

# Capstone 2 – Milestone Report 1

## Conducting NLP on Amazon Store Review Data

Zach Palamara



## Project Aim and Background

### Problem Statement

The goal of this project was to perform sentiment analysis on product reviews in the Amazon store. The target audience for this project are sellers on the Amazon store. They will be able to use this model to analyze consumer behavior and their response to products. This information will be useful in deciding on what new products to bring to the market and how to improve existing products. Additionally, this data can be used to help recommend products to consumers based on their behavior tendencies.

### Dataset and Source

I am using a data set from the following website: [Amazon Review Data \(2018\)](#). The data is in a .json format and can be easily read into a Pandas df for EDA and pre-processing. Each review is a single record in the dataset containing the following features:

- reviewerID - ID of the reviewer, e.g. A2SUAM1J3GNN3B
- asin - ID of the product, e.g. 0000013714
- reviewerName - name of the reviewer
- vote - helpful votes of the review

- style - a dictionary of the product metadata, e.g., "Format" is "Hardcover"
- reviewText - text of the review
- overall - rating of the product
- summary - summary of the review
- unixReviewTime - time of the review (unix time)
- reviewTime - time of the review (raw)
- image - images that users post after they have received the product

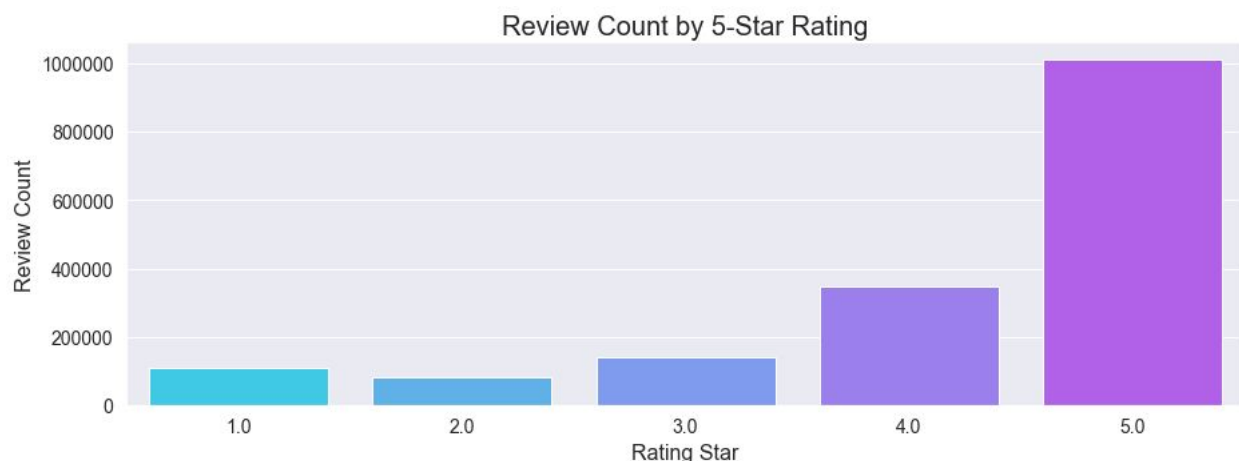
## Project Goals

I used Natural Language Processing (NLP) to solve this problem. More specifically, I used Fine-Grained Sentiment Analysis which is a common method for determining sentiment in 5-Star reviews. After performing EDA I tested out a variety of methods for feature engineering. Some of the methods I used include text vectorization using a bag-of-words or bag-of-ngrams. After the features were properly cleaned and engineered, I used three different classifiers to construct my machine learning models for text and numerical features.

