Yelp Dataset Analysis

STATS 140XP Final Project Annaka Schone Damien Ha Sylvia Deng Joyce Xu

Caitlin Ree Meiyi Ye Wenhong Sun

Introduction

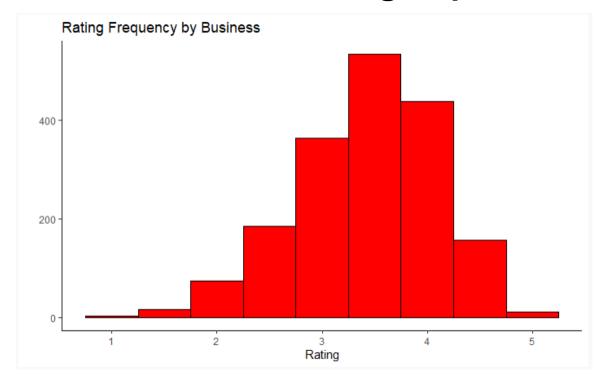
Yelp is a platform that crowd-sources reviews about businesses, with the majority being restaurants. The Yelp dataset being examined in this study was assembled for the Yelp Dataset Challenge in 2018. It contains Yelp business data, reviews and user data.

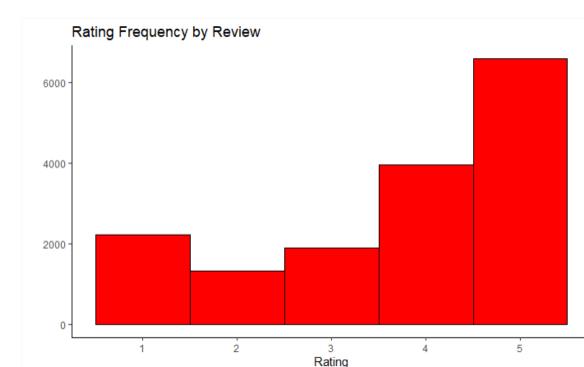
For the purpose of this study, the dataset has been reduced to business data and reviews of restaurants in the Chinese Food category, for the purpose of mitigating the effects of "Restaurant Category" in our analysis while keeping the dataset as large as possible. This study attempts to answer the following questions:

- 1. Do a restaurant's Yelp reviews accurately reflect its star rating?
- 2. What are the factors that contributed to the overall star rating of Chinese restaurants? Which factor has the greatest effect?

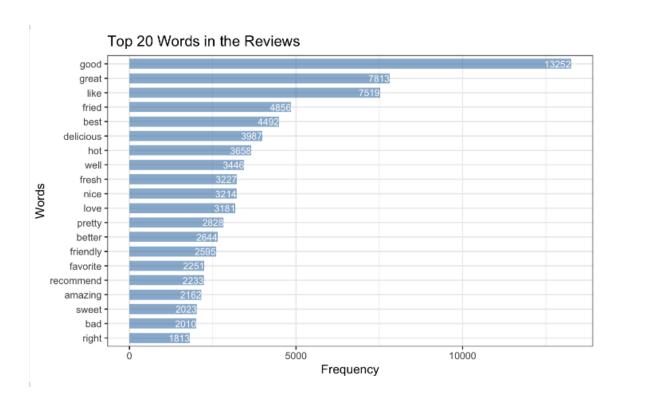
Exploratory Data Analysis

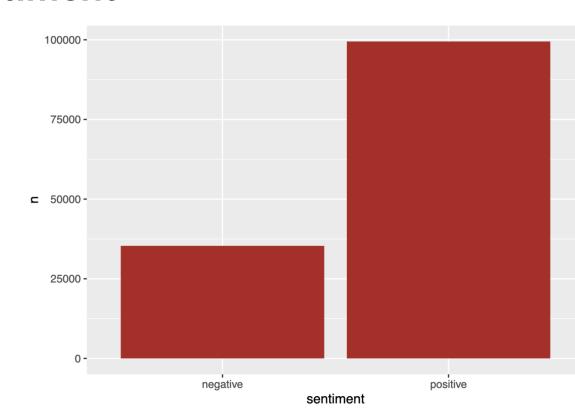
Distribution of Ratings by Business and Review





Distribution of Review Words and Sentiment





Word Cloud: Words that most commonly appear in the Yelp Reviews



Is Yelp ratings a

good way to

determine

restaurant

quality?

It probably 15!

Do a restaurant's Yelp reviews accurately reflect its star rating?

