

# YELP REVIEW ANALYSIS

Ziren Xia & Huihan Nie & Jie Gu

Jiongnan Liu & Yan Gao

## CONTENTS

1	Introduction	2
2	Data Exploration	3
3	Technology	3
4	Rating Distribution	3
4.1	Rating Star and Review Amounts Distribution . . . . .	3
4.2	Common Characteristics . . . . .	5
5	Geographic Distribution	5
5.1	World Map US Map . . . . .	5
5.2	Atlanta Distribution Map . . . . .	6
6	Regional Analysis (Atlanta)	7
6.1	Food trend discovery . . . . .	8
6.2	Review analysis . . . . .	8
6.3	Business value of the regional analysis . . . . .	9
7	Sentimental model	9
7.1	Background . . . . .	9
7.2	The Dataset . . . . .	10
7.3	features . . . . .	10
7.4	Implementation . . . . .	10
7.5	Improvement . . . . .	11
7.6	Business Use Case . . . . .	12
8	Text Analysis Trigram Counts	14
8.1	Analysis and Results . . . . .	14
8.2	Business Insight . . . . .	15
9	Conclusion	15
9.1	Recommendations for business . . . . .	15
9.2	Recommendations for Yelp . . . . .	15

## LIST OF FIGURES

Figure 1	Top Rated Restaurant Based on Sorting by Rating and Review Count . . . . .	3
Figure 2	High Customer Flow Based on Restaurant Rating Count Distribution . . . . .	4
Figure 3	Yelp Restaurant Rating Distribution by State by Star . . . . .	4
Figure 4	Yelp Five Star Restaurant Common Characteristics . . . . .	5
Figure 5	Restaurant Distribution on World Map . . . . .	5
Figure 6	Restaurant Distribution on US Map . . . . .	6
Figure 7	Top 5 cities with restaurant count . . . . .	6

Figure 8	Restaurant count choropleth Atlanta map plotted with restaurants that have stars over 4.5 . . . . .	6
Figure 9	Choropleth Atlanta map with respect to population density . . . . .	7
Figure 10	Choropleth Atlanta map with respect to average income . . . . .	7
Figure 11	Atlanta Businesses Distribution . . . . .	7
Figure 12	Most popular highly rated restaurants in Atlanta . . . . .	8
Figure 13	Atlanta cuisine food type rank . . . . .	8
Figure 14	Common words used in reviewer ratings . . . . .	9
Figure 15	Dataset Example . . . . .	10
Figure 16	Left: Review Count by Each Star; Right: Review Stats by Each Star . . . . .	10
Figure 17	Model Accuracy and Model Mean Absolute Error by Number of Features . . . . .	11
Figure 18	True Positive Rate by Each Label . . . . .	12
Figure 19	Regrouping 1 and 2-star, 4 and 5-star; Left: Model Accuracy, Right: Model MAE; (blue line: regrouping, red line: without regrouping . . . . .	12
Figure 20	Regrouping 2, 3, and 4-star; Left: Model Accuracy, Right: Model MAE; (blue line: regrouping, red line: without regrouping . . . . .	13
Figure 21	Blue Line: Regrouping 1 and 2-star, 4 and 5-star; Red Line: Regrouping 2, 3, and 4-star . . . . .	13
Figure 22	Trigram Count Result from One Star Rated Comments . . . . .	14
Figure 23	Text Model for Restaurants of Different Price Range . . . . .	15

## 1 INTRODUCTION

Yelp is a well known mobile app which publishes crowd-sourced reviews about businesses. The company was founded in 2004 by former Paypal employees, and it grew in usage and raised several rounds of funding in the following years. As of today, it has over 199 million business and restaurant reviews worldwide. The customer reviews on Yelp are a great source of feedback and offer tremendous insights into what customers like and dislike about the product or service. It also has impacts on the purchasing decisions of the future customers as positive and negative reviews could affect the reputation of a business. Here, we are utilizing the power of Big Data techniques to do an in-depth analysis of restaurant business reviews on Yelp. By Rating Distribution, we were able to analyze how to help users choose a perfect restaurant, how to help businesses find a good location to start a new restaurant and how to boost a restaurant's reputation. By Geographic Distribution, we could visually see the top rated restaurants' distribution in the USA and Canada. We chose Atlanta which offered the best food and service as our target city for evaluation and found the most highly rated popular cuisine and food. We investigated the top rated restaurants in Atlanta to learn about the taste of our customers. Additionally, we predicted the ratings based on text-based reviews via logistic regression, and we did some improvements to increase the accuracy of detecting positive, negative sentiment and also neutral reviews. From text analysis, we found the most substantial factors that impact the ratings. By exploring patterns in the positive and negative reviews, we are able to analyze the major merits and complaints of businesses. We gave the business insights and provided Yelp some recommendations for optimizing user experience.

## 2 DATA EXPLORATION

This dataset is a subset of Yelp's businesses, reviews and user data which includes information about businesses across 16 metropolitan areas in two countries. It was originally put together for the Yelp Dataset Challenge which is a chance for students to conduct research or analysis on Yelp's data and share their discoveries. This dataset contains five JSON files with the size of 11 GB. In total, there are 5,200,000 user reviews, information on 174,000 businesses and the data spans 16 metropolitan areas in the USA and Canada.

