

# NLP For Drugs.com Data Set

## Packages Import

In [3]:

```
### Visualization
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
sns.color_palette("Blues", as_cmap=True)

### Standard Packages
import numpy as np
import warnings
import nltk
import re
import pandas as pd
pd.set_option('display.max_colwidth', None)
warnings.filterwarnings("ignore")

### NLTK
from nltk.tokenize import TweetTokenizer
from nltk.corpus import stopwords, wordnet
from nltk import pos_tag
from nltk.stem import WordNetLemmatizer
from nltk.sentiment.vader import SentimentIntensityAnalyzer
from nltk.tokenize import RegexpTokenizer
nltk.download('wordnet')
nltk.download('vader_lexicon')
import contractions

### Scikit-Learn
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import MultinomialNB
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.model_selection import GridSearchCV
from sklearn.compose import ColumnTransformer
from sklearn.svm import SVC
from sklearn.metrics import ConfusionMatrixDisplay, confusion_matrix, accuracy_score, f1_score, recall_score, precision_score

### ImbLearn
from imblearn.pipeline import Pipeline as imbpipeline
from imblearn.over_sampling import SMOTE
```

```
[nltk_data] Downloading package wordnet to
[nltk_data] /Users/albertcc/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package vader_lexicon to
```

```
[nltk_data] /Users/albertcc/nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
```

## Bringing in two .tsv files as test and train

```
In [4]: # Load in the test and train datasets provided in the data file
data_test = pd.read_csv('data/drugsComTest_raw.tsv', sep='\t')
data_train = pd.read_csv('data/drugsComTrain_raw.tsv', sep='\t')
```

```
In [5]: data_test.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53766 entries, 0 to 53765
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   Unnamed: 0      53766 non-null  int64
 1   drugName        53766 non-null  object
 2   condition       53471 non-null  object
 3   review          53766 non-null  object
 4   rating          53766 non-null  float64
 5   date            53766 non-null  object
 6   usefulCount     53766 non-null  int64
dtypes: float64(1), int64(2), object(4)
memory usage: 2.9+ MB
```

```
In [6]: data_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 161297 entries, 0 to 161296
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   Unnamed: 0      161297 non-null  int64
 1   drugName        161297 non-null  object
 2   condition       160398 non-null  object
 3   review          161297 non-null  object
 4   rating          161297 non-null  float64
 5   date            161297 non-null  object
 6   usefulCount     161297 non-null  int64
dtypes: float64(1), int64(2), object(4)
memory usage: 8.6+ MB
```

## Merge Test and Train dataframes

- The data provided is already split into test and train tsv files. I would like to combine these to not only have more data to work with, but any cleaning could be applied to the merged dataset before splitting into a training and testing set.

```
In [7]: # Combining both tsv files
merged_df = pd.concat([data_test, data_train], axis=0)
```

```
In [8]: merged_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 215063 entries, 0 to 161296
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   Unnamed: 0      215063 non-null int64
 1   drugName        215063 non-null object
 2   condition       213869 non-null object
 3   review          215063 non-null object
 4   rating          215063 non-null float64
 5   date            215063 non-null object
 6   usefulCount     215063 non-null int64
dtypes: float64(1), int64(2), object(4)
memory usage: 13.1+ MB
```

```
In [9]: # Drop first column since these appear to be entry numbers
merged_df = merged_df.drop(merged_df.columns[0],axis=1)
```

```
In [10]: # Sanity check that the Unnamed column has been dropped
merged_df.head()
```

Out[10]:

	drugName	condition	review	rating	date	usefulCount
0	Mirtazapine	Depression	"I've tried a few antidepressants over the years (citalopram, fluoxetine, amitriptyline), but none of those helped with my depression, insomnia & anxiety. My doctor suggested and changed me onto 45mg mirtazapine and this medicine has saved my life. Thankfully I have had no side effects especially the most common - weight gain, I've actually lost alot of weight. I still have suicidal thoughts but mirtazapine has saved me."	10.0	February 28, 2012	22

1	Mesalamine	Crohn's Disease, Maintenance	"My son has Crohn's disease and has done very well on the Asacol. He has no complaints and shows no side effects. He has taken as many as nine tablets per day at one time. I've been very happy with the results, reducing his bouts of diarrhea drastically."	8.0	May 17, 2009	17
2	Bactrim	Urinary Tract Infection	"Quick reduction of symptoms"	9.0	September 29, 2017	3
3	Contrave	Weight Loss	"Contrave combines drugs that were used for alcohol, smoking, and opioid cessation. People lose weight on it because it also helps control over-eating. I have no doubt that most obesity is caused from sugar/carb addiction, which is just as powerful as any drug. I have been taking it for five days, and the good news is, it seems to go to work immediately. I feel hungry before I want food now. I really don't care to eat; it's just to fill my stomach. Since I have only been on it a few days, I don't know if I've lost weight (I don't have a scale), but my clothes do feel a little looser, so maybe a pound or two. I'm hoping that after a few months on this medication, I will develop healthier habits that I can continue without the aid of Contrave."	9.0	March 5, 2017	35
4	Cyclafem 1 / 35	Birth Control	"I have been on this birth control for one cycle. After reading some of the reviews on this type and similar birth controls I was a bit apprehensive to start. Im giving this birth control a 9 out of 10 as I have not been on it long enough for a 10. So far I love this birth control! My side effects have been so minimal its like Im not even on birth control! I have experienced mild headaches here and there and some nausea but other than that ive been feeling great! I got my period on cue on the third day of the inactive pills and I had no idea it was coming because I had zero pms! My period was very light and I barely had any cramping! I had unprotected sex the first month and obviously didn't get pregnant so I'm very pleased! Highly recommend"	9.0	October 22, 2015	4

- Noticed how 'condition' has some missing values, but other columns are fine

```
In [11]: # Drop null values that are in 'condition'

merged_df = merged_df.dropna(subset=['condition'])
```

```
In [12]: merged_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 213869 entries, 0 to 161296
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   drugName        213869 non-null object
1   condition       213869 non-null object
2   review         213869 non-null object
3   rating         213869 non-null float64
4   date           213869 non-null object
5   usefulCount    213869 non-null int64
dtypes: float64(1), int64(1), object(4)
memory usage: 11.4+ MB
```

- Reduces the dataset to 213,869 values

```
In [13]: # Want to see which conditions are the most frequent

merged_df['condition'].value_counts()
```

```
Out[13]: Birth Control          38436
Depression                    12164
Pain                          8245
Anxiety                       7812
Acne                          7435
...
Amyotrophic Lateral Sclerosis      1
Hyperuricemia Secondary to Chemotherapy  1
146 users found this comment helpful.  1
Microscopic polyangiitis           1
Systemic Candidiasis               1
Name: condition, Length: 916, dtype: int64
```

## Looking at unique drugs under 'Birth Control' condition

```
In [14]: # Focus will be seeing which birth control drugs there are
merged_df['drugName'][merged_df['condition'] == 'Birth Control'].nunique
```

Out[14]: 181

```
In [15]: # Want to find which kinds of birth control are the most common
merged_df['drugName'][merged_df['condition'] == 'Birth Control'].value_
```

```
Out[15]: Etonogestrel          4394
Ethinyl estradiol / norethindrone 3081
Levonorgestrel          2884
Nexplanon              2883
Ethinyl estradiol / levonorgestrel 2107
...
Briellyn                1
Lillow                  1
Larin 24 Fe             1
Dasetta 7 / 7 / 7       1
Philith                 1
Name: drugName, Length: 181, dtype: int64
```

```
In [16]: # Select the the top 7 birth controls since we want to focus on these
bc_drugs = ['Etonogestrel', 'Ethinyl estradiol / norethindrone', 'Nexp
           'Ethinyl estradiol / norgestimate', 'Implanon']

bc_data = merged_df[merged_df['drugName'].isin(bc_drugs)]
```

In [17]: `bc_data.head()`

Out[17]:

	drugName	condition	review	rating	date	usefulCount
			"I want to share my experience to possibly ease some of the ladies out there.\n\nThe incident occur in Nov 18 (the condom broke) and it was 4 days after my period ended and approximately 4 days before ovulation. This means I was fertile. \n\nwhen we noticed that the condom broke, we immediately bought Plan B One Step and was taken about 45 minutes or an hour later. \n\nThe sooner you take the greater the chances it will work.\n\nSo I waited anxiously, praying I get my period. I felt side effects such as fatigue, bloated,			
42	Levonorgestrel	Emergency Contraception		10.0	December 28, 2013	20

In [18]: `bc_data['drugName'].value_counts()`

Out[18]:

Levonorgestrel	4896
Etonogestrel	4402
Ethinyl estradiol / norethindrone	3619
Nexplanon	2892
Ethinyl estradiol / norgestimate	2682
Ethinyl estradiol / levonorgestrel	2400
Implanon	1506

