# Visualizing the Yelp Dataset

Ashutosh Verma, Karthik Balasubramanian, Md Khaled Hassan, Suneil Shrivatsav and Tianji Zhou
Department of Computer Science, Georgia Institute of Technology
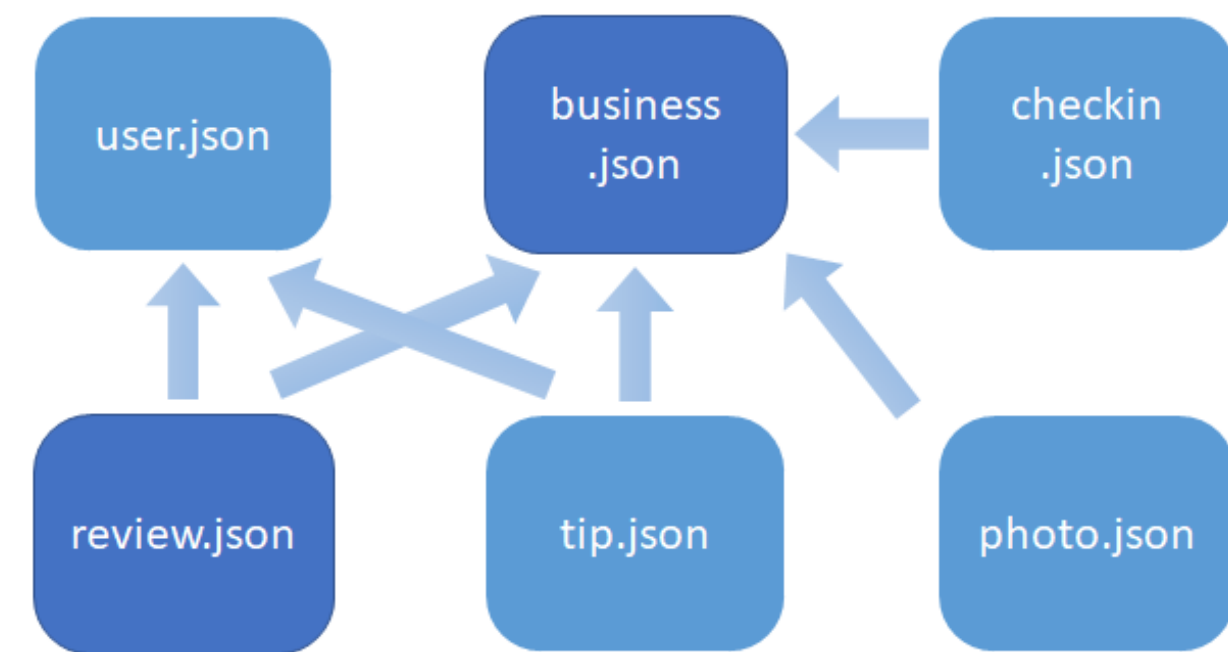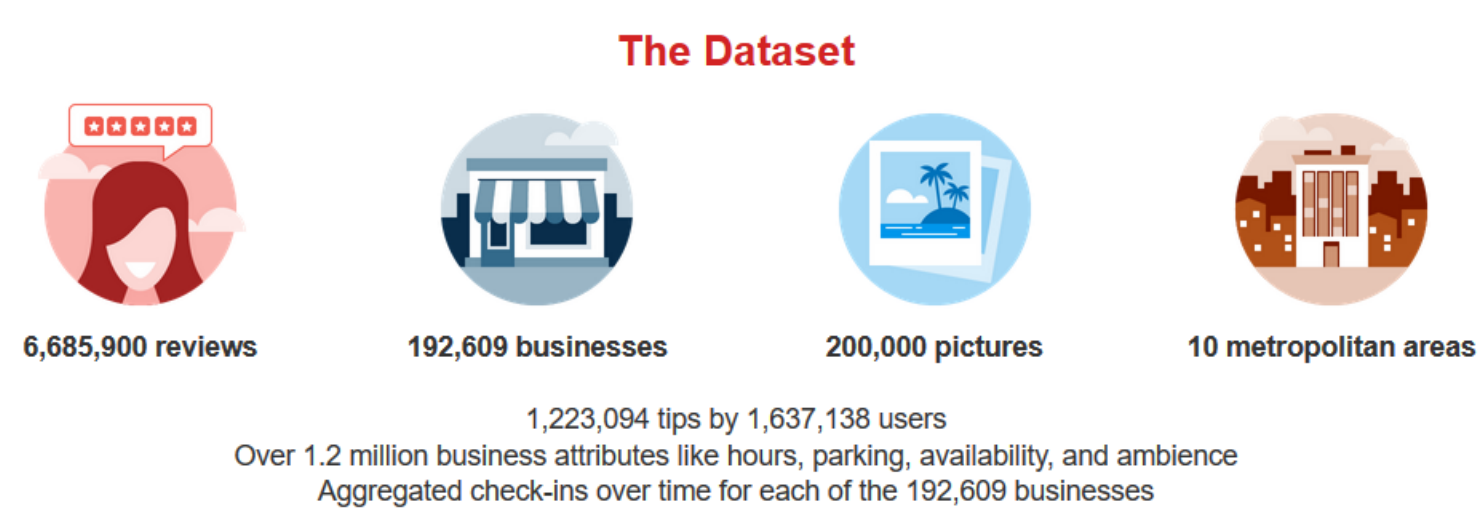
## I. Introduction

