

Github Link

https://github.com/aakhunaizi/IT9002_202306801

Dataset Link

<https://snap.stanford.edu/data/web-FineFoods.htm>

✓ Task 1 - Problem Statement Formulation and Definition

- ✓ Motivation

Working as an Associate Product Manager for a leading FoodTech company in the region, Calo, has given me direct exposure to the constantly growing and evolving landscape of online commerce and how critical it is for businesses to have a better understanding of their customers than ever before.

Being a highly competitive market, a customer-obsessed company can have an edge over its competitors by investing time and resources into analysing its customers' reviews as they are an invaluable source of insights.

Analyzing the sentiment of customer reviews helps in identifying the customers' preferences and pain points when using a company's products and services, making it easier to make informed decisions on how resources should be allocated and what efforts should be prioritized to give customers the best experience.

Although this project does not use Calo customer reviews, it aims to serve as a proof of concept on the usage of sentiment analysis on customer reviews to gain a competitive advantage in the market.

Problem Statement

Analyzing the sentiment of customer reviews, if done in a manual and traditional manner, has proven to be a challenging and extremely time-consuming process, especially for regional or international companies that are scaling up at a rapid pace.

Customer reviews can be written in multiple languages, are diverse in nature, and can include the usage of sarcasm, informal language, slang, or even different writing styles. In addition to these factors, human bias and subjectivity make accurately understanding customer sentiment and gaining valuable insights an even more challenging process.

Expected Result

Implementing a customer review sentiment analysis solution powered by natural language processing (NLP) and machine learning (ML) techniques can help companies overcome the limitations of manual, traditional methods.

Some of its benefits include but are not limited to:

- Automating the customer review analysis process
- Freeing up valuable time and resources
- Gathering deeper customer insights by uncovering hidden trends in the data as well as predicting future customer behavior
- Improving the product by gaining a better understanding of the customers' pain points and prioritizing them based on their impact.
- Enhancing the customer experience by proactively tackling negative trends before they escalate and cause more damage.

The scope of this project will primarily focus on accurately predicting whether a customer's text review has a positive or negative sentiment.

Double-click (or enter) to edit

✓ Task 2 - Selection of an Appropriate Data Set (Data Collection)

Source of Dataset

The source of the dataset is the Stanford Network Analysis Project (SNAP) and consists of reviews of fine foods from amazon.

The dataset was selected based on the following factors:

- Reputable source of data (SNAP)
- Sufficient data size (~500,000 reviews)
- Quality of data (Reliable with minimal skew)
- Reviews originating from a large number of users (256,059 users)
- Large number of products being reviewed (74,258 products)
- Large number of users with over 50 reviews (260 users)
- Sufficient number of words per review (Median of 56 words)
- Duration of data gathering (October 1999 to October 2012)

- ✓ Initial Data Visualization and Summary of Data

```
1 #Importing pandas library
2 import pandas as pd
3
4 #Creating a dataframe df and reading the csv dataset
5 df = pd.read_csv(r"content/drive/MyDrive/Reviews.csv", nrows=50000)
6
7 #Displaying the first five rows of the dataset
8 df.head()
```

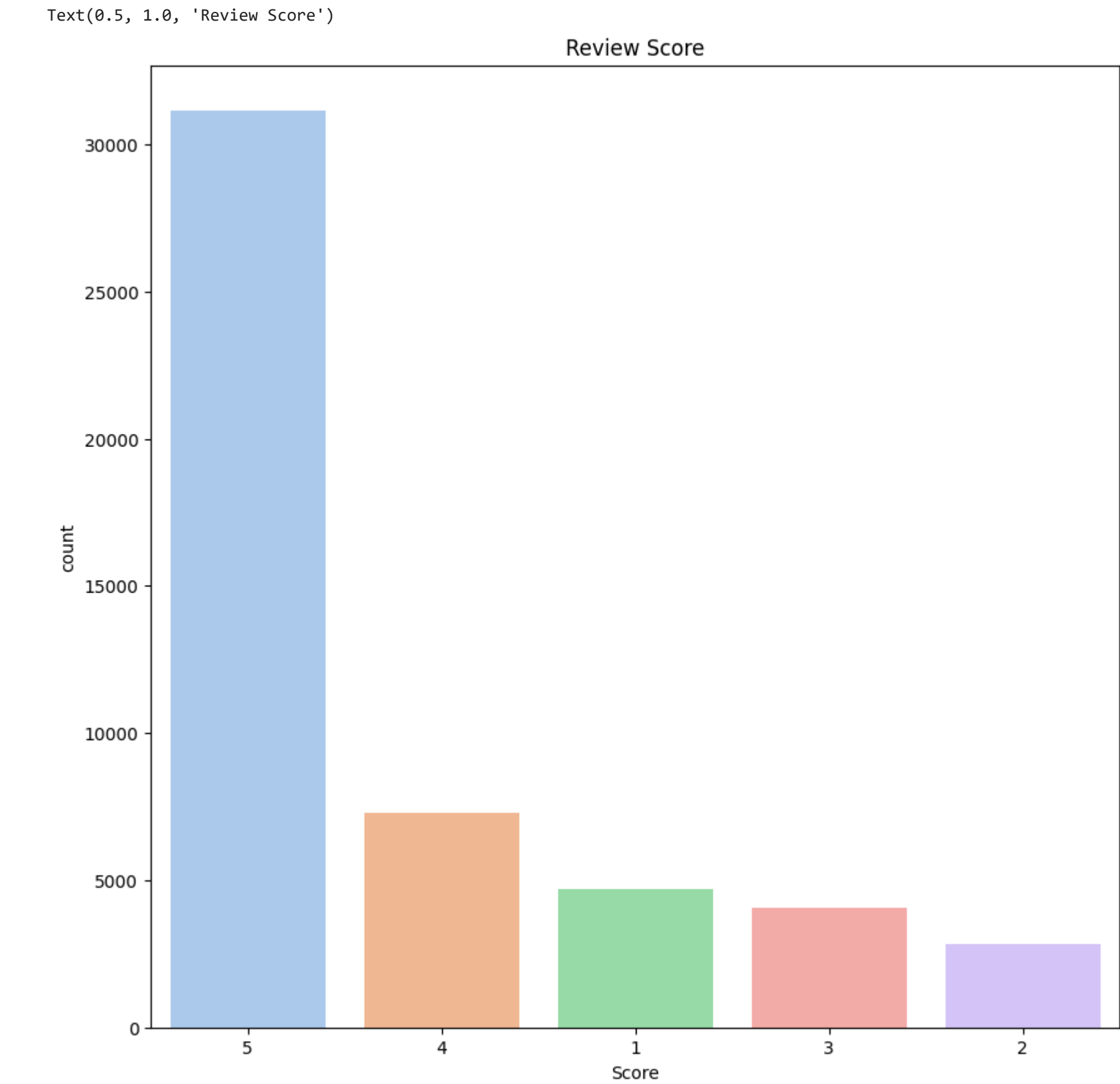
	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	B001E4KF6G	A3SGH71AUHU8GW	delmarian	1	1	1	1353862400	Good Quality Dog Food	I have bought several of the Vitality canned...
1	2	B00813GRG4	A1D87F6ZC0VENK	dill pa	0	1	1346897600	Not as Advertised	Product arrived labeled as Jumbo Salted Peann...	
2	3	B000LQOCH0	ABXLMMJXXJAN	Natalia Corres "Natalia Corres"	1	4	1219017600	"Delight" says it all	This is a confection that has been around a fe...	
3	4	B000UA0AQJ	A39SBORC6FGXVU	Karl	3	2	19332000	Cough Medicine	If you are looking for the secret ingredient I...	
4	5	B006KZZZ7K	A1UQRSCLF8W1T	Michael D. Bigham "M. Wassa"	0	0	1350777600	Great taffy	Great taffy at a great price. There was a wid...	

```
1 #Description of the dataset
2 df.info()
```

```
> class('pandas.core.frame.DataFrame')
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Id                     50000 non-null    int64
1   ProductId             50000 non-null    object
2   UserId                50000 non-null    object
3   ProfName              49997 non-null    object
4   HelpfulnessNumerator  50000 non-null    int64
5   HelpfulnessDenominator 50000 non-null    int64
6   Score                 50000 non-null    int64
7   Time                  50000 non-null    int64
8   Summary               49998 non-null    object
9   Text                  50000 non-null    object

dtypes: int64(5), object(5)
memory usage: 3.84 MB
```

```
1 #Review score distribution
2
3 #Importing seaborn and pyplot
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 #Plotting a countplot of the review scores (1-5)
8 sns.countplot(x='Score', data=df, palette="pastel", order = df['Score'].value_counts().index)
9 fig = plt.gcf()
10 fig.set_size_inches(10,10)
11 plt.title('Review Score')
```



