



By: Brian Jankowitz

Yelp Data to Predict Affluency



Brian Jankowitz

Data Analyst

- Presented in front of audiences 300+ people
- Background data analyzation

Agenda

1. What is Yelp?
2. Problem Statement
3. Data Gathering
4. Analyzing Data
5. Modeling
6. Model Selection
7. Conclusion
8. Recommendations



What is Yelp?

- A place for business reviews
- Crowd-sourced



See All 420

Camellia Claimed

4.5 stars - 140 reviews

\$\$ · Ramen, Gelato

[Write a Review](#) [Add Photo](#) [Share](#) [Save](#)

Popular Dishes

Salt and Pepper Fried ... 13 Photos · 14 Reviews

Spicy Miso Ramen 23 Photos · 14 Reviews

Fried Crispy Pork Belly 10 Photos · 12 Reviews

Order Food

Delivery Takeout

Delivery Address

Enter delivery address

Start Order

camelliany.com

(212) 228-2070

[Get Directions](#)

[Message the Business](#)

Problem Statement

- How can we target people in non affluent area to help provide them health insurance
 - Find non-affluent area in a non-traditional way
 - Use of big data related to commercial activity and cost of product and services as an indicator for affluency.
 - Increase healthcare coverage
 - Market to those people
- Facts (US Census Bureau)
 - In 2018, 28 million people without health insurance
 - 72.2% of people 100% or more below that poverty line are uninsured

Data Gathering- Yelp

- Scrapped Yelp
- Gathered information on businesses in all 5 boroughs of New York City
- Was this information enough by itself? No
- Get more data
 - Round 1: 5000 places to only 1200 places
 - Round 2: 32000 places to only 12000 places