**Problem Statement:**
The aim of this project is to use very large dataset from Yelp(~8.6GB) and provide meaningful visualizations/dashboard to business analysts, restaurant owners & economic policy makers. The dataset contains a multitude of indicators like user reviews, ratings, co-ordinates, review counts, categories etc. which we intend to correlate and produce such visualizations. We also intend to use the poverty/population dataset per county provided in homework-2 for such visualizations.

**Motivations:**
Our study will provide detailed insights for restaurant owners, customers, suppliers, economists and business analysts.

## II. Dataset

The Yelp Open Dataset (https://www.yelp.com/dataset) is a subset of Yelp's business, reviews and users data.
It includes 192609 business in 10 metropolitan areas. The dataset contains a 8GB core data set and a supplementary data set for user uploaded images.

The Dataset

6,685,900 reviews    192,609 businesses    200,000 pictures    10 metropolitan areas

Over 1.2 million business attributes like hours, parking, availability, and ambience
Aggregated check-ins over time for each of the 192,609 businesses

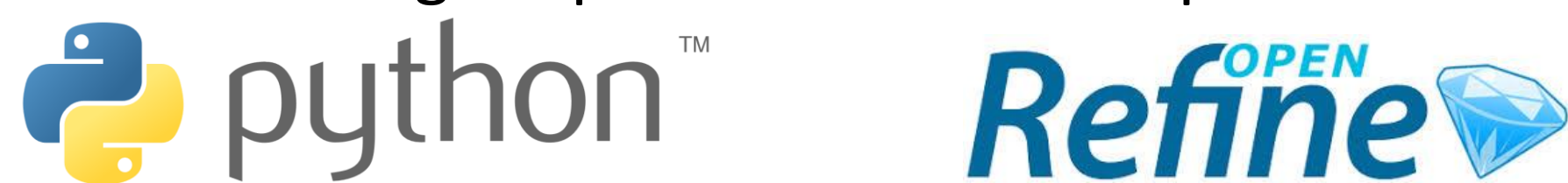1,223,094 tips by 1,637,138 users

The core data set is a combination of six different json files: business, review, user, tip, check-in and photo. Each item in the business, user and review data files is given an ID, by which these data files can be connected.

To investigate the relation between Yelp dataset business with county economy and population, we will also explore two other data file from the United States Department of Agriculture (https://www.ers.usda.gov/), county_poverty.csv and county_detail.csv.

## III. Approaches

The proposed methodology broadly consists of three parts:
(1). **Data cleaning and refinement**: Data present in Business.json was not directly usable in Tableau as it required some Data cleanup and Data refining. We use Python and Google OpenRefine to cleanup and filter information.