- ✓ Shape of the Dataset

```
1 #Displaying data shape
2 display(df.shape)

(50000, 10)
```

Dataset Labels

Score is the only label within the dataset and it is used to extract the sentiment values.

```
1 #Dropping irrelevant columns
2 df.drop(['Id', 'ProfileName', 'HelpfulnessNumerator', 'HelpfulnessDenominator', 'Time', 'Summary'], axis=1, inplace=True)
```

```
1 #Displaying the first five rows of the dataset
2 df.head()
```

	ProductId	UserId	Score	Text
0	B001E4KFG0	A3SGXH7AUHU8GW	5	I have bought several of the Vitality canned d...
1	B00813GRG4	A1D87F6ZCVE5NK	1	Product arrived labeled as Jumbo Salted Peanut...
2	B000LQOCH0	ABXLMWJXXAIN	4	This is a confection that has been around a fe...
3	B000UA0QIQ	A395BORC8FGVXV	2	If you are looking for the secret ingredient i...
4	B006K2ZZ7K	A1UQRSCLF8GW1T	5	Great taffy at a great price. There was a wid...

```
1 #Description of the dataset
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0   ProductId   50000 non-null   object
1   UserId      50000 non-null   object
2   Score       50000 non-null   int64
3   Text        50000 non-null   object
dtypes: int64(1), object(3)
memory usage: 1.5+ MB
```

```
1 #Displaying data shape
2 display(df.shape)
```

```
(50000, 4)
```

```
1 #Finding duplicate rows
2 df.duplicated().sum()
```

```
101
```

```
1 #Dropping duplicate rows
2 df.drop_duplicates(inplace=True)
```

```
1 #Sum of duplicate rows
2 df.duplicated().sum()
```

```
0
```

```
1 #Sum of null value rows
2 df.isnull().sum()
```

```
ProductId      0
UserId         0
Score          0
Text           0
dtype: int64
```

```
1 #Assigning reviews with score > 3 as positive sentiment
2 #Assigning reviews with score < 3 as negative sentiment
3 #Removing reviews with score = 3
4 df = df[df.Score != 3]
5 df['Sentiment'] = df['Score'].apply(lambda rating : 1 if rating > 3 else 0)
```

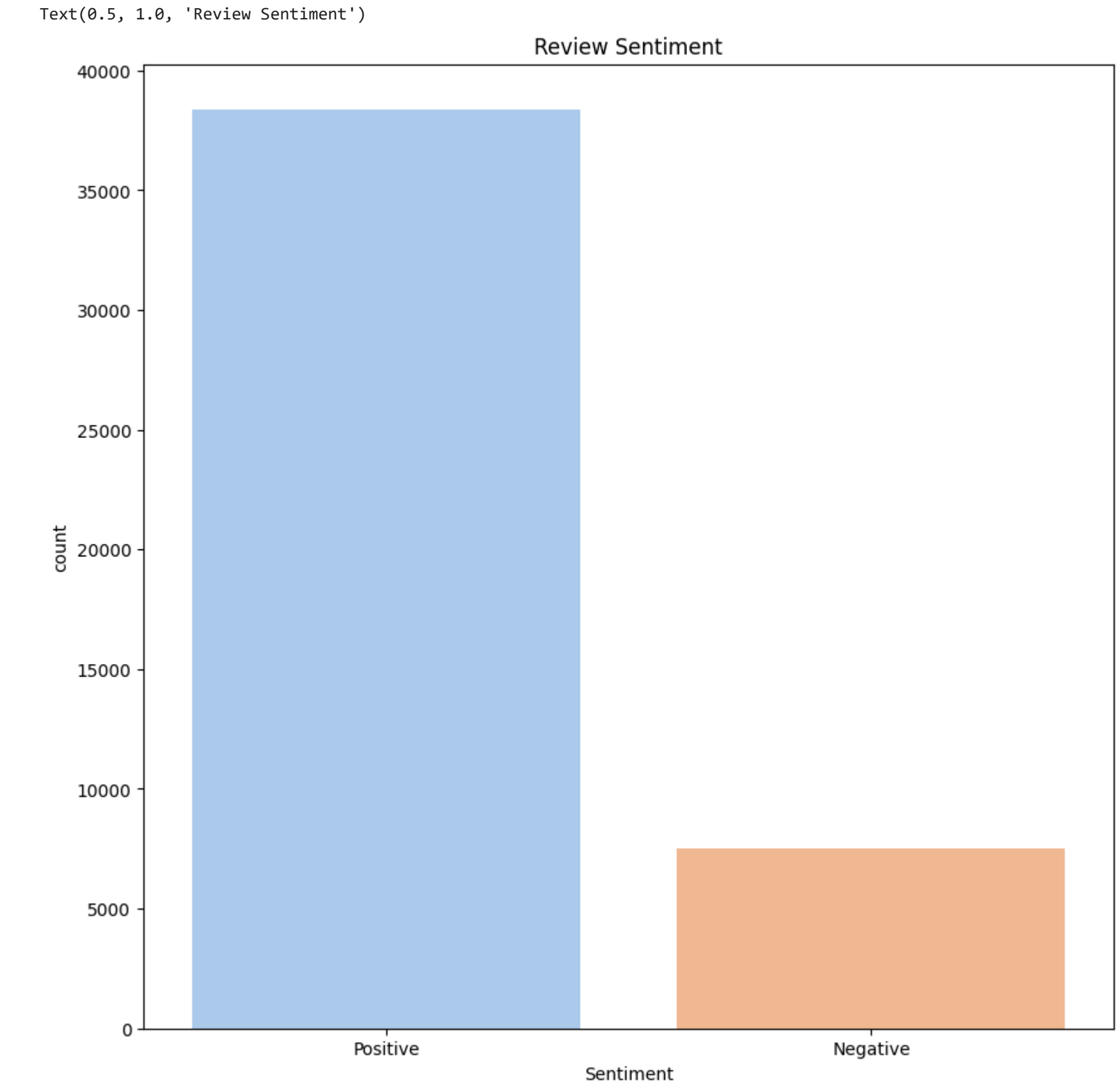
```
1 #Displaying the first five rows of the dataset
2 df.head()
```