**TAKE
THE
NEXT
STEP...**

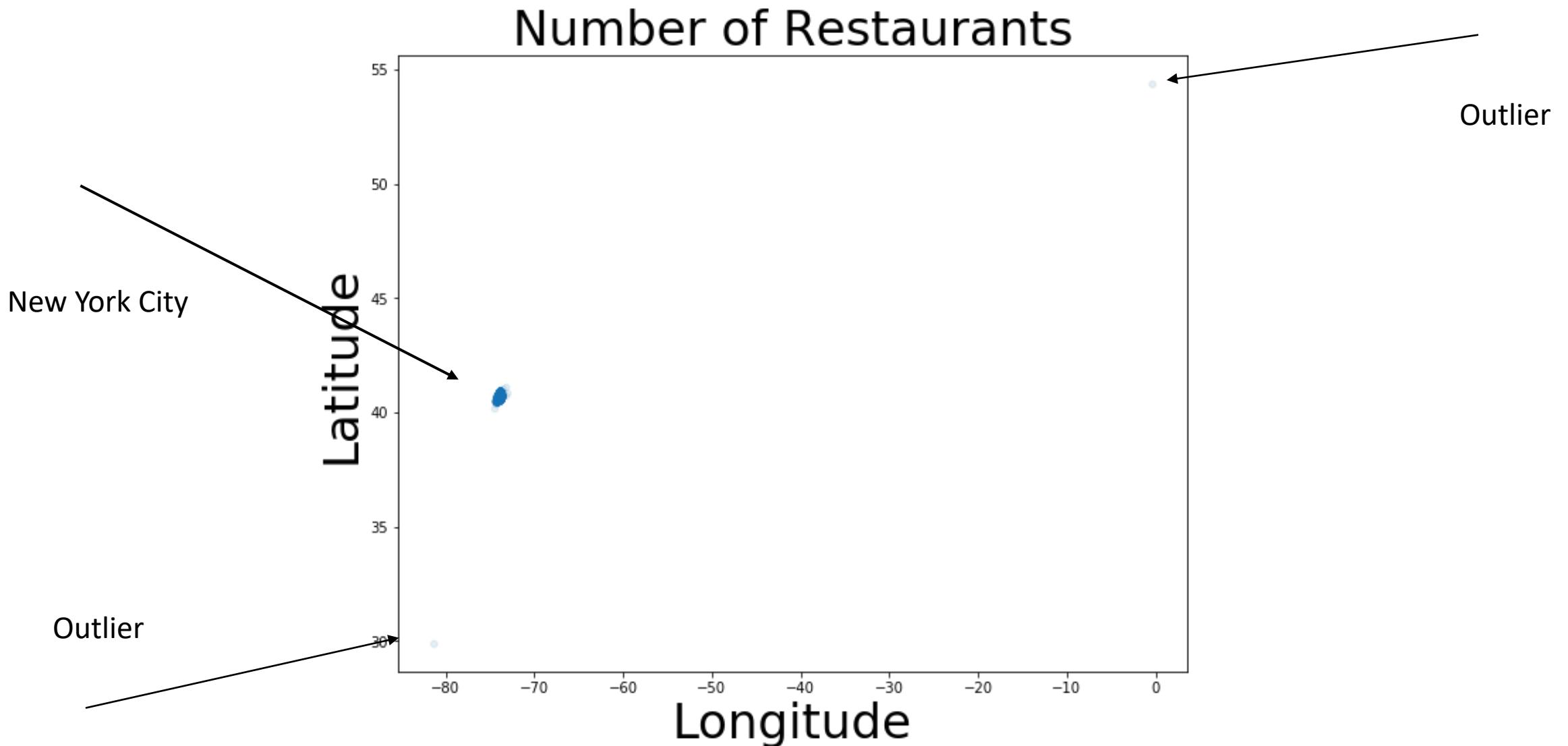


Data Gathering- Income

- IRS data
- Last year available 2017
- Number of returns and amount for AGI
- Why AGI and not just income?
 - To adjust for other areas
 - Each person pays different taxes
- Affluenct > 25=

