

---

# Yelp Dataset Analysis Report

---

**Shihao Lin**

Pitzer College

1050 N. Mills. Ave, CA 91711

shilin@pitzer.students.edu

**Ziyuan Shang**

Scripps College

Claremont, CA, 91711

zshang9431@scrippscollege.edu

## Abstract

Rating prediction is one of the most important application of machine learning. In this project, we explore the Yelp dataset and are looking for appropriate machine learning algorithms to predict the star users will assign to a business through the content of review they wrote. Currently, we use the basic topic modeling processed the data and use classification algorithms like Naive Bayes(NB) algorithm, MultinomialNB, BernoulliNB, Logistic Regression and Linear Support Vector Machine to do binary classification (positive and negative) on the dataset.

## 1 Introduction

Yelp is an American multinational corporation headquartered in San Francisco, California. It was founded in 2004, and dedicated to helping people find a great local business like restaurants and hair salons. The company allows users to give star ratings and write reviews for the business. More than 155 millions of reviews have been written by Yelp users by the end of Q1 2018. Among these reviews, myriads of hidden information can be extracted, and this information can be beneficial to business owners as well as customers. By studying these reviews, owners can be more aware of things that they need to improve, and customers will be able to find better matches for their own tastes. However, how to extract the most important information has remained a big challenge for a long time. The project specifically focuses on building a model that can extract useful features from millions of reviews. Processing these reviews using sentiment analysis, finding the most featured words and predicting the proper ratings solely based on review text.

## 2 Data description

Yelp's datasets include six files, which are business.json, review.json, user.json, checkin.json, tip.json, and photos. In our midterm modeling, we used the business.json and review.json. Later, we will further explore the review.json for our final modeling. In the following subsections, we are going to illustrate the datasets and its format.

In this dataset, the total organizations are 146702, the total users are 1326101, the total reviews are 1493480, the total negative reviews (from rating scale 0 to 4) are 504771, and total positive reviews (from rating scale 4 to 5) are 988709.

The above figure 1 illustrates the review star distribution. Based on figure 1, we can notice that the total star reviews are mostly made up by positive reviews.

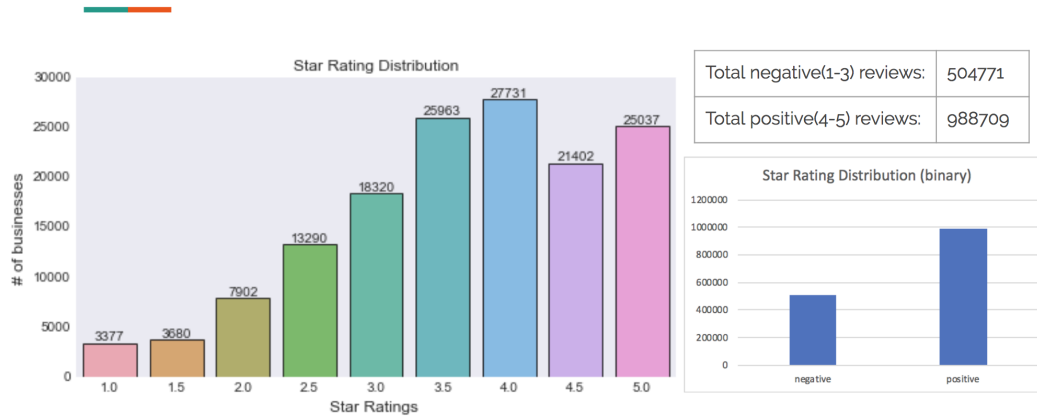


Figure 1: Star Rating Distribution

## 2.1 Files Description

The file Business.json contains business data including location data, attributes, and categories. More specific, this file includes information of city, state, latitude, and longitude of the store, stars, review counts, business parking, ambience, price level, Good for meal and etc. In the later section, we would set up the model to illustrate the correlations between all of those attributes and business rating. In other words, we are going to build up a model to predict the business rating applying some of the features to different types of existing algorithms to train the computer.

The file Review.json contains full review text data including the user\_id that wrote the review and the business\_id the review is written for. We will not use this for the midterm report. But we will use this file for the final modeling to predict a solo review star, which is based on its review text.

