# Coursera Data Science Final Project

## Density and Specializations of Boston's Restaurants

By Alexander Hahn - 8/1/2020

# Introduction

- We want to discover a good place to build a new Mexican Restaurant in Boston

- We want to find a place with high density of restaurants but a low density of Mexican restaurants
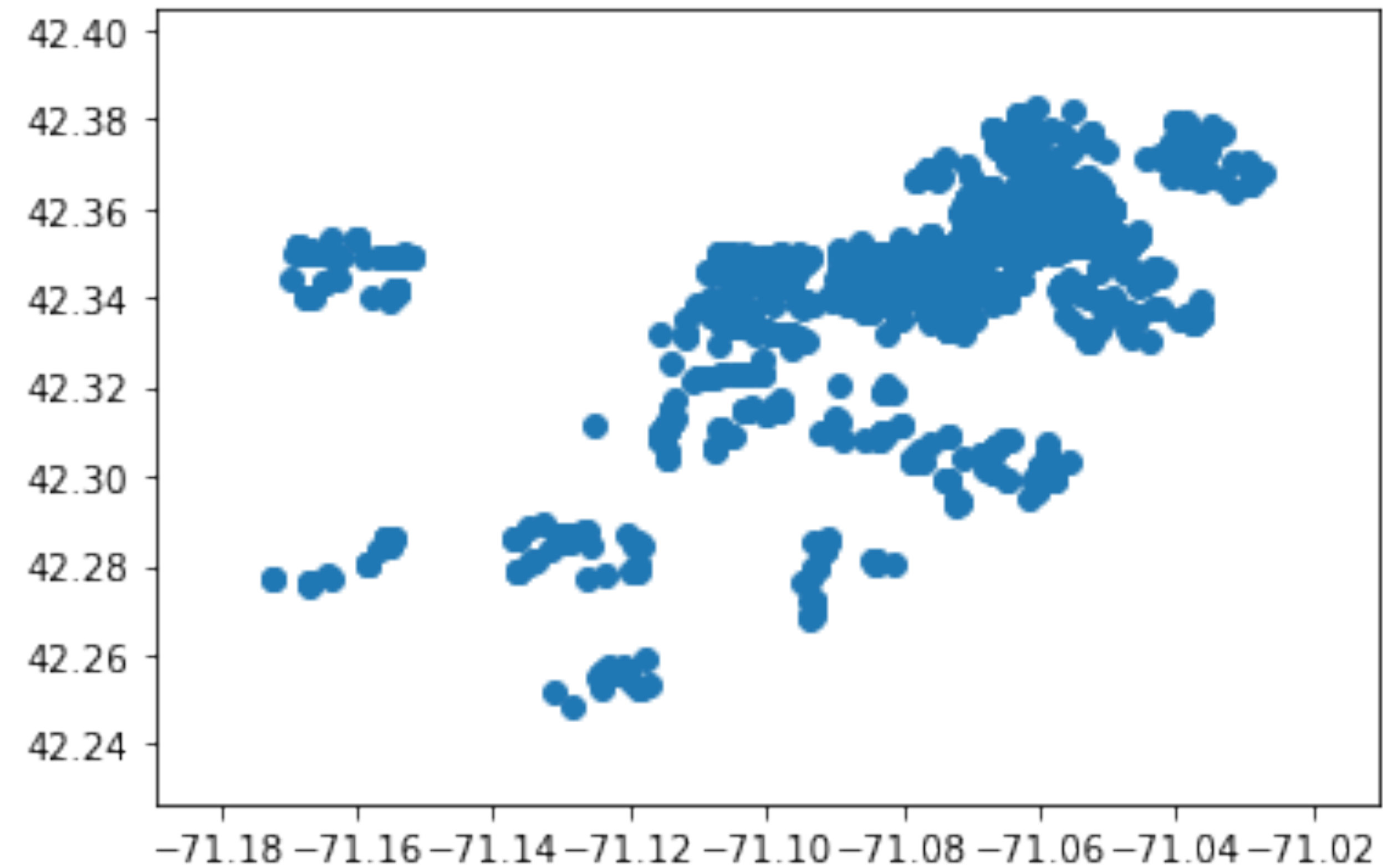
# Data

## Foursquare and District Data

| African | Thai | German | Mediterranean |
|---|---|---|---|
| **American** | BBQ | Greek | Middle eastern |
| **Asian** | Burger | Indian | Pizza |
| **Chinese** | Burger | Italian | Southern |
| **Japanese** | Fish | Latin | Spanish |
| **Korean** | French | Mexican | Steak |

