

Sentiment Analysis Yelp

Springboard Capstone Project 2
Serik Omarov



Background

- **Yelp**
- **Online Social network services**
- **Consumer provides reviews and ratings**
- **1 to 5 star- rating**
- **Customer feedback on services, quality and location etc**



Business Problem

- ❖ Will the customer reviews help in the indication of the provided rating?
- ❖ How Can restaurants access their success and faults based on reviews?
- ❖ What aspects of the business are correlated between positive and negative sentiments?



Outline

☐ Prepare Data for Supervised ML

- Filtering and Joining, Preprocessing, Normalization and Labeling

☐ Descriptive Analytics

- Explore data and visualize and understand the attributes

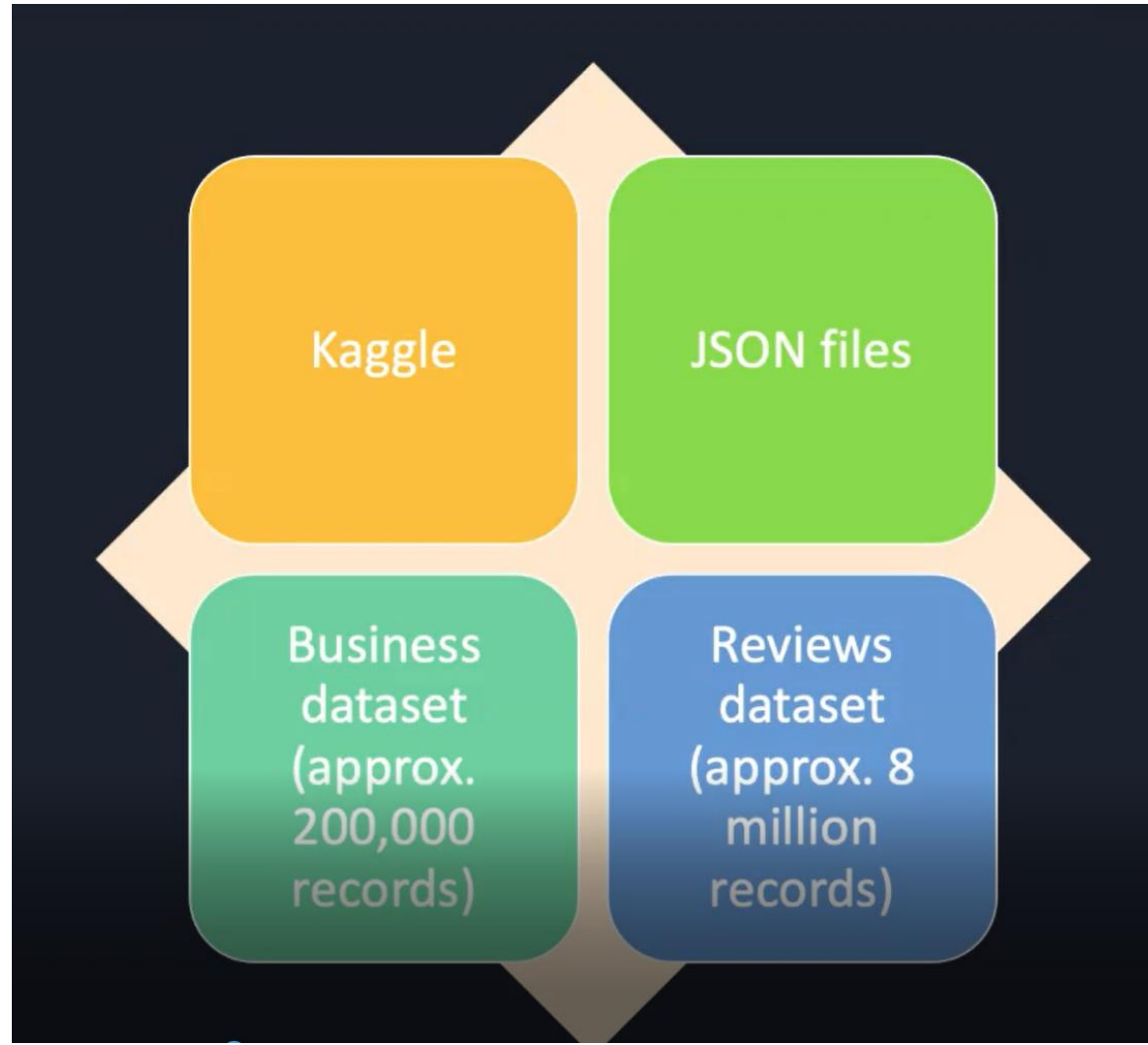
☐ Build Model

- Build and optimize a model that predicts sentiment given customer review

☐ Extract Results

- Extract features importance that classify customer sentiments

Data Sources



Project Approach

- **Latent Semantic Analysis (LCA) and Singular Value Decomposition (SVD)**
 - ✓ Topic Modeling
 - ✓ Dimensionality Reduction
 - ✓ Relationship between documents and Terms
- **Logistic Regression**
 - ✓ Predictive Binary Classification
 - ✓ Distinguish between two classes
 - ✓ BOW and TF-IDF to extract terms and coefficients for feature importance