Name: drugName, dtype: int64

- Wanted to only include in our dataset the top 7 drugs with condition = birth control, however when we filtered for these drugs we see additional conditions were selected



```
In [19]: bc_data['condition'].value_counts()
```

```
Out[19]: Birth Control          18942
Emergency Contraception      1651
Abnormal Uterine Bleeding    812
Acne                        439
Endometriosis              178
Menstrual Disorders         124
Ovarian Cysts              106
Polycystic Ovary Syndrome    89
Not Listed / Othe           19
Premenstrual Syndrome       11
0</span> users found this comment helpful. 8
Postmenopausal Symptoms     6
2</span> users found this comment helpful. 4
8</span> users found this comment helpful. 2
1</span> users found this comment helpful. 2
4</span> users found this comment helpful. 1
Gonadotropin Inhibition     1
3</span> users found this comment helpful. 1
9</span> users found this comment helpful. 1
Name: condition, dtype: int64
```

```
In [20]: # Let's try to get rid of these conditions that took in the 'Useful' r
bc_data = bc_data[~bc_data['condition'].str.contains('comment')]
```

```
In [21]: bc_data['condition'].value_counts()
```

```
Out[21]: Birth Control          18942
Emergency Contraception      1651
Abnormal Uterine Bleeding    812
Acne                        439
Endometriosis              178
Menstrual Disorders         124
Ovarian Cysts              106
Polycystic Ovary Syndrome    89
Not Listed / Othe           19
Premenstrual Syndrome       11
Postmenopausal Symptoms     6
Gonadotropin Inhibition     1
Name: condition, dtype: int64
```

```
In [22]: bc_data['drugName'].value_counts()
```

```
Out[22]: Levonorgestrel          4896
          Etonogestrel          4402
          Ethinyl estradiol / norethindrone 3619
          Nexplanon             2883
          Ethinyl estradiol / norgestimate 2682
          Ethinyl estradiol / levonorgestrel 2400
          Implanon              1496
          Name: drugName, dtype: int64
```

```
In [23]: bc_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 22378 entries, 42 to 161273
Data columns (total 6 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   drugName        22378 non-null  object
 1   condition       22378 non-null  object
 2   review          22378 non-null  object
 3   rating          22378 non-null  float64
 4   date            22378 non-null  object
 5   usefulCount     22378 non-null  int64
dtypes: float64(1), int64(1), object(4)
memory usage: 1.2+ MB
```

```
In [24]: # Create new column called 'sentiment' that will have the target variable
bc_data['sentiment'] = ['Positive' if x > 7.0 else 'Negative' for x in bc_data['rating']]
```

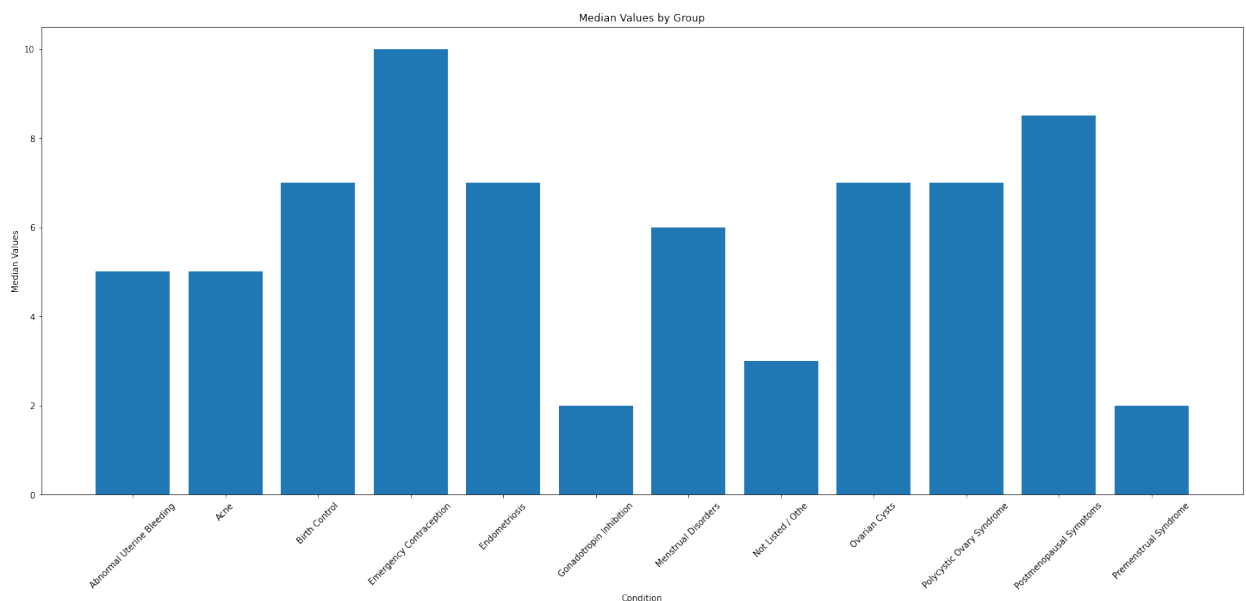
## Analyze the median rating based on condition

```
In [25]: median_values = bc_data.groupby('condition')['rating'].median()

# Create bar chart
fig, ax = plt.subplots(figsize = (25, 10))
ax.bar(median_values.index, median_values.values)
plt.xticks(rotation=45)

# Set axis labels and title
ax.set_xlabel('Condition')
ax.set_ylabel('Median Values')
ax.set_title('Median Values by Group')

# Show plot
plt.show()
```



```
In [26]: bc_data.loc[bc_data['condition'] == 'Emergency Contraception'].rating.
```

```
Out[26]: 10.0    1025
          9.0     216
          1.0     155
          8.0     106
          7.0      43
          5.0      43
          6.0      27
          4.0      14
          3.0      13
          2.0       9
          Name: rating, dtype: int64
```

## Create columns that count emphasis and capital letters in text, as this could express sentiment

```
In [27]: # Creating a 'punc_emphasis' column that scores how many exclamation p
bc_data['punc_emphasis'] = bc_data['review'].apply(lambda x: sum([1 fo
# Creating a 'capt_emphasis' column that scores how many capitalized w
bc_data['capt_emphasis'] = bc_data['review'].apply(lambda x: sum([1 fo
```

```
In [28]: bc_data.head()
```

Out[28]:

drugName	condition	review	rating	date	usefulCount	sentiment	pu
		"I want to share my experience to possibly ease some of the ladies out there.\n\nThe incident occur in Nov 18 (the condom broke) and it was 4 days after my period ended and approximately 4 days before ovulation. _ . .					

```
In [29]: bc_data['sentiment'].value_counts(normalize=True)
```

```
Out[29]: Negative    0.535526
Positive    0.464474
Name: sentiment, dtype: float64
```

## Sentiment Analysis Against Condition

- Within the conditions we have selected, how do the reviews look pertaining to each condition?

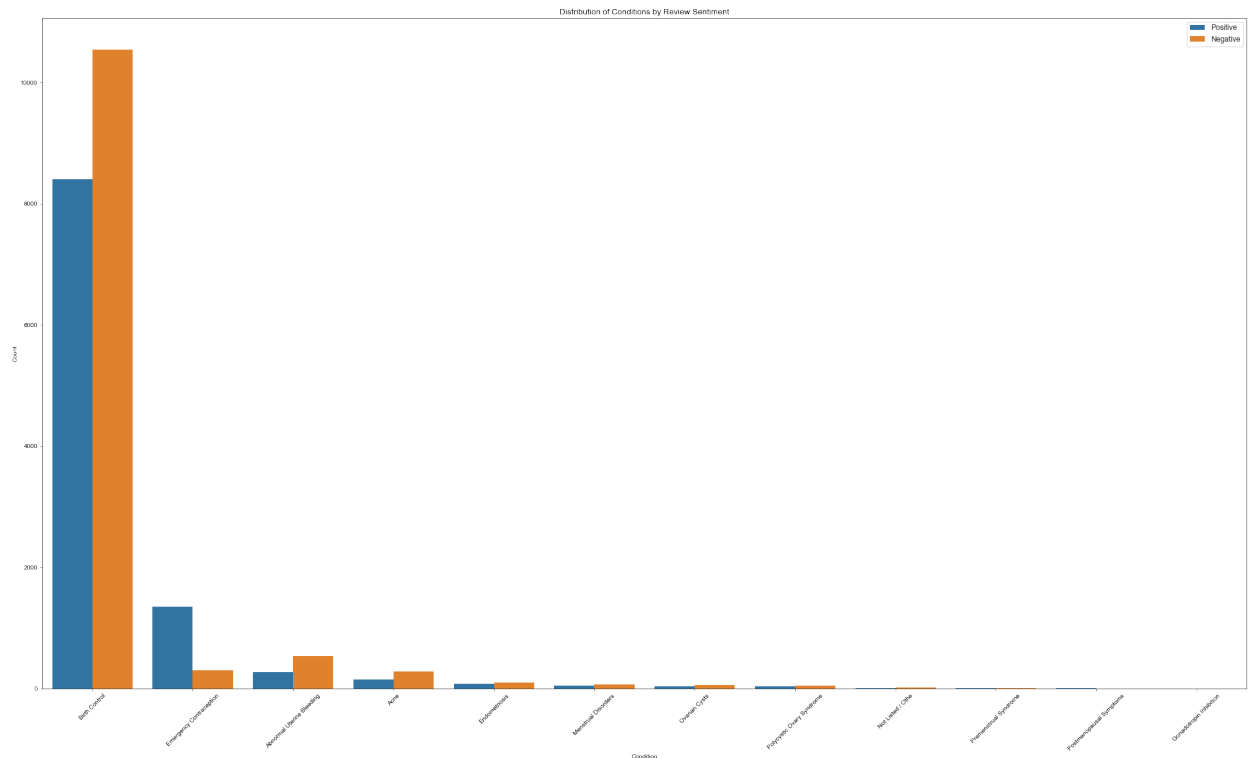
```
In [30]: # Let's try plotting sentiment against groups

fig = plt.figure(figsize = (35, 20))

hue_order = ['Positive', 'Negative']
sns.countplot(x='condition', hue='sentiment', data=bc_data, hue_order=
plt.xticks(rotation=45)

sns.set(style='white', font_scale=1.1)

plt.legend(loc='upper right')
plt.xlabel('Condition')
plt.ylabel('Count')
plt.title('Distribution of Conditions by Review Sentiment');
```