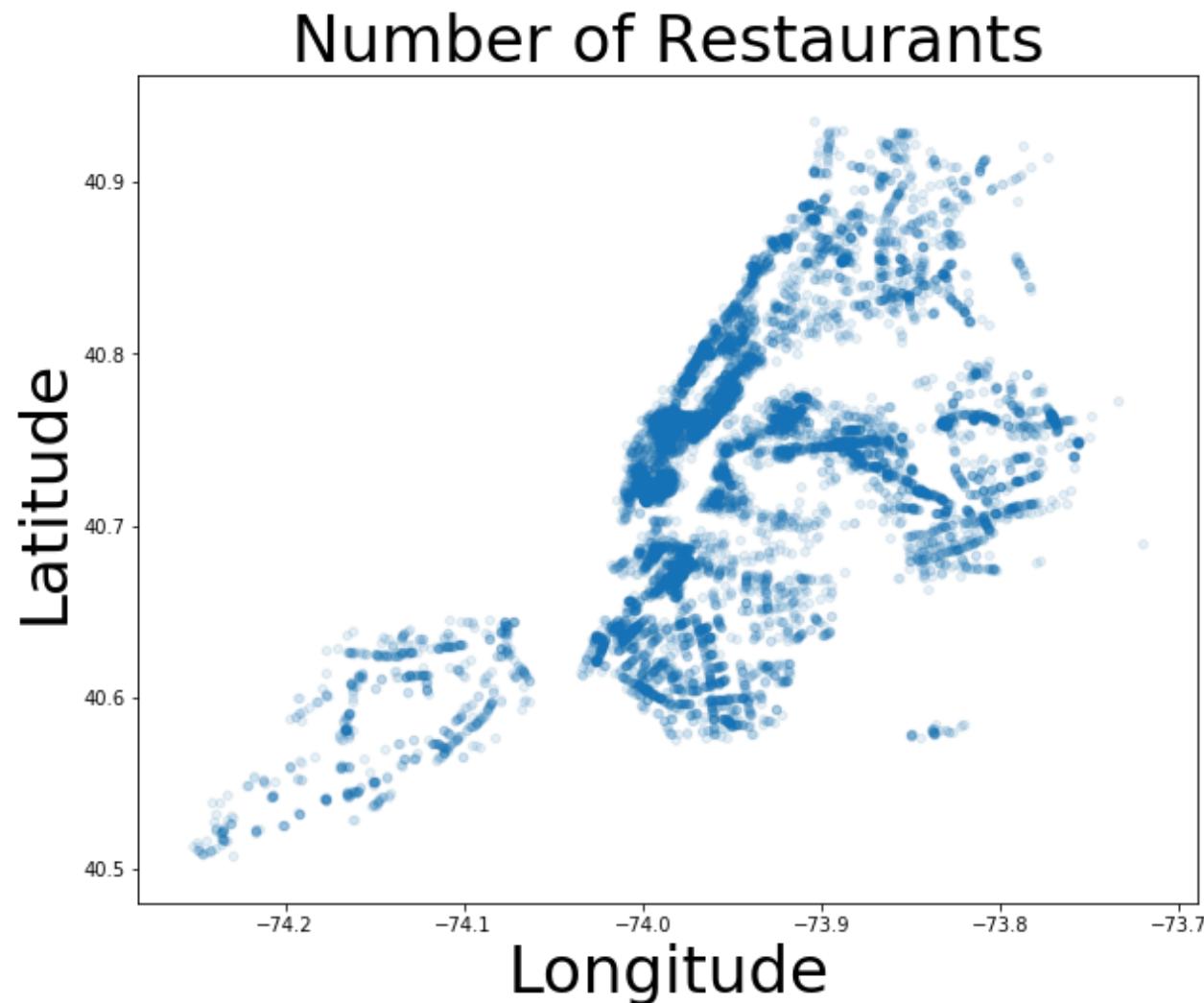


Problems with Data- Outliers



Coverage of New York City

- Dark vs light colors
- Business density