(2). **Data visualization and analysis**: The current data visualization is looking at the business data table. We use Tableau too visualize these on a map to allow more intuitive understanding of food diversity data and its intersection with poverty.
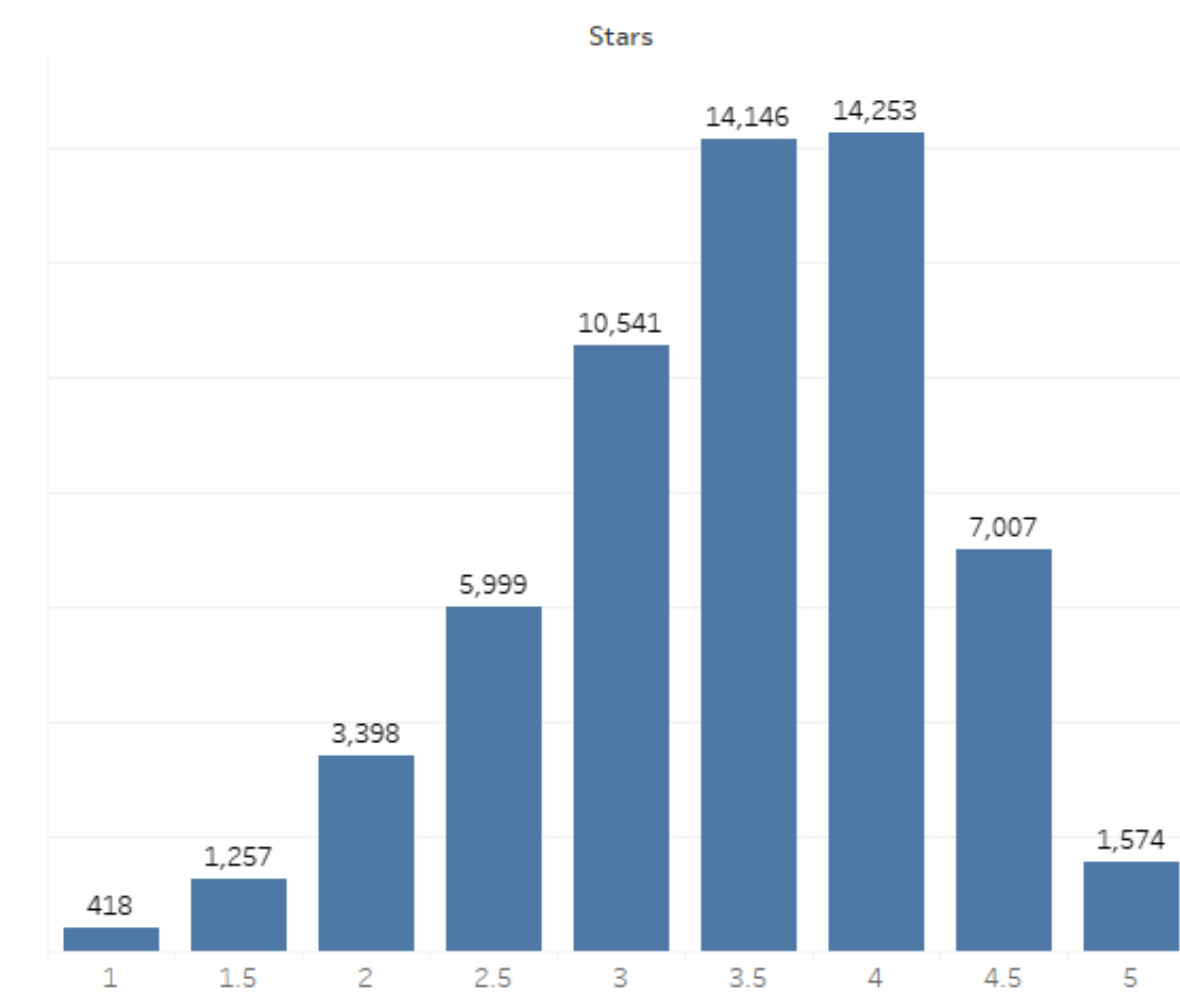
(3). **Machine learning and text mining**: The machine