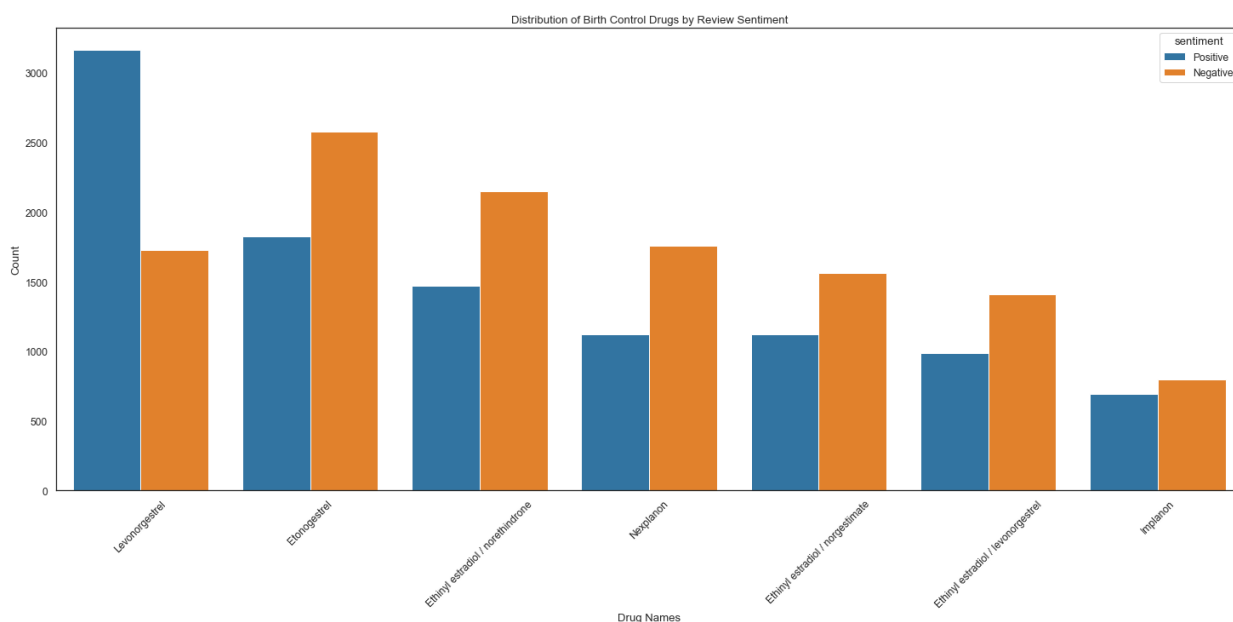
```
In [31]: # Let's try plotting sentiment against groups

fig = plt.figure(figsize = (25, 10))

sns.countplot(x='drugName', hue='sentiment', palette='tab10', data=bc_)
plt.xticks(rotation=45)

sns.set(style='white', font_scale=1.1)

plt.xlabel('Drug Names')
plt.ylabel('Count')
plt.title('Distribution of Birth Control Drugs by Review Sentiment');
```



## Can we do anything with 'UsefulCount'?

```
In [32]: bc_data['usefulCount'].value_counts(bins=3)
```

```
Out[32]: (-1.248, 415.667]    22371
          (415.667, 831.333]      6
          (831.333, 1247.0]       1
          Name: usefulCount, dtype: int64
```

```
In [33]: bc_data['usefulCount'].min()
```

```
Out[33]: 0
```

```
In [34]: bc_data['usefulCount'].describe()
```

```
Out[34]: count      22378.000000
         mean         8.008535
         std        18.501084
         min         0.000000
         25%         2.000000
         50%         4.000000
         75%         9.000000
         max        1247.000000
         Name: usefulCount, dtype: float64
```

- Not sure if this is too useful of a feature, maybe we could filter the reviews that were found useful above a certain threshold to take in user input.

## After research on birth controls, wanted to read what reviews are saying

```
In [35]: # Look at the reviews for Levonogestrel
         bc_data[bc_data['drugName'] == 'Levonorgestrel']
```

```
Out[35]:
```

drugName	condition	review	rating	date	usefulCount	sentiment
		"I want to share my experience to possibly ease some of the ladies out there.\n\nThe incident occur in Nov 18 (the condom broke) and it was 4 days after my period ended and approximately 4 days before ovulation. This means I was fertile. \n\nwhen we noticed that ..				

## Binning the years these reviews were written into two groups to see if there's a difference over time

```
In [36]: # Convert date to date time
bc_data['date'] = pd.to_datetime(bc_data['date'])
```

```
In [37]: # Look at the minimum and maximum dates
bc_data['date'].describe()
```

```
Out[37]: count                22378
         unique                3014
         top      2015-10-21 00:00:00
         freq                   38
         first    2008-02-27 00:00:00
         last     2017-12-12 00:00:00
         Name: date, dtype: object
```

```
In [38]: bc_data['date'].value_counts(bins=2)
```

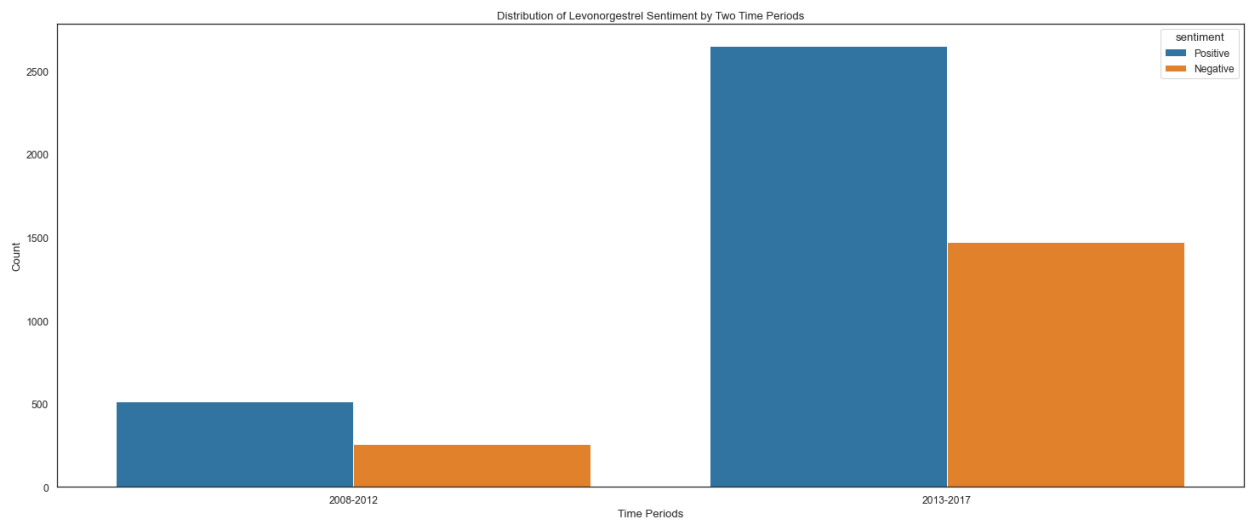
```
Out[38]: (2013-01-19, 2017-12-12]                18297
         (2008-02-23 10:10:33.599999999, 2013-01-19]    4081
         Name: date, dtype: int64
```

```
In [39]: # Create new column called 'date_column' that will have grouped time r
bc_data['date_column'] = ['2013-2017' if x.year > 2013 else '2008-2012']
```

**Can we see the difference in reviews of these drugs over time?**



```
In [40]: # Let's try plotting sentiment of Levonorgestrel against the two assign  
  
fig = plt.figure(figsize = (25, 10))  
  
sns.countplot(x='date_column', hue='sentiment', data=bc_data[bc_data['  
  
# plt.xticks(rotation=45)  
  
plt.xlabel('Time Periods')  
plt.ylabel('Count')  
plt.title('Distribution of Levonorgestrel Sentiment by Two Time Periods')
```

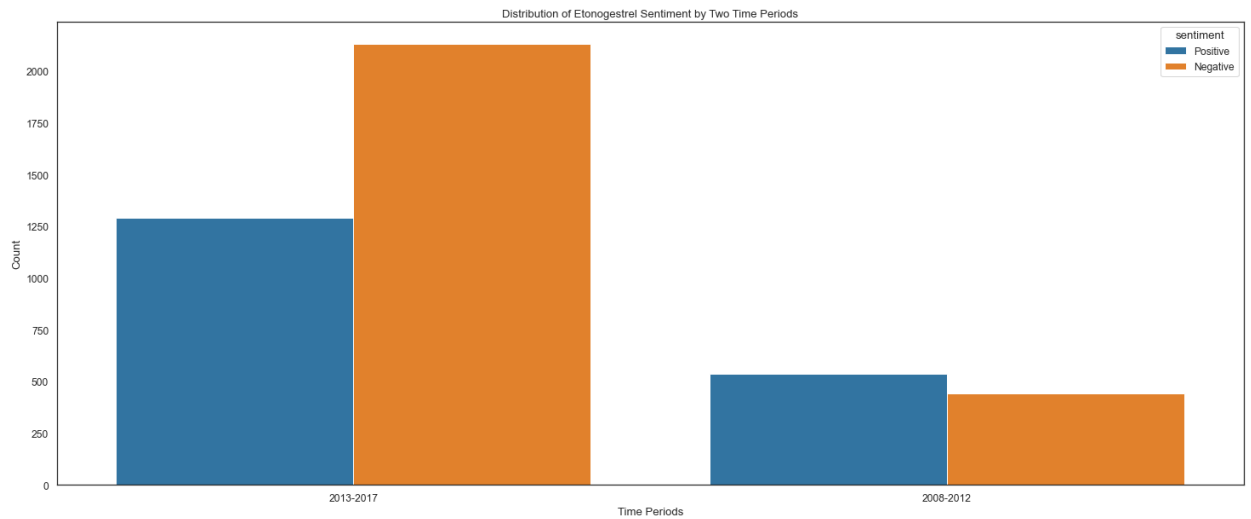


```
In [41]: # Let's try plotting sentiment of Etonogestrel against the two assigned time periods

fig = plt.figure(figsize = (25, 10))

sns.countplot(x='date_column', hue='sentiment', data=bc_data[bc_data['date_column'].isin(['2013-2017', '2008-2012'])])

# plt.xticks(rotation=45)
plt.xlabel('Time Periods')
plt.ylabel('Count')
plt.title('Distribution of Etonogestrel Sentiment by Two Time Periods')
```