## 3 Data Analysis

Before we actually start analyzing the data, we transform the format from json to csv, which is convenient for us to analyze the data. In the following subsections, we generate various graphs to help visualize the relationship between different feature.

### 3.1 Location v.s Review counts

In the first part, we link the review counts to states. We find out that Arizona has the most reviews, which is a bit surprising. The topmost reviews states are Arizona, Nevada, Oregon, North Carolina, and Ohio. Surprisingly, we find that California is not on this list. We also link the review counts to the city. We find that Las Vegas has the most reviews, which is 22113 and is also the main component of the total reviews of Nevada. The following figure 2 illustrates the topmost reviews city in the U.S. By using the latitudes and longitudes of business organizations and its review counts, we can plot out a graph about that information, just like the following figure 3 about Las Vegas and Phoenix. Each review will represent a small black dot. The more reviews in the same location, the dot will be darker. Therefore, we can use this graph to predict where is the most popular place in a city.

### 3.2 Review Count v.s. rating

In figure 4 "Review count v.s. Rating", we label the number of review in one single business organization into horizontal axis, and the probability to get a positive review into vertical axis. We found that if a business has fewer reviews, the probability of it to be rated as positive is around 50%. The business organization with a medium amount of review tends to have high variance on its probability to be rated as positive. The organization with large reviews often tend to have the probability to be rated as positive.

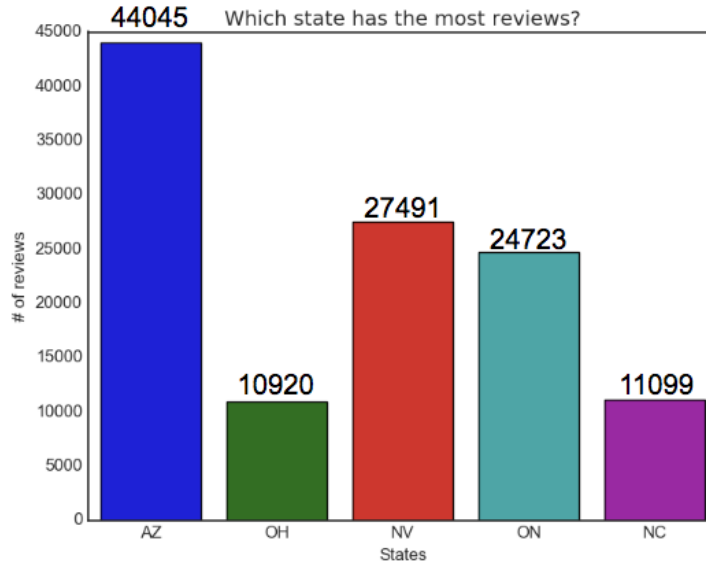


Figure 2: What state has the most reviews?

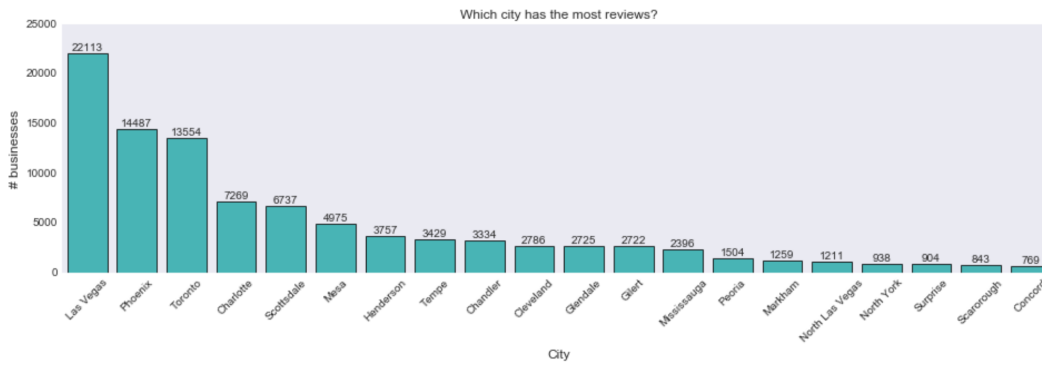


Figure 3: Which city has the most reviews

### 3.3 Reservation v.s. Rating

We want to discover the relationship between whether an organization accepts reservation would influence it to be rated as positive. Therefore, we analyzed the data and put it into a histogram as presented in figure 5. On the horizontal axis, 0 represents that reservation service is unavailable, and 1 represents that reservation service is available. The vertical axis shows the probability for each type of organization to be rated as positive. As a result, we found the probability of being rate as positive for both types are very similar and centering around 50 percentage. Thus, we conclude that there is no such correlation between the service of reservation and business rating.