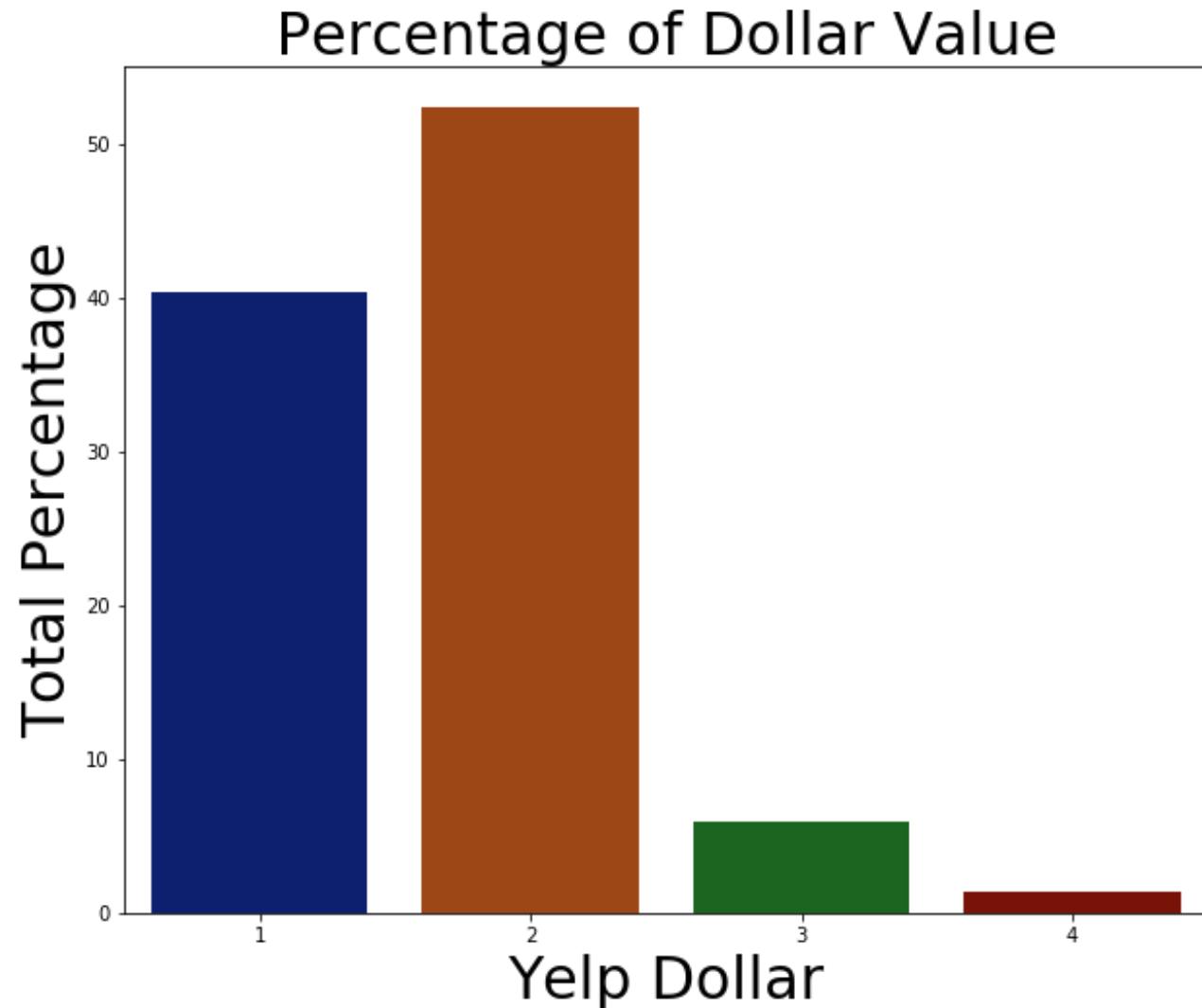


Borough	Population (2018)
The Bronx	1,432,132
Brooklyn	2,582,830
Manhattan	2,278,906
Queens	2,278,906
Staten Island	476,179