Data Pre-Processing

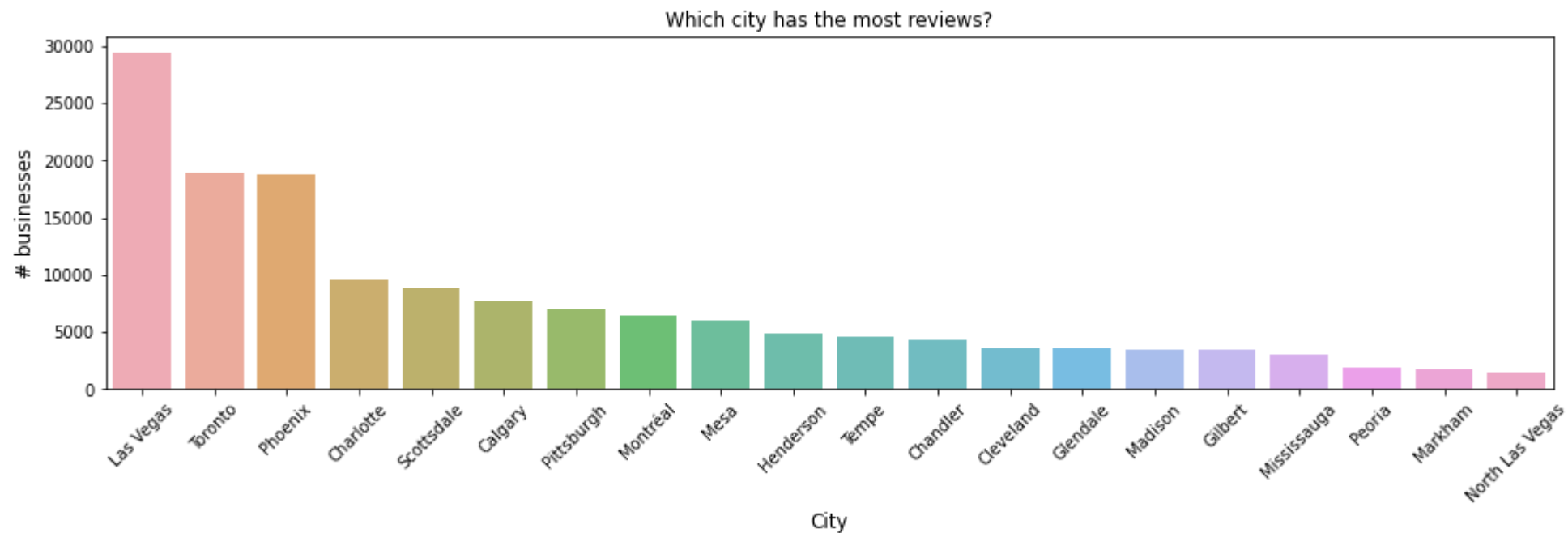
- **Step 1- Cleaning**

- ✓ Dropping columns and NA Values
- ✓ Filtering Data and narrowing dataset (City Las Vegas and Japanese Restaurants)
- ✓ Joining All business data with every review
- ✓ Labeling based on review stars (4, 5 stars positive review and 3 and under negative)

