_finder = BigramCollocationFinder.from_words(tokens)
6
7 # Applying a frequency filter - bigram must appear at least 500 times
8 tokens_pmi_finder.apply_freq_filter(500)
9
10 # Scoring the tokens based on our bigram measures
11 tokens_pmi_scored = tokens_pmi_finder.score_ngrams(bigram_measures.pmi)
12
13 # Show scores
14 tokens_pmi_scored
```

executed in 19.7s, finished 14:28:29 2021-06-23

```
Out[24]: [('sci', 'fi'), 12.014985625405476),
          ('th', 'century'), 10.941059050317797),
          ('production', 'values'), 9.938295540111962),
          ('low', 'budget'), 9.617255903440476),
          ('special', 'effects'), 9.442463595883456),
          ('new', 'york'), 9.437849429863778),
          ('high', 'school'), 8.821587200156959),
          ('highly', 'recommend'), 8.735826357910884),
          ('years', 'ago'), 8.705289579685402),
          ('black', 'white'), 8.56233814091226),
          ('supporting', 'cast'), 8.237654151722769),
          ('half', 'hour'), 8.223432319941598),
          ('year', 'old'), 7.977964402482492),
          ('character', 'development'), 7.918511003378143),
          ('read', 'book'), 7.718890452656954),
          ('takes', 'place'), 7.602421569135796),
          ('writer', 'director'), 7.505539139359406),
          ('worth', 'seeing'), 7.297108484893386),
          ('camera', 'work'), 7.150295316173729),
          ('', ''), 7.116256222401521)
```

1 All of the pairs make sense within the context of our dataset. While many didn't appear in the list of top ngrams, it isn't surprising to see phrases like "sci fi," "production values," and "low budget" as having a high likelihood of appearing alongside one another.

All of the pairs make sense within the context of our dataset. While many didn't appear in the list of top ngrams, it isn't surprising to see phrases like "sci fi," "production values," and "low budget" as having a high likelihood of appearing alongside one another.

6 Modeling

Now that we have a better understanding of our dataset, we'll move on to our machine learning models. For the models, we'll be using two different vectorizers - Count Vectorization and TF-IDF Vectorization.

6.1 Count Vectorizer

The Count Vectorizer will convert all of our tokens into frequency representations, similar to what we did with our frequency distribution. While it's a good start, count vectorization has a couple of flaws:

- It doesn't distinguish between more and less important tokens.
- It only will consider the most frequent terms as the most statistically important. Words like "movie" and "film" appear frequently for both positive and negative reviews, and these will both be considered statistically important for both classifiers.

Below, we will set up our vectorizer using our original review set, since the vectorizer has the capability of tokenizing for us. We will also separately specify our stop words list. Creating our vectorizers this way will allow us flexibility if we wish to use a different tokenizer or stop words list in another iteration of this project.

```
In [25]: 1 # Instantiating count vectorizer using my_tokenizer and stop_words_list
2 cv = CountVectorizer(strip_accents='unicode',
3                       tokenizer=my_tokenizer,
4                       stop_words=stop_words_list
5                       )
6
7 # Assigning X and y for train test split
8 X = df['review']
9 y = df['sentiment']
10
11 # Apply train test split
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
13                                                    random_state=8)
14
15 # Fit both x train and test for CV
16 X_train_cv = cv.fit_transform(X_train)
17 X_test_cv = cv.transform(X_test)
18
19 #Check X_train_cv
20 X_train_cv
```

executed in 18.6s, finished 14:28:48 2021-06-23

```
Out[25]: <34707x93802 sparse matrix of type '<class 'numpy.int64'>'
         with 3392481 stored elements in Compressed Sparse Row format>
```

6.2 TF-IDF Vectorizer

The Term Frequency - Inverse Document Frequency (TF-IDF) Vectorizer makes up for some of the shortcomings in the Count Vectorizer by penalizing words that are too abundant or too rare. In other words, if a token is found frequently in review A but not frequently in review B, then that token will be considered more significant as it might provide additional insight to the rest of the corpus.

```
In [26]: 1 # Instantiating TF-IDF vectorizer using my_tokenizer and stop_words_list
2 tfidf = TfidfVectorizer(strip_accents='unicode',
3                       tokenizer=my_tokenizer,
4                       stop_words=stop_words_list
5                       )
6
7 # Fit both x train and test for TF-IDF
8 X_train_tfidf = tfidf.fit_transform(X_train)
9 X_test_tfidf = tfidf.transform(X_test)
10
11 # Check X_train_tfidf
12 X_train_tfidf
```

executed in 18.8s, finished 14:29:07 2021-06-23

Out[26]: <34707x93802 sparse matrix of type '<class 'numpy.float64'>' with 3392481 stored elements in Compressed Sparse Row format>

▼ 6.3 Model Evaluation Function

Before running our models, we'll create our model evaluation function. This will be used to compare the performances of the various machine learning models.