## 3 TECHNOLOGY

Data Analytics: Apache Spark(PySpark) Spark Machine Learning Data Visualization: JavaScript, the D3.js Emerging Big Data Technologies: Docker

## 4 RATING DISTRIBUTION

### 4.1 Rating Star and Review Amounts Distribution

Create PySpark DataFrame from Yelp business dataset. Extract only restaurant business data and do a depth analysis of the distribution of rating stars, rating review quantities and top rated restaurant attributes. Through the distribution of data, we could help users to find a perfect restaurant, we could help businesses to find a high volume location to open a restaurant and find the five star restaurant common characteristics for helping startups choose sensible strategies to start the business. We could help the restaurant to boost its reputation by some strategies. And also we could provide Yelp some suggestions to optimize App user experience. As we know, Yelp can sort results by "recommended", "distance", "rating", "most reviewed", but sorting by only one option each time. Enhance sorting filter in more options, such as sorting by both rating stars(rating) and review amounts(most reviewed), sorting by both rating stars and distance.

state	city	stars	review_count	name
ABE	Vancouver	4.5	14	Kitanoya Guu Garlic
BC	BURNABY	3.5	59	IHOP
BC	BURNABY	2.5	31	IHOP
BC	Bowen Island	5.0	8	Shika Provisions
BC	Bowen Island	4.5	26	Cocoa West Chocol...
BC	Bowen Island	4.5	13	Barcelona Tapas &...
BC	Bowen Island	4.0	52	Tuscany Restaurant
BC	Bowen Island	4.0	40	The Snug Cafe
BC	Bowen Island	4.0	30	Artisan Eats Cafe...
BC	Bowen Island	4.0	24	Bowen Island Pub
BC	Bowen Island	4.0	22	Rustique Bistro
BC	Bowen Island	4.0	16	Miksa Restaurant
BC	Bowen Island	4.0	9	Branch on Bowen
BC	Bowen Island	4.0	6	Sushi to Go
BC	Bowen Island	3.5	51	Doc Morgan's Pub ...
BC	Bowen Island	3.5	7	Lime and Moon Pie...
BC	Bowen Island	3.5	6	Leftbank
BC	Burnaby	5.0	7	Active Body Burnaby
BC	Burnaby	5.0	7	Five Loaves Two Fish
BC	Burnaby	5.0	6	SOCRATES in the H...

Figure 1: Top Rated Restaurant Based on Sorting by Rating and Review Count

Users can choose a location, a state, a city and even a specific postal code, sorting by stars, then the review count and it's easy to find the perfect restaurant with the highest stars and highest review count. As a customer, we hope to find an ideal restaurant with both high rating and most reviewed. And we could see that in British Columbia, if we want to go to a good restaurant in Burnaby, "Five Loaves Two Fish" will be a good choice. As a consequence, if Yelp could enhance the

sorting filter in both rating and review count, it will help the customers to find the ideal restaurant.

state	city	postal_code	total_review_count
TX	Austin	78704	124492
FL	Orlando	32819	121451
TX	Austin	78701	120005
OR	Portland	97214	104339
MA	Boston	02116	93835
OR	Portland	97209	74662
OR	Portland	97204	65349
MA	Boston	02113	64217
TX	Austin	78702	63891
MA	Cambridge	02138	58856
MA	Cambridge	02139	58794
OR	Portland	97205	57987
GA	Atlanta	30309	54481
OR	Portland	97202	53466
GA	Atlanta	30308	51272
GA	Atlanta	30318	50933
OH	Columbus	43215	49479
FL	Orlando	32803	49345
TX	Austin	78758	48169
MA	Boston	02111	46514

Figure 2: High Customer Flow Based on Restaurant Rating Count Distribution

Every business owner must figure out how location will or not contribute to the success of the business and choose a spot accordingly. There are many issues to consider when you're looking for a space to house your business. Location is very important for a restaurant business. A high customer flow location will increase your restaurant's customer volume. From the data we can see that if you want to open a restaurant in Austin, the location in postal code "78704" will be a good location because it has the highest total restaurant review quantities which means this district has a high volume of customers. And you may assume that this location may have lots of pedestrian traffic nearby and a good neighborhood.

Yelp Rating Distribution By State By Star

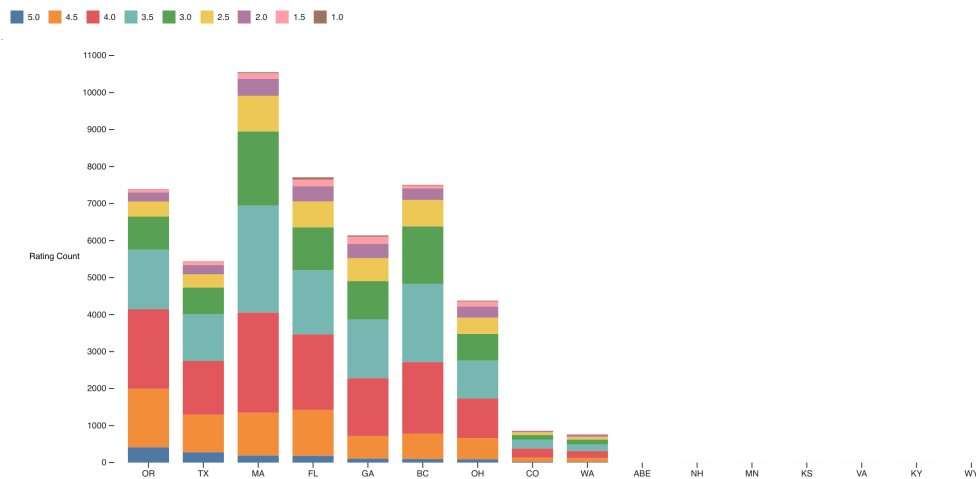


Figure 3: Yelp Restaurant Rating Distribution by State by Star

Extract rating distribution csv file from PySpark Dataframe, do data visualization via JavaScript, the D3.js. From "Yelp Rating Distribution By State By Star", we can see that the top 3 states with the largest number of top rated restaurants(5.0 rating star) are Oregon, Texas and Massachusetts. Oregon, Massachusetts and Florida have the largest number of high rated restaurants(above 4.0 star). And Massachusetts seems to have the largest number of restaurant's reviews in 16 states in the USA and Canada. If you are interested in opening chain restaurants. Massachusetts, Florida, Oregon and British Columbia are all good places to.