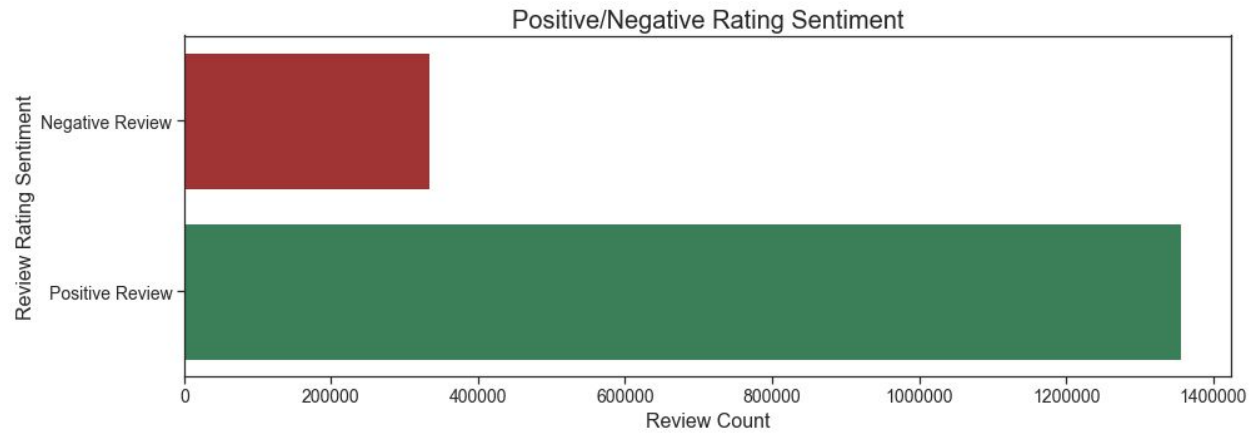
## Data Cleaning and EDA

### Initial Findings

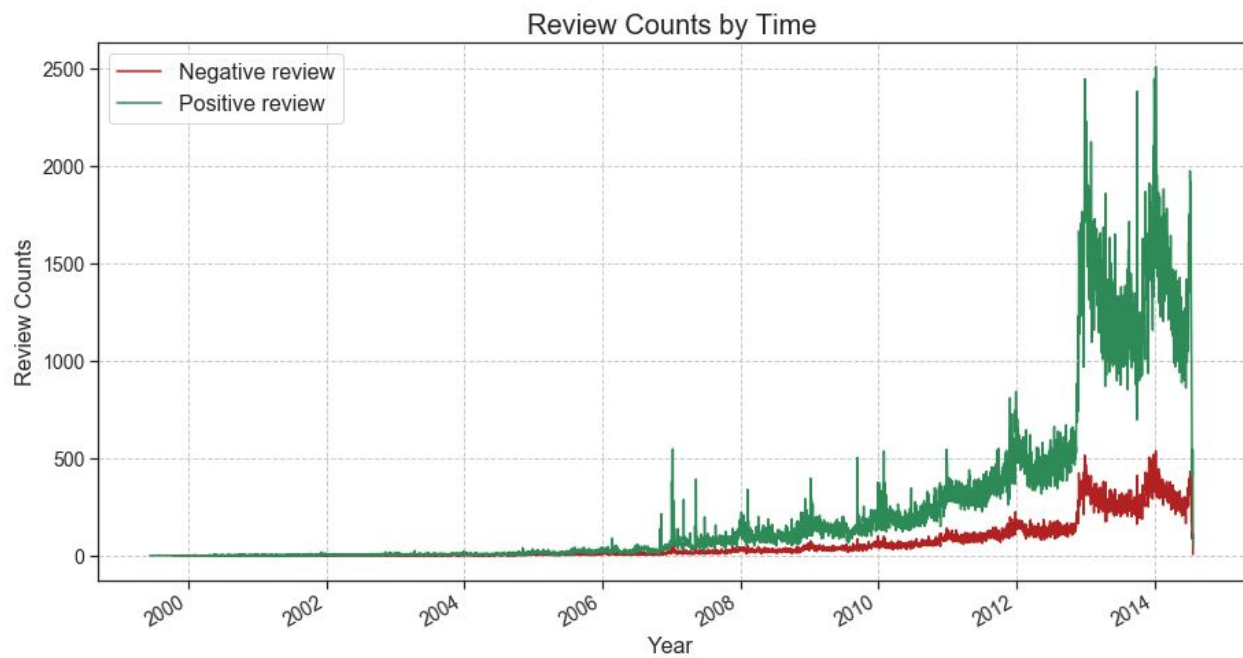
The first step I took in examining the data was to look at a simple count plot of reviews broken out by their respective 5-star rating.