- We pulled data from an external CSV with Boston District Data

- We use Foursquare to find information on the various venues. Venue name, Restaurant Category and Longitude/Latitude will all be collected.
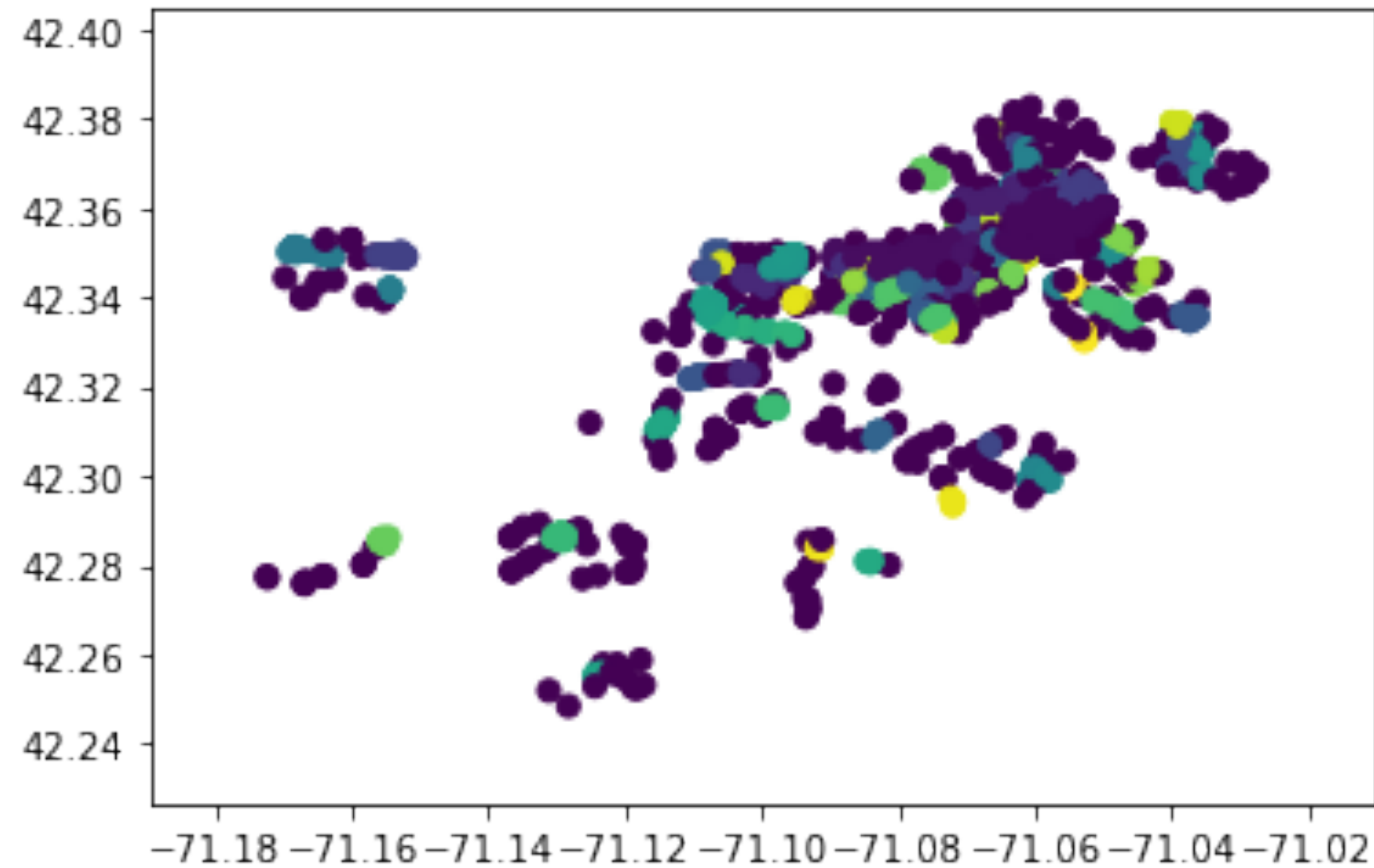
# Exploratory Data Analysis

## Foursquare location data

- We pulled a total amount of 57200 venue's from foursquare with a net amount of 3956 venues loaded in

- The graph on the right shows the information loaded in from foursquare showing the X/Y location data from each of the venues.

# Predictive Modeling

## Using DB Scan clustering



- We used DB Scan to auto-generate clusters finding 92 unique clusters using an epsilon of 0.001 and a min cluster value of 5

- Used the formula show below to calculate the value level of each district based on total number of venues and the number of unique instances in each cluster

- Value = (# of Restaurants in Cluster/Total Restaurants) * (# of Unique Restaurants/ Total # of Restaurants in cluster)

# Conclusion

- The top candidate we found from the analysis was cluster 1 with a value of 0.0195.

- The next three candidates we see a sharp drop in value but still look promising.

- The 1st cluster has a Longitude and Latitude value of -71.059 and 42.356 respectively.

| Cluster | Total Value | Cluster | Total Value |
|--------:|-------------|--------:|-------------|
| 1 | 0.019473950429944400 | 6 | 0.000505816894 2842690 |
| 2 | 0.003287809812847 7500 | 12 | 0.000505816894 2842690 |
| 32 | 0.003034901365705 6100 | 0 | 0.000505816894 2842690 |
| 8 | 0.002529084471421350 | 7 | 0.000505816894 2842690 |
| 16 | 0.001517450682852 8100 | 73 | 0.000505816894 2842690 |
| 47 | 0.001011633788568 5400 | 10 | 0.000505816894 2842690 |
| 20 | 0.001011633788568 5400 | 55 | 0.000252908447 14213500 |
| 3 | 0.000758725341426 4040 | 58 | 0.000252908447 14213500 |
| 33 | 0.000758725341426 4040 | 64 | 0.000252908447 14213500 |
| 57 | 0.000505816894 2842690 | 22 | 0.000252908447 14213500 |
| 52 | 0.000505816894 2842690 | 21 | 0.000252908447 14213500 |
| 70 | 0.000505816894 2842690 | 75 | 0.000252908447 14213500 |

# Future Directions

- Increased the number of samples pulled

- Introduce additional data based on population of the districts and density of people in the area.

- Could pull data in about affluence of each region to see if the area we are looking at could sustain the expense level of the restaurant we have in mind.

- Preform additional analysis combine the cluster data we have above with the sources listed above.