## 4.2 Common Characteristics

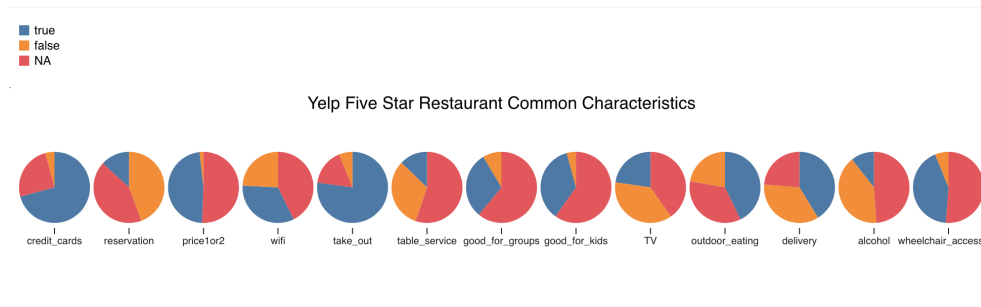


Figure 4: Yelp Five Star Restaurant Common Characteristics

Extract top rated restaurant attributes csv file from PySpark Dataframe, do data visualization via JavaScript, the D3.js. Finding Yelp Five Star Restaurant Common Characteristics could help businesses to set their perfect restaurant model. From the data, we can see that most of the five-star restaurants do accept credit cards, don't need reservations, have wifi, do provide food take out, food delivery and outdoor seating, are good for kids and groups, and are wheelchair accessible. It seems the table services and alcoholic beverages and TV are not very important for being a critical element to enhance the restaurant's reputation. And it seems that the price of most of the high rated restaurants are in range 1 and 2 which means perhaps the low end restaurants are not expected to be as high as the high end restaurants for most of the consumers.

## 5 GEOGRAPHIC DISTRIBUTION

### 5.1 World Map US Map

Although the yelp data set contains enormous amount of geographic data, it is hard for us to directly get the information out of plain table. Thus, it is important for us to show the geographic data using map chart. Firstly, in order to get a general idea of the data distribution world wide, we collect the longitude and latitude data of each restaurant and plotted them in the world map.



Figure 5: Restaurant Distribution on World Map

From the red scatter points we can tell that the restaurant data of our data set primarily located in the United States. Thus, to get a closer look at the pattern, we need to zoom in to the US since the data set is not collected worldwide. The image below shows the restaurant distribution in the US. From the US map as shown in figure 6, we can find out that the data in our data set is incomplete in national wide and are more city focused. Thus, to filter useful information for geographic analysis,

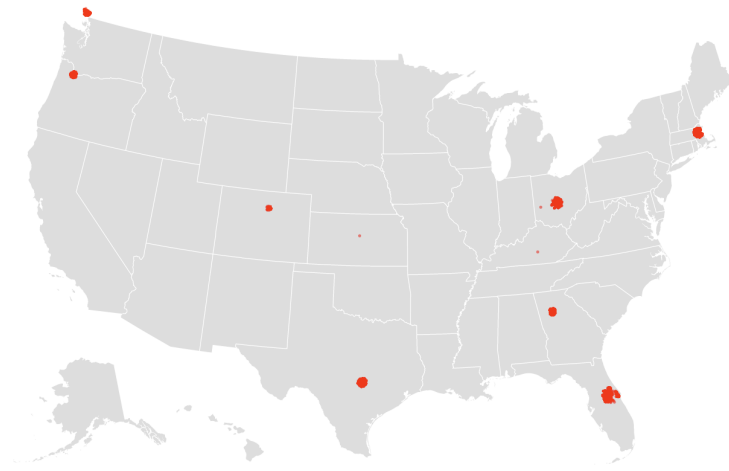


Figure 6: Restaurant Distribution on US Map

we have to take city map as the base map and visualize the restaurant location in that city. We use pyspark to rank the cities in our data by the number of restaurants they have and slice out the top 5 cities because we want our data to be as much as possible.

city	count
Portland	5737
Austin	4965
Atlanta	4180
Orlando	3749
Boston	2846

Figure 7: Top 5 cities with restaurant count

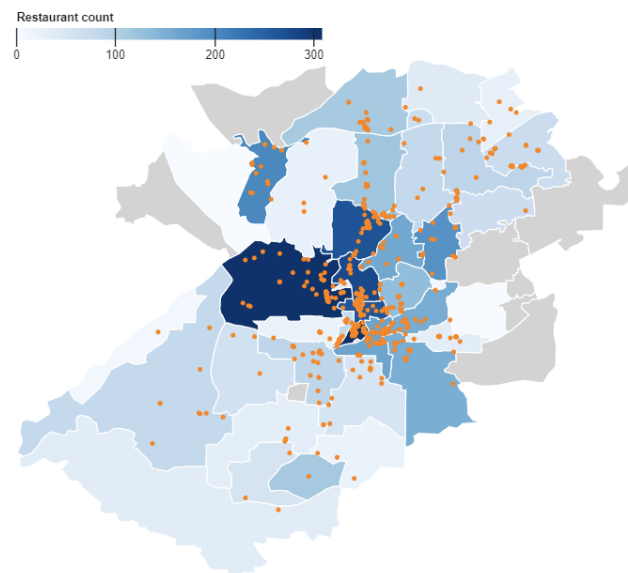


Figure 8: Restaurant count choropleth Atlanta map plotted with restaurants that have stars over 4.5

## 5.2 Atlanta Distribution Map

The figure 7 shows that the top 5 cities are: Portland, Austin, Atlanta, Orlando, Boston. Take Atlanta as the example map, the oranges dots plotted on the choropleth map of figure 8 of Atlanta are the restaurants whose rating are equal or higher than 4.5, which is a threshold that show enough data yet not cluttered too much in the graph. The intensity of the blue shows the number of restaurant in that area. As we can see that the restaurant tend to gather at the center of the city which is intuitively correct since usually the center of the city has the largest number of population. In order to prove that, we also configured the choropleth map with respect to population and average income.