This figure shows that most of the reviews are positive and overwhelmingly 5-stars. Additionally, I created a figure that depicts the count of negative and positive reviews. For this, I created a new binary target feature, which represents 1,2 and 3 star reviews as 0 or "Negative" and 4 and 5 star reviews as 1 or "Positive".



When viewing positive and negative review counts in time-series, there's a clear jump in the volume of reviews around the start of 2013. Also, it appears that positive sentiment seems to diverge from negative sentiment with time.



# Text Feature Engineering

## Removing Stop Words

One of the most important steps in any NLP problem is to remove “stop words” within a corpus. Stop words are words such as “the”, “a” and “are”, which are common words that do not really provide any significant meaning to our features. It is important to remove these words, because they take up unnecessary space in memory and processing power. I utilized the Natural Language Toolkit (NLTK) framework to remove stop words and punctuation in this dataset.

## Tokenization

Once the stop words were removed, each review for each data record needed to be tokenized. Tokenization is a way of separating a large text document into tokens, which in this case would be singular words. I utilized the `word_tokenize()` method provided by the NLTK framework.

## Lemmatization

The next step in engineering the text features involved lemmatization. In its simplest form, lemmatization is a process in which different forms of the same word are mapped to a singular form of that word. This is also very important in trying to reduce the size and memory space that your data set takes up. For example, the words “running”, “ran” and “run” all get mapped to a singular form of “run”. This was executed using the `PorterStemmer` class from the NLTK framework.

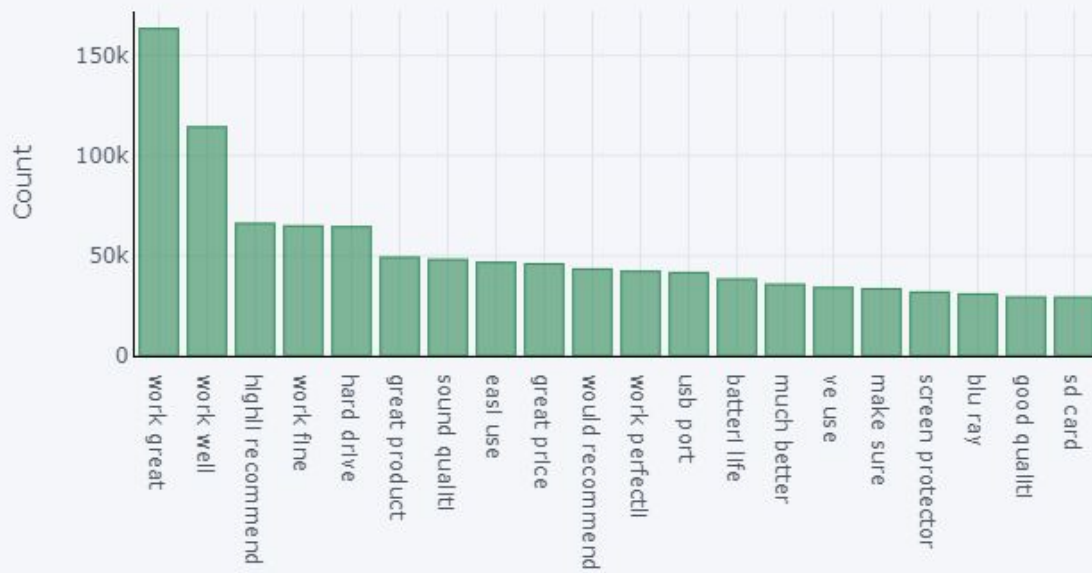
## Stemming

The final step I took in engineering the review data involved stemming. Stemming simply removes prefixes and suffixes from tokenized words. Similar to Lemmatization, this is also done in an effort to reduce the size and complexity of the data. This was also executed using the `WordNetLemmatizer` class from the NLTK framework.

