

Yelp Data Analysis

Chao Chang, Shuyang Chen, Kunning Wang, Youhui Ye

STAT 628

October 31, 2019

Contents

- 1 Data Overview
- 2 Data Cleaning
- 3 Preliminary Analysis
- 4 Future Work

Data Overview

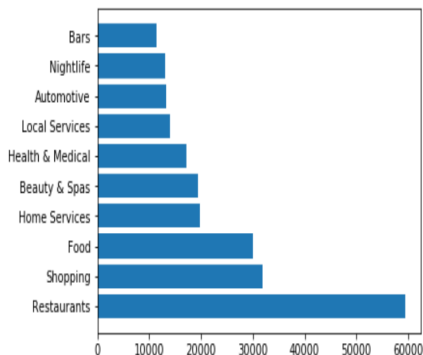


Figure: All Categories

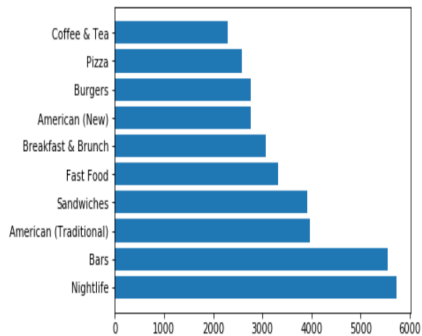
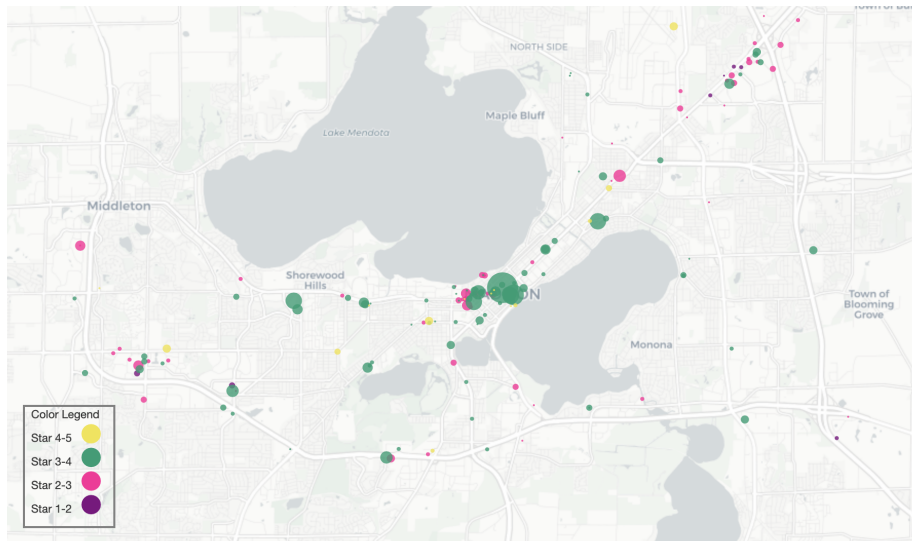


Figure: Restaurants

American Traditional Restaurants in Madison



Purpose

- Simplify and merge all data sets with interested features
- Convert a text to a vector
- Select recent reviews and interested features(name, review text, stars, locations and attributes)
- Deal with the review text

Dealing with the review text

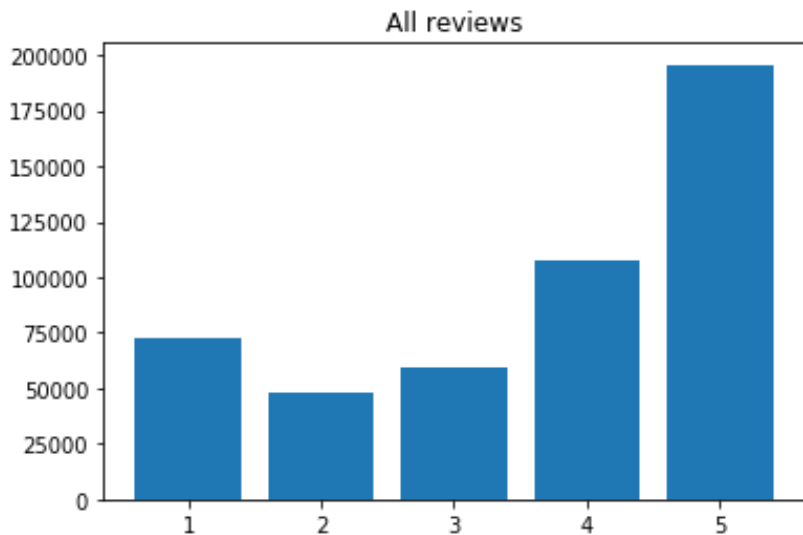
Original	Lowercase, tokenize and expand abbreviation	Remove stop words and add neg tags	Remove all punctuations
I don't like cheeseburger, though my kid does.	['i', 'not', 'like', 'cheeseburger', ,', 'though', 'my', 'kid', 'does', '.']	['not', 'neg_like', 'neg_cheeseburger', ,', 'though', 'kid', '.']	['not', 'neg_like', 'neg_cheeseburger', 'though', 'kid']

Table: Dealing Process

Chao C., Shuyang C., Kunming W., Youhui Y.



Distribution Plots



High frequency words

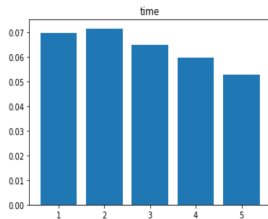


Figure: time

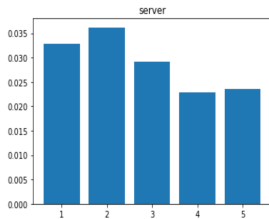


Figure: server

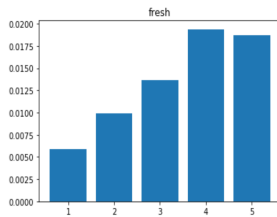


Figure: fresh

Different kinds of foods and drinks

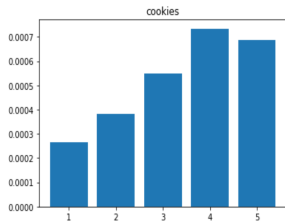


Figure: cookies

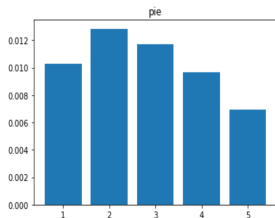


Figure: pies

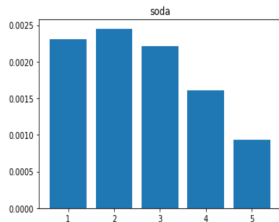


Figure: soda

Insights from the attributes

- **Common Attributes:** 10 most common attributes such as Goodforkids, Wifi, Alcohol.
- **Transform string label to integer label:**
Eg: Outdoorseating: True to 1, None to 0, False to -1
- **Suggestion based on linear regression:** Having outdoor seats, wifi and delivery contributes to higher stars while allowing dogs and being noisy affect in a negative way.

- **Extract noun phrases from the review text**
"Credit Card" is more informative than "credit" or "card".
- **do stem extracting and lemmatization:**
"cats" → "cat" and "loving" → "love"
- **Using tf-idf to select useful words**
- **Do linear regression and significance test**
- **Explore attributes and locations**
- **Give suggestions to American traditional restaurants**

Thank you