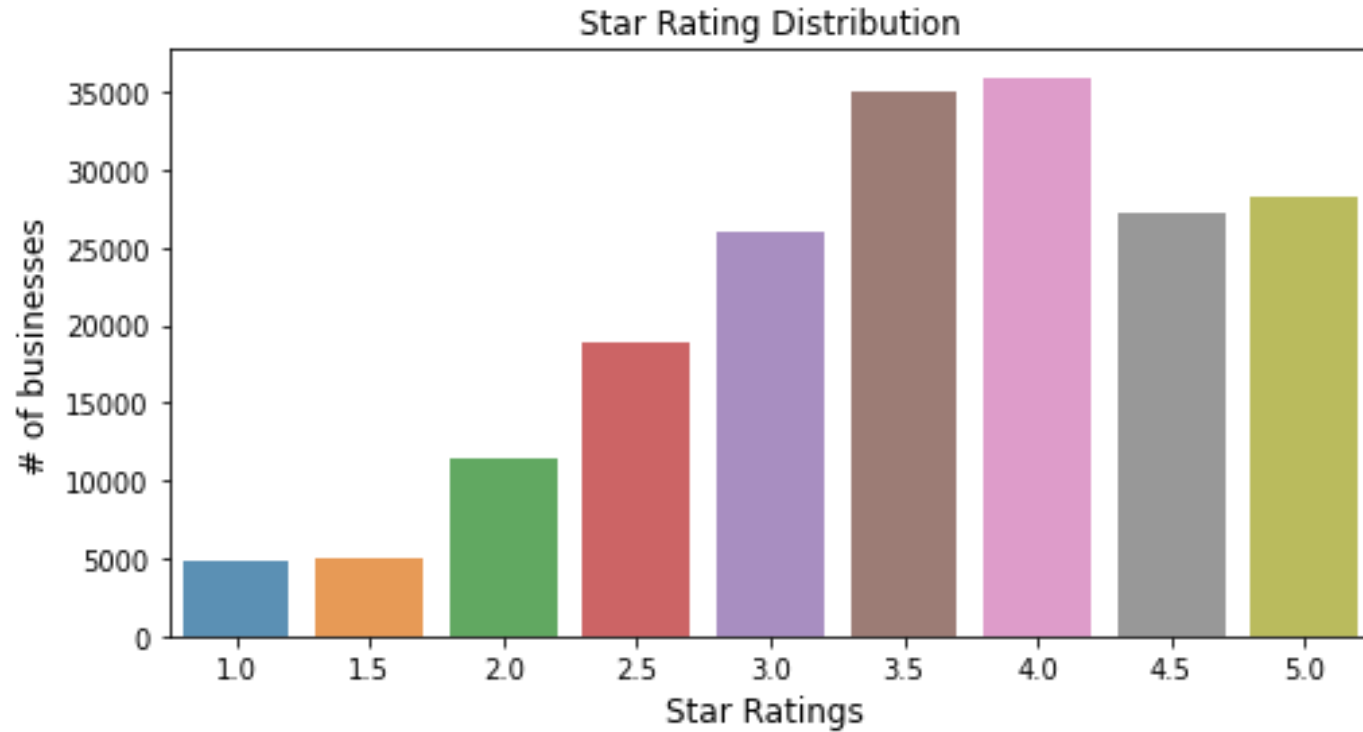
- **Step2 Normalization**

- ✓ Regex (removing white spacing and special characters)
- ✓ Porter Stemmer
- ✓ Vectorize