```
In [27]: 1 def evaluate_model(model, X_train, X_test, y_train=y_train,
2                   y_test=y_test, cmap='Greens', normalize=None,
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 14ms, finished 14:29:07 2021-06-23

6.4 Logistic Regression

We'll start off by running one of the more basic machine learning models, logistic regression. We will perform this both on our CV and TF-IDF vectorizers.

6.4.1 Logistic Regression CV

```
In [28]: 1 # Logistic Regression with Count Vectoriser
2 cv_log = LogisticRegression(random_state=8, n_jobs=-1)
3
4 # Fit X_train_cv
5 cv_log.fit(X_train_cv, y_train)
```

executed in 7.05s, finished 14:29:14 2021-06-23

Out[28]: LogisticRegression(n_jobs=-1, random_state=8)

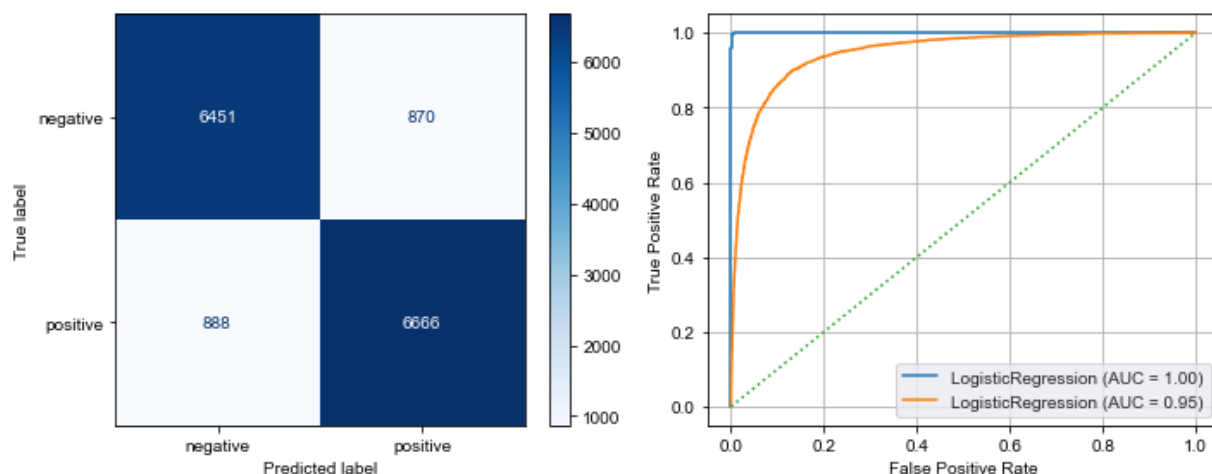
```
In [29]: 1 # Evaluate model performance
2 evaluate_model(cv_log, X_train_cv, X_test_cv, cmap='Blues')
```

executed in 1.93s, finished 14:29:16 2021-06-23

Training Accuracy: 99.74%

Test Accuracy: 88.18%

	precision	recall	f1-score	support
negative	0.88	0.88	0.88	7321
positive	0.88	0.88	0.88	7554
accuracy			0.88	14875
macro avg	0.88	0.88	0.88	14875
weighted avg	0.88	0.88	0.88	14875



6.4.2 Logistic Regression TF-IDF

```
In [49]: 1 # Logistic Regression with TF-IDF Vectoriser
2 tfidf_log = LogisticRegression(penalty='l2',C=100)
3 tfidf_log.fit(X_train_tfidf, y_train)
```

executed in 5.09s, finished 14:39:10 2021-06-23

Out[49]: LogisticRegression(C=100)

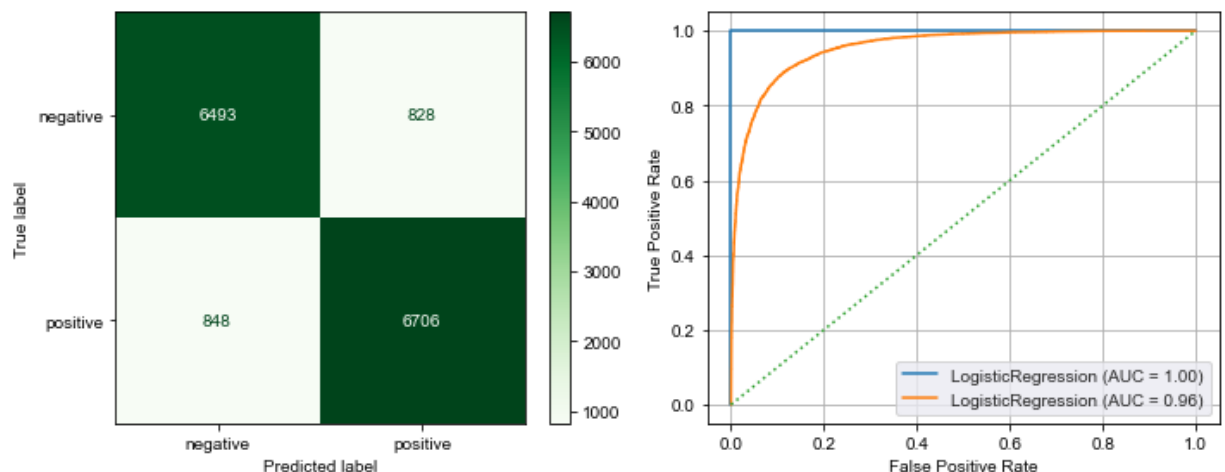
```
In [50]: 1 evaluate_model(tfidf_log, X_train_tfidf, X_test_tfidf, cmap='Greens')
```

