

Capstone 3:

Customer Feedback With Natural Language Processing

Problem Statement

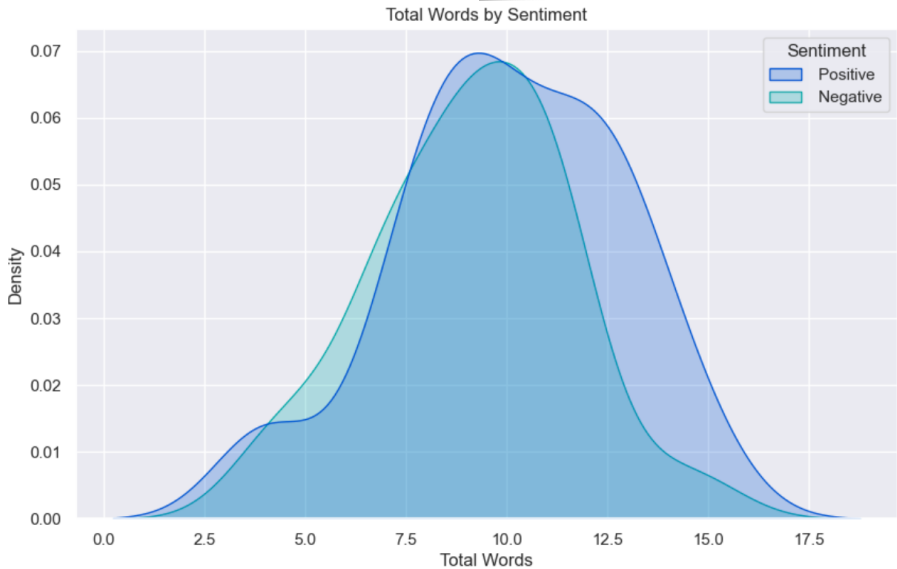
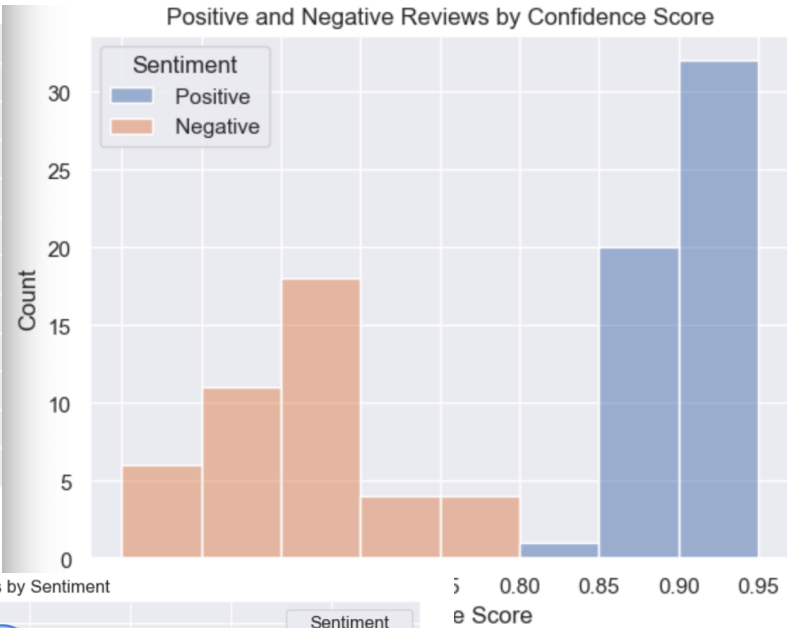
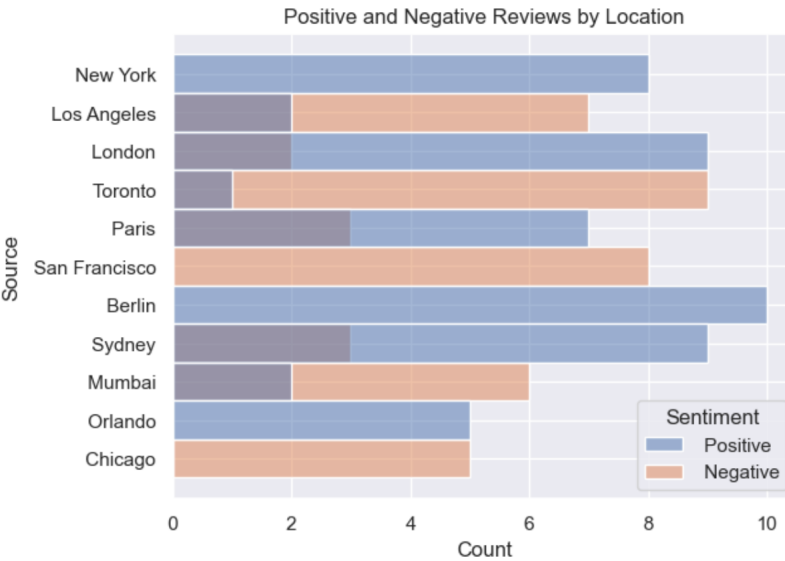
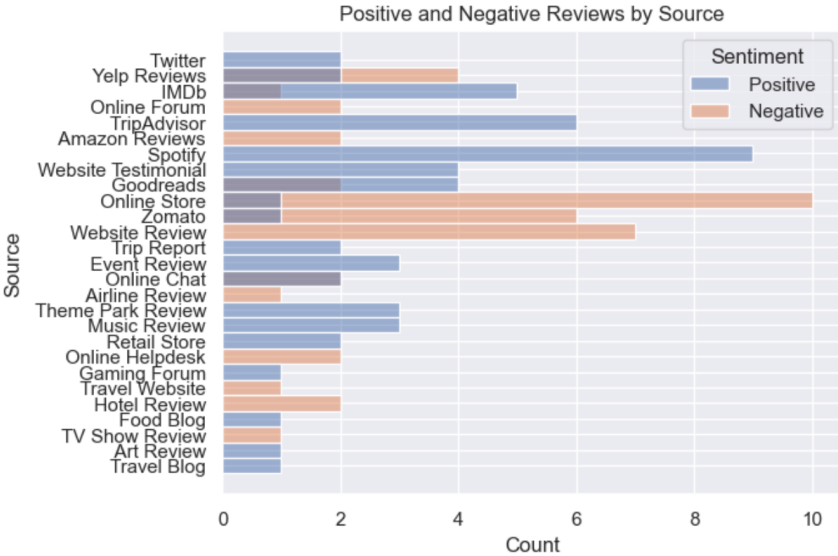
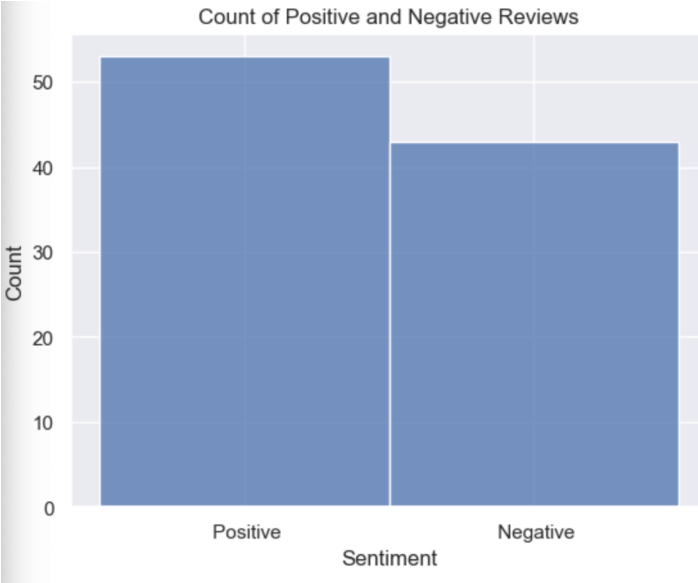
Customer feedback is an important tool for many companies and currently there are many different platforms and ways to receive that information. Analyzing text data can be very challenging due to the qualitative nature of the data. In order to analyze feedback, natural language processing (NLP) combined with machine learning can be used on such data to determine the sentiment. Utilizing NLP allows for efficient processing of text without human intervention and can determine if the sentiment of the text is positive, negative, or neutral.