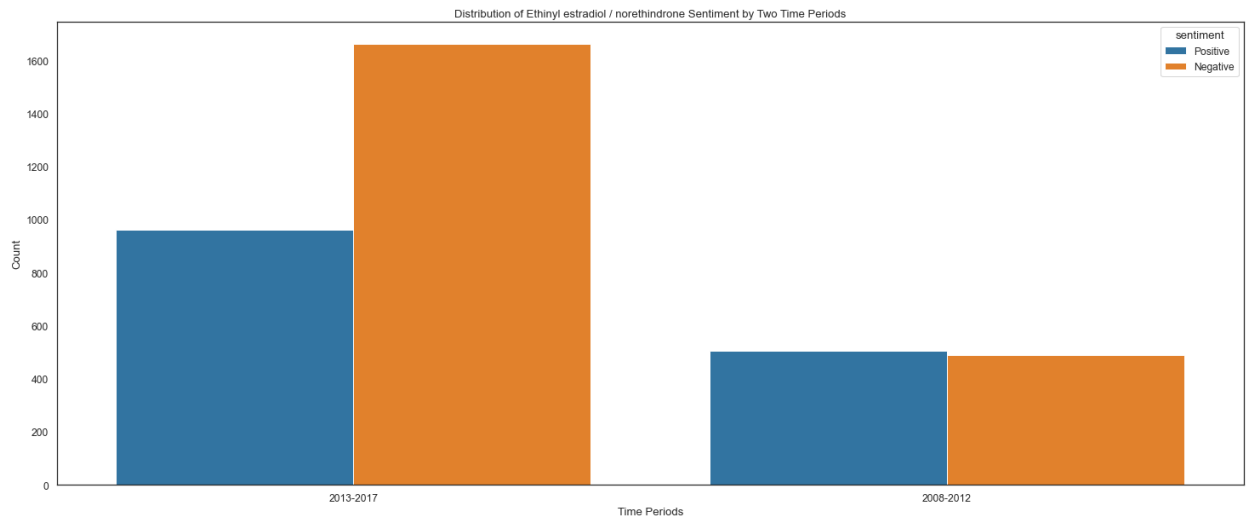


```
In [42]: # Let's try plotting sentiment of Ethinyl estradiol / norethindrone ag

fig = plt.figure(figsize = (25, 10))

sns.countplot(x='date_column', hue='sentiment', data=bc_data[bc_data['

# plt.xticks(rotation=45)
plt.xlabel('Time Periods')
plt.ylabel('Count')
plt.title('Distribution of Ethinyl estradiol / norethindrone Sentiment
```



## Cleaning Text Reviews

```
In [43]: # Create function that will lowercase the text

def lower_case(text):
    text = text.lower()
    return text

# Create function to remove the html apostrophes in the text

def apostrophe(text):
    text = text.replace('&#039;', '\')
    return text

# Want to expand the contractions so we can see if these words have in

def fixcontractions(text):
    text = contractions.fix(text)
    return text

# Create a function that uses a regex tokenizer to remove punctuation
```

```
def remove_punctuation(text):
    tokenizer = RegexpTokenizer(r'\w+\'?\\w+')
    text = tokenizer.tokenize(text)
    text = ' '.join(text)
    return text

# Remove stopwords from the reviews

def remove_stopwords(text, stop_words_list = set(stopwords.words('english'))):
    text = text.split()
    text = [word for word in text if word not in stop_words_list]
    text = ' '.join(text)
    return text

# Create a function that lemmatizes words

def lemmatize(text):
    lemmatizer = WordNetLemmatizer()
    text = text.split()
    text = [lemmatizer.lemmatize(word) for word in text]
    text = ' '.join(text)
    return text

def clean_text(text):
    text = lower_case(text)
    text = apostrophe(text)
    text = fixcontractions(text)
    text = remove_punctuation(text)
    text = remove_stopwords(text)
    text = lemmatize(text)
    return text
```

```
In [44]: # Manually testing the contractions.fix function
contractions.fix("I've aren't Tim's got a lovely bunch of coconuts")
```

```
Out[44]: "I have are not Tim's got a lovely bunch of coconuts"
```

```
In [45]: # Original review text location 6
bc_data['review'][14]
```

```
Out[45]: '"Started Nexplanon 2 months ago because I have a minimal amount of c
ontraception's I can take due to my inability to take the hormon
e that is used in most birth controls. I'm trying to give it tim
e because it is one of my only options right now. But honestly if I h
ad options I'd get it removed.\r\nI've never had acne probl
ems in my life, and immediately broke out after getting it implanted.
Sex drive is completely gone, and I used to have sex with my boyfrien
d a few days a week, now its completely forced and not even fun for m
e anymore. I mean I'm on birth control because I like having sex
but don't want to get pregnant, why take a birth control that ta
kes away sex? Very unhappy and hope that I get it back with time or I
'm getting it removed.'"
```

```
In [46]: # Testing one of the reviews to see what it is doing to the text, as a
clean_text(bc_data['review'][14])
```

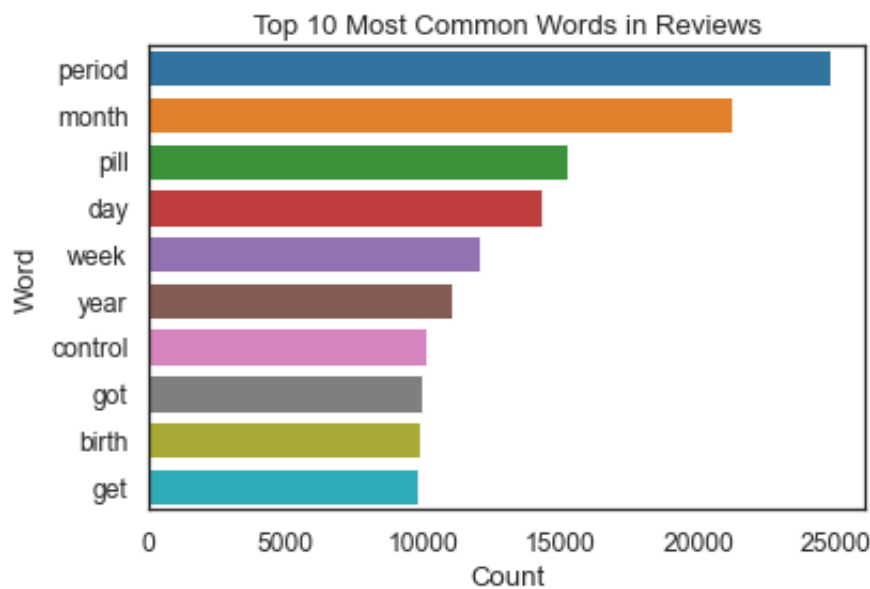
```
Out[46]: "started nexplanon month ago minimal amount contraception's take due
inability take hormone used birth control trying give time one option
right honestly option would get removed never acne problem life immed
iately broke getting implanted sex drive completely gone used sex boy
friend day week completely forced even fun anymore mean birth control
like sex want get pregnant take birth control take away sex unhappy h
ope get back time getting removed"
```

In [47]: *### Plotting the top 10 most common words in the 'text' column in an s*

```
text = ' '.join(bc_data['review'])
text = clean_text(text)
text = text.split()

freq = pd.Series(text).value_counts()[:10]
freq = freq.to_frame()
freq = freq.reset_index()
freq.columns = ['word', 'count']
freq = freq.sort_values(by='count', ascending=False)

fig = plt.figure(figsize=(6,4))
sns.barplot(x='count', y='word', data=freq, palette='tab10')
plt.xlabel('Count')
plt.ylabel('Word')
plt.title('Top 10 Most Common Words in Reviews')
plt.show()
```



## Modeling

In [48]: `bc_data.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 22378 entries, 42 to 161273
Data columns (total 10 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   drugName            22378 non-null  object
 1   condition           22378 non-null  object
 2   review              22378 non-null  object
 3   rating              22378 non-null  float64
 4   date                22378 non-null  datetime64[ns]
 5   usefulCount         22378 non-null  int64
 6   sentiment           22378 non-null  object
 7   punc_emphasis       22378 non-null  int64
 8   capt_emphasis       22378 non-null  int64
 9   date_column         22378 non-null  object
dtypes: datetime64[ns](1), float64(1), int64(3), object(5)
memory usage: 2.5+ MB
```

In [49]: `# How is the distribution of sentiment amongst the Birth Control`  
`bc_data['sentiment'].value_counts(normalize=True)`

Out[49]: Negative 0.535526  
Positive 0.464474  
Name: sentiment, dtype: float64

## First Simple Model - Count Vectorizer / Logistic Regression / No Features

In [50]: *# Using review text against sentiment analysis*