executed in 1.50s, finished 14:39:12 2021-06-23

Training Accuracy: 100.00%

Test Accuracy: 88.73%

	precision	recall	f1-score	support
negative	0.88	0.89	0.89	7321
positive	0.89	0.89	0.89	7554
accuracy			0.89	14875
macro avg	0.89	0.89	0.89	14875
weighted avg	0.89	0.89	0.89	14875



Accuracy is surprisingly high for our first modeling attempt. Both CV and TF-IDF score above 88% accuracy.

Coefficient Analysis

Given the strong performance of the logistic regression model using TF-IDF, it is worth taking a look at the coefficients of the model. In this case, the coefficients will be the individual tokens in our corpus. We will be able to tell which words influenced the model in determining whether a review is positive or negative.

6.4.3 Coefficient Analysis

Given the strong performance of the logistic regression model using TF-IDF, it is worth taking a look at the coefficients of the model. In this case, the coefficients will be the individual tokens in our corpus. We will be able to tell which words influenced the model in determining whether a review is positive or negative.

In [32]:

```

1 # Create list of feature names
2 feature_names = tfidf.get_feature_names()
3
4 # Create empty dataframe to store coefficients
5 df_coef = pd.DataFrame()
6
7 # Populate columns with feature names and corresponding coefficients
8 df_coef['features'] = feature_names
9 df_coef['coefficients'] = tfidf_log.coef_.flatten()
10
11 # Sort coefficients by descending values
12 df_coef = df_coef.sort_values(by='coefficients', ascending=False)
13
14 # Create another dataframe with only the top 10 and bottom 10 coefficients
15 df_top_bottom = df_coef.iloc[np.r_[0:10, -10:0]]
16
17 # View top 10 and bottom 10 coefficients
18 df_top_bottom

```

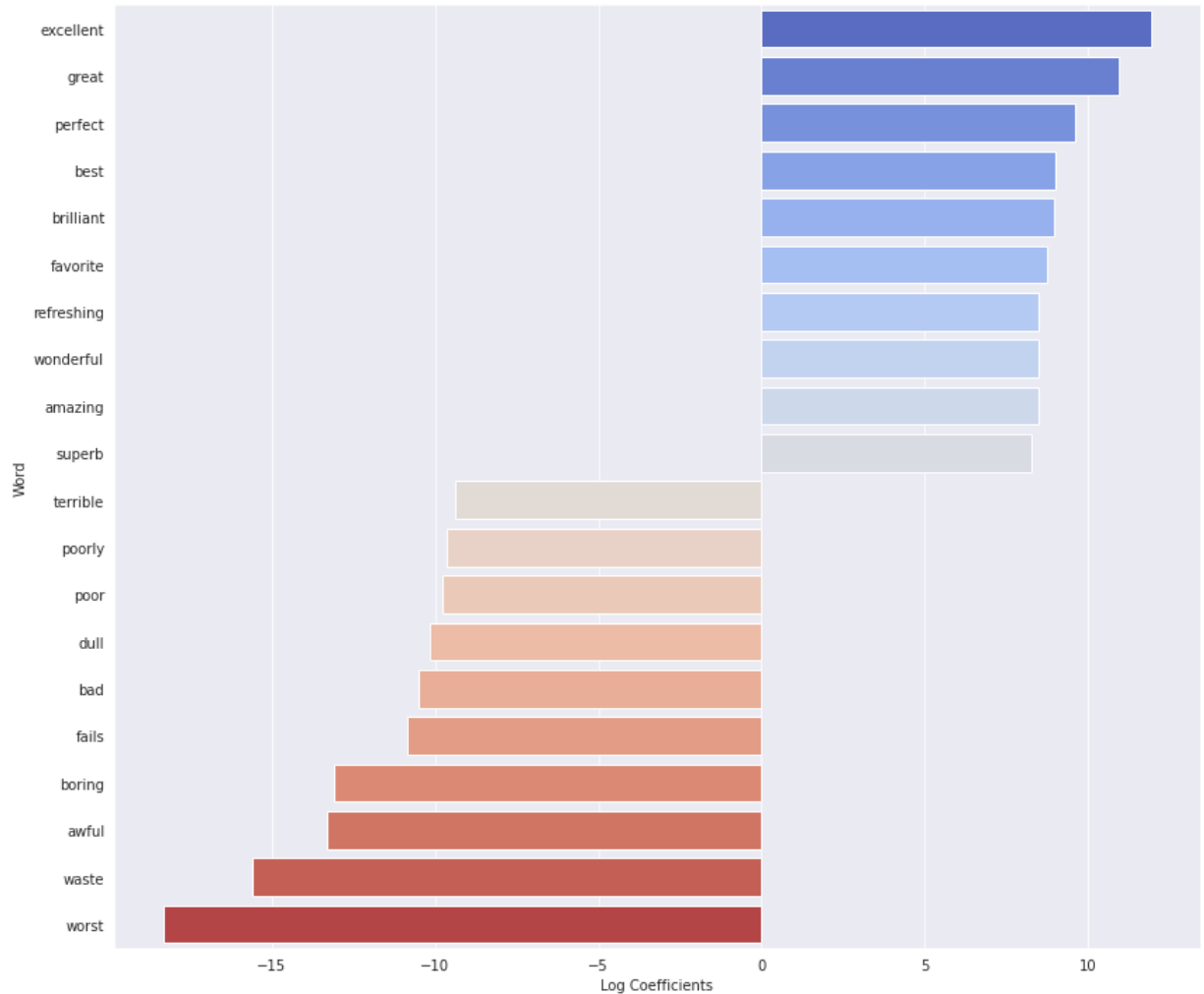
executed in 215ms, finished 14:29:22 2021-06-23