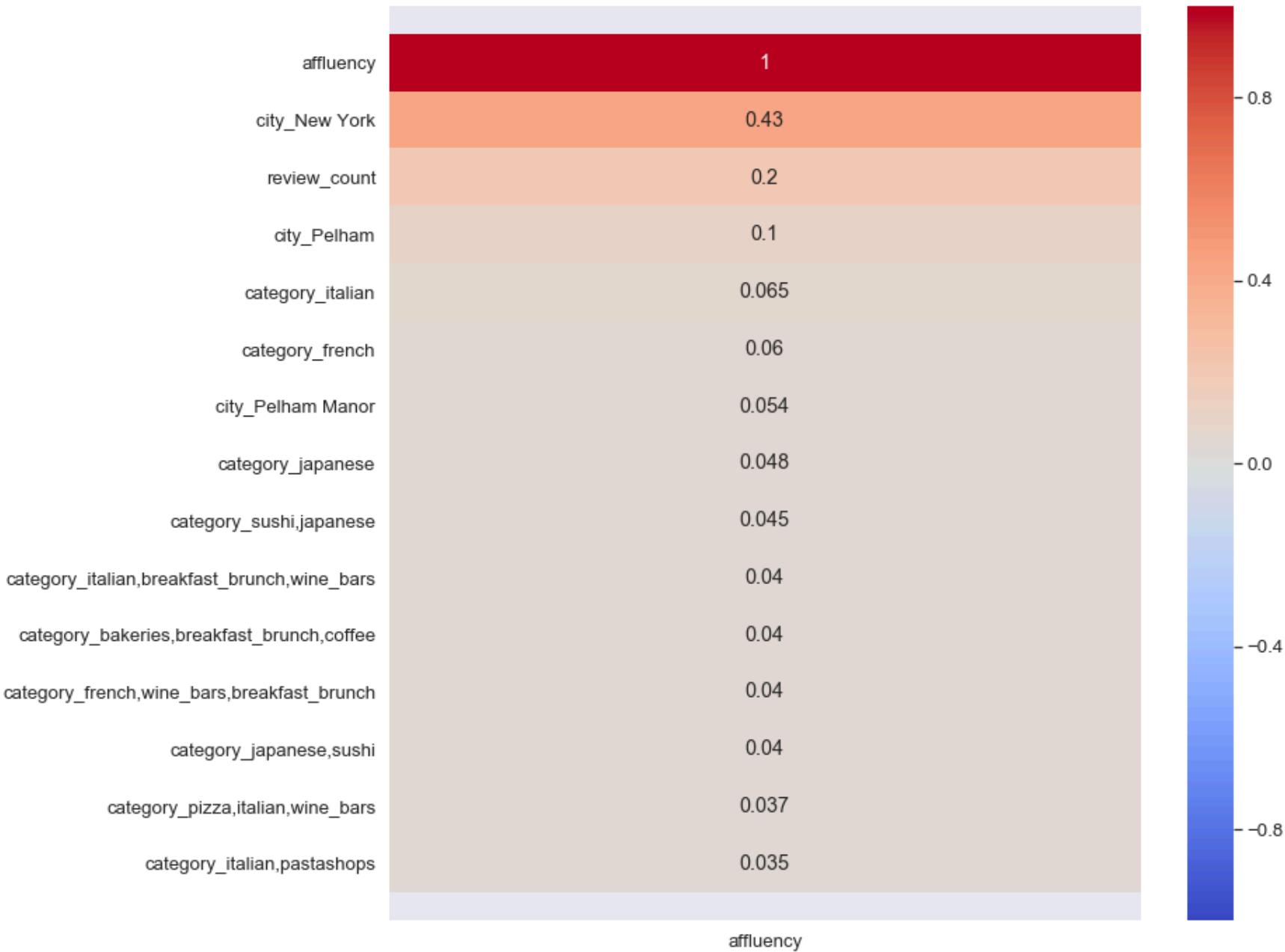


EDA

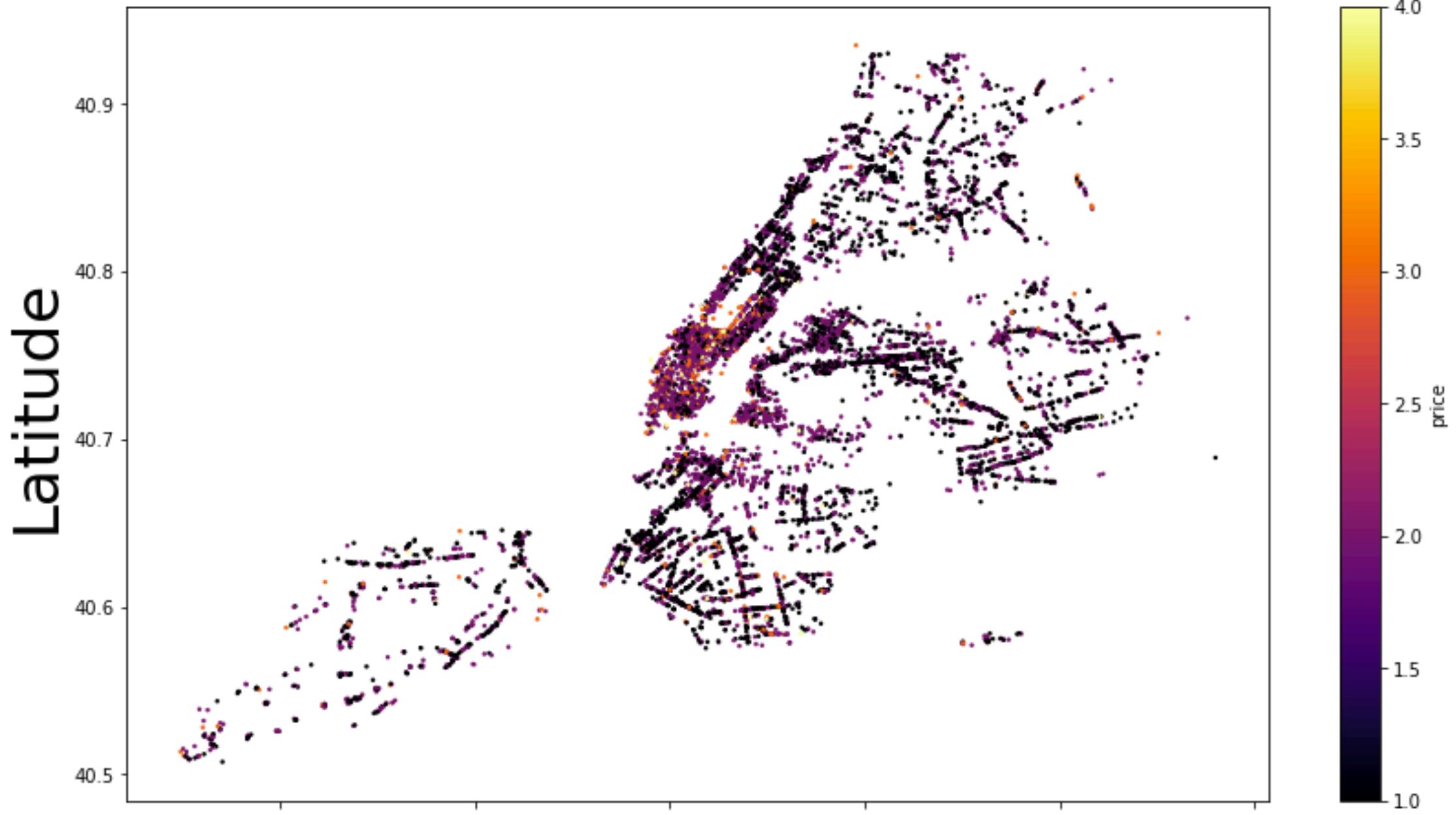
- 3(\$\$\$) and 4(\$\$\$\$) less than 10%
- 1(\$) and 2(\$\$) more than 90%