```
X1 = bc_data['review']
y1 = bc_data['sentiment']

X_train_1, X_test_1, y_train_1, y_test_1 = train_test_split(X1, y1, te

# For Train Set, apply clean_text function

X_train_1 = X_train_1.apply(clean_text)

### Train - Tokenize the training data with a simple split of words, a

X_train_1 = X_train_1.apply(lambda x: x.split())
X_train_1 = X_train_1.map(' '.join)

### Train - Vectorize the training data using CountVectorizer

cv = CountVectorizer()
X_train_1 = cv.fit_transform(X_train_1)

### Train - Fit training data to Logistic Regression Model

logit = LogisticRegression()
logit.fit(X_train_1, y_train_1)

### VALIDATION - Perform a cross validation on the logistic regression

scores = cross_val_score(logit, X_train_1, y_train_1, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())
```

```
Cross Validation Scores: [0.82658475 0.83105278 0.82513966 0.8184357
5 0.82234637]
Mean Cross Validation Score: 0.8247118638250607
```

In [51]: *# Logistic Regression Test Set Preprocessing*

```
X_test_1_logit = X_test_1.apply(clean_text)
X_test_1_logit = X_test_1.apply(lambda x: x.split())
X_test_1_logit = X_test_1.map(' '.join)
X_test_1_logit = cv.transform(X_test_1)
```

In [52]: logit\_pred = logit.predict(X\_test\_1\_logit)



```
In [53]: print('Logistic Regression Accuracy: ', accuracy_score(y_test_1, logit_pred))
print('Logistic Regression F1 Score: ', f1_score(y_test_1, logit_pred))
print('Logistic Regression Precision Score: ', precision_score(y_test_1, logit_pred))
print('Logistic Regression Recall Score: ', recall_score(y_test_1, logit_pred))
```

```
Logistic Regression Accuracy: 0.7888739946380697
Logistic Regression F1 Score: 0.7882641149643028
Logistic Regression Precision Score: 0.7979275466120083
Logistic Regression Recall Score: 0.7888739946380697
```

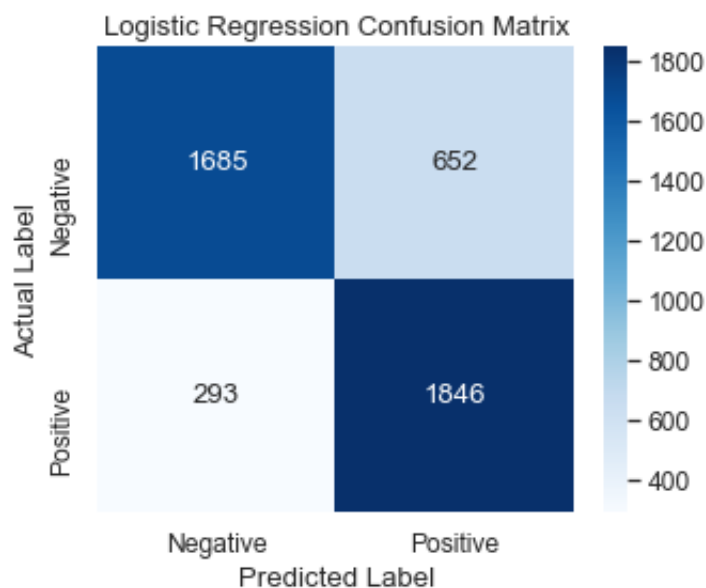
```
In [54]: # Get scores for the sentiments individually to see

# Quick look at decision matrix for our first model:

# Logistic Regression Confusion Matrix

cm = confusion_matrix(y_test_1, logit_pred)
cm_df = pd.DataFrame(cm, index=['Negative', 'Positive'], columns=['Negative', 'Positive'])

fig_cm1 = plt.figure(figsize=(5,4))
sns.heatmap(cm_df, annot=True, fmt='g', cmap='Blues')
plt.title('Logistic Regression Confusion Matrix')
plt.ylabel('Actual Label')
plt.xlabel('Predicted Label')
plt.show()
fig_cm1.savefig('images/logistic_regression_confusion_matrix.png', dpi=100)
```



```
In [55]: # Quick math check to see if these add up to the train split. It does!
1685 + 1846 + 293 + 652
```

```
Out[55]: 4476
```

## Second Model - TFIDF Vectorizer / Logistic Regression / No Features

- Want to test if TFIDF Vectorizer makes a difference compared to Count Vectorizer

```
In [56]: X1_tfidf = bc_data['review']
y1_tfidf = bc_data['sentiment']

X_train_1_tfidf, X_test_1_tfidf, y_train_1_tfidf, y_test_1_tfidf = train_test_split(X1_tfidf, y1_tfidf,
# For Train Set, apply clean_text function

X_train_1_tfidf = X_train_1_tfidf.apply(clean_text)

### Train - Tokenize the training data with a simple split of words, and a split of sentences

X_train_1_tfidf = X_train_1_tfidf.apply(lambda x: x.split())
X_train_1_tfidf = X_train_1_tfidf.map(' '.join)

### Train - Vectorize the training data using TfidfVectorizer

tfidf = TfidfVectorizer()
X_train_1_tfidf = tfidf.fit_transform(X_train_1_tfidf)

### Train - Fit training data to Logistic Regression Model

logit_tfidf = LogisticRegression()
logit_tfidf.fit(X_train_1_tfidf, y_train_1_tfidf)

### VALIDATION - Perform a cross validation on the decision tree classifier

scores = cross_val_score(logit_tfidf, X_train_1_tfidf, y_train_1_tfidf, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())
```

```
Cross Validation Scores: [0.83216978 0.83524155 0.83463687 0.82206704 0.8273743 ]
```

```
Mean Cross Validation Score: 0.8302979099811388
```

In [57]: *# Logistic Regression Test (TFIDF) Set Preprocessing*

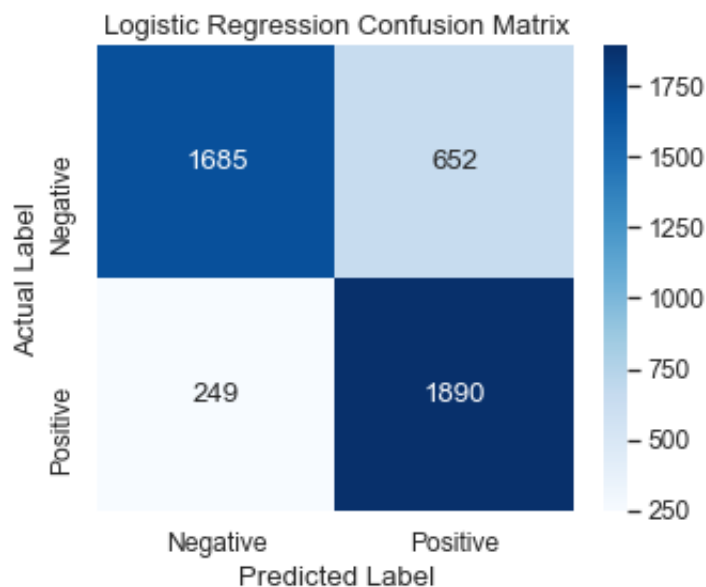
```
X_test_1_tfidf = X_test_1_tfidf.apply(clean_text)
X_test_1_tfidf = X_test_1_tfidf.apply(lambda x: x.split())
X_test_1_tfidf = X_test_1_tfidf.map(' '.join)
X_test_1_tfidf = cv.transform(X_test_1_tfidf)
```

In [58]: `logit_pred_tfidf = logit_tfidf.predict(X_test_1_tfidf)`

In [59]: *### Logistic Regression Confusion Matrix w/ TFIDF*

```
cm = confusion_matrix(y_test_1_tfidf, logit_pred_tfidf)
cm_df = pd.DataFrame(cm, index=['Negative', 'Positive'], columns=['Neg', 'Pos'])

fig_cm1b = plt.figure(figsize=(5,4))
sns.heatmap(cm_df, annot=True, fmt='g', cmap='Blues')
plt.title('Logistic Regression Confusion Matrix')
plt.ylabel('Actual Label')
plt.xlabel('Predicted Label')
plt.show()
fig_cm1b.savefig('images/logistic_regression_confusion_matrix_tfidf.png')
```



```
In [60]: print('Logistic Regression (TFIDF) Accuracy: ', accuracy_score(y_test_1, y_test_1_tfidf))
print('Logistic Regression (TFIDF) F1 Score: ', f1_score(y_test_1_tfidf, y_test_1_tfidf))
print('Logistic Regression (TFIDF) Precision Score: ', precision_score(y_test_1_tfidf, y_test_1_tfidf))
print('Logistic Regression (TFIDF) Recall Score: ', recall_score(y_test_1_tfidf, y_test_1_tfidf))
```

```
Logistic Regression (TFIDF) Accuracy: 0.798704200178731
Logistic Regression (TFIDF) F1 Score: 0.7978723877450785
Logistic Regression (TFIDF) Precision Score: 0.8102055684985682
Logistic Regression (TFIDF) Recall Score: 0.798704200178731
```

- We see here that our Logistic regression model using TFIDF has a higher accuracy than the model with Count Vector.
- Originally this was a lower accuracy due to the fact that we tried using three classes to predict sentiment that was imbalanced between the classes.
- In this version we made only a positive and negative sentiment for Birth Control and it was more balanced.

### Grid Search for best Logistic Regression parameters