	ProductId	UserId	Score	Text	Sentiment
0	B001E4KFG0	A3SGXH7AUHU8GW	5	I have bought several of the Vitality canned d...	1
1	B00813GRG4	A1D87F6ZCVE5NK	1	Product arrived labeled as Jumbo Salted Peanut...	0
2	B000LQOCH0	ABXLMWJXXAIN	4	This is a confection that has been around a fe...	1
3	B000UA0QIQ	A395BORC8FGVXV	2	If you are looking for the secret ingredient i...	0
4	B006K2ZZ7K	A1UQRSCLF8GW1T	5	Great taffy at a great price. There was a wid...	1

```
1 #Checking if any reviews with score = 3 exist
2 df.loc[df['Score'] == 3]
```

```
ProductId  UserId  Score  Text  Sentiment
```

```
1 #Distribution of positive and negative sentiment reviews
2 sns.countplot(x='Sentiment', data=df, order=df.Sentiment.value_counts().index, palette="pastel")
3 fig = plt.gcf()
4 fig.set_size_inches(10,10)
5 plt.xticks([0, 1], ['Positive', 'Negative'])
6 plt.xlabel("Sentiment")
7 plt.title("Review Sentiment")
```



```
1 #Total positive and negative sentiments
2 df.Sentiment.value_counts()
3 #1 is positive sentiment
4 #0 is negative sentiment
```

```
1      38348
0       7515
Name: Sentiment, dtype: int64
```

Task 3 - Text Preprocessing

This section includes a series of processes to clean, simplify, and standardize text to make it easier for natural language processing (NLP) algorithms to understand and extract the meaning of text.

Converting Text to Lowercase

```
1 #Creating a function to convert text to lowercase
2 def text_lowercase(text):
3     return text.lower()
4
5 #Converting all words in the 'Text' column to lowercase
6 df['Text'] = df['Text'].apply(text_lowercase)
```

Removing URLs from Text

```
1 #Importing regular expressions
2 import re
3
4 #Creating a function to remove URLs from text
5 def remove_url(text):
6     text = re.sub(r'http[s+]', '', text)
7     return text
8
9 #Removing any existing URLs in the rows of the 'Text' column
10 df['Text'] = df['Text'].apply(remove_url)
```

Removing HTML Tags from Text

```
1 #Creating a function to remove HTML tags from text
2 def remove_html(text):
3     pattern = re.compile('<.*>')
4     text = re.sub(pattern, '', text)
5     return text
6
7 #Removing any existing HTML tags in the rows of the 'Text' column
8 df['Text'] = df['Text'].apply(remove_html)
```

Removing Punctuation from Text

```
1 #Importing string
2 import string
3
4 #Creating a function to remove punctuation from text
5 def remove_punctuation(text):
6     translator = str.maketrans('', '', string.punctuation)
7     return text.translate(translator)
8
9 #Removing any existing punctuation in the rows of the 'Text' column
10 df['Text'] = df['Text'].apply(remove_punctuation)
```

Removing Whitespaces from Text

```
1 #Creating a function to remove whitespaces from text
2 def remove_whitespace(text):
3     return " ".join(text.split())
4
5 #Removing any existing whitespaces in the rows of the 'Text' column
6 df['Text'] = df['Text'].apply(remove_whitespace)
```

Removing Stopwords


```
1 #Importing the corpus of stopwords
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
4
5 #Creating a function to remove stopwords
6 def remove_stopwords(text):
7     stop_words = set(stopwords.words("english"))
8     word_tokens = word_tokenize(text)
9     filtered_text = [word for word in word_tokens if word not in stop_words]
10    return ''.join(filtered_text)
11
12 #Removing stopwords from all rows in the 'Text' column
13 df['Text'] = df['Text'].apply(remove_stopwords)
```

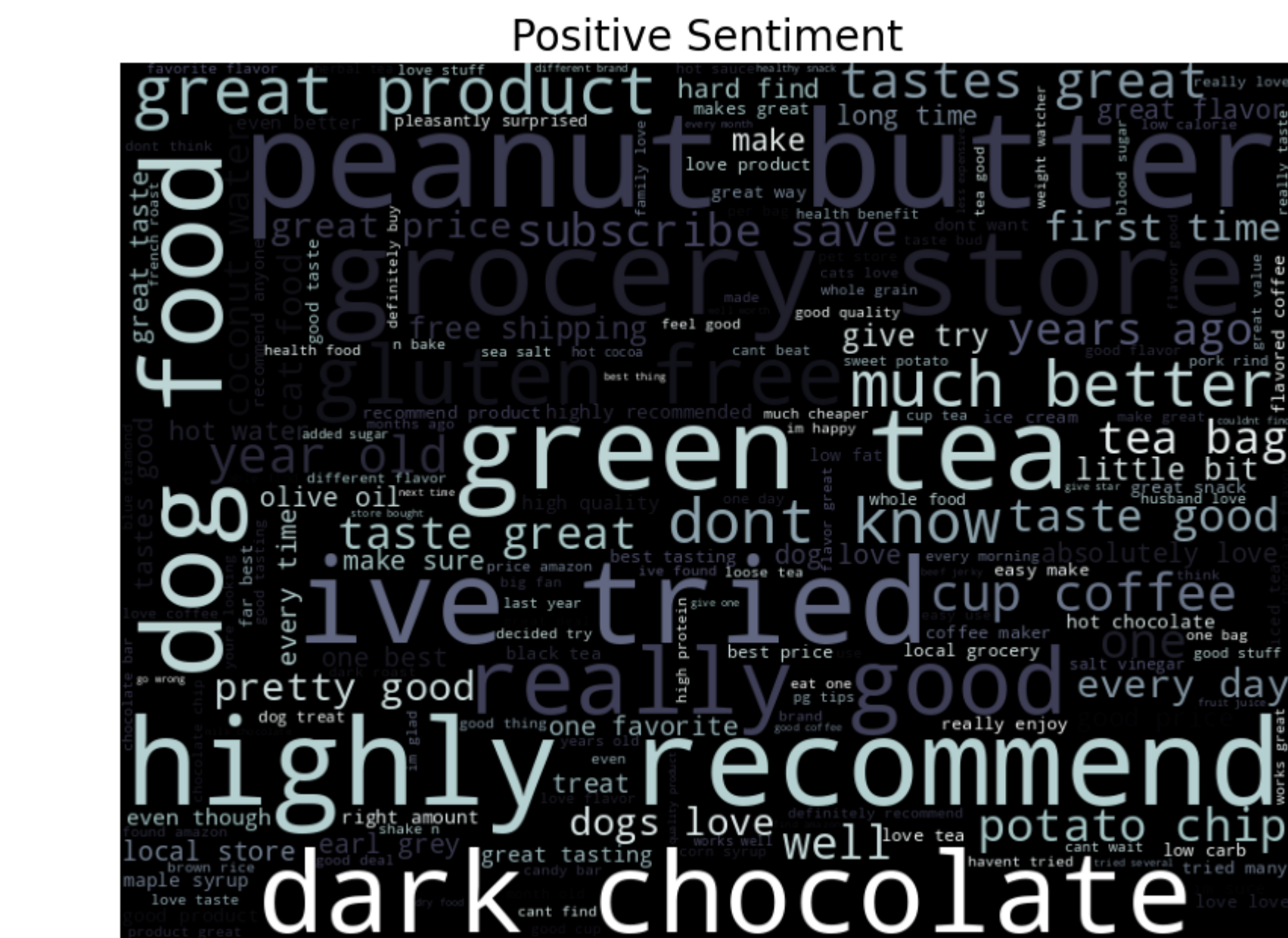
```
1 #Positive sentiment
2 positive = df[df['Sentiment'] == 1].copy()
3 #Negative sentiment
4 negative = df[df['Sentiment'] == 0].copy()
5
6 from sklearn.utils import resample
7
8 #Resampling
9 positive = resample(positive, replace=False, n_samples=len(negative), random_state=123)
10 df = pd.concat([positive, negative])
11 print(df.Sentiment.value_counts())
```

```

1 import wordcloud
2 from wordcloud import WordCloud
3 def show_wordcloud(df, title):
4     wordcloud_text = ',.join(df["Text"].astype(str))
5     stopwords = set(wordcloud.STOPWORDS)
6
7     fig_wordcloud = WordCloud(wordcloud_text=stopwords, background_color='black',
8                               colormap='bone', collocation_threshold = 3, width=886, height=688).generate(wordcloud_text)
9
10    plt.figure(figsize=(10,10), frameon=True)
11    plt.imshow(fig_wordcloud)
12    plt.axis('off')
13    plt.title(title, fontsize=20 )
14    plt.show()