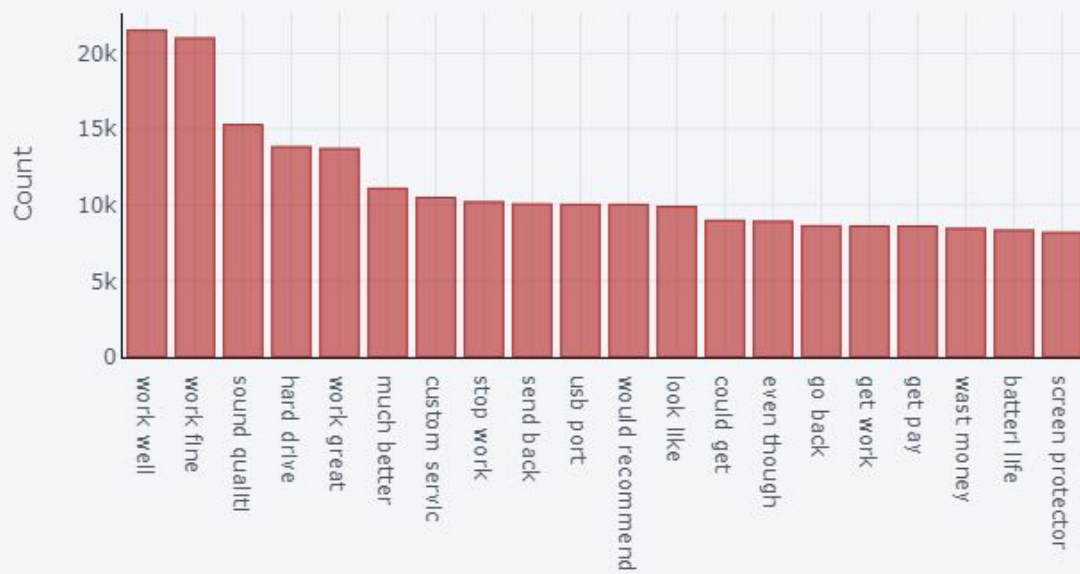
## Extracting N-Gram Features

In the simplest of terms, N-grams are a statistical representation of the likelihood that a series of words appear in a text document. The letter “N” in N-grams represents the number of words in the series. For example, I created lists of the top 20 Bigrams and Trigrams of both positive and negative reviews. The results can be seen in the figures below.