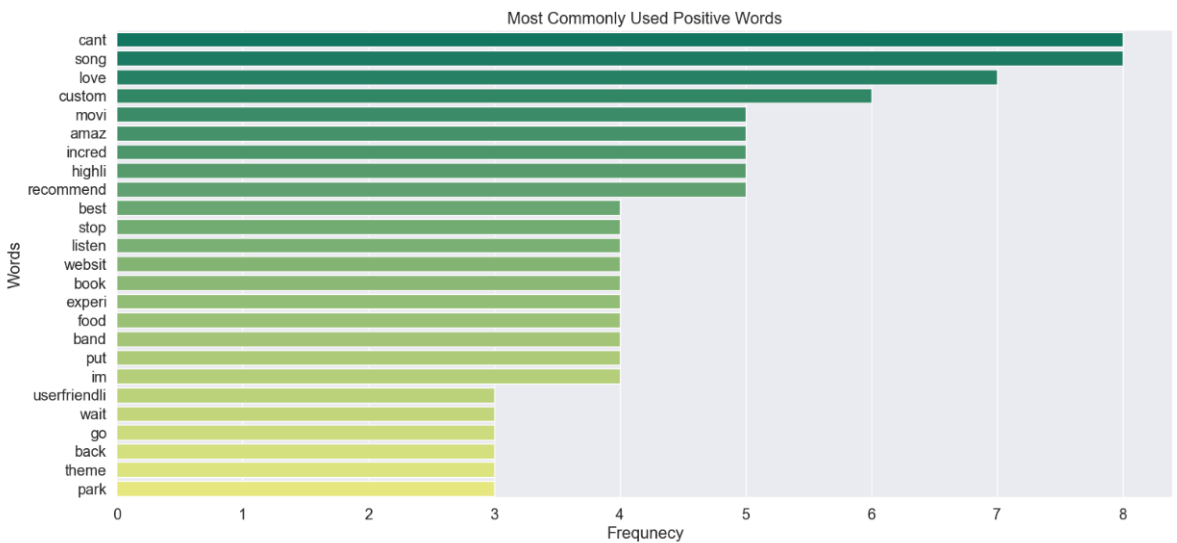
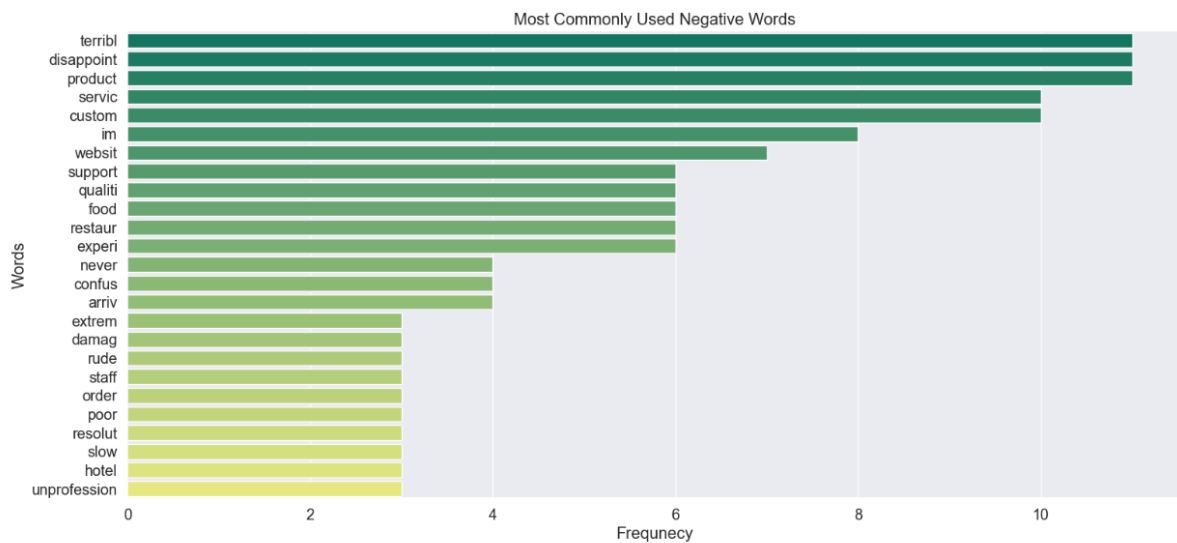
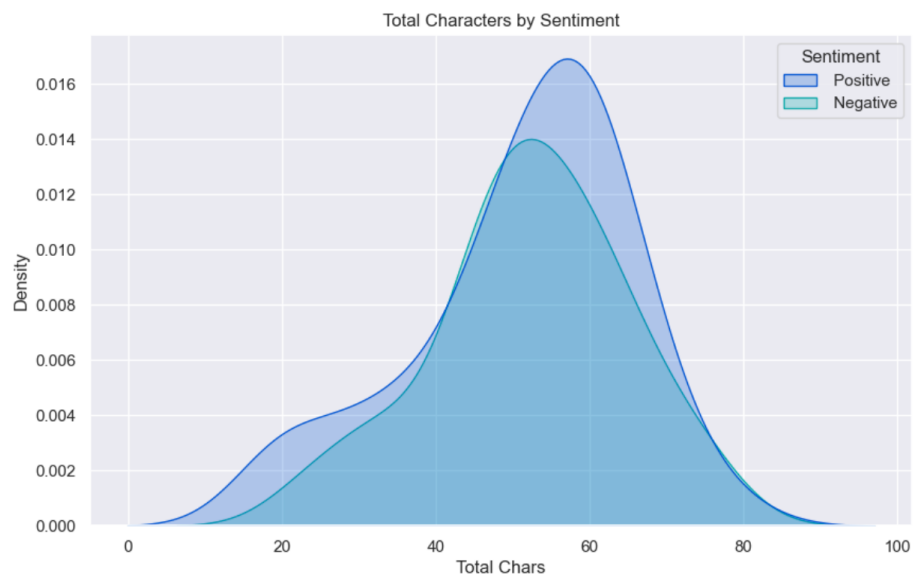
Data Wrangling