Exploratory Data Analysis

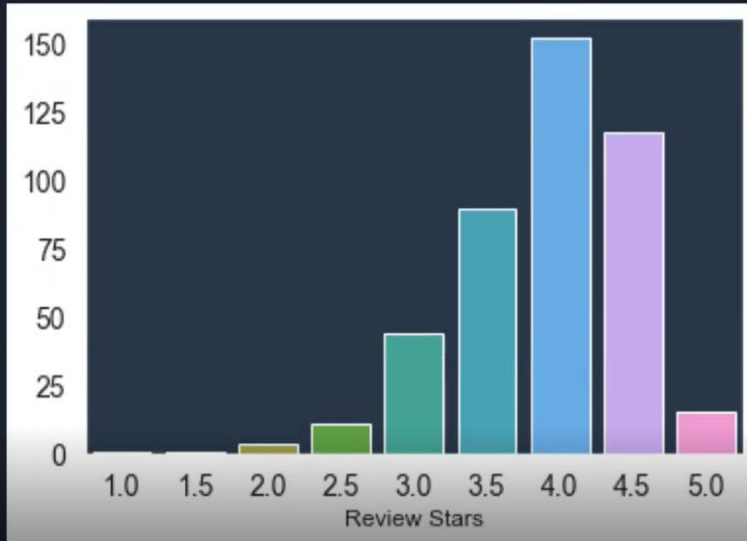


Exploratory Data Analysis

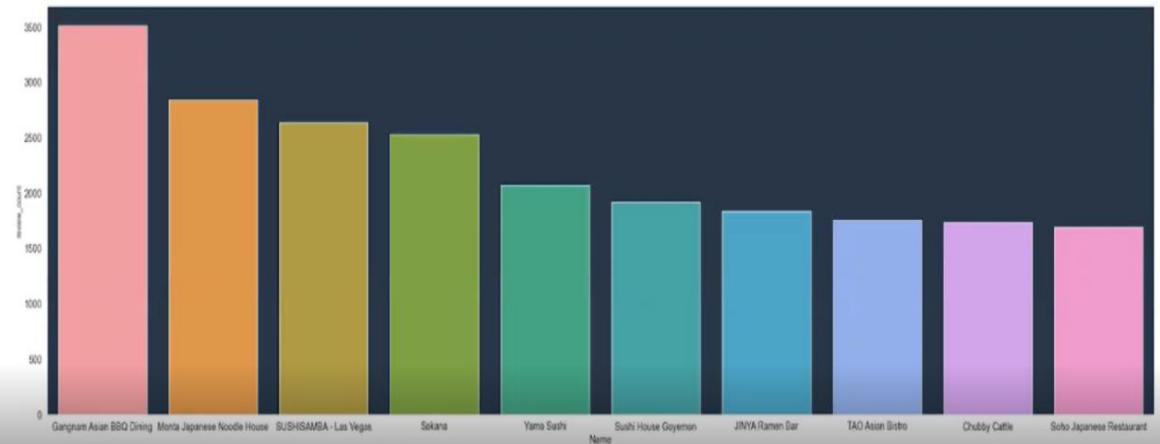


Exploratory Data Analysis: Las Vegas

Distribution of Star Ratings

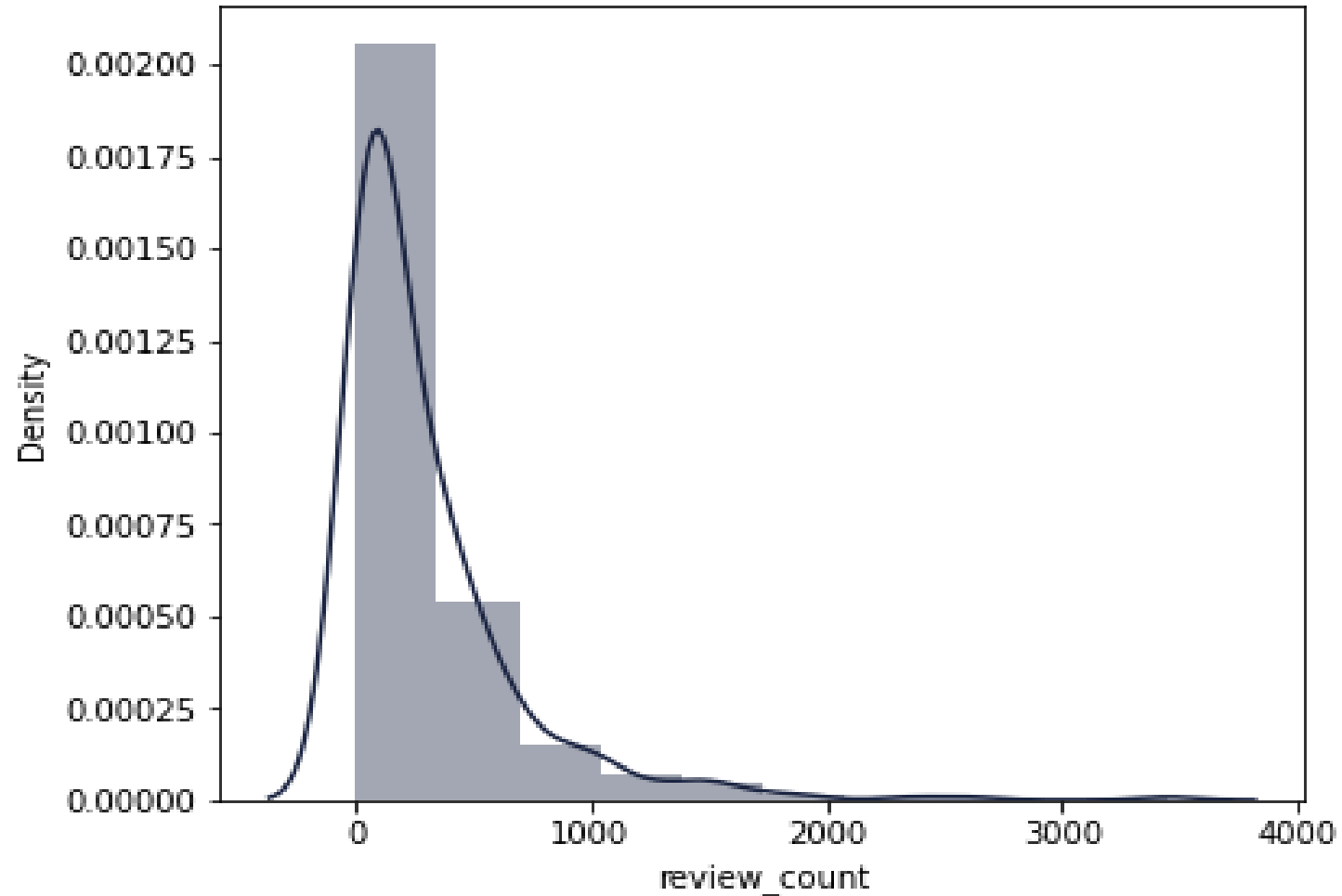


Top 10 Restaurants Based on Review Count