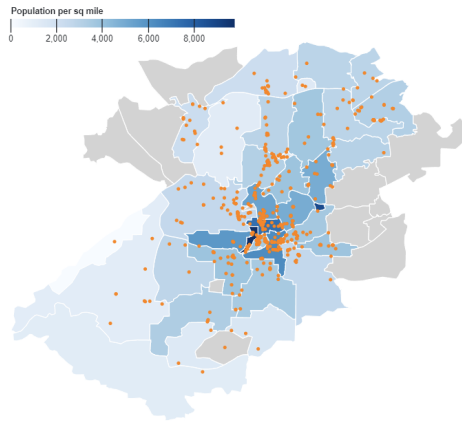


Figure 9: Choropleth Atlanta map with respect to population density

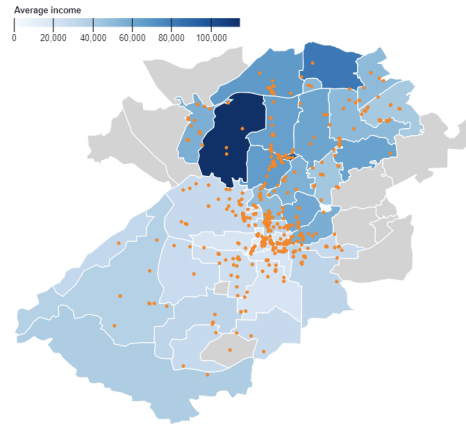


Figure 10: Choropleth Atlanta map with respect to average income

As we can see from figure 9 and 10, the restaurants do tend to clutter at the location that has the highest population density. However, it is not the case for the income choropleth map. This is because people with higher income tends to live in a certain neighborhood which usually located far away from the city downtown. With low population density, the location will have less restaurant.

## 6 REGIONAL ANALYSIS (ATLANTA)

In this section, a deeper understanding of the restaurant industry in Atlanta, GA is presented. Below are the top 18 business categories in the city, with more than 4000 restaurants ranked first and around 800 estate businesses ranked last. In the figure, there could be some overlaps in each bar due to the synonyms in the category terms, and some multi-category businesses. However, we can still conclude that the most of the businesses are restaurant related.

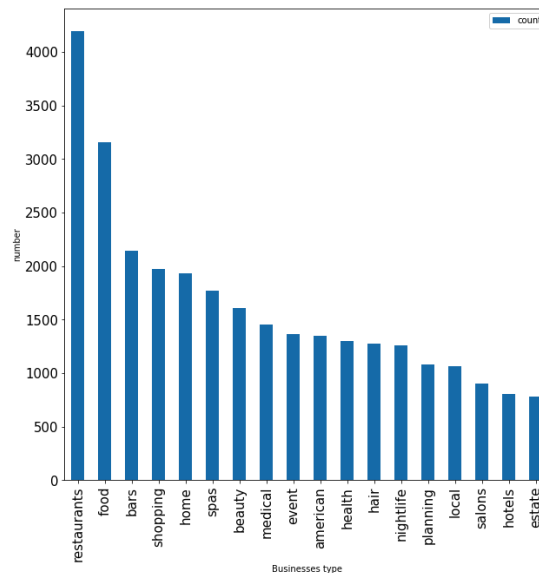


Figure 11: Top eighteen Atlanta businesses categories

## 6.1 Food trend discovery

By taking a closer look at the restaurants category, the top 10 most popular highly rated restaurants in Atlanta are revealed below. The threshold of restaurants on this

	name	review_count	stars	categories
0	Mary Mac's Tea Room	3861	4.0	Venues & Event Spaces, Event Planning & Services, Comfort Food, <b>Southern</b> , Restaurants, Nightlife, Tea Rooms, Food, Bars
1	Fox Bros. Bar-B-Q	3825	4.5	<b>Southern</b> , Restaurants, Barbeque, <b>American</b> (Traditional)
2	Atlanta Breakfast Club	3761	4.5	Restaurants, Breakfast & Brunch, <b>American</b> (Traditional), <b>Southern</b> , Coffee & Tea, <b>American</b> (New), Food
3	Poor Calvin's	3379	4.5	<b>Southern</b> , Nightlife, Thai, Food, <b>American</b> (New), Ethnic Food, Restaurants, Beer, Wine & Spirits, Asian Fusion, Bars, Comfort Food, Specialty Food, Seafood
4	Antico Pizza	3135	4.0	Italian, Pizza, Restaurants
5	Two Urban Licks	2730	4.0	Bars, <b>American</b> (New), Nightlife, <b>Southern</b> , Restaurants
6	South City Kitchen Midtown	2618	4.5	Restaurants, <b>American</b> (New), <b>Southern</b> , Breakfast & Brunch, Gluten-Free
7	Fat Matt's Rib Shack	2165	4.0	Restaurants, Barbeque, Fast Food, Nightlife, <b>Southern</b> , Beer Bar, Bars, <b>American</b> (Traditional)
8	FLIP burger boutique	1909	4.0	Burgers, Specialty Food, Restaurants, Barbeque, Food, Nightlife, Bars, <b>American</b> (New)
9	Canoe	1844	4.5	Seafood, Restaurants, Breakfast & Brunch, <b>American</b> (New)

Figure 12: Most popular highly rated restaurants in Atlanta(>4.0 stars)

list is 4.0 average customer rating, and they are in descending order by the number of reviews. In this table, 7 out of 10 restaurants offer southern food, and 8 out of 10 are classified as American(New/Traditional) restaurants. This is an interesting pattern as we later discovered these two are the most popular types of cuisine in the highly rated restaurants in Atlanta. Unfortunately, Yelp does not specify cuisine and

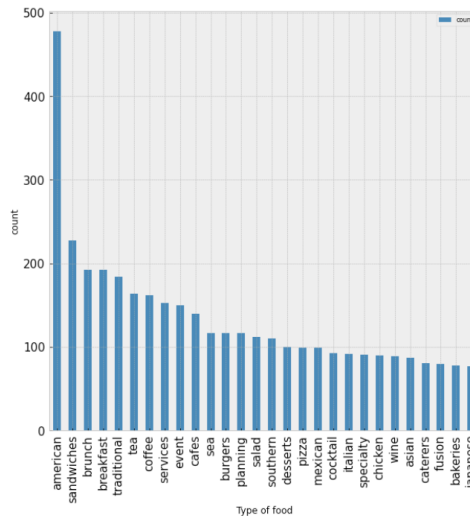


Figure 13: Atlanta food type rank among highly rated(>4.0 stars) restaurants

food types, but it does not affect the food trend to be discovered. The figure above reveals the food type rank in the highly rated (>4 stars) restaurants in Atlanta. As we can observe, the most popular cuisines are American, traditional, and southern. The most popular types of food are sandwiches, breakfast brunch, and burgers.