```

```
1 #Positive sentiment wordcloud
2 show_wordcloud(positive, "Positive Sentiment")
```



```
1 #Negative sentiment wordcloud
2 show_wordcloud(negative, "Negative Sentiment")
```



This section includes a series of processes to convert raw text data into a format that can be easily understood and analyzed by computers to identify and extract relevant information from the text.

```
1 #Creating the TF-IDF model
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 cv = TfidfVectorizer(max_features = 10000)
4 X = cv.fit_transform(df['Text'])

1 X.toarray()

array([[0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       ...,
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.]])

1 #Feature names
2 cv.get_feature_names_out()

array(['08', '08oz', '0999', ..., 'zoo', 'zucchini', 'zukes'], dtype=object)
```

<https://colab.research.google.com/drive/19TVsj9UyaDP-jK7khFTGPHoJ5tKLW/?authuser=2#scrollTo=y9MQmjVAnmk&printMode=true>


```
1 #Importing NLTK and downloading the necessary packages
2 import nltk
3 nltk.download('punkt')
4 nltk.download('averaged_perceptron_tagger')

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /root/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
True

1 #Importing the necessary packages
2 from nltk.tokenize import word_tokenize
3 from nltk import pos_tag
4
5 #Creating a function to convert text into word tokens with their tags
6 def pos_tagging(text):
7     word_tokens = word_tokenize(text)
8     return pos_tag(word_tokens)
9
10 #Applying PoS tagging
11 df['Text'].apply(pos_tagging)

12895 [(2, CD), (dogs, NNS), (product, NN), (used, V...
37498 [(vegetarian, JJ), (runner, NN), (needed, VBD)...
29665 [(price, NN), (fair, NN), (works, VBD), (fine,...
32618 [(read, JJ), (reviews, NGS), (coconut, VBP), (...
17089 [(reading, VBG), (book, NN), (food, NN), (pets...
...
49968 [(oily, RB), (like, IN), (flavor, NN), (combin...
49973 [(bought, VBD), (3, CD), (flavors, NNS), (bars...
49975 [(local, JJ), (grocer, NN), (recently, RB), (d...
49977 [(review, NN), (hazelnut, NN), (chocolate, NN)...
49986 [(saw, JJ), (local, JJ), (health, NN), (food, ...
Name: Text, Length: 15030, dtype: object
```

Word Embeddings (Word2Vec)

```
1 #Importing the necessary packages
2 from gensim.models import Word2Vec

1 # Preparing the dataset
2 sentences = df['Text'].apply(word_tokenize)

1 sentences[3]

['looking',
 'secret',
 'ingredient',
 'robitussin',
 'believe',
 'found',
 'got',
 'addition',
 'root',
 'beer',
 'extract',
 'ordered',
 'good',
 'made',
 'cherry',
 'soda',
 'flavor',
 'medicinal']

1 # Training the Word2Vec model
2 model = Word2Vec(sentences, min_count=1)

1 # Finding Word Vectors
2 vector = model.wv['nice']

1 vector

array([-0.590732 ,  0.92159253,  0.1848151 ,  1.5201386 , -0.943574 ,
        -1.415056 ,  0.80786003,  1.1470503 ,  0.9573283 ,  0.04966137,
         0.33011842, -2.3189325 ,  0.35450727,  0.7481503 , -0.40901548,
        -0.76155436,  0.9590922 , -0.45441058, -0.04096768, -1.1891639 ,
         0.24980803,  0.97609363,  0.33616623,  0.00270797, -0.14561011,
         1.1671588 ,  0.59211636,  0.40352037, -0.3215095 ,  0.5768783 ,
         1.4682668 ,  0.7994466 , -0.39168414,  0.23299551, -0.633959 ,
         1.2858373 ,  0.19383235,  0.18010634, -0.11365937, -0.3674978 ,
         0.12718546,  1.15359319 ,  0.39793825, -1.1267297 ,  0.60801446,
        -0.25949377, -0.58693475, -0.46380773, -0.0125768 ,  0.56996644,
         0.64050915,  0.23959985, -0.57818973, -0.7628686 ,  0.06949516,
        -0.6013539 ,  1.3061968 ,  0.29499927, -0.8568597 ,  0.22514129,
        -0.37452996,  1.0736032 ,  0.2740991 ,  0.34485506, -0.9845416 ,
         0.26510638, -0.65449625,  0.12075504, -0.92173046,  1.2225282 ,
        -0.5074568 , -0.09389084,  0.6696838 , -0.32956028,  0.10613009,
        -0.13416097, -0.23970069 , -1.0207237 , -0.60013 ,  0.11266125,
        -0.02002736, -0.04463236, -0.8999008 ,  0.71046346, -0.8675138 ,
        -1.2307748 , -0.36060684,  1.453321 ,  0.66359425, -0.13236582,
        -0.04292237,  0.57859665,  0.2707664 ,  0.43542433,  1.6847239 ,
         0.45083925,  0.73301335, -0.15457162,  0.30170008,  0.792024485],
      dtype=float32)

1 #Most similar words
2 similar = model.wv.most_similar('nice')

1 similar

[('light', 0.9425188302993774),
 ('subtle', 0.9390808343887329),
 ('pleasant', 0.9384437203407288),
 ('overpowered', 0.9327387213706977),
 ('bitterness', 0.9298967123031616),
 ('overpowering', 0.926109254360199),
 ('rich', 0.9244909286490823),
 ('pungent', 0.9240808073695374),
 ('mild', 0.92403444166755676),
 ('lingering', 0.9233019351959229)]
```