learning and text mining analysis are done using Python and scikit-learn package, for the purpose of predicting ratings or sentiment analysis based on the given attributes.
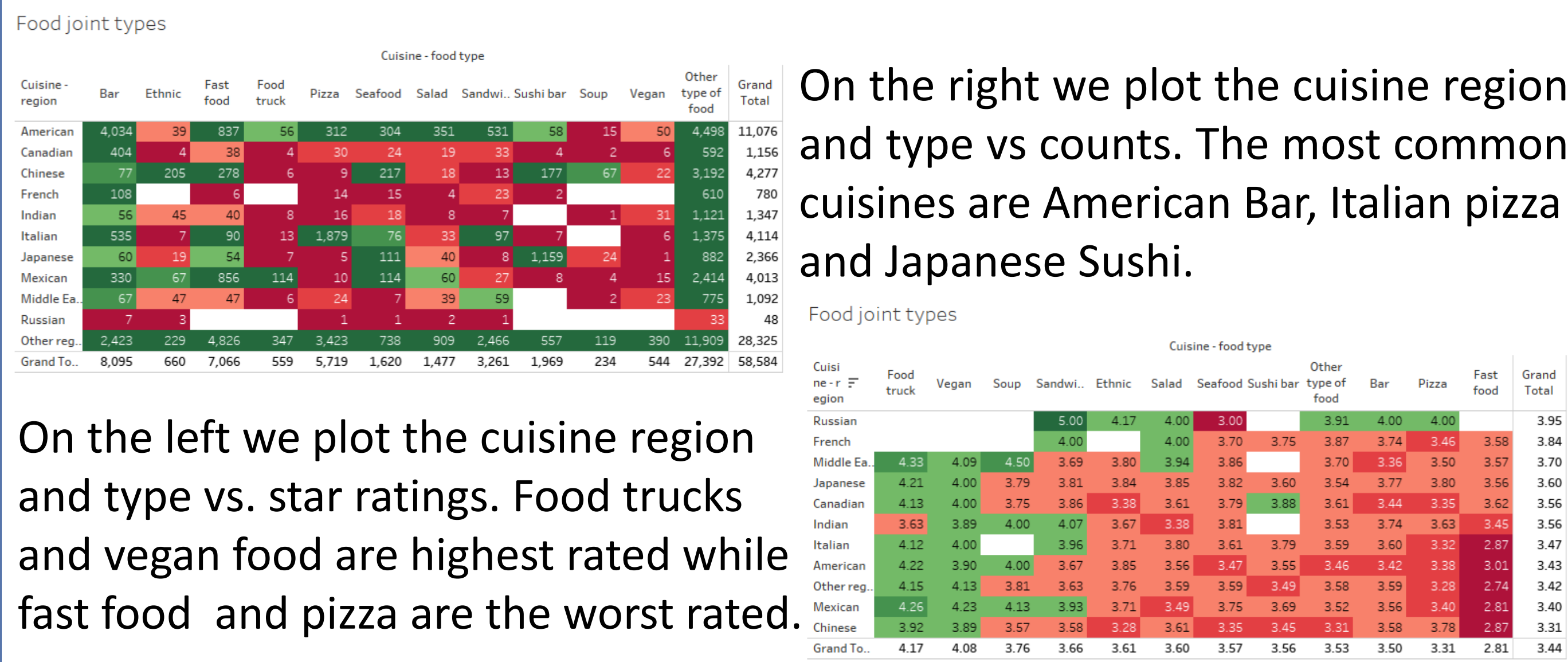
**List of Innovations:**
(1). Exploring Yelp's data through a lens of food diversity and poverty.
(2). Geographical data blend.