### 3.4 Price v.s. Rating

The price level of a restaurant is always a big factor to influence people's decision of selecting a restaurant. Sometime, people would think an expensive restaurant tend to taste better than a cheap restaurant. In this paper, we are going to present our finding of the correlation between price level and business rating. By plugging the price level data and business rating data into "Seaborn" analyzer, we get a histogram in figure 6. Each number on the horizontal axis represents different price levels. "1.0" means very cheap, and "4.0" means very expensive. The vertical axis shows the probability for different price levels of an organization to be rated as positive. As a result, we found the probability of all four different price level are very similar and close to 50 percentage, and the probability of price

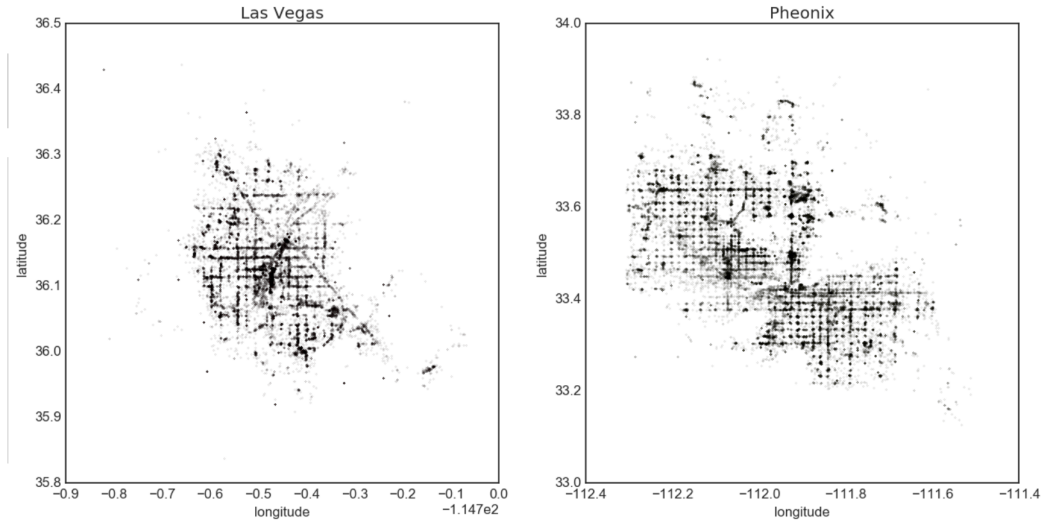


Figure 4: Business popularity on city map

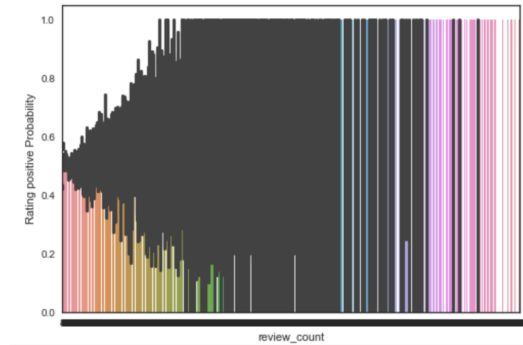


Figure 5: Review Count v.s. Rating

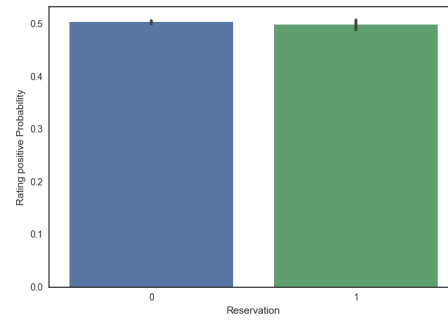


Figure 6: Reservation v.s Rating

level 4 is slightly lower than the price level. Thus, the result indicates the price level may not have much influence on Business Rating. A very expensive restaurant may not always taste very good.