Top 20 bigrams in postive reviews after processing

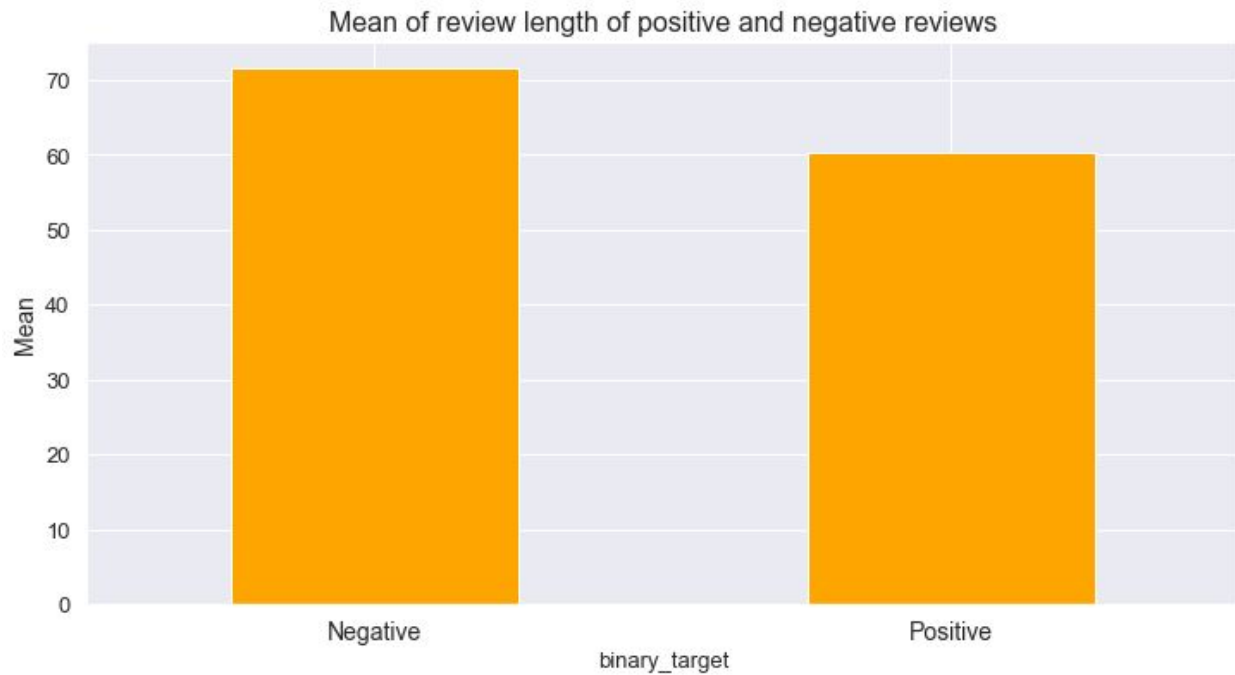


Top 20 bigrams in negative reviews after processing

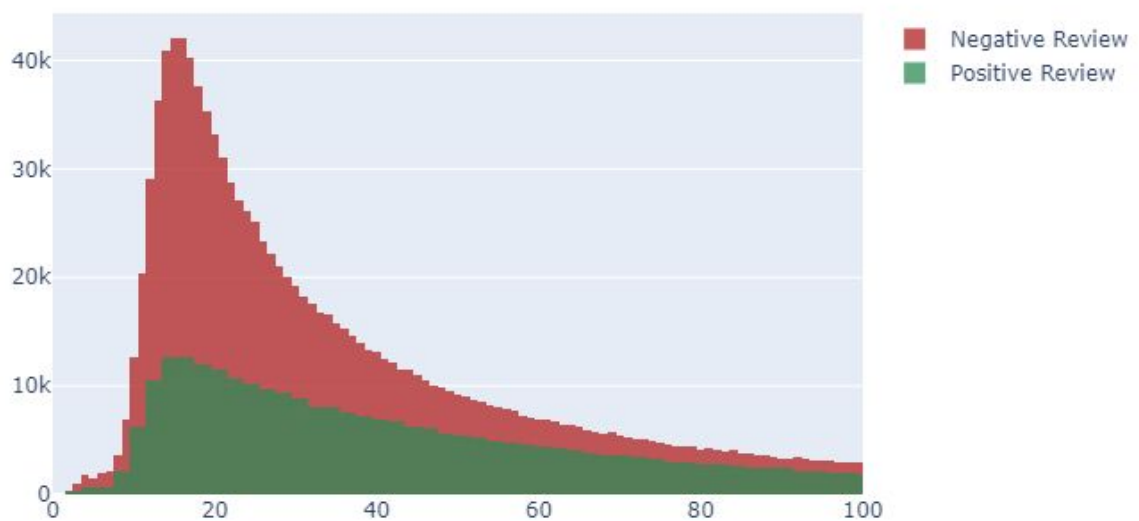


## Review Length

On average, negative reviews seem to be longer than positive reviews. Intuitively, it makes sense that people tend to be more vocal expressive about something that displeases them. Additionally, both positive and negative reviews seem to have right tail distributions.



## Distributions of Review Lengths



## Word Cloud

One of the most popular ways to visualize common words in an NLP project is through a word cloud. Pictured below is a word cloud containing the 20 most common bigrams found in positive reviews.