## IV. Business Data Analysis

We first look at the distribution of rating: how many Business ids receiving a particular star rating. The plot shows that most of the ratings are centered around 3-4.5 stars on a scale of 5. However, the number of ratings of 1-2 stars and 5 stars are low indicating a skewed dataset. The dataset shows an average rating of 3.44.
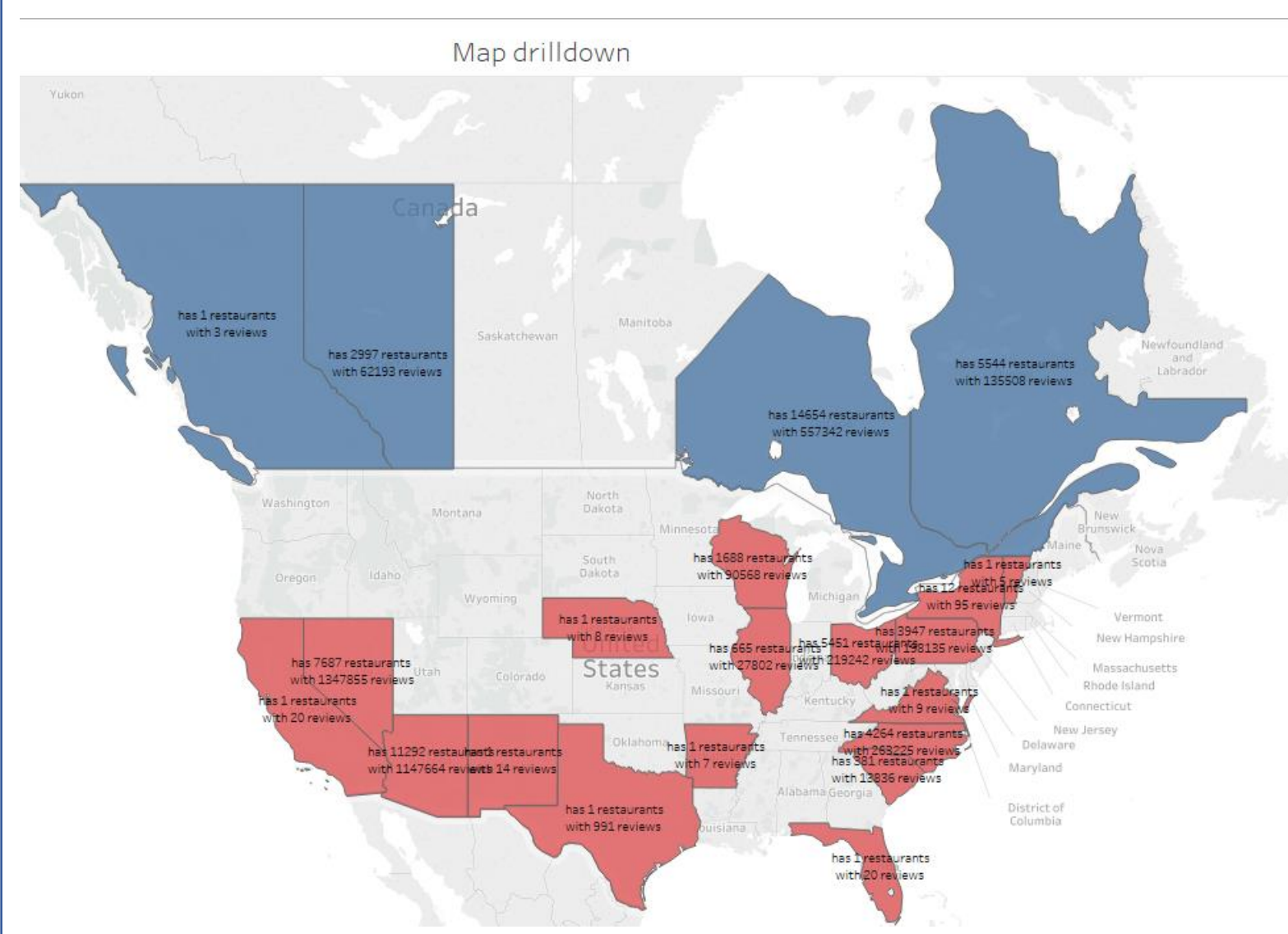
We then look at the restaurant types, especially the cuisine type and region, and plot them vs. the restaurant type counts and star ratings in heat maps.

On the right we plot the cuisine region and type vs counts. The most common cuisines are American Bar, Italian pizza and Japanese Sushi.