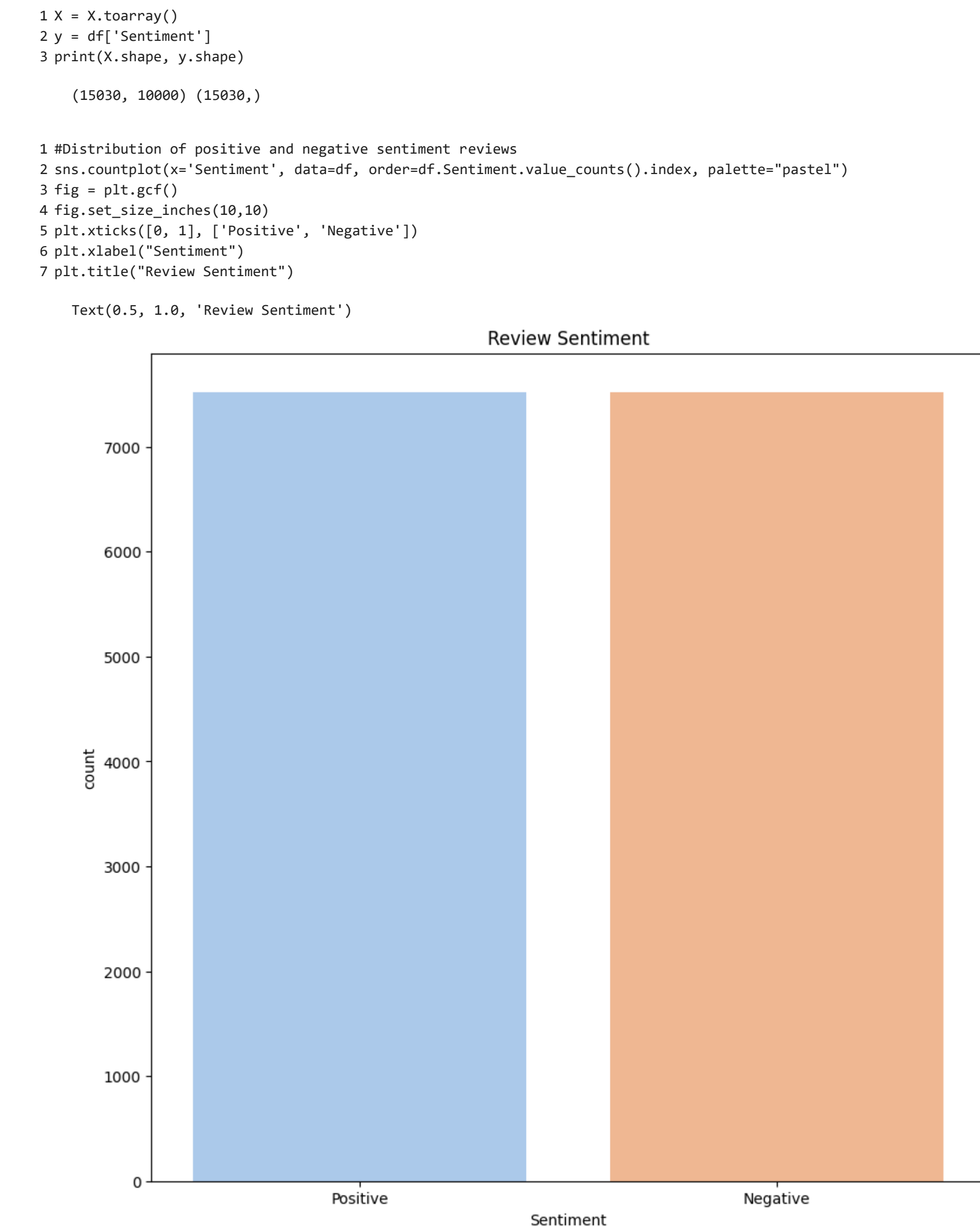
Task 5 - Text Classification/Prediction and Task 6 - Evaluation, Inferences, Recommendation and Reflection

Justification for Method Selection

Naive Bayes, Multinomial Naive Bayes, and Logistic Regression were selected as the classification methods due to their effectiveness for multi-class problems. They perform well at differentiating between multiple text categories and are a great choice for a use case such as sentiment analysis.

Moreover, they are highly accurate and scalable, being able to handle massive collections of text without any issues.

Although more classifiers were tested, these three classifiers were selected as there was a noticeable variation between their performances, facilitating for a clearer and more streamlined evaluation process.



```
1 #Train-Test Split
2 from sklearn.model_selection import train_test_split
3
4 SEED=123
5 X_train,X_test,y_train,y_test=train_test_split(X, y, test_size=0.3, random_state=SEED, stratify=y)
6 print(X_train.shape, y_train.shape)
7 print(X_test.shape, y_test.shape)

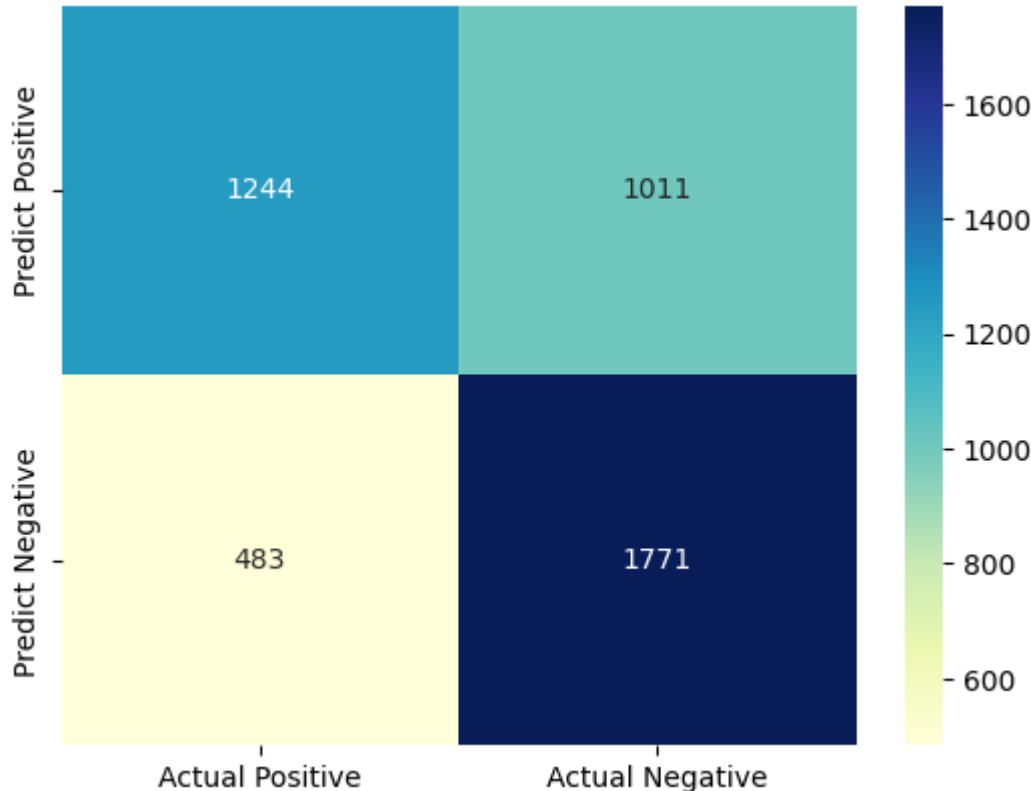
(10521, 10000) (10521,)
(4509, 10000) (4509,)
```

Gaussian Naive Bayes

```
1 #Gaussian Naive Bayes
2 from sklearn.naive_bayes import GaussianNB
3 from sklearn.metrics import accuracy_score
4
5 gnb = GaussianNB()
6 gnb.fit(X_train, y_train)
7 y_pred_train = gnb.predict(X_train)
8 y_pred_test = gnb.predict(X_test)
9 print("\nTraining Accuracy score:",accuracy_score(y_train, y_pred_train))
10 print("Testing Accuracy score:",accuracy_score(y_test, y_pred_test))

Training Accuracy score: 0.8557171371542629
Testing Accuracy score: 0.6686626746506986
```