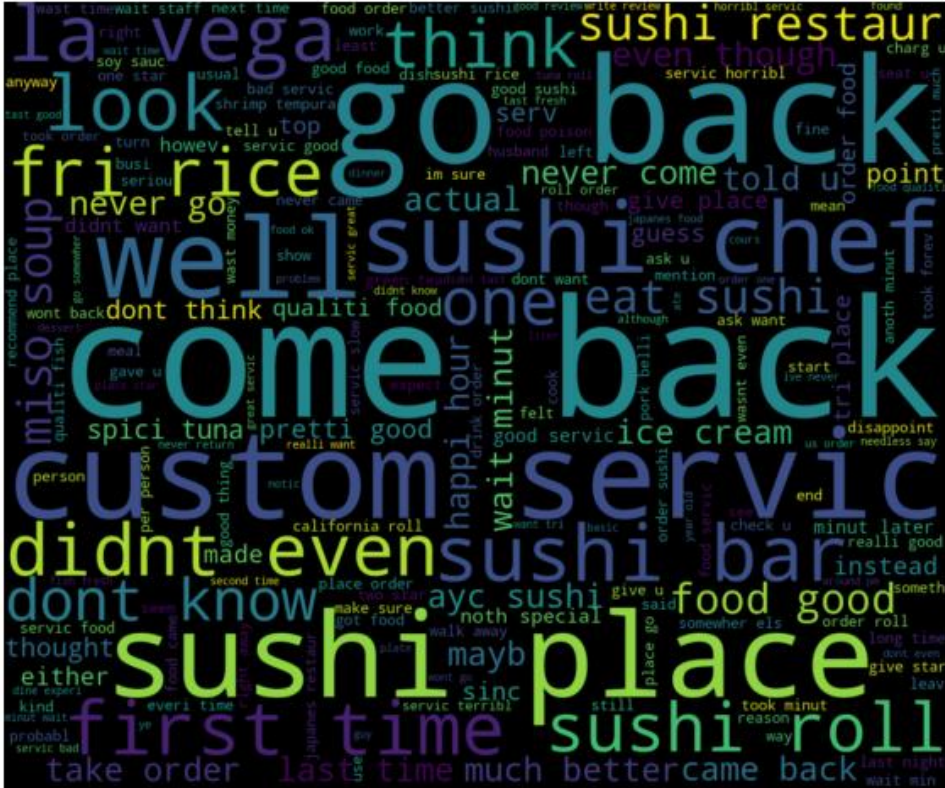
Exploratory Data Analysis: Las Vegas

Distribution of Count of Reviews



Exploratory Data Analysis

Negative Word Cloud



- Customer Service
- Fried Rice

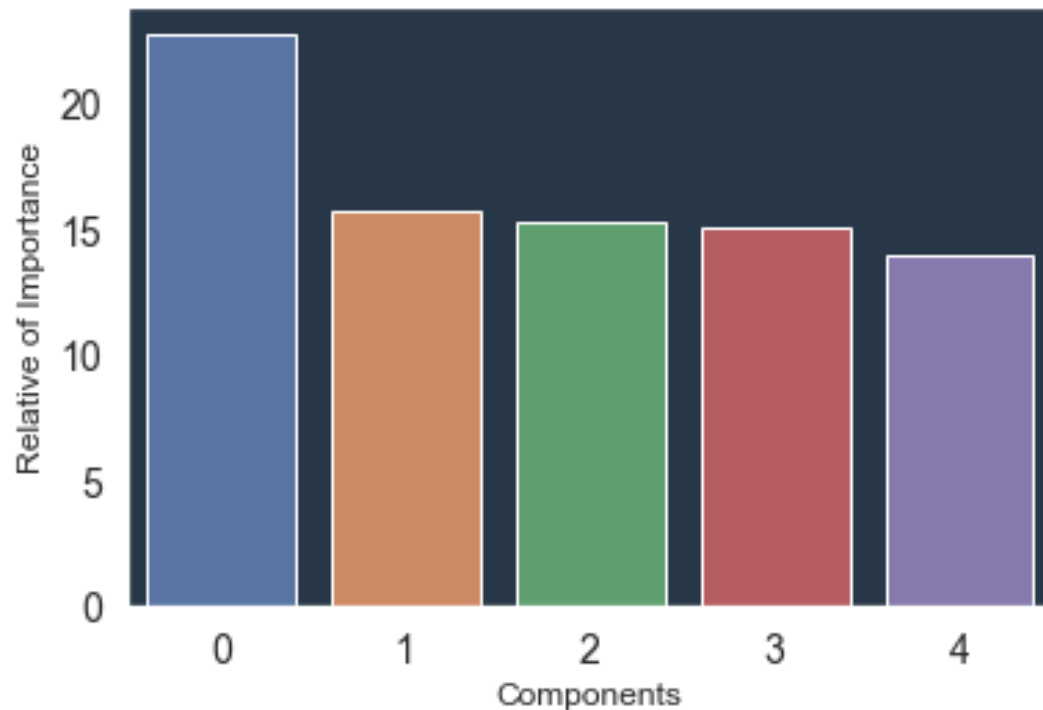
Positive Word Cloud



- Sushi Place
- Happy Hour
- Ice Cream

Latent Semantic Analysis (LCA)