### 3.5 Ambience v.s. Rating

Different type ambience of a restaurant tends to have different chance to be rated as a positive good restaurant. We want to discover the correlation between ambience and rating. In the yelp dataset, it provides 10 categories for ambience, "None", "casual", "divey", "intimate", "romantic", "trendy", "classy", "touristy", "hipster", and "upscale". Then we plugged all of those data with its business rating into "Seaborn" analyzer. We get a bar diagram in figure 7. Different categories are labeled in the horizontal axis. The vertical axis shows the probability for different categories of an organization to be rated as positive. As we observed, "upscale" ambience restaurant tend to have higher chance to be rated as positive. However, in general, since the positive rating probability for all categories is all laid between 40 percentage and 60 percentage, the ambience does not seem to have much correlation with the business rating.

### 3.6 City v.s. Rating

We are curious about whether the overall quality of food in some city is higher than other cities. Therefore, we input all organizations' city location and corresponding business rating into "Seaborn" analyzer and get a histogram as shown in figure 8. The cities' names are labeled in a horizontal axis, and the corresponding positive rating probability is presented by a vertical axis. The different color

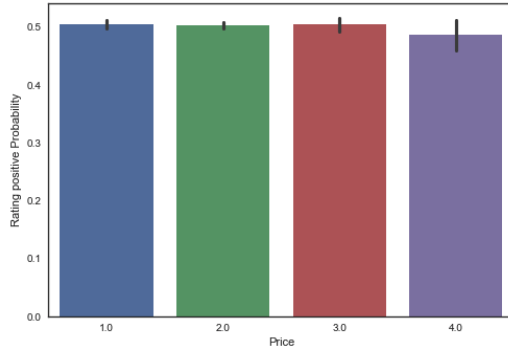


Figure 7: Price v.s. Rating

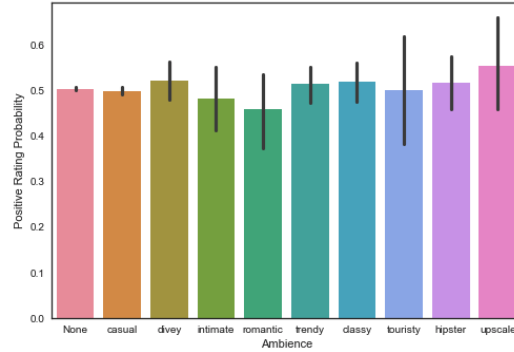


Figure 8: Ambience v.s. Rating

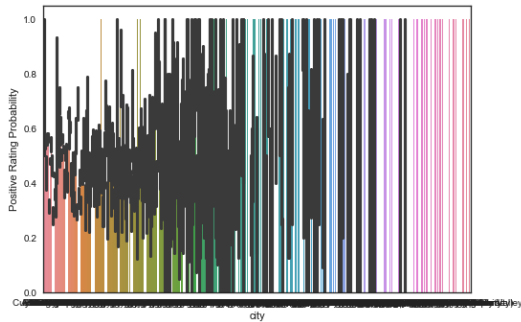


Figure 9: City v.s. Rating

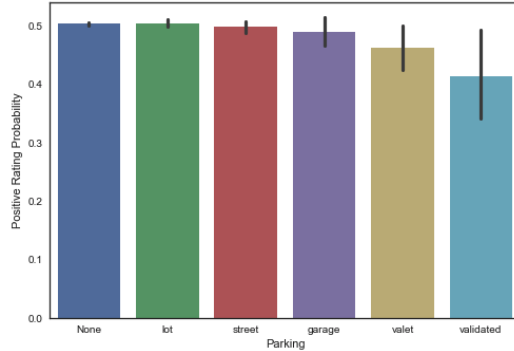


Figure 10: Parking v.s. Rating

bar represents the different city, and the black bar is the variance of the probability. By observation, we found some city definitely have higher overall quality restaurant than others.