Out[32]:

	features	coefficients
27604	excellent	11.915176
34781	great	10.955981
61686	perfect	9.589469
7567	best	8.993472
10249	brilliant	8.936870
28941	favorite	8.737998
67983	refreshing	8.482725
92076	wonderful	8.476920
2324	amazing	8.462058
80600	superb	8.286361
82796	terrible	-9.405800
63726	poorly	-9.644778
63722	poor	-9.785211
24488	dull	-10.177873
5496	bad	-10.512454
28381	fails	-10.835542
9417	boring	-13.084709
5181	awful	-13.297779
90379	waste	-15.609118
92326	worst	-18.315653

```
In [33]: 1 # Define plot coefficient function
2 def plot_coefficients(df, filename='image', cmap=cmap):
3
4     # Create coefficients plot
5     with sns.axes_style("darkgrid"):
6         plt.figure(figsize=(12, 10))
7         ax = sns.barplot(data=df, x='coefficients', y='features', palette='c
8         ax.set(xlabel='Log Coefficients', ylabel='Word')
9
10    # View coefficient plot
11    plot_coefficients(df_top_bottom)
```

executed in 693ms, finished 14:29:23 2021-06-23



1 Based on this plot of the top 10 influential words in the positive and negative direction, it looks like the model latched onto polar descriptive adjectives that we would normally associate with positive or negative sentiment.

Based on this plot of the top 10 influential words in the positive and negative direction, it looks like the model latched onto polar descriptive adjectives that we would normally associate with positive or negative sentiment.

6.5 Random Forest

A different type of machine learning model we can try is random forests. We'll attempt this model with both the count and TF-IDF vectorizer.

6.5.1 Random Forest CV

```
In [34]: 1 # Initiate a random forest model for CV
2 rf_cv = RandomForestClassifier(random_state=8, n_jobs=-1)
3
4 # Fit to X_train_cv and y_train
5 rf_cv.fit(X_train_cv, y_train)
```

executed in 1m 8.01s, finished 14:30:31 2021-06-23

Out[34]: RandomForestClassifier(n_jobs=-1, random_state=8)

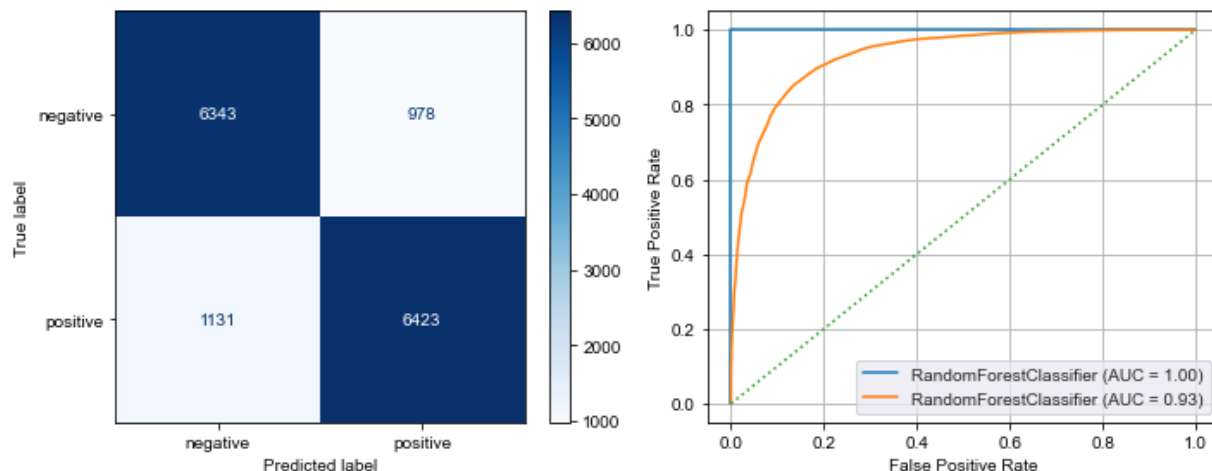
```
In [35]: 1 evaluate_model(rf_cv, X_train_cv, X_test_cv, cmap='Blues')
```