Correlation Matrix



Price of Restaurants

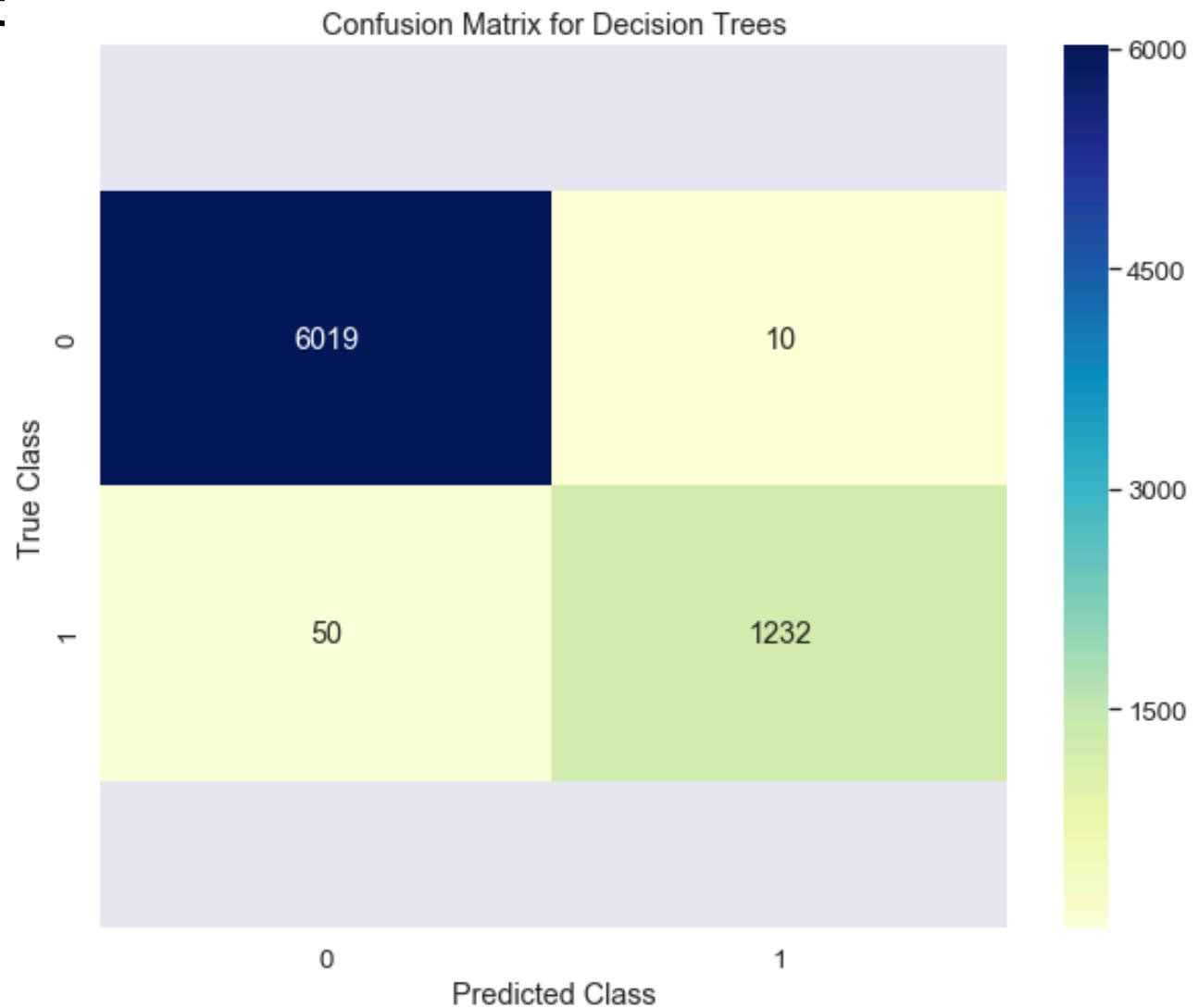


Modeling

Model	Train	Test	Sensitivity	Specificity
Logistic Regression	83.1%	83%	22%	96%
KNN	97.8%	82%	96.4%	96.1%
Decision Trees	100%	99.9%	21.9%	96.1%
Bagging	99.9%	99%	21.9%	96.1%
Random Forest	99.9%	99.9%	21.9%	96.1%

Baseline Train	Baseline Test
82.3%	82.3%

Model Selection



Conclusion

- We can predict the affluency of an area
 - High true positives
- There are some questions about the data
 - How accurate is this?
 - Tribeca
 - Richest zip code in New York City (source: Business Insider)
 - Don Bella Pizza rated \$
 - More research needs to be done

Recommendations

- More time should be allocated to this
- Look at more expensive areas
 - Times Square and World Trade Center
 - Plain slice of pizza \$2, Times Square \$4+
- Use Big data in other ways
 - Apartment prices
 - Express trains
 - How close to train (only certain areas)
- Use big group to collaborate, see if \$\$ can be used to supplement another model for predicting prices