## 6.2 Review analysis

We have also done a review analysis in both 5 stars and 1 star Atlanta restaurant review. Below are the most common words used in each division. On the left side, we have the most common words used in 5 stars reviewer ratings. Other than the subject that customers appreciate for, some complimentary words are also used, such as "great", "good", "amazing", "delicious", "best", etc. On the other hand, people tend to use less descriptive terms in the 1 star reviewer ratings. Instead, almost all of the words are the related subjects that customers complained about.



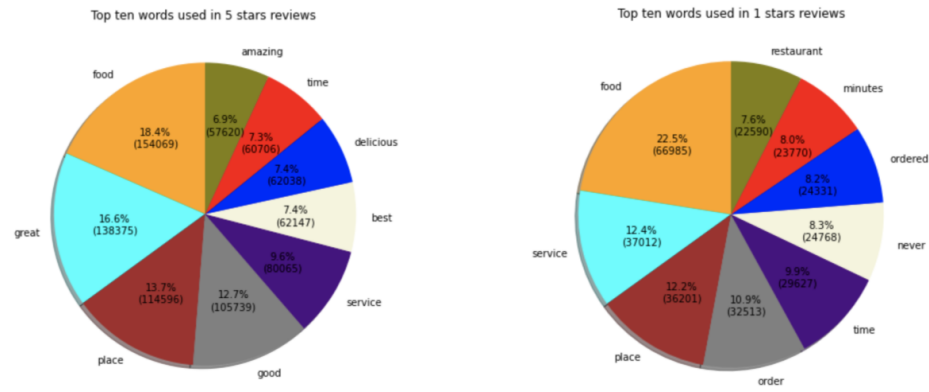


Figure 14: Most common words used in 5stars and 1star reviewer ratings for Atlanta restaurants

### 6.3 Business value of the regional analysis

In this section, some interesting data such as business distribution, food trend and reviewer text use in high and low rating stars are reported. These results should be valuable for both local restaurant owners in Atlanta and Yelp developers. For the former, they can use the data such as food type rank to analyze the popularity of their food. This data should be helpful if the owner wants to expand their menu or anyone wants to open a restaurant in Atlanta. For Yelp developers, the food trend and the most popular highly rated restaurants can appear on the recommendation list for any user who is located in Atlanta. The common words for different reviewer ratings can be generated as tags, so the customers can click on the tags to submit reviews instead of typing. This could encourage more users to write reviews as the review submission is more convenient. Also, this can eliminate the problem in which customers leave blank content reviews (star only reviews), as other users might be curious about the experience of the reviewer.

## 7 SENTIMENTAL MODEL

### 7.1 Background

Given a large number of subjective text reviews of a product or restaurant on the web today, we want to derive a lot of benefits by detecting and interpreting the users' opinions and sentiment in their reviews. For example, sentiment analysis can quantify a review, provide succinct review summaries to businesses or readers, or build an automatic recommendations system for both business owners and customers.

We focus on sentiment detection and classification, which takes a text and attempts to classify it based on the classifier's detected sentiment. In this project, we applied the natural language processing technique and multiclass logistic regression to classify a set of reviews based on the ratings and a set of the predefined selected feature.

Specifically: 1. Pre-process each text review, converting them to a collection of unigrams, bigrams, and trigrams. 2. Pre-select a set of n-grams as features set based on their frequency in all the n-grams we created. 3. Vectorize each label point, set the value to 1.0 if a feature existed in the review and 0.0 if not existed in the review. 4. Build a classifier from the training data to classify each text review from 1-star to 5-star. 5. Experiment with different features set with varying proportions of n-grams. 6. Compare the evaluation model accuracy and true positive rate, on label and model levels.

stars	text
4.0	Apparently Prides...
4.0	This store is pre...
5.0	I called WVM on t...
2.0	I've stayed at ma...
4.0	The food is alway...

only showing top 5 rows

Figure 15: Dataset Example

stars	count	stars	std_dev	skew	kurtosis
1.0	1262800	1.0	711.22	2.41	7.85
2.0	711378	2.0	648.4	2.25	7.44
3.0	926656	3.0	589.97	2.29	8.19
4.0	1920037	4.0	533.98	2.44	9.64
5.0	3814532	5.0	457.63	2.95	14.25

Figure 16: Left: Review Count by Each Star; Right: Review Stats by Each Star

## 7.2 The Dataset

Our dataset consists of nearly 860,000 restaurant reviews from yelp.com. For the purpose of this project and the limited memory resources available to us, we randomly select a 0.1 percent of these reviews when implementing and testing our models, and use 70 percent as training dataset and 30 percent as testing dataset. Through further analysis of the sample data, we found that there are more positive reviews (5-star) and negative reviews (1-star) than neutral reviews and the distribution of length of strong reviews (1-star and 5-star) tend to be more skewed than that of neutral reviews, as shown in Figure 16

In terms of preprocessing the text-based review, we applied NLTK (natural language processing tool kit) stopwords list to remove the unnecessary stopwords and used NLTK PorterStemmer to simplify the English words so that we can have a more general collection of n-grams, which can be later used for the word counting process.

## 7.3 features

In dealing with pre-selecting a set of features before vectorizing the data points, we decided to choose an arbitrary number as a threshold. Any n-gram whose frequency is greater than that threshold will be selected into the features set. By having a higher threshold, we will have a smaller size of features set, and vice versa. We decided to use length as one of the features to train our model. In order to properly assign a category to each review based on the length of the review, we build a bucket based on the standard deviation of the distribution of the length of all text reviews. The category of length that each review is assigned to is determined by how many standard deviations away from the mean of the length the review is.