### 3.7 Parking v.s. Rating

The transportation in most of the place in the U.S. is driving. In average, every household has around 2 cars in the U.S. Thus, we think whether a restaurant would provide parking service is essential to its rating. Then we plug the data about its restaurant parking situation and corresponding business rating into "Seaborn" analyzer and get a histogram. As in figure 9, we observe that the positive rating probability for all categories is under 50 percent. Thus, the service of parking does not seem to have much influence on the business rating.

### 3.8 Good For Meal v.s. Rating

At first, we consider that different meal type may have an influence on the business rating. For example, restaurant, which focuses on making breakfast, in general, may have higher business rating than those focuses on making brunch. In yelp data, it provides 7 categories for "Good For Meal", which are "None", "late night", "dessert", "breakfast", "dinner", "lunch", and "brunch". While we plugged all of those data and corresponding business rating, we get a bar diagram, as shown in figure 10. The meal categories are labeled in a horizontal axis, and the corresponding positive rating probability is presented by a vertical axis. By observation, we found that the positive rating probability for all categories are very similar and around 50 percentage. Thus, distinguishing meal specialty of a restaurant may not help to distinguish its business rating.

## 4 Business Analysis

Before actually applying the data into the algorithm, we decide to filter out some features of the dataset. In the end, we remain the three features, review\_count Reservation, and Parking, to predict

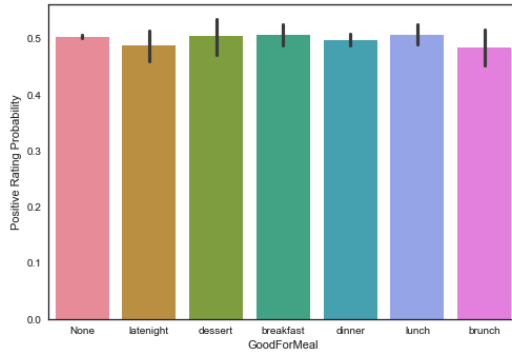


Figure 11: Good For Meal v.s Rating

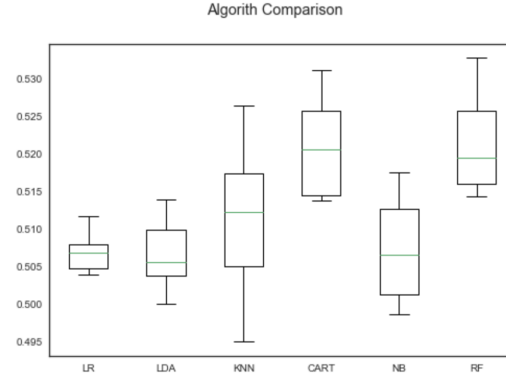


Figure 12: Algorithm.png

the business rating. We set 75 percentage dataset as the training dataset, and remain 25 percentage dataset as the testing dataset. Then we purposely select six different methods to predict the outcome, which is "Logistic Regression", "Linear Discriminant Analysis", "K-Neighbors Classifier", "Decision Tree Classifier", "Gaussian Naive Bayes", and "Random Forest Classifier" from *sklearn* package. All of six methods could be used to analyze a binary outcome.

In the previous section, although the features in the dataset are independence to each other, we notice that most of the features in this dataset do not have much correlation with the business rating. Therefore, we assume that the accuracy for those algorithms to have a right prediction is low. The overall result is present in the following table, and the visualized version is in figure 12.

By observation, the overall accuracy is very low and around 50 percentage. We observed that "Logistic regression" method has the lowest variance, "K-nearest Neighbors classifier" has the highest variance, and "Random Forest classifier" brings up the highest accuracy, which is yielding roughly 52.11%. Although "Decision tree classifier" is the underlying concept for "Random Forest classifier", the accuracy performance does not have much different. In conclusion, the correlation between the feature and its outcome is essential for generating a high accuracy prediction.

Algorithm	Accuracy Percentage
Logistic_Regression_classifier	0.506369 (0.003203)
Linear_Discriminant_Analysis	0.506685 (0.004129)
K-Neighbors_classifier	0.511014 (0.009426)
Decision_Tree_classifier	0.520832 (0.006517)
Gaussian_Naive_Bayes	0.506978 (0.006656)
Random_Forest_classifier	0.521114 (0.006188)

## 5 Review Analysis

In our review analysis, we mainly explored the relationship between review content (what words does user use) and rating of the corresponding business. We first attempted our reviews analysis on a smaller sample. Form our reviews dataset, we generated a random sample of 3000 reviews. We will scale the code to perform analysis for the entire dataset in the future.