executed in 9.60s, finished 14:30:40 2021-06-23

Training Accuracy: 100.00%

Test Accuracy: 85.82%

	precision	recall	f1-score	support
negative	0.85	0.87	0.86	7321
positive	0.87	0.85	0.86	7554
accuracy			0.86	14875
macro avg	0.86	0.86	0.86	14875
weighted avg	0.86	0.86	0.86	14875



6.5.2 Random Forest TF-IDF

```
In [36]: 1 # Initiate a random forest model for TF-IDF
2 rf_tfidf = RandomForestClassifier(random_state=8, n_jobs=-1)
3
4 # Fit to X_train for TF-IDF
5 rf_tfidf.fit(X_train_tfidf, y_train)
```

executed in 58.1s, finished 14:31:38 2021-06-23

Out[36]: RandomForestClassifier(n_jobs=-1, random_state=8)

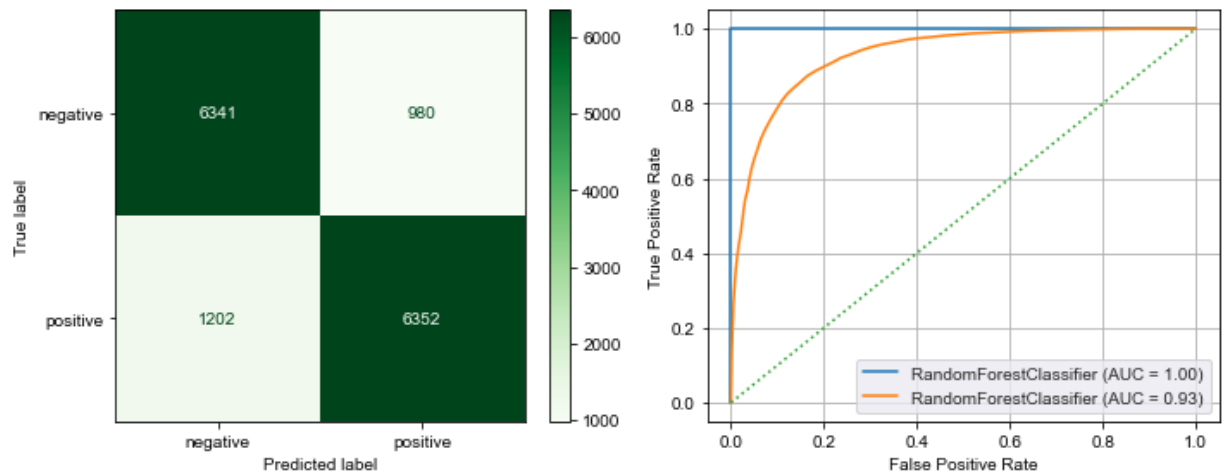
```
In [37]: 1 evaluate_model(rf_tfidf, X_train_tfidf, X_test_tfidf, cmap='Greens')
```

executed in 9.63s, finished 14:31:48 2021-06-23

Training Accuracy: 100.00%

Test Accuracy: 85.33%

	precision	recall	f1-score	support
negative	0.84	0.87	0.85	7321
positive	0.87	0.84	0.85	7554
accuracy			0.85	14875
macro avg	0.85	0.85	0.85	14875
weighted avg	0.85	0.85	0.85	14875



The random forest models were less accurate than our logistic regression model with TF-IDF.

▼ 6.6 XGBoost

Finally, we'll try XGBoost, which is a more advanced machine learning model. Hopefully this will give us a performance boost of some kind.

▼ 6.6.1 XGBoost CV

```
In [38]: 1 # Instantiate XGB classifier for CV
2 xgb_cv = XGBClassifier(random_state=8, n_jobs=-1)
3
4 # Fit to X_train_cv and y_train
5 xgb_cv.fit(X_train_cv, y_train)
```

executed in 29.3s, finished 14:32:17 2021-06-23

```
Out[38]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
      colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
      importance_type='gain', interaction_constraints='',
      learning_rate=0.300000012, max_delta_step=0, max_depth=6,
      min_child_weight=1, missing=nan, monotone_constraints='()',
      n_estimators=100, n_jobs=-1, num_parallel_tree=1, random_state=8,
      reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
      tree_method='exact', validate_parameters=1, verbosity=None)
```