## 7.4 Implementation

After finishing features selection, we vectorized the data point by checking if each feature, an n-gram, in the features set exists in each review, and set the value to 1.0

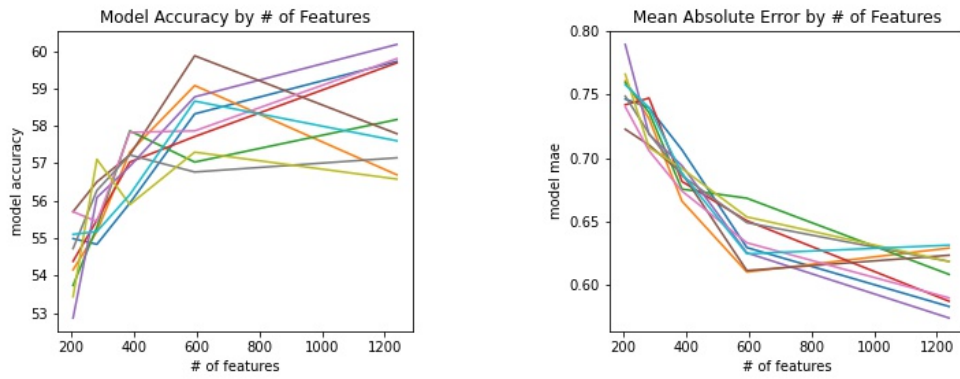


Figure 17: Model Accuracy and Model Mean Absolute Error by Number of Features

if it exists and 0.0 if it doesn't. We also decided to choose the length of each text review as one of the features in addition to the n-grams.

We chose a multiclass logistic regression model provided from the PySpark machine learning library. We ran the model against the training dataset and evaluated the model accuracy against the testing dataset. In order to make the result more robust, we ran ten times for both the training and evaluation parts. In each run, we change the threshold and see how model accuracy changes in response to the change in the threshold. As a result, we can see that the model accuracy improved as the size of the features set increases due to a lower threshold, and the model means absolute error decreases as the number of features increases, as shown in Figure 17. MAE is calculated as follows:

$$\text{MAE} = \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{n}$$

, where

$$\hat{y}_i$$

is the predicted rating by our classifier and  $y$  is the actual rating

However, we found out that our model has higher accuracy in detecting strong positive, 5-star, and negative sentiment, 1-star, and lower accuracy in detecting neutral sentiments, which span from 2-star to 4-star. We think the reason for this is that sample data have more 1-star and 5-star reviews than 2-star/3-star/4-star reviews, and the strong negative and positive ratings tend to come along with strong words and phrases that are easier to study and detect than the more neutral phrases are.

## 7.5 Improvement

In order to improve model accuracy, we tried to regroup the neutral reviews by regrouping either 1-star and 2-star together and 4-star and 5-star together or 2-star, 3-star, and 4-star together. The former approach would build a model suitable for detecting general positive and negative sentiment but failing to distinguish strong sentiment from slightly strong sentiment. The latter way would give us a model that is best suitable for detecting strong sentiments. First way to improve our model is to reduce the number of classes by regrouping the slightly neutral ratings. For instance, if we care more about general negative sentiment and positive sentiment, we can relabel 2-star to 1-star, and 4-star to 5-star, because 2-star reviews are arguably negative sentiment, and 4-star are arguably positive sentiment. By doing so, our model accuracy in detecting positive and negative sentiment improved compared to the original model. Alternatively, the second way is to relabel both 2-star and 4-star to 3-star, meaning that we only care about the ability to detect the extreme negative (1-star) and extreme positive sentiment (5-star). By doing so, our accuracy