```
In [61]: param_grid = {'C': [1, 10, 100, 1000, 10000], 'penalty': ['none', 'l1', 'l2']}
grid = GridSearchCV(logit, param_grid, cv=5, scoring='accuracy', n_jobs=-1)
grid.fit(X_train_1, y_train_1)
print('Best Parameters: ', grid.best_params_)
```

```
Best Parameters: {'C': 1, 'penalty': 'l2'}
```

- Looks that the default parameters are the best parameters for this model

```
In [62]: param_grid = {'C': [1, 10, 100, 1000, 10000], 'penalty': ['none', 'l1', 'l2']}
grid = GridSearchCV(logit_tfidf, param_grid, cv=5, scoring='accuracy', n_jobs=-1)
grid.fit(X_train_1_tfidf, y_train_1_tfidf)
print('Best Parameters: ', grid.best_params_)
```

```
Best Parameters: {'C': 1, 'penalty': 'l2'}
```

- Same as above

## Third Model: Count Vectorizer / Multinomial Bayes / No Features

```
In [63]: # Using the same X_train and y_train from our first models but using a
X1_nb = bc_data['review']
y1_nb = bc_data['sentiment']

X_train_1_nb, X_test_1_nb, y_train_1_nb, y_test_1_nb = train_test_split(X1_nb, y1_nb,
                                test_size=0.2, random_state=42)

# For Train Set, apply clean_text function
X_train_1_nb = X_train_1_nb.apply(clean_text)

### Train - Tokenize the training data with a simple split of words, a
X_train_1_nb = X_train_1_nb.apply(lambda x: x.split())
X_train_1_nb = X_train_1_nb.map(' '.join)

### Train - Vectorize the training data using CountVectorizer
cv = CountVectorizer()
X_train_1_nb = cv.fit_transform(X_train_1_nb)

# Train - fitting the training data to a Naive Bayes Classifier
nb = MultinomialNB()
nb.fit(X_train_1_nb, y_train_1_nb)

# Validation - Performing cross validation on the Naive Bayes Classifier
scores = cross_val_score(nb, X_train_1_nb, y_train_1_nb, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())
```

```
Cross Validation Scores: [0.82016197 0.82435074 0.81564246 0.8047486
0.81061453]
Mean Cross Validation Score: 0.8151036585080476
```

### Grid Search for best conditions with NB

```
In [64]: param_grid = {'alpha': [0.1, 0.5, 1, 2, 5, 20, 50]}
grid = GridSearchCV(nb, param_grid, cv=5, scoring='accuracy', n_jobs=-1)
grid.fit(X_train_1_nb, y_train_1_nb)
print('Best Parameters: ', grid.best_params_)
```

Best Parameters: {'alpha': 0.5}

- Default parameter is 1.0 for alpha, we can try seeing what the results are with 0.5

```
In [65]: # Naive Bayes Test Set Preprocessing

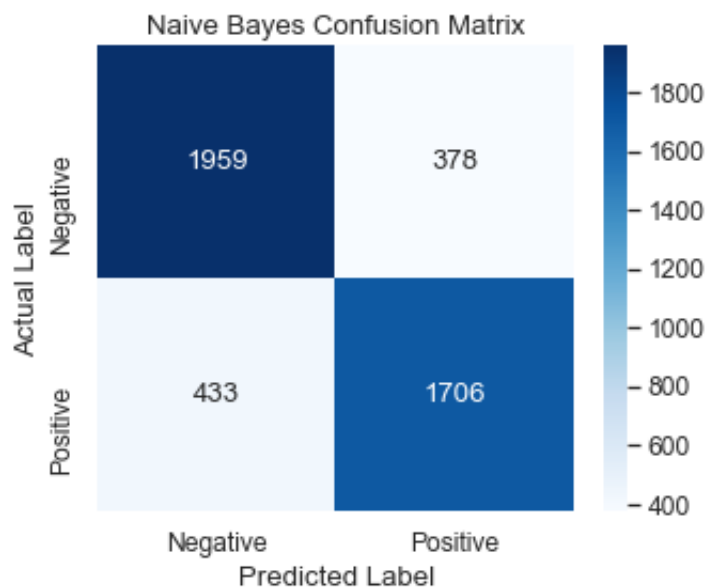
X_test_1_nb = X_test_1_nb.apply(clean_text)
X_test_1_nb = X_test_1_nb.apply(lambda x: x.split())
X_test_1_nb = X_test_1_nb.map(' '.join)
X_test_1_nb = cv.transform(X_test_1_nb)
```

```
In [66]: nb_pred_cv = nb.predict(X_test_1_nb)
```

In [67]: *### Naive Bayes Confusion Matrix w/ CountVectorizer*

```
cm = confusion_matrix(y_test_1_nb, nb_pred_cv)
cm_df = pd.DataFrame(cm, index=['Negative', 'Positive'], columns=['Neg

fig_cm3 = plt.figure(figsize=(5,4))
sns.heatmap(cm_df, annot=True, fmt='g', cmap='Blues')
plt.title('Naive Bayes Confusion Matrix')
plt.ylabel('Actual Label')
plt.xlabel('Predicted Label')
plt.show()
fig_cm3.savefig('images/mnb_confusion_matrix.png', dpi=300)
```



In [68]:

```
print('Naive Bayes Accuracy: ', accuracy_score(y_test_1_nb, nb_pred_cv))
print('Naive Bayes F1 Score: ', f1_score(y_test_1_nb, nb_pred_cv, average='micro'))
print('Naive Bayes Precision Score: ', precision_score(y_test_1_nb, nb_pred_cv, average='micro'))
print('Naive Bayes Recall Score: ', recall_score(y_test_1_nb, nb_pred_cv, average='micro'))
```

```
Naive Bayes Accuracy: 0.8188114387846291
Naive Bayes F1 Score: 0.8186851910506858
Naive Bayes Precision Score: 0.8188069919574142
Naive Bayes Recall Score: 0.8188114387846291
```

## 4th Model: Count Vectorizer / Logistic Regression / All Features

```

In [69]: # Additional features to see if it helps improve the model
X2 = bc_data[['review', 'drugName', 'condition', 'punc_emphasis', 'cap
y2 = bc_data['sentiment']

X_train_2, X_test_2, y_train_2, y_test_2 = train_test_split(X2, y2, te

# For Train Set, apply clean_text function

X_train_2['review'] = X_train_2['review'].apply(clean_text)

### Train - Tokenize the training data with a simple split of words, a

X_train_2['review'] = X_train_2['review'].apply(lambda x: x.split())
X_train_2['review'] = X_train_2['review'].map(' '.join)

### Train - Vectorize the training data using CountVectorizer

cv = CountVectorizer()
X_train_2 = cv.fit_transform(X_train_2['review'])

### Train - Fit training data to Logistic Regression Model

logit = LogisticRegression()
logit.fit(X_train_2, y_train_2)

### VALIDATION - Perform a cross validation on the logistic regression

scores = cross_val_score(logit, X_train_2, y_train_2, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())

```

Cross Validation Scores: [0.82658475 0.83105278 0.82513966 0.8184357  
5 0.82234637]

Mean Cross Validation Score: 0.8247118638250607

```

In [70]: logit.coef_

```

```

Out[70]: array([[ 0.43082765, -0.53249771,  0.15085222, ..., -0.28975559,
                -0.15145837, -0.39983239]])

```

```

In [71]: # Logistic Regression Test (All Features) Preprocessing

X_test_2['review'] = X_test_2['review'].apply(clean_text)
X_test_2['review'] = X_test_2['review'].apply(lambda x: x.split())
X_test_2['review'] = X_test_2['review'].map(' '.join)
X_test_2 = cv.transform(X_test_2['review'])

```

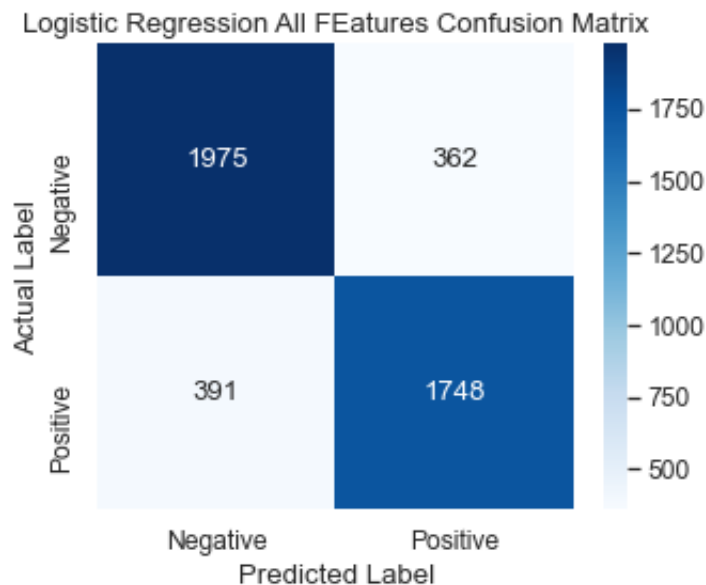