### 5.1 Data preparation

The data used in this analysis comes from the preprocessed review.csv mentioned in Section 3. The only useful columns are *text* - the actual review content given by the user in a string format, and *review\_stars* - the star rating associated with the review. A third column named *Rating* is created based on *review\_stars* - reviews with 1, 2 or 3 stars are considered negative sentiment; reviews with 4 or 5 stars are considered positive.

### 5.1.1 Text pre-processing

First, we performed several necessary pre-processing on review text, including tokenizing, removing stop words and stemming.

Tokenizing refers to the act of separating a sentence into individual word tokens. In our analysis, we used the default *word\_tokenizer* in the *nlk* package under Python. As an example, if we pass the following sentence: "The python programmer named pythoner is pythoning a game pythonly." into the *word\_tokenizer*, the returned result will be

```
'The', 'python', 'programmer', 'named', 'pythoner', 'is', 'pythoning', 'a', 'game', 'pythonly'
```

After tokenization, we need to remove the common stop words in English. Stop words are words that appear extremely common but do not convey much meaning or sentiment of bodies of text. We used the bag of stop words gathered by the *nlk* package and eliminated the words that falls into this bag from each piece of text. Using the above example, after removing the stop words, the returned result is

```
'The', 'python', 'programmer', 'named', 'pythoner', 'pythoning', 'game', 'pythonly'
```

Observed that the words 'is' and 'a' are removed from the list, whereas the words that conveying meanings are kept.

Then we stemmed the individual words in each piece of text using Porter Stemmer. Stemming refers to "the process of reducing inflected (or sometimes derived) words to their word stem, base or root from generally a written word form". In particular, Porter Stemmer is the de facto algorithm developed in 1980, and it was widely used to stem English words. In our analysis, we used the *nlk* implementation of Porter Stemmer to bring words from text to its stem form. Here, we assumed that words with the same stem convey roughly the same meaning, and therefore are weighted equally in sentiment analysis.

After passing the above list to a Porter Stemmer, the list returned is in the following form:

```
'the', 'python', 'programm', 'name', 'python', 'python', 'game', 'pythonli'
```

Observed that all python-related words, except 'pythonly', are being turned into its stemmed for 'python'. The single exception of the adverb might be an example of under stemming.

### 5.1.2 Feature selection

After text processing, another important area that we explored is the feature selection for sentiment prediction. Feature selection for sentiment prediction is a popular topic with many on-going researched. Popular methods for feature selection include information gain and HMM-LDA. For our purposes, we use a rather simplistic feature: individual words.

## 5.2 Model training and evaluation

### 5.2.1 Training

For sentiment predictions, we have tested a few machine learning algorithms that can make binary (positive and negative) classifications using different sets of features. We have explored from the *nlk* package, trained on sparse matrix build by only individual words as features. We also used several more advanced algorithms. Including MNB, BernoulliNB, Logistic Regression, Linear SVC from the *sklearn* package.

### 5.2.2 Evaluation

Below is the graph comparing the accuracy of different algorithms we've tried. We can observe that the Naive Bayes Classifier performed poorly, yielding roughly 50.87%. Then, we used several more advanced algorithms, including MNB, BernoulliNB, Logistic Regression, Linear SVC, the test accuracy are as follows: After text processing, another important area that we explored is the feature selection for sentiment prediction. Feature selection for sentiment prediction is a popular topic with many on-going researched. Popular methods for feature selection include information gain and HMM-LDA. For our purposes, we use a rather simplistic feature: individual words.

We see that the selected algorithms are producing results with around 71% of accuracy. We will use tf-idf and other methods to improve the score in the future.

Classifier	Accuracy Percentage
MNB_classifier	71.8%
BernoulliNB_classifier	64.93%
LogisticRegression_classifier	73.07%
LinearSVC_classifier	71.47%

## 190 **References**

191 [1] <https://www.yelp.com/dataset-challenge>.