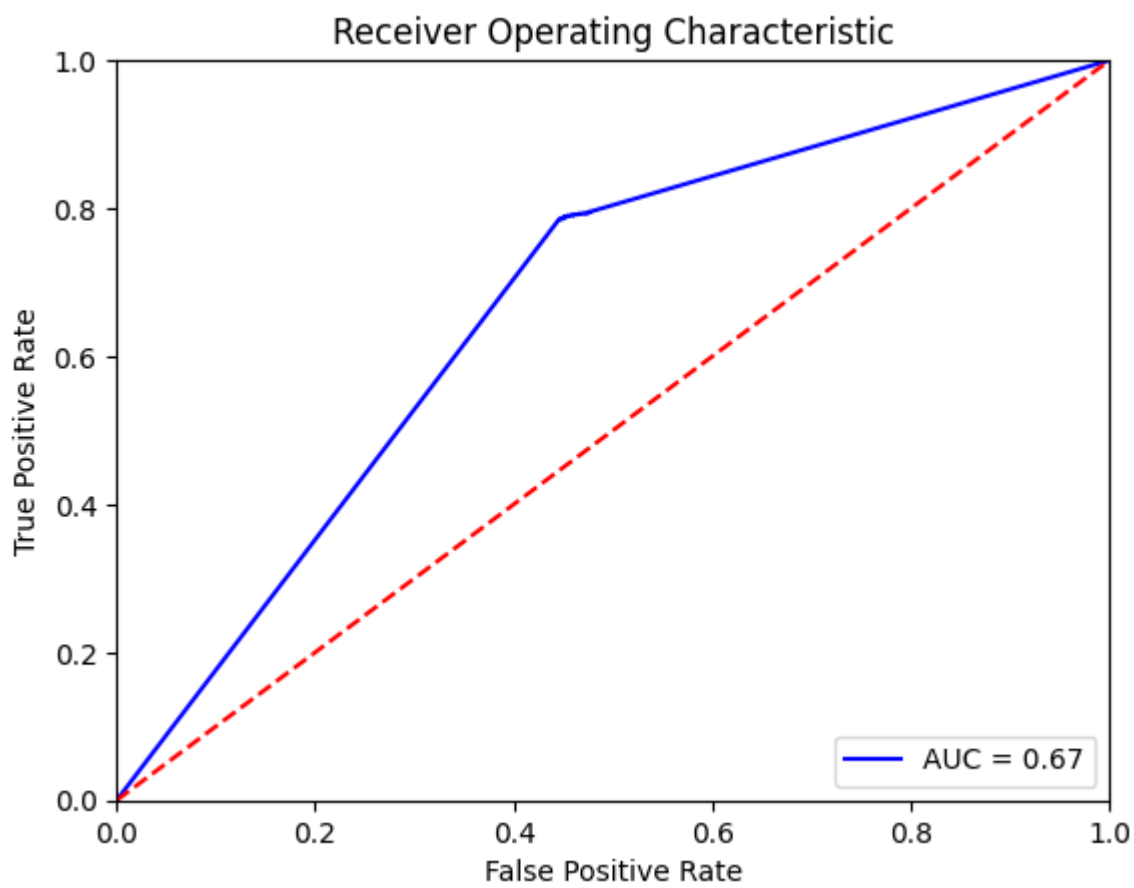
```
1 #Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3
4 cm = confusion_matrix(y_test, y_pred_test)
5
6 cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive', 'Actual Negative'],
7                             index=['Predict Positive', 'Predict Negative'])
8 sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='YlGnBu')
9 plt.show()
```



```
1 #Classification report
2 from sklearn.metrics import classification_report
3 print(classification_report(y_test, y_pred_test, target_names=['negative', 'positive']))
```

	precision	recall	f1-score	support
negative	0.72	0.55	0.62	2255
positive	0.64	0.79	0.70	2254
accuracy			0.67	4589
macro avg	0.68	0.67	0.66	4589
weighted avg	0.68	0.67	0.66	4589

```
1 #AUC-ROC Curve
2 from sklearn import metrics
3 probs = gnb.predict_proba(X_test)
4 preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
6 roc_auc = metrics.auc(fpr, tpr)
7
8 plt.title('Receiver Operating Characteristic')
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right')
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.ylabel('True Positive Rate')
15 plt.xlabel('False Positive Rate')
16 plt.show()
```

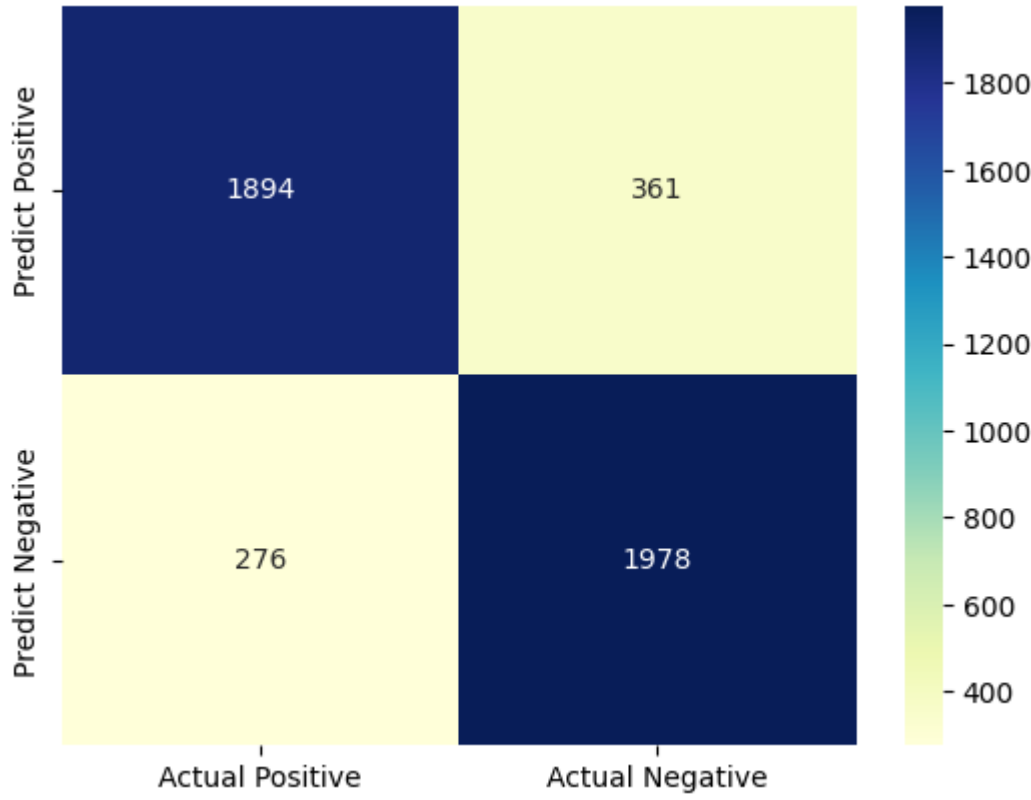


▼ Multinomial Naive Bayes

```
1 from sklearn.naive_bayes import MultinomialNB
2
3 mnb = MultinomialNB()
4 mnb.fit(X_train, y_train)
5
6 y_pred_train = mnb.predict(X_train)
7 y_pred_test = mnb.predict(X_test)
8 print("\nTraining Accuracy score:", accuracy_score(y_train, y_pred_train))
9 print("Testing Accuracy score:", accuracy_score(y_test, y_pred_test))
```

Training Accuracy score: 0.9124607927903137
Testing Accuracy score: 0.8587269904635174