```
In [72]: logit_pred_all = logit.predict(X_test_2)
```

```
In [73]: ### Logistic Regression All Features Confusion Matrix w/ CountVectorizer

cm = confusion_matrix(y_test_2, logit_pred_all)
cm_df = pd.DataFrame(cm, index=['Negative', 'Positive'], columns=['Neg

fig_cm4 = plt.figure(figsize=(5,4))
sns.heatmap(cm_df, annot=True, fmt='g', cmap='Blues')
plt.title('Logistic Regression All FEatures Confusion Matrix')
plt.ylabel('Actual Label')
plt.xlabel('Predicted Label')
plt.show()
fig_cm4.savefig('images/logit_all_confusion_matrix.png', dpi=300)
```



```
In [74]: print('Logistic Regression All Features Accuracy: ', accuracy_score(y_
print('Logistic Regression All Features F1 Score: ', f1_score(y_test_2
print('Logistic Regression All Features Precision Score: ', precision_
print('Logistic Regression All Features Recall Score: ', recall_score(
```

```
Logistic Regression All Features Accuracy: 0.8317694369973191
Logistic Regression All Features F1 Score: 0.8317140169306872
Logistic Regression All Features Precision Score: 0.8317285793797539
Logistic Regression All Features Recall Score: 0.8317694369973191
```

**5th Model: Want to see training on Levonogestrel on its own**

- The previous models are evaluated on all Birth Control types, but want to see how the model performs on just one of the drugs

```
In [75]: data_lev = bc_data[bc_data['drugName'] == 'Levonorgestrel']  
data_lev.head()
```

```
very irritable  
but I&#039;m  
assuming it  
was because I  
was hurting  
and hungry all  
day. Other than  
that I&#039;ve  
been great!"  
  
"I went in to  
have my Skyla  
placed  
yesterday  
morning. After  
reading all of  
these reviews I  
was  
hyperventilating  
and crying on  
my way there,  
bc I did not  
want to
```

```
In [76]: data_lev['sentiment'].value_counts()
```

```
Out[76]: Positive      3168  
Negative      1728  
Name: sentiment, dtype: int64
```

In [77]: data\_lev.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4896 entries, 42 to 161257
Data columns (total 10 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   drugName            4896 non-null   object
 1   condition           4896 non-null   object
 2   review              4896 non-null   object
 3   rating              4896 non-null   float64
 4   date                4896 non-null   datetime64[ns]
 5   usefulCount         4896 non-null   int64
 6   sentiment           4896 non-null   object
 7   punc_emphasis       4896 non-null   int64
 8   capt_emphasis       4896 non-null   int64
 9   date_column         4896 non-null   object
dtypes: datetime64[ns](1), float64(1), int64(3), object(5)
memory usage: 420.8+ KB
```

```
In [78]: # Loading in review vs sentiment for only the Levonogestrel drug

X3 = data_lev['review']
y3 = data_lev['sentiment']

X_train_3, X_test_3, y_train_3, y_test_3 = train_test_split(X3, y3, te

# For Train Set, apply clean_text function

X_train_3 = X_train_3.apply(clean_text)

### Train - Tokenize the training data with a simple split of words, a

X_train_3 = X_train_3.apply(lambda x: x.split())
X_train_3 = X_train_3.map(' '.join)

### Train - Vectorize the training data using CountVectorizer

cv = CountVectorizer()
X_train_3 = cv.fit_transform(X_train_3)

### Train - Fit training data to Logistic Regression Model

logit_lev = LogisticRegression()
logit_lev.fit(X_train_3, y_train_3)

### VALIDATION - Perform a cross validation on the logistic regression

scores = cross_val_score(logit_lev, X_train_3, y_train_3, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())

Cross Validation Scores: [0.78061224 0.79182631 0.79438059 0.7879948
9 0.81481481]
Mean Cross Validation Score: 0.7939257695415332
```

```
In [79]: # Logistic Regression Test Set Preprocessing

X_test_3_logit = X_test_3.apply(clean_text)
X_test_3_logit = X_test_3.apply(lambda x: x.split())
X_test_3_logit = X_test_3.map(' '.join)
X_test_3_logit = cv.transform(X_test_3)
```

```
In [80]: logit_pred_3 = logit_lev.predict(X_test_3_logit)
```

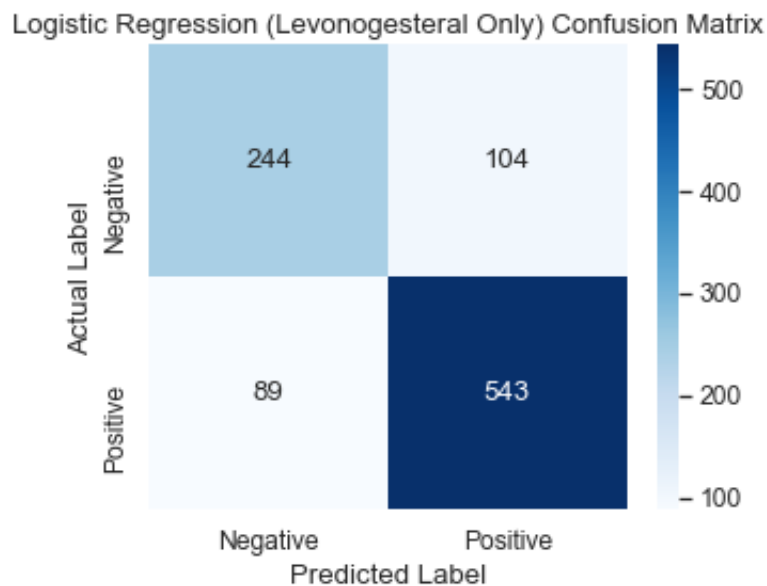
```
In [81]: ### Get scores for the sentiments individually to see

# Quick look at decision matrix for our fifth model:

### Logistic Regression Confusion Matrix

cm = confusion_matrix(y_test_3, logit_pred_3)
cm_df = pd.DataFrame(cm, index=['Negative', 'Positive'], columns=['Neg', 'Pos'])

fig_cm3 = plt.figure(figsize=(5,4))
sns.heatmap(cm_df, annot=True, fmt='g', cmap='Blues')
plt.title('Logistic Regression (Levonogesteral Only) Confusion Matrix')
plt.ylabel('Actual Label')
plt.xlabel('Predicted Label')
plt.show()
```



```
In [82]: print('Logistic Regression (Levonogesteral Only) Accuracy: ', accuracy)
print('Logistic Regression (Levonogesteral Only) F1 Score: ', f1_score)
print('Logistic Regression (Levonogesteral Only) Precision Score: ', precision)
print('Logistic Regression (Levonogesteral Only) Recall Score: ', recall)
```

```
Logistic Regression (Levonogesteral Only) Accuracy: 0.8030612244897959
Logistic Regression (Levonogesteral Only) F1 Score: 0.8020471356008727
Logistic Regression (Levonogesteral Only) Precision Score: 0.80143073395447
Logistic Regression (Levonogesteral Only) Recall Score: 0.8030612244897959
```

## 6th Model - TFIDF / Random Forest Classifier / All Features

```
In [83]: # Additional features to see if it helps improve the model

X2 = bc_data[['review', 'drugName', 'condition', 'punc_emphasis', 'cap
y2 = bc_data['sentiment']

X_train_4, X_test_4, y_train_4, y_test_4 = train_test_split(X2, y2, te

# For Train Set, apply clean_text function

X_train_4['review'] = X_train_4['review'].apply(clean_text)

### Train - Tokenize the training data with a simple split of words, a

X_train_4['review'] = X_train_4['review'].apply(lambda x: x.split())
X_train_4['review'] = X_train_4['review'].map(' '.join)

### Train - Vectorize the training data using CountVectorizer

tfidf = TfidfVectorizer()
X_train_4 = tfidf.fit_transform(X_train_4['review'])

### Train - Fit training data to Logistic Regression Model

rfc = RandomForestClassifier()
rfc.fit(X_train_4, y_train_4)

### VALIDATION - Perform a cross validation on the logistic regression

scores = cross_val_score(rfc, X_train_4, y_train_4, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())
```

```
Cross Validation Scores: [0.83998883 0.84389835 0.84357542 0.8438547
5 0.84189944]
Mean Cross Validation Score: 0.8426433582579692
```

### Grid Search for Best Conditions for RFC

```
In [84]: param_grid = {'n_estimators': [200, 300, 400, 500],
                        'max_features': ['auto', 'sqrt', 'log2'],
                        'max_depth' : [4,5,6,7,8],
                        'criterion': ['gini', 'entropy']}
grid = GridSearchCV(rfc, param_grid, cv=5, scoring='accuracy', n_jobs=
grid.fit(X_train_4, y_train_4)
print('Best Parameters: ', grid.best_params_)
```

Best Parameters: {'criterion': 'gini', 'max\_depth': 8, 'max\_features': 'sqrt', 'n\_estimators': 200}

```
In [85]: X_test_4['review'] = X_test_4['review'].apply(clean_text)
X_test_4['review'] = X_test_4['review'].apply(lambda x: x.split())
X_test_4['review'] = X_test_4['review'].map(' '.join)
X_test_4 = tfidf.transform(X_test_4['review'])
```

## Results of Grid Search Tuning of Hyper Parameters

Logit Results = {'C': 1, 'penalty': 'l2'}

NB Results = {'alpha': 0.5}

RFC results = {'criterion': 'gini', 'max\_depth': 8, 'max\_features': 'auto', 'n\_estimators': 300}