We first used sentiment analysis to Calculate the sentiment of each review based off lexicon and then created 2 indicators based off the sentiment scores: overall sentiment per business (positive/negative) and sentiment percentage. Finally, we built a linear regression model to compare the businesses' original average stars and sentiment reflected through reviews.

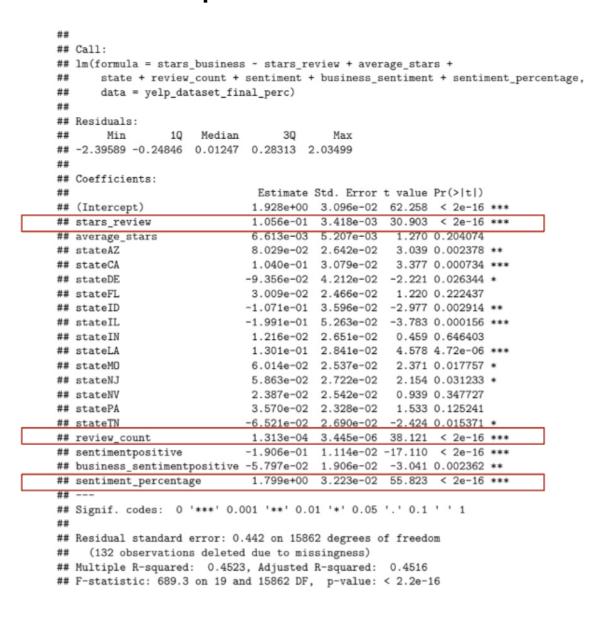
```
## lm(formula = stars_business ~ business_sentiment + sentiment_percentage
```

The results indicate that:

- 1. A business' percentage of positive review sentiments is positively correlated with its average stars.
- 2. Sentiment percentage is a more accurate reflection of star rating than overall business sentiment.

Which aspects of a restaurant have the greatest effect on its Yelp star rating?

We applied both linear regression and logistic regression to test the effect of various potential variables on Yelp star ratings.



| ## | 0-33 | | | | | |
|-----|--------------------|--------------|------------|----------|------------|--------------------------------|
| | Call: | | | | | |
| | glm(formula = bina | - | | _ | | |
| ## | review_count - | e sentiment, | ramily = ' | binomial | .", data = | <pre>= yelp_dataset_fina</pre> |
| | Coefficients: | | | | | |
| ## | coefficients. | Fetimate 9 | Std. Error | z value | Pr(> z) | |
| | (Intercept) | -0.6240284 | | | , | ** |
| | stars_review | 0.2113596 | | | | |
| | average_stars | 0.0932458 | | | | |
| | - | -0.6344344 | 0.2231646 | -2.843 | 0.00447 | ** |
| ## | stateCA | -1.4146165 | 0.2681666 | -5.275 | 1.33e-07 | *** |
| ## | stateDE | -0.1395539 | 0.3559806 | -0.392 | 0.69504 | |
| ## | stateFL | -0.3902056 | 0.2007986 | -1.943 | 0.05198 | |
| ## | stateID | -0.5804108 | 0.2862261 | -2.028 | 0.04258 | * |
| ## | stateIL | 1.1888855 | 0.6319747 | 1.881 | 0.05994 | |
| ## | stateIN | -0.2298430 | 0.2299165 | -1.000 | 0.31747 | |
| | | -0.1959170 | 0.2645303 | | | |
| | stateMO | -1.0312063 | 0.2003593 | | | *** |
| | | -0.1744060 | 0.2235313 | | | |
| | | | 0.2345429 | | | |
| | | -0.7446485 | 0.1855902 | | | |
| | | -1.3417298 | 0.2065399 | | | |
| | review_count | 0.0149429 | | | | |
| | sentimentpositive | 1.0210404 | 0.0877576 | 10.472 | < 2e-16 | *** |
| | Signif. codes: 0 | '***' 0 001 | '**' 0 01 | '*' 0 05 | 0 1 | |
| *** | Dignii. Codes. O | 0.001 | 0.01 | . 0.00 | . 0.1 | |

Linear Regression Model

- 1. Most significant variables: stars on particular reviews, review count, sentiment percentage
- 2. Again, sentiment percentage is positively correlated with average stars, as does stars for each particular review

Logistic Regression Model

- 1. Similar to linear regression, stars on particular reviews and review count were significant variables
- 2. Location does not seem to have as much of an impact

Future Research

Future research likely involves re-running this analysis on other restaurant category subsets to see how if the results are consistent.