- ❑ *5 concept of the top 10 words*
- ❑ *Understand relationship between document and terms*
- ❑ *Average conceptual idea of consumer's reviews and experiences*

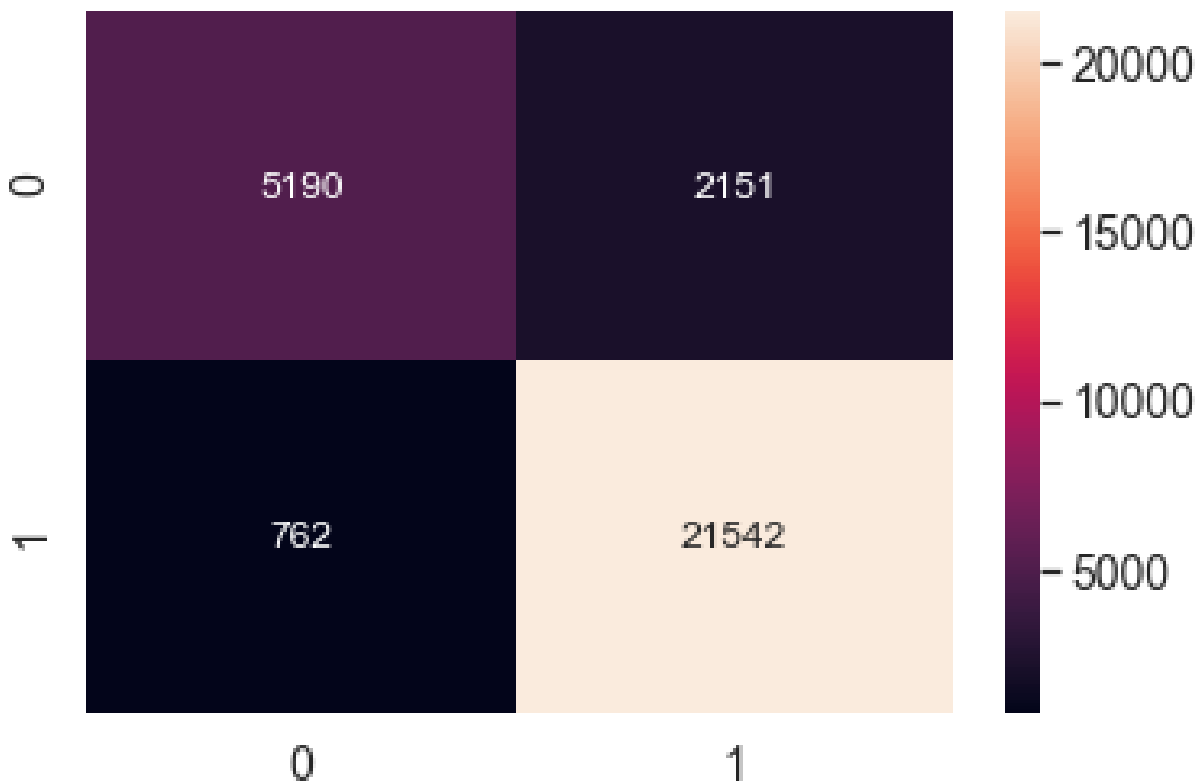


```
{ 'Concept 0': [('come back', 0.2431157699795417),  
  ('sushi place', 0.17997107780816396),  
  ('happi hour', 0.17804732633872644),  
  ('la vega', 0.16806065163897932),  
  ('great servic', 0.1601527456538924),  
  ('great food', 0.1536805335141426),  
  ('servic great', 0.1356425593111377),  
  ('food great', 0.13150955476848292),  
  ('first time', 0.12967232600828385),  
  ('highli recommend', 0.12603940331979532)],  
  'Concept 1': [('great food', 0.5218811933526182),  
  ('great servic', 0.49262986808609666),  
  ('food great', 0.4280575951443961),  
  ('servic great', 0.24755979873048756),  
  ('great price', 0.08176026154274607),  
  ('place great', 0.05714059964624758),  
  ('love place', 0.05361362620297335),  
  ('food servic', 0.05290096731960285),  
  ('definit come', 0.05038352348897965),  
  ('great atmospher', 0.04896926979287501)],  
  'Concept 2': [('happi hour', 0.6312600412663628),  
  ('hour menu', 0.11258051002344965),  
  ('sushi place', 0.1124378037487351),  
  ('great servic', 0.09678187490235811),  
  ('la vega', 0.09524375998037303),  
  ('best sushi', 0.09243326025487358),  
  ('great food', 0.08591520988129625),  
  ('food great', 0.06385241634253805),  
  ('great happi', 0.06063360733451816),  
  ('favorit sushi', 0.05834041062886958)],  
  'Concept 3': [('happi hour', 0.6548203854641816),  
  ('come back', 0.38109426694757226),  
  ('definit come', 0.23598214267045062),  
  ('hour menu', 0.11871321779958796),  
  ('would definit', 0.06106450893741144),  
  ('great happi', 0.05941249895675935),  
  ('hour price', 0.05189120393828811),  
  ('late night', 0.05106092400534962),  
  ('back tri', 0.04495207156130668),  
  ('first time', 0.04375899004490264)],  
  'Concept 4': [('come back', 0.38007304367868544),  
  ('sushi place', 0.34655698703219584),  
  ('best sushi', 0.31087854882458366),  
  ('definit come', 0.2601014523884045),  
  ('la vega', 0.2152988359892752),  
  ('happi hour', 0.1397904583508747),  
  ('favorit sushi', 0.13953106840489732),  
  ('one best', 0.08178624597495925),  
  ('sushi restaur', 0.0720665027759246),  
  ('ayc sushi', 0.06836812706133981)] }
```