## Apply Hyper Parameters Tuning to Models and Evaluate Results

```
In [86]: # Applying the hyperparameters from the grid search to the Logistic Regression
logit = LogisticRegression()
logit.fit(X_train_1_tfidf, y_train_1_tfidf)

### Applying the hyperparameters from the grid search to the Naive Bayes
nb = MultinomialNB(alpha=0.5)
nb.fit(X_train_1_nb, y_train_1_nb)

### Applying the hyperparameters from the grid search to the Random Forest
rfc = RandomForestClassifier(criterion='gini', max_depth=8, max_features='sqrt')
rfc.fit(X_train_4, y_train_4)
```

```
Out[86]: RandomForestClassifier(max_depth=8, n_estimators=300)
```

## Running Models on Test Data and Evaluating Results

```
In [87]: # Logistic Regression Test Set Predictions
logit_pred = logit.predict(X_test_1_tfidf)

### ^ This is only using text vs sentiment, need to update this to include sentiment

# Multinomial Bayes Test Set Predictions
nb_pred = nb.predict(X_test_1_nb)

# Random Forest Classifier Test Set Predictions
rfc_pred = rfc.predict(X_test_4)
```

```
In [180]: X_train_1_nb_, X_test_1_nb_, y_train_1_nb_, y_test_1_nb_ = train_test_split(X_train_1_nb, y_train_1_nb, X_train_4, y_train_4, test_size=0.2, random_state=42)
```

```
In [181]: type(X_test_1_nb_)
```

```
Out[181]: pandas.core.series.Series
```



```
In [190]: X_test_1_nb_.reset_index().drop('index', axis=1)['review'][1]
```

```
Out[190]: 'birth control ever gotten cycle regulated long period time usually s  
witch medicine regulated month everything go haywire luter a breakthro  
ugh bleeding spotting headache body ache weight gain love birth contr  
ol anemic abnormal period bleeding week month even take iron suppleme  
nt also pregnant work good'
```

```
In [169]: nb.predict(X_test_1_nb[1])
```

```
Out[169]: array(['Negative'], dtype='<U8')
```

```
In [88]: ### Logistic Regression Test Set Evaluation
```

```
print('Logistic Regression Accuracy: ', accuracy_score(y_test_1_tfidf,  
print('Logistic Regression F1 Score: ', f1_score(y_test_1_tfidf, logit  
print('Logistic Regression Precision Score: ', precision_score(y_test_  
print('Logistic Regression Recall Score: ', recall_score(y_test_1_tfidf
```

```
Logistic Regression Accuracy:  0.798704200178731  
Logistic Regression F1 Score:  0.7978723877450785  
Logistic Regression Precision Score:  0.8102055684985682  
Logistic Regression Recall Score:  0.798704200178731
```

```
In [89]: ### Naive Bayes Test Set Evaluation
```

```
print('Naive Bayes Accuracy: ', accuracy_score(y_test_2, nb_pred))  
print('Naive Bayes F1 Score: ', f1_score(y_test_2, nb_pred, average='w  
print('Naive Bayes Precision Score: ', precision_score(y_test_2, nb_pr  
print('Naive Bayes Recall Score: ', recall_score(y_test_2, nb_pred, av
```

```
Naive Bayes Accuracy:  0.8197050938337802  
Naive Bayes F1 Score:  0.8195849999935735  
Naive Bayes Precision Score:  0.8196954748629014  
Naive Bayes Recall Score:  0.8197050938337802
```

In [90]: *### Random Forest Test Set Evaluation*

```
print('Random Forest Accuracy: ', accuracy_score(y_test_4, rfc_pred))
print('Random Forest F1 Score: ', f1_score(y_test_4, rfc_pred, average
print('Random Forest Precision Score: ', precision_score(y_test_4, rfc
print('Random Forest Recall Score: ', recall_score(y_test_4, rfc_pred,
```

```
Random Forest Accuracy:  0.711349419124218
Random Forest F1 Score:  0.6874723745521364
Random Forest Precision Score:  0.7799721388607994
Random Forest Recall Score:  0.711349419124218
```

## Deployment of Model: Making Word Clouds and Visualizations of Predicted Sentiment

In [115]: *# Create function that takes in our bc\_data that contains the review a*

```
def create_word_clouds(bc_data):
    bc_data['review'] = bc_data['review'].apply(clean_text)
    negative = bc_data.loc[bc_data['sentiment'] == 'Negative']
    positive = bc_data.loc[bc_data['sentiment'] == 'Positive']
    negative_text = ' '.join(negative['review'])
    positive_text = ' '.join(positive['review'])
    negative_wordcloud = WordCloud(width=600, height=400, collocations
                                   background_color='white').generate(n
    positive_wordcloud = WordCloud(width=600, height=400, collocations
                                   background_color='white').generate(p
    fig, ax = plt.subplots(1, 2, figsize=(20, 10))
    fig.tight_layout(pad=-1.0)
    ax[0].imshow(negative_wordcloud, interpolation='bilinear')
    ax[0].set_title('Negative Sentiment \n', fontsize=28)
    ax[0].axis('off')
    ax[1].imshow(positive_wordcloud, interpolation='bilinear')
    ax[1].set_title('Positive Sentiment \n', fontsize=28)
    ax[1].axis('off')
    plt.show()
```

```
create_word_clouds(bc_data)
```



## Deployment of Model: Create Function for Testing Data with Model

```
In [193]: def multinomial_bayes(data):

    # Load in data, if it is the same format we have the tab separation
    df = pd.read_csv(data, sep='\t')

    # Data cleaning

    df = df.dropna(subset=['condition'])

    # Function for picking the top 7 birth controls

    bc_drugs_nb = ['Etonogestrel', 'Ethinyl estradiol / norethindrone',
                   'Ethinyl estradiol / norgestimate', 'Implanon']

    bc_data_nb = df[df['drugName'].isin(bc_drugs_nb)]

    # Further cleaning of conditions that were not correct

    bc_data_nb = bc_data_nb[~bc_data_nb['condition'].str.contains('com

    # Creating sentiment column to categorize positive or negative review

    bc_data_nb['sentiment'] = ['Positive' if x > 7.0 else 'Negative' for x in bc_data_nb['review']]

    # Creating a 'punc_emphasis' column that scores how many exclamation marks are in the review

    bc_data_nb['punc_emphasis'] = bc_data_nb['review'].apply(lambda x: x.count('!'))
```

```
# Creating a 'capt_emphasis' column that scores how many capitaliz
bc_data_nb['capt_emphasis'] = bc_data_nb['review'].apply(lambda x:

# Loading in the features for model analysis

X5 = bc_data_nb[['review', 'drugName', 'condition', 'punc_emphasis
y5 = bc_data_nb['sentiment']

# Train test split

X_train_5_nb, X_test_5_nb, y_train_5_nb, y_test_5_nb = train_test_

# Apply clean_text function for Train Set

X_train_5_nb['review'] = X_train_5_nb['review'].apply(clean_text)

# Tokenize the training data with a simple split of words

X_train_5_nb['review'] = X_train_5_nb['review'].apply(lambda x: x.
X_train_5_nb['review'] = X_train_5_nb['review'].map(' '.join)

# Instantiate CountVectorizer()

cv = CountVectorizer()
X_train_5_nb = cv.fit_transform(X_train_5_nb['review'])

# Train - fitting the training data to a Naive Bayes Classifier

nb = MultinomialNB(alpha=0.5)
nb.fit(X_train_5_nb, y_train_5_nb)

# Validation - Cross Validation on the Naive Bayes Classifier Mode

scores = cross_val_score(nb, X_train_5_nb, y_train_5_nb, cv=5)
print('Cross Validation Scores: ', scores)
print('Mean Cross Validation Score: ', scores.mean())

# Naive Bayes Classifier Test Preprocessing

X_test_5_nb['review'] = X_test_5_nb['review'].apply(clean_text)
X_test_5_nb['review'] = X_test_5_nb['review'].apply(lambda x: x.sp
X_test_5_nb['review'] = X_test_5_nb['review'].map(' '.join)
X_test_5_nb = cv.transform(X_test_5_nb['review'])

# Naive Bayes Classifier Predictions

nb_pred_final = nb.predict(X_test_5_nb)

### Naive Bayes Test Set Evaluation
```

```
print('Naive Bayes Accuracy: ', accuracy_score(y_test_5_nb, nb_pred_fina
print('Naive Bayes F1 Score: ', f1_score(y_test_5_nb, nb_pred_fina
print('Naive Bayes Precision Score: ', precision_score(y_test_5_nb, nb_p
print('Naive Bayes Recall Score: ', recall_score(y_test_5_nb, nb_p

return nb
```

```
In [194]: multinomial_bayes('data/drugsComTrain_raw.tsv')
```

```
Cross Validation Scores: [0.81533879 0.81310499 0.80156366 0.8023082
7 0.82390171]
Mean Cross Validation Score: 0.8112434847356663
Naive Bayes Accuracy: 0.819237641453246
Naive Bayes F1 Score: 0.8190684386952239
Naive Bayes Precision Score: 0.8190585424161556
Naive Bayes Recall Score: 0.819237641453246
```

```
Out[194]: MultinomialNB(alpha=0.5)
```

```
In [195]: model = multinomial_bayes('data/drugsComTrain_raw.tsv')
```

```
Cross Validation Scores: [0.81533879 0.81310499 0.80156366 0.8023082
7 0.82390171]
Mean Cross Validation Score: 0.8112434847356663
Naive Bayes Accuracy: 0.819237641453246
Naive Bayes F1 Score: 0.8190684386952239
Naive Bayes Precision Score: 0.8190585424161556
Naive Bayes Recall Score: 0.819237641453246
```

```
In [ ]: # Figure out how to save this as a pickle file
```