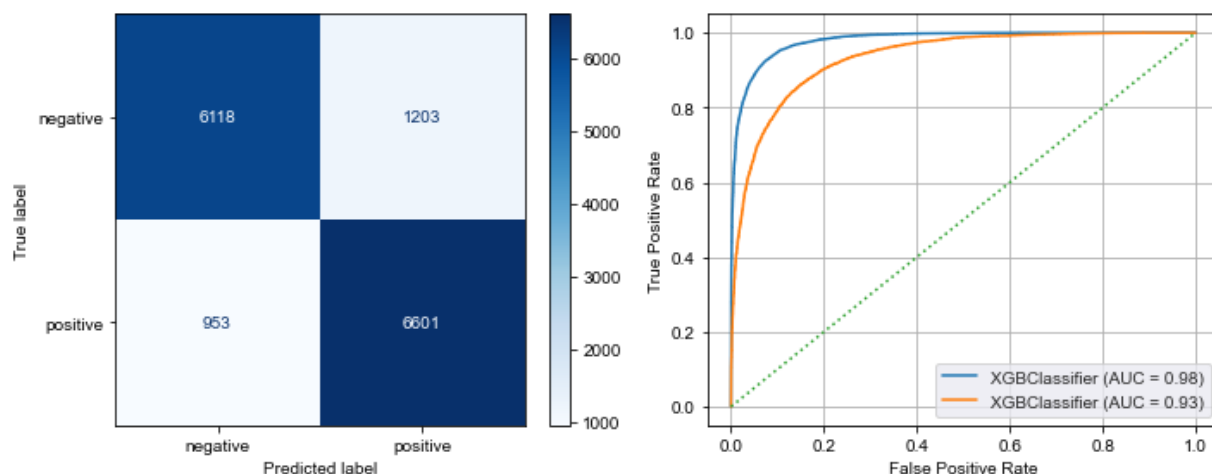
```
In [39]: 1 evaluate_model(xgb_cv, X_train_cv, X_test_cv, cmap='Blues')
```

executed in 7.52s, finished 14:32:25 2021-06-23

Training Accuracy: 92.43%

Test Accuracy: 85.51%

	precision	recall	f1-score	support
negative	0.87	0.84	0.85	7321
positive	0.85	0.87	0.86	7554
accuracy			0.86	14875
macro avg	0.86	0.85	0.85	14875
weighted avg	0.86	0.86	0.85	14875



6.6.2 XGBoost TF-IDF

```
In [40]: 1 # Instantiate XGB model for TF-IDF
2 xgb_tfidf = XGBClassifier(random_state=8, n_jobs=-1)
3
4 # Fit to X_train_tfidf
5 xgb_tfidf.fit(X_train_tfidf, y_train)
```

executed in 1m 17.3s, finished 14:33:42 2021-06-23

```
Out[40]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
importance_type='gain', interaction_constraints='',
learning_rate=0.300000012, max_delta_step=0, max_depth=6,
min_child_weight=1, missing=nan, monotone_constraints='()',
n_estimators=100, n_jobs=-1, num_parallel_tree=1, random_state=8,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
tree_method='exact', validate_parameters=1, verbosity=None)
```

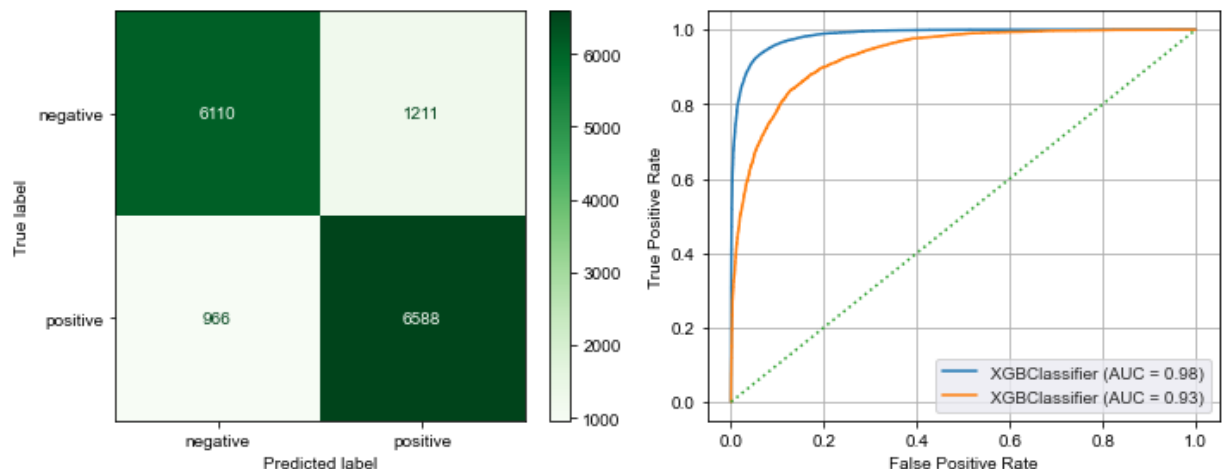
```
In [41]: 1 evaluate_model(xgb_tfidf, X_train_tfidf, X_test_tfidf, cmap='Greens')
```

executed in 7.19s, finished 14:33:49 2021-06-23

Training Accuracy: 93.37%

Test Accuracy: 85.36%

	precision	recall	f1-score	support
negative	0.86	0.83	0.85	7321
positive	0.84	0.87	0.86	7554
accuracy			0.85	14875
macro avg	0.85	0.85	0.85	14875
weighted avg	0.85	0.85	0.85	14875



Similar to the random forest model, XGBoost did not provide us with any additional accuracy. It looks like logistic regression will be our best bet for now.