On the left we plot the cuisine region and type vs. star ratings. Food trucks and vegan food are highest rated while fast food and pizza are the worst rated.
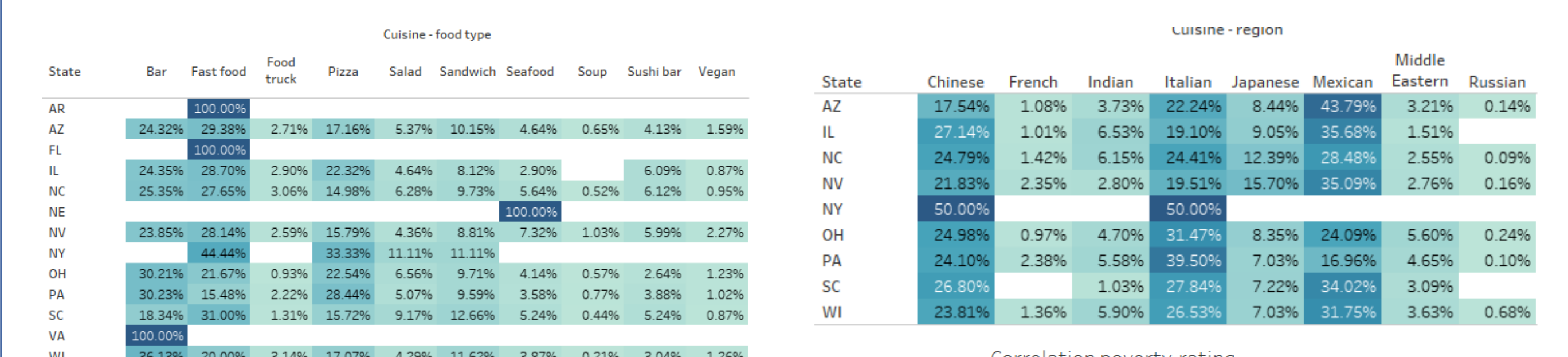
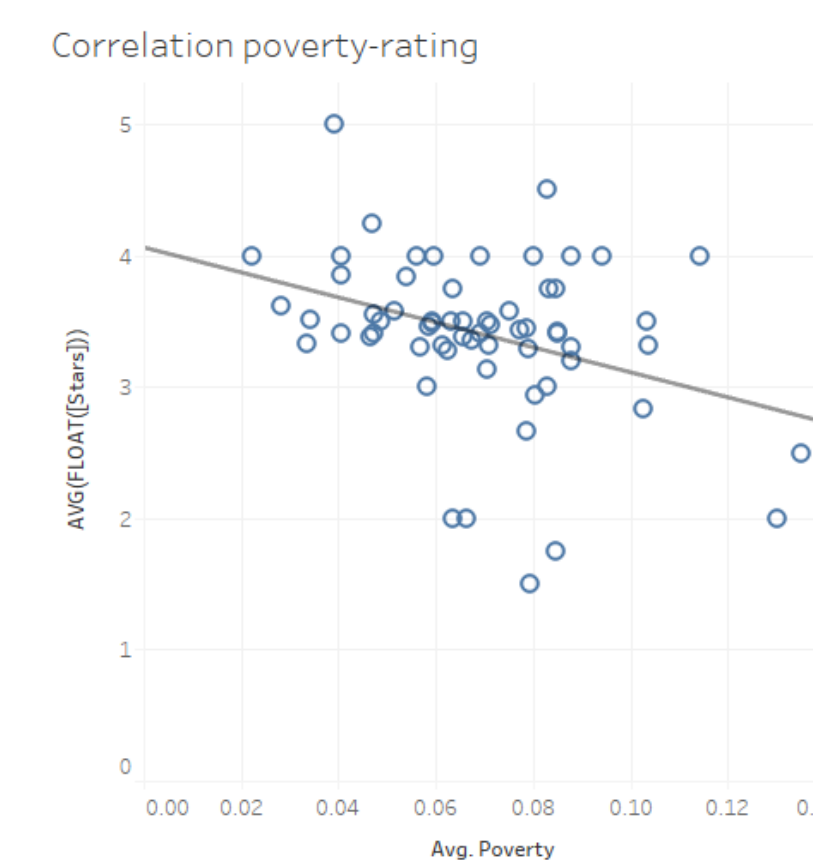## V. Geographical Data and County Economy

On the left is a choropleth map showing the number of restaurants in the different states which gives us an idea of the income levels which can be correlated to number of restaurants across the states.

The following two tables show type of restaurants in each state where data is available.
Arizona, Illinois, Ohio and Pennsylvania all have a robust combination of food by region and by cuisine type. On the other hand, New York is the state with the least food diversity in terms of region of origin of the cuisine.
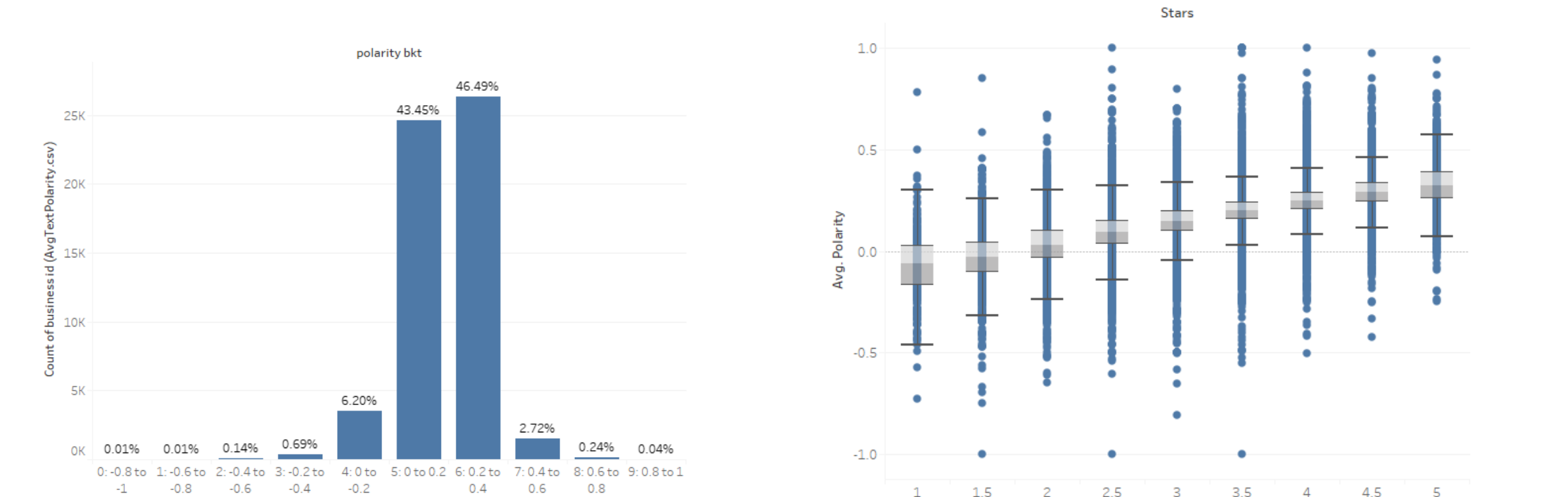
The plot on the right shows a statistical analysis of the dataset fitting the poverty level and average rating of the restaurants. The plot indicates that with increasing poverty, the quality of the restaurants goes down.