Logistic Regression Implementation

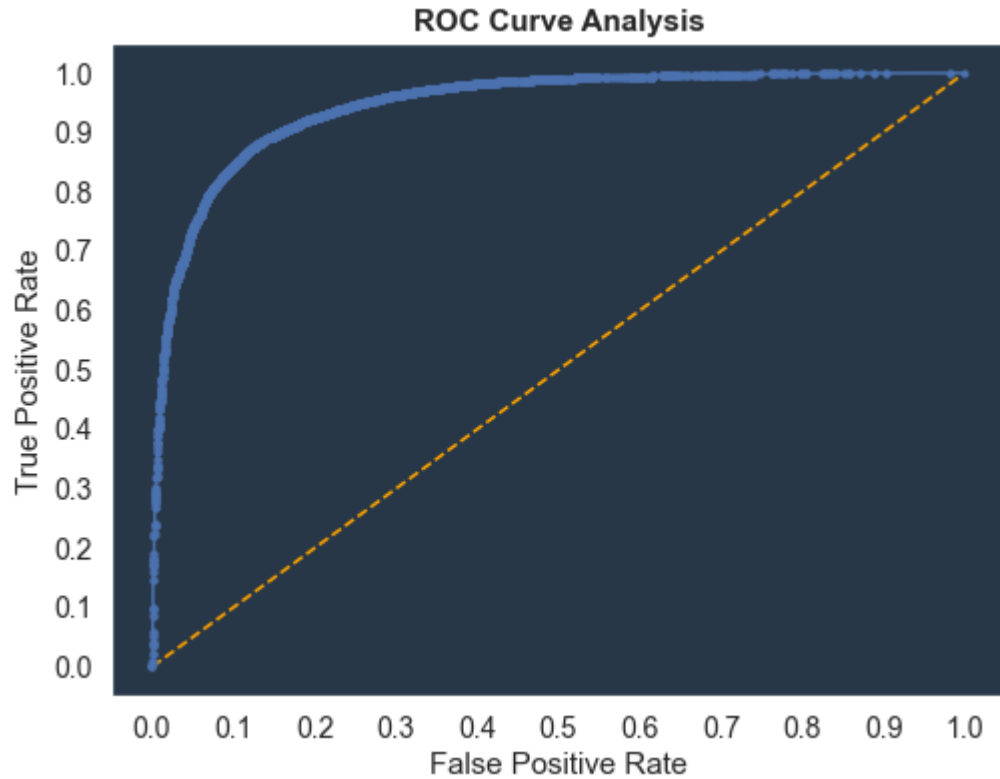
- ✓ N-Grams (n_gram_grange) : Uni, Bi, Tri
- ✓ Feature Selection(max_feature)
- ✓ Word Frequency Exclusion (min/max_df)