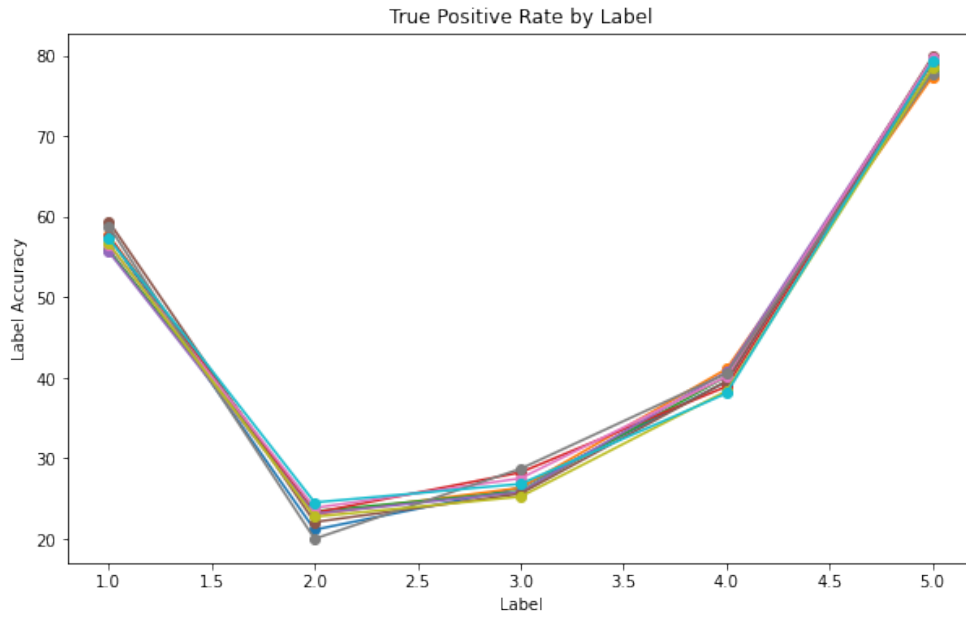


Figure 18: True Positive Rate by Each Label

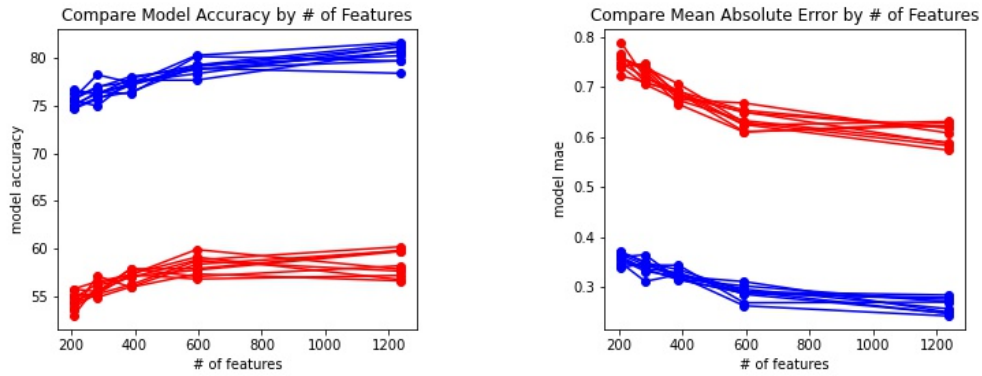


Figure 19: Regrouping 1 and 2-star, 4 and 5-star; Left: Model Accuracy, Right: Model MAE; (blue line: regrouping, red line: without regrouping)

in detecting positive and negative sentiment also improved. However, both models achieved higher accuracy at the expense of predicting the power of detecting slightly neutral sentiment.

From both Figure 19 and Figure 20 we can see that, by regrouping, the model has higher accuracy (the blue line) than the one without regrouping (the red line), and the mean absolute error is also lower than the original ones. However, from Figure 21, we can see that regrouping 1 and 2-star, and 4 and 5-star does not improve the true positive rate of detecting neutral sentiment, label 2 (blue line)

Which regrouping approach to choose depends on whether we care more about detecting the general positive and negative sentiments or just strong sentiments. If we just want to detect the general positive and negative sentiment, then the first approach is preferred. If we care more about the detecting extreme sentiment, then the second approach is preferred.

## 7.6 Business Use Case

Yelp can help detect users' sentiment for food products or restaurants and help build a better recommendation system for restaurants, generating potential ad rev-

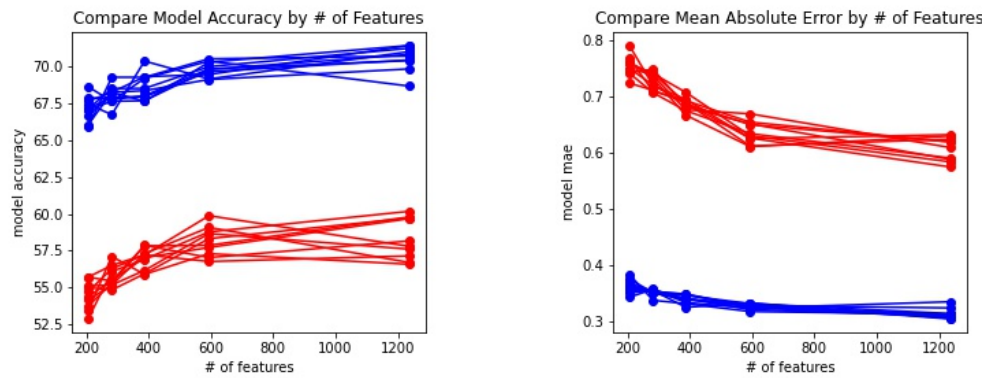


Figure 20: Regrouping 2, 3, and 4-star; Left: Model Accuracy, Right: Model MAE; (blue line: regrouping, red line: without regrouping)

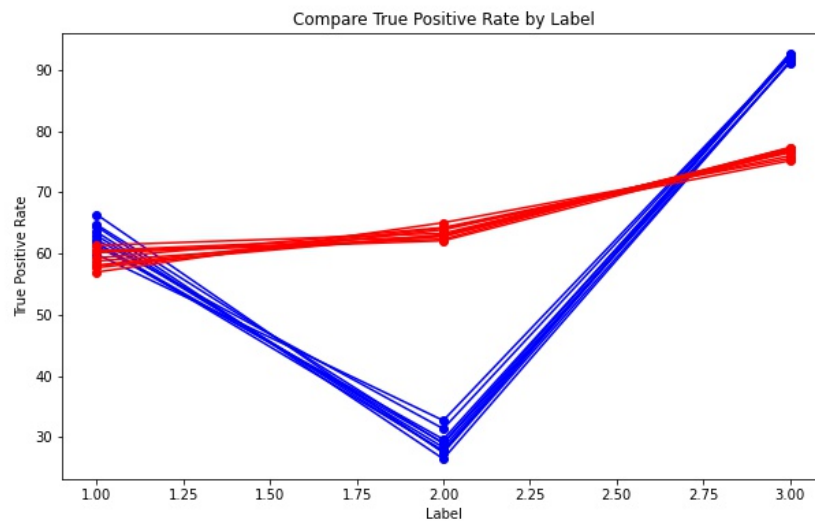


Figure 21: Blue Line: Regrouping 1 and 2-star, 4 and 5-star; Red Line: Regrouping 2, 3, and 4-star

enue on some famous social media platforms where ratings are unavailable to users, such as Twitter, Facebook, or Instagram.

When a national or global franchise restaurant launches a new food product or an event, and the restaurant owner wants to get well-quantified feedback from text-based reviews posted by users online, in real-time, Yelp can be the solution provider that helps detect the sentiments using this model by partnering with the social media platforms mentioned above.