6.7 Logistic Regression with Grid Search

As a last step, we'll expand upon our logistic regression model with TF-IDF by performing a grid search. This process will iterate through multiple parameters and return a logistic regression model that optimizes for accuracy.

```
In [52]: 1 # Initiate new model and perform grid search
2 tfidf_log_gs = LogisticRegression()
3
4 # If run = True, code will perform full grid search
5 # If run = False, code will use previously calculated best parameters
6 run = False
7
8 if run == True:
9
10 # Define lists of parameters to compare
11     params = {'C':[0.01,0.1,1,10,100],
12               'penalty':['l1','l2','elastic_net'],
13               'solver':['liblinear', 'newton-cg', 'lbfgs', 'sag','saga']}
14
15
16 else:
17
18     # Previously calculated best parameters
19     params = {'C':[10],
20               'penalty':['l2'],
21               'solver':['newton-cg']}
22
23
24 # Run the grid search with a focus on accuracy
25 log_grid_search = GridSearchCV(tfidf_log_gs,params,scoring='accuracy',
26                                verbose=100,
27                                n_jobs=-1)
28
29 # Fit grid search to training data and display best parameters
30 log_grid_search.fit(X_train_tfidf, y_train)
31
32 # Print best parameters
33 log_grid_search.best_params_
```

executed in 20.2s, finished 14:53:21 2021-06-23

```
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: 11.6s
[Parallel(n_jobs=-1)]: Done 2 out of 5 | elapsed: 11.8s remaining: 17.7
s
[Parallel(n_jobs=-1)]: Done 3 out of 5 | elapsed: 11.9s remaining: 7.9
s
[Parallel(n_jobs=-1)]: Done 5 out of 5 | elapsed: 15.3s remaining: 0.0
s
[Parallel(n_jobs=-1)]: Done 5 out of 5 | elapsed: 15.3s finished
```

Out[52]: {'C': 10, 'penalty': 'l2', 'solver': 'newton-cg'}

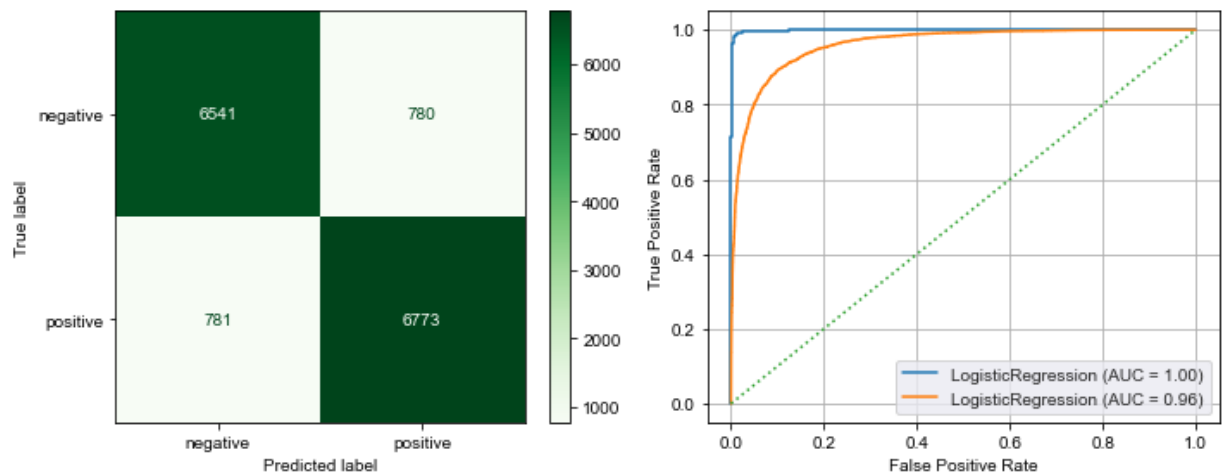
```
In [53]: 1 evaluate_model(log_grid_search.best_estimator_, X_train_tfidf, X_test_tfidf,
2                  cmap='Greens')
```

executed in 1.34s, finished 14:53:22 2021-06-23

Training Accuracy: 98.82%

Test Accuracy: 89.51%

	precision	recall	f1-score	support
negative	0.89	0.89	0.89	7321
positive	0.90	0.90	0.90	7554
accuracy			0.90	14875
macro avg	0.90	0.90	0.90	14875
weighted avg	0.90	0.90	0.90	14875



Using grid search, we were able to eek out just a bit more performance. Our baseline logistic regression model with TF-IDF scored 88.73%, and our grid search model improved accuracy to 89.51%.