## VI. Text Mining: Review Text Sentiment Polarity

We carried out a sentiment analysis on more than 1.8 million reviews using TextBlob library of python. We extracted the polarity of each of the reviews in a range of [-1.0 1.0], where, -1 being an extremely negative and +1 being an extremely positive review.
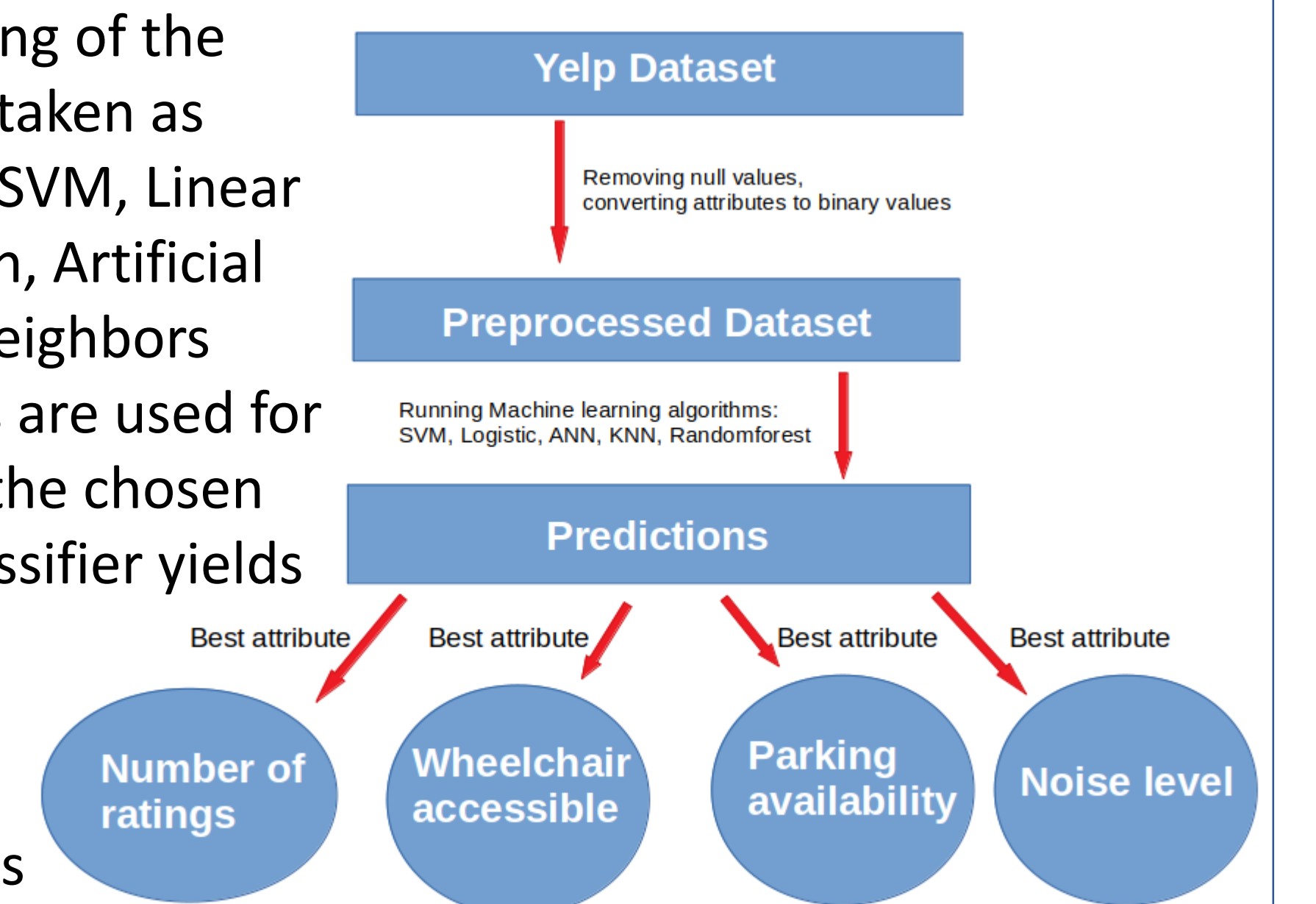
On the left is the distribution of the polarity of the reviews given to restaurants. Most reviews are in the range of [0,0.4]. On the right we explore the correlation between polarity and star ratings. It appears the higher star ratings have slightly higher polarity.