The data used for this project was obtained from kaggle. The dataset contains a collection of customer sentiments from various sources and platforms. A few examples of where the data was collected is Twitter, Facebook, and review platforms. The [source code](#) can be reviewed for a more detailed view of all of the exploratory data analysis and modeling. The data came with 98 rows and 9 columns. The data includes text, sentiment (positive or negative), source of the text, time/date of when the text was published, user ID, location, and confidence score. There were a couple rows with missing data that were dropped and the time column was split into additional features to separate out month, day, and hour for each row. New features including, total words, total characters, and total words after transformation were created for visual analysis.

Exploratory Data Analysis

The features of the dataset were visualized with histograms to see if there were any correlations related to the sentiment outcome. Sentiment was observed by, source, location, confidence score, month, day, and hour. Of the 96 rows, 53 of the responses were recorded as positive and 43 as negative. A density plot of total words and total characters were observed by sentiment, and most common positive and negative words were identified. Overall there was a large discrepancy between sentiment among all features. This may be attributed to the size of the dataset and the diversity of the source of the sentiment.

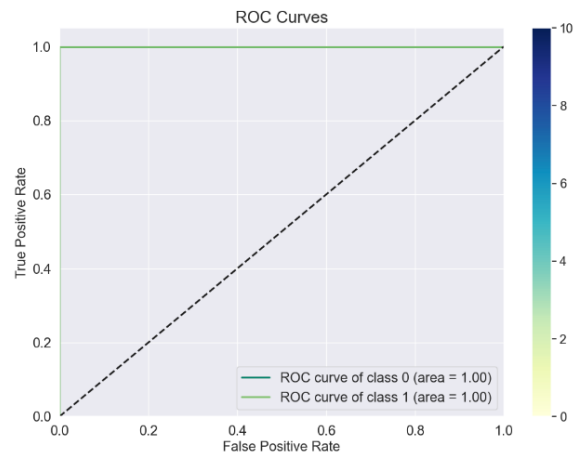
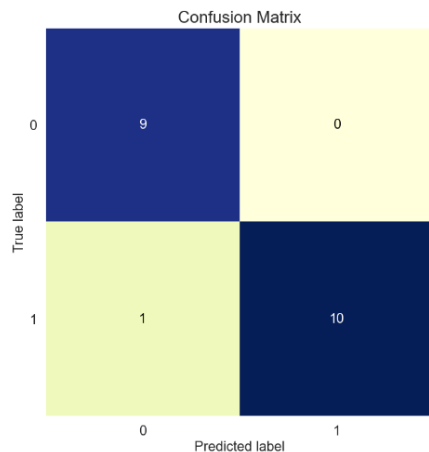




Modeling

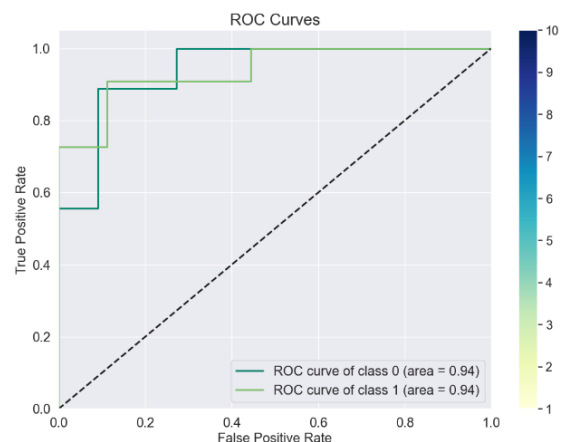
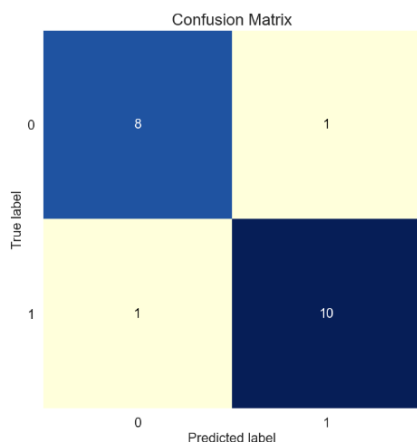
The data was prepared for modeling by converting the positive and negative sentiment responses to integers. Naive Bayes, Random Forest, and LSTM models were chosen for this dataset, and were conducted on both sentiment with and without text. For sentiment classification with text, Naive Bayes performed the best with an accuracy score of 0.95.

Accuracy of the model: 0.95
Precision Score of the model: 1.0
Recall Score of the model: 0.909



The data was then transformed to create dummy features for classification of sentiment without text. Naive Bayes performed the best with an accuracy score of 0.9.

Accuracy of the model: 0.9
Precision Score of the model: 0.909
Recall Score of the model: 0.909



Future work with a larger dataset could be explored to compare performance of the models.