7 Interpretation

Most of our models performed reasonably well. Perhaps the method of using decision trees in Random Forest and XGBoost isn't as useful when working with a dataset with this many features. All models suffered from overfitting, and it might be a good step in the future to examine the causes of this more closely.

7.1 Best Model - Logistic Regression with Grid Search using TF-IDF Vectorizer

Despite our implementation of more complex machine learning models, the tried and true logistic regression model with grid search was the best performer at 89.51% accuracy.

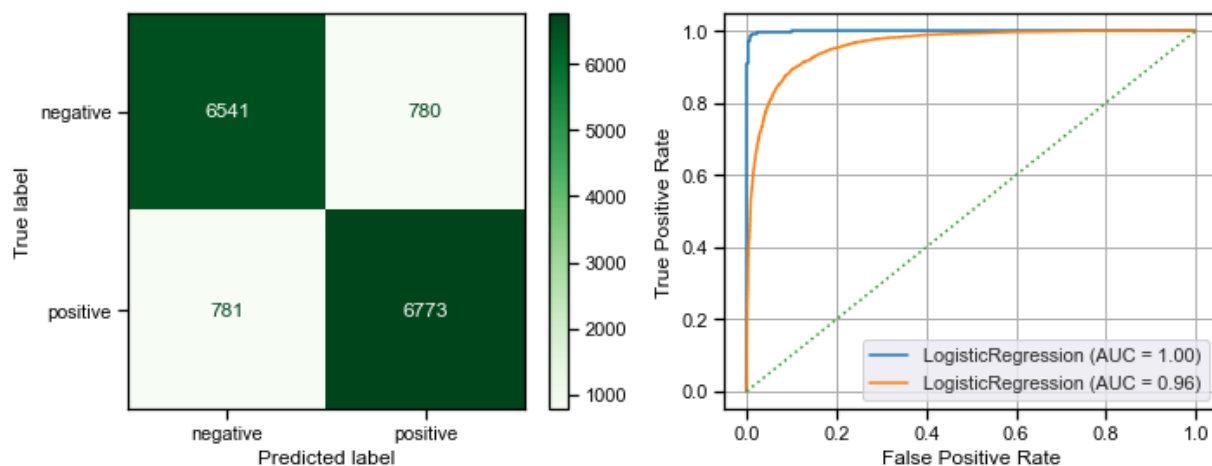

```
In [57]: 1 sns.set_context('notebook')
2
3 evaluate_model(log_grid_search.best_estimator_, X_train_tfidf, X_test_tfidf,
4               cmap='Greens')
```

executed in 1.37s, finished 15:30:07 2021-06-23

Training Accuracy: 98.82%

Test Accuracy: 89.51%

	precision	recall	f1-score	support
negative	0.89	0.89	0.89	7321
positive	0.90	0.90	0.90	7554
accuracy			0.90	14875
macro avg	0.90	0.90	0.90	14875
weighted avg	0.90	0.90	0.90	14875



8 Conclusions and Recommendations

The goal of this project was to determine whether or not it would be possible to distinguish between positive and negative reviews using only the content of the review and the 'positive' or 'negative' classifiers. We accomplished this with reasonable success with grid search logistic regression and the TF-IDF vectorizer scoring 89.51% accuracy.

We also learned from our EDA and visual exploration of coefficients that the machine learning models are latching onto strong positive and strong negative adjectives to make its decisions. This does not fall out of line with what we might expect.

Using this model, we might be able to calculate a new score for movies on IMDB that more closely reflects what is found on the Rotten Tomatoes Critic Score. A metric like this could help avoid issues where IMDB users felt very similarly about a movie, but gave it different numerical scores based on their own interpretation of the 1-10 scale.

The benefits of this new type of score might not be immediately tangible, but running this model and scoring all movies in IMDB's database might return an alternate top 250 movies list that would at the very least be interesting to users.

8.1 Next Steps

We aren't entirely sure how each of the 50,000 reviews were deemed 'positive' or 'negative.' According to the data source on Kaggle, the reviews were "highly polar," but we don't understand how this informed the data collection process.

It might be worthwhile to pull 50,000 reviews from IMDB in a bucket of scores from 1-3, 4-7, and 8-10. This would allow us to perform a multiclass sentiment analysis and see if we can determine whether reviews are positive, negative, or neutral. While computationally expensive and possibly a fool's errand, the same logic could be applied to bucketing every 1-10 score in buckets of 1, 2, 3, etc.

Also, while we explored three models for this project, there are many more machine learning models that we could have attempted, such as support vector machines, k nearest neighbors, and a variety of deep learning / neural network models. Exploring more models might lead to a higher accuracy and could inform future experimentation.