## VII. Machine Learning: Best Features Extraction

The business dataset consisting of the entries for the restaurants is taken as the dataset for explorations. SVM, Linear regression, Logistic regression, Artificial neural networks, k-nearest neighbors and RandomForest classifiers are used for predicting the ratings. From the chosen classifiers, RandomForest classifier yields the best results.

The best features which contribute most to the ratings include "Number of ratings", "Wheelchair accessible", "Availability of parking" and "Noise level".

Yelp Dataset → Removing null values, converting attributes to binary values → Preprocessed Dataset → Running Machine learning algorithms: SVM, Logistic, ANN, KNN, Randomforest → Predictions

Best attribute: Number of ratings    Best attribute: Wheelchair accessible    Best attribute: Parking availability    Best attribute: Noise level

## VIII. Conclusions

The key points that we could analyze from these visualizations are:
(1). High poverty areas have lower rating restaurant.
(2). A few businesses are predominant in a specific area.
(3). Food Trucks and Vegan Food are the highest rated ones while Fast Food and Pizza places are the worst rated.
(4). The sentiment polarity is well correlated with average review rating of the restaurant.
Based on all the above data points a business analyst or an investor could make a decision on which sector he should invest in and where for a successful business .

### References

[1]. Reviews, Reputation, and Revenue: The Case of Yelp.com, L. Michael, Harvard Business School NOM Unit Working Paper No. 12-016 (2016).
[2] Fake It Till You Make It: Reputation, Competition, and Yelp Review Fraud, Luca, Michael & Zervas, Georgios. (2013). SSRN Electronic Journal. 10.2139/ssrn.2293164.
[3]. Improving Restaurants by Extracting Subtopics from Yelp Reviews, Huang, J.; Rogers, S.; Joo, E. Improving Restaurants by Extracting Subtopics from Yelp Reviews. in iConference Social Media Expo (2014).
[4]. Restaurant Setup Business Analysis Using Yelp Dataset , Sindhu B Hegde, Supriya Satyappanavar, Shankar Setty, Advances in Computing Communications and Informatics (ICACCI) 2018 International Conference on, pp. 1455-1462, (2018).
[5]. The Geography of Taste: Using Yelp to Study Urban Culture, Rahimi, S.; Mottahedi, S.; Liu, X. ISPRS Int. J. Geo-Inf. 7, 376(2018).
[6]. Spatial analysis of users-generated ratings of yelp venues, Sun, Y. & Paule, J.D.G. Open geospatial data, softw. stand. https://doi.org/10.1186/s40965-017-0020-9 (2017).
[7]. Nowcasting the Local Economy: Using Yelp Data to Measure Economic Activity E. L. Glaeser, H. Kim, M. Luca, National Bureau of Economic Research,Working Paper No. 24010 (2017).
[8]. Identifying Restaurant Features via Sentiment Analysis on Yelp Reviews, Yu B., Zhou J., Zhang Y., Cao Y., arXiv e-prints, arXiv:1709.08698 (2017).
[9-15]. Learning from data. Chapters 1-9 Abu-Mostafa, Yaser S., Malik Magdon-Ismail, and Hsuan-Tien Lin. New York, NY, USA:: AMLBook, 2012.
[16]. Restaurants Review Star Prediction for Yelp Dataset, M. Yu, M. Xue, W. Ouyang, pdfs.semanticscholar.org (2015).