## 8 TEXT ANALYSIS TRIGRAM COUNTS

### 8.1 Analysis and Results

We conducted the text analysis by applying trigram counts on the review comments. Because the original data-set is enormous we only take 1 percent of original data as our sample data. The reason we utilized trigram count rather than unigram or bigram is that it provides proper context for the descriptive words so we can better classify positive and negative comments. For example, phrases such as “great customer service “ and “best food ever” clearly indicate the reasons behind the high ratings. Subsequently, we designed a model that consists of 4 different dimensions and for each dimension, there are a few indicators as shown in Figure 23. And we aggregated those trigram counts under the descriptive word category. As we can see from Figure 22, obvious phrases such as “horrible customer service”, “worst customer service” can simply be categorized into employee attitude under service quality dimension. And “wait 10 mins” or “wait 20 mins” imply long wait time, therefore we total these into wait time and order under service quality dimension. Once we have the aggregation results, we started comparing the one-star-rated comments to five-star-rated comments for restaurants of the same price range. Here is what we have found: for low-end restaurants with one dollar sign, service quality and Food quality are the most substantial factors that impact the ratings. And words related to atmosphere and price-value perception are barely present in comments. On the other hand, for premium high-end restaurants with four dollars signs, Service quality and value proposition are the most substantial factors that impact the ratings. And the atmosphere and food quality don’t contribute to negative comments. Even in one-star-rated comments, customers mentioned that the place was clean and the atmosphere was great.

```

In [11]: processed_data
Out[11]: [('never go back', 53),
          ('horribl custom servic', 28),
          ('never come back', 28),
          ('worst custom servic', 28),
          ('poor custom servic', 23),
          ('terribl custom servic', 21),
          ('worst servic ever', 20),
          ('give zero star', 18),
          ('go somewher els', 16),
          ('get money back', 15),
          ('everi time go', 14),
          ('bad custom servic', 14),
          ('wish could give', 14),
          ('wait 10 minut', 14),
          ('wait 20 minut', 14),
          ('custom servic skill', 13),
          ('0 star would', 13),
          ('ask speak manag', 13),
          ('wast time money', 13),
          ('worst custom servic', 13)]

```

Figure 22: Trigram Count Result from One Star Rated Comments

★ ★ ★ ★ ★

Dimensions	Indicators	Descriptive words
Food quality	Taste	Derogatory: sub-par (6), really bad (5)
Service quality	Employee attitude	Derogatory: horrible service (28), poor service (23), worst service (20), rude service (6)
	Wait time and order	Derogatory: long wait time (102), wrong order (20)
Atmosphere	Environment	Commendatory: Derogatory:
Price and value	Price	Commendatory: Derogatory:
	Price-quality ratio	Commendatory: Derogatory:

★ ★ ★ ★ ★

Dimensions	Indicators	Descriptive words
Food quality	Taste	Commendatory: great food (36), authentic Mexican food (25), good food (21)
	Health	Commendatory: fresh food (10)
Service quality	Employee attitude	Commendatory: great service (92), super friendly (83), excellent service (18)
	Wait time and order	Commendatory: serve fast (25)
Atmosphere	Environment	Commendatory: Derogatory:
Price and value	Price	Commendatory: reasonable price (18) good price (20)
	Price-quality ratio	Commendatory: Derogatory:

Text Model for Low End Restaurants

★ ★ ★ ★ ★

Dimensions	Indicators	Descriptive words
Food quality	Taste	Commendatory: Derogatory:
	Health	Commendatory: Derogatory:
Service quality	Employee attitude	Derogatory: poor service (30) terrible service (20)
	Wait time and order	Derogatory: long wait time (30)
Atmosphere	Environment	Commendatory: great atmosphere (25) clean (35)
Price and value	Price	Derogatory: too expensive (45)
	Price-quality ratio	Derogatory: not worth the price (50)

★ ★ ★ ★ ★

Dimensions	Indicators	Descriptive words
Food quality	Taste	Commendatory: best meal (40), best food (30), best dessert (20)
	Health	Commendatory: Derogatory:
Service quality	Employee attitude	Commendatory: great service (60), top-notch (30)
	Wait time and order	Commendatory: Derogatory:
Atmosphere	Environment	Commendatory: Derogatory:
Price and value	Price	Commendatory: Derogatory:
	Price-quality ratio	Commendatory: worth every penny (30), well worth the price (45)

Text Model for High End Restaurants

Figure 23: Text Model for Restaurants of Different Price Range

## 8.2 Business Insight

In order to improve restaurants' ratings in Yelp's review section, we can strategize based on restaurants' price range. For low end restaurants, they should focus on food quality and service quality since service quality and food quality are the most substantial factors that impact the ratings. For premium and high end restaurants, they should prioritize service quality and extra added values since service quality and value proposition are the most substantial factors that impact the ratings.

## 9 CONCLUSION

After a comprehensive analysis of Yelp's reviews and user data which includes information about businesses across 16 metropolitan areas in two countries. We have come to conclusion on two major aspects' Yelp app. One is the recommendations for business, primarily restaurants in terms of how they can improve their business performance based on the feedback from Yelp app. The other one is the recommendations for Yelp itself in terms of how to improve the functionality of Yelp app.

### 9.1 Recommendations for business

Business can use the data such as food type rank to analyze the popularity of their food. This data should be helpful if the owner wants to expand their menu or anyone wants to open a restaurant in Atlanta.

Also business can strategize based on restaurants' price range. For low end restaurants, they should focus on food quality and service quality since service quality and food quality plays a significant role in impacting the ratings. For premium and high end restaurants, they should prioritize service quality and extra added values since service quality and value proposition are what customers care the most after spending a lot of money on premium charge.

### 9.2 Recommendations for Yelp

For Yelp developers, the food trend and the most popular highly rated restaurants can appear on the recommendation list for any user who is located in Atlanta. The common words for different reviewer ratings can be generated as tags, so the cus-

tomers can click on the tags to submit reviews instead of typing. This could encourage more users to write reviews as the review submission is more convenient. Also, this can eliminate the problem in which customers leave blank content reviews (star only reviews), as other users might be curious about the experience of the reviewer.

Yelp can help detect users' sentiment for food products or restaurants and help build a better recommendation system for restaurants, generating potential ad revenue on some famous social media platforms where ratings are unavailable to users, such as Twitter, Facebook, or Instagram.

When a national or global franchise restaurant launches a new food product or an event, and the restaurant owner wants to get well-quantified feedback from text-based reviews posted by users online, in real-time, Yelp can be the solution provider that helps detect the sentiments using this model.

We also conducted a qualitative analysis which is based on the usage of Yelp from our own focused group. Our goal for this part of analysis is to make recommendations on how to improve Yelp's functionality to optimize user experiences. And our recommendation is that Yelp can enhance sorting filters in more options, such as sorting by both rating stars and review amounts (most reviewed), sorting by both rating stars and distance.