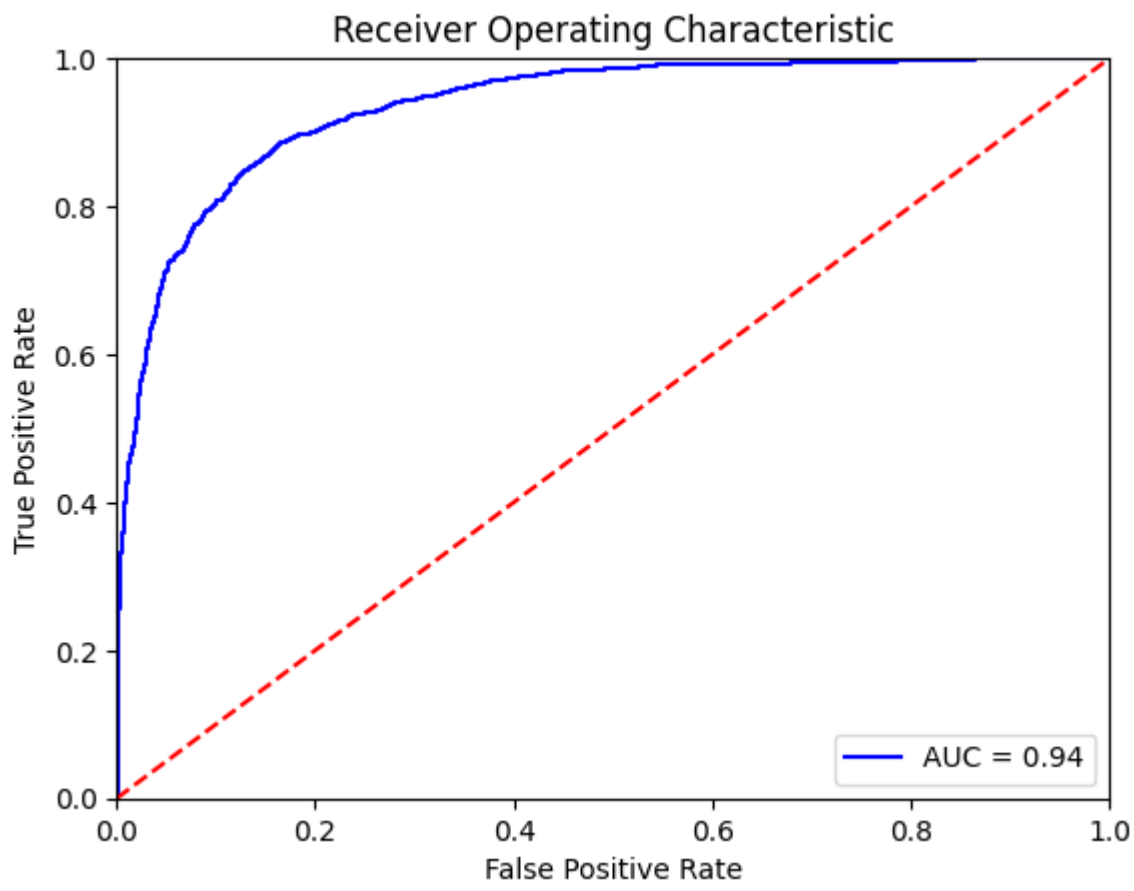
```
1 #Confusion Matrix
2 cm = confusion_matrix(y_test, y_pred_test)
3
4 cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive', 'Actual Negative'],
5                             index=['Predict Positive', 'Predict Negative'])
6 sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='YlGnBu')
7 plt.show()
```



```
1 #Classification report
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 print(classification_report(y_test, y_pred_test, target_names=['negative', 'positive']))
```

	precision	recall	f1-score	support
negative	0.87	0.84	0.86	2255
positive	0.85	0.88	0.86	2254
accuracy			0.86	4589
macro avg	0.86	0.86	0.86	4589
weighted avg	0.86	0.86	0.86	4589

```
1 #AUC-ROC Curve
2 probs = mnb.predict_proba(X_test)
3 preds = probs[:,1]
4 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
5 roc_auc = metrics.auc(fpr, tpr)
6
7 plt.title('Receiver Operating Characteristic')
8 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
9 plt.legend(loc = 'lower right')
10 plt.plot([0, 1], [0, 1], 'r--')
11 plt.xlim([0, 1])
12 plt.ylim([0, 1])
13 plt.ylabel('True Positive Rate')
14 plt.xlabel('False Positive Rate')
15 plt.show()
```