Model Evaluation



Measures	Score
Accuracy	90%
Recall	97%
F-1	94%
Precision	91%

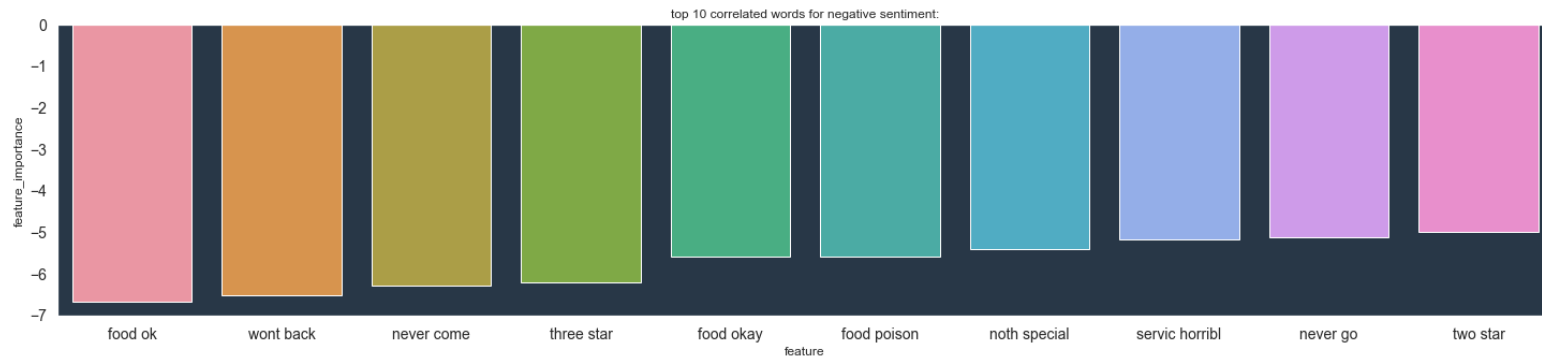
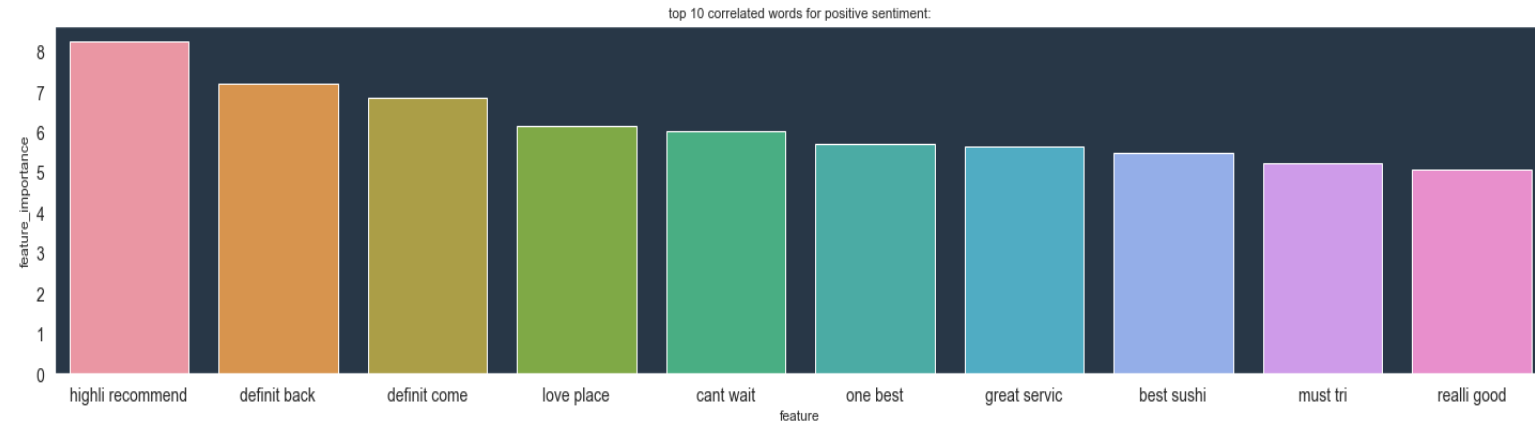
Logistic Regression



Evaluation Metrics

- **Accuracy:89.25%**
- **Recall:96.55 %**
- **F-1: 92.99 %**
- **Precision: 89.69 %**

Feature Importance



Conclusion

- Throughout this analysis I found that sentiment words were more positive than negative of Yelp users 'experiences at the designated restaurants. In relation to the positive sentiments of users, there is a positive correlation between that of positive reviews with high ratings, and negative reviews with lower ratings. By segmenting the area to Las Vegas and categorizing Japanese restaurants, i was able to gain insight on how they operate.
- Each of the individual users provide their opinions throughout their reviews; as the positive outweighs the negative, Japanese restaurants are providing great dining services and food to their customers, which increases the positivity disclosed in their review, as well as an input of a higher rating. Although, through this project, Japanese restaurant owners can also view aspects that drive more negative sentiments which they can take initiative on and remedy over time.

Recommendation

- Apply sarcasm and joke detection using NLP to separate negative and positive reviews
- There is an additional dataset from the Kaggle data competition that includes checkin, tip and user
 - Future work will focus on NLP of the other attributes in the business dataset and predicting overall star rating based on those attributes.