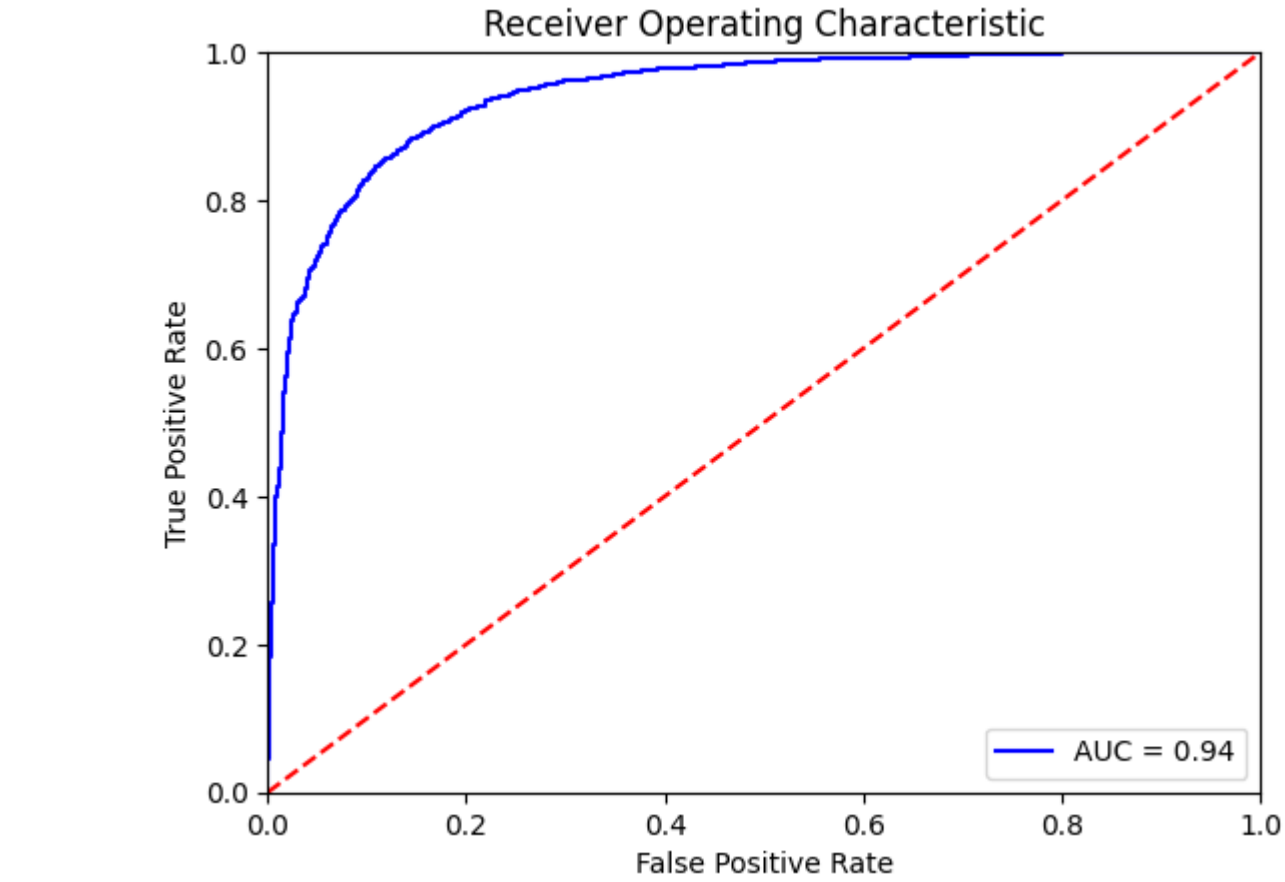
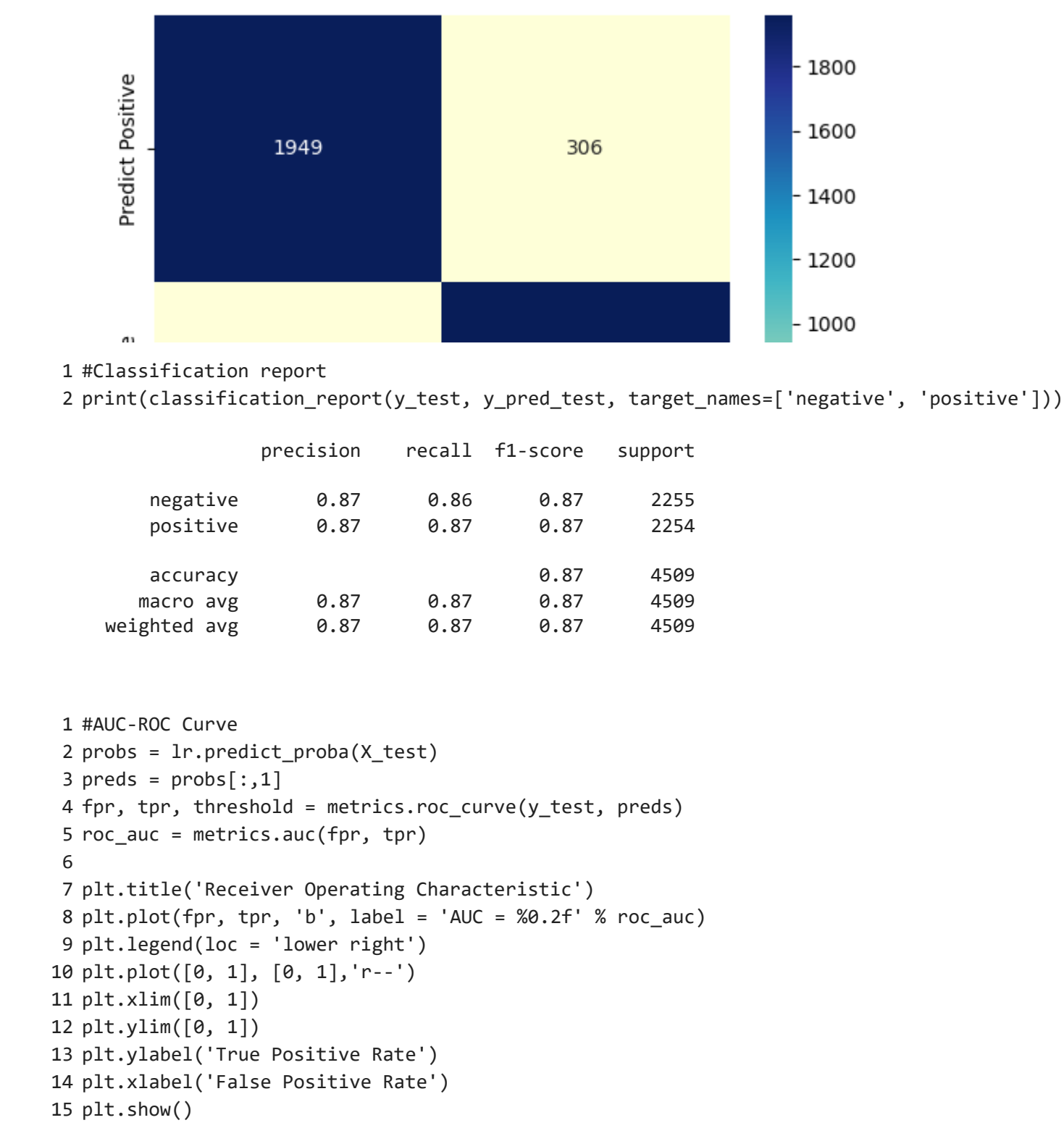


▼ Logistic Regression Classifier

```
1 from sklearn.linear_model import LogisticRegression
2
3 lr = LogisticRegression(random_state=5650)
4 lr.fit(X_train, y_train)
5
6 y_pred_train = lr.predict(X_train)
7 y_pred_test = lr.predict(X_test)
8 print("\nTraining Accuracy score:", accuracy_score(y_train, y_pred_train))
9 print("Testing Accuracy score:", accuracy_score(y_test, y_pred_test))
```

Training Accuracy score: 0.9196844486425245
Testing Accuracy score: 0.86715457972943

```
1 #Confusion Matrix
2 cm = confusion_matrix(y_test, y_pred_test)
3
4 cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive', 'Actual Negative'],
5                             index=['Predict Positive', 'Predict Negative'])
6 sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='YlGnBu')
7 plt.show()
```

Evaluation and Inference

It is evident from the plots that the AUC for the Logistic Regression ROC curve is higher than that for the Gaussian Naive Bayes ROC curve, but relatively similar to the Multinomial Naive Bayes ROC curve. Therefore, we can say that Logistic Regression clearly did a better job of classifying the positive sentiment than Gaussian Naive Bayes, but had only a slightly better performance than Multinomial Naive Bayes.

The classifiers rank as follows (best to worst):

1. Logistic Regression
2. Multinomial Naive Bayes
3. Gaussian Naive Bayes

Recommendation

Based on the evaluation and inference of the various metrics discussed above, Logistic Regression and Multinomial Naive Bayes are both great classifier choices for projects revolving around sentiment analysis.

Reflection

Undertaking this project was a wonderful journey of exploring the intersection between code and linguistics. Converting a mere set of endless rows of data into a beautiful set of visualizations and insights was nothing short of fascinating.

Similar to many aspects of life, this journey also had its fair share of challenges and hardships.

To begin, being excited to work on a project of this nature for the first time overloaded me with ideas of topics to explore, and it took me a decent amount of time to settle on a topic that was not too complex to be completed within the deadline set for the project, but also accommodated for future work and enhancements related to my field of work.

The dataset selection process was not straightforward either as it was initially difficult to determine the qualities and factors that make up a good dataset. The challenges persisted with text preprocessing, text representation, and text classification as these processes required lots of research and code customization, and countless trials and errors to reach a satisfactory result.

A major learning that revealed itself towards the end of the project was that the size of the dataset, containing over 500,000 rows, was the root cause of many unanticipated issues such as long code execution times and crashes due to insufficient computational resources. This, however, was rectified by limiting the number of rows in the dataframe to 50,000.

In addition, the reviews were skewed heavily towards the positive sentiment, causing a massive drop in the accuracy and overall performance of the classifiers. This was combated by resampling the data and having an even split between reviews of positive sentiment and negative sentiment.

Beyond the scope of this course, I plan to revisit this project with a different dataset consisting of food-related customer reviews submitted by my organization's customers, and combining it with automated flows that proactively award loyal customers or compensate customers that had subpar experiences with their meals.

References

1. <https://snap.stanford.edu/data/web-FineFoods.html>
2. <https://www.kaggle.com/datasets/snap/amazon-fine-food-reviews>
3. <https://www.kaggle.com/code/shashankai/text-preprocessing-using-python>
4. <https://www.kaggle.com/code/sonalisingh1411/nlp-part-1-amazon-fine-food-sentiment-analysis>
5. <https://www.kaggle.com/code/kshitiimohan/sentiment-analysis-universal-sentence-encoder-91>