# Estimating Neighborhood Affluence with Yelp Data

By: Bernard Kurka, Thomas Ludlow, Brittany Allen

NYC DSI 6 Project 4

1/18/2019

# **Problem Overview**

◉  How can we develop a tool leveraging more agile data sources (like Yelp) that will help us measure local economic activity?

# Existing Methods

◉ Traditional wealth estimation is based on demographic characteristics (e.g. income or county business activity (CBP))

◉ Reporting lag with government surveys and statistics

yelp.com homepage

- Yelp's filters and results page

- Using Yelp's Fusion API we collected the first 100 "best match", results by all categories

# Data Snapshot

| zipcode | pr_1 ($) | pr_2 ($$) | pr_3 ($$$) | pr_4 ($$$$) |
|---|---|---|---|---|
| 10179 | 15 | 38 | 28 | 13 |
| 10012 | 31 | 63 | 0 | 0 |
| 10019 | 13 | 54 | 17 | 8 |
| 11235 | 28 | 57 | 7 | 4 |

# Feature Engineering

- Weighted price counts
  - i.e. for $pr\_2$=6 ⟹ $pr\_2w$=12

- Sum of all weighted price counts
  - $pr\_totw = pr\_1w + pr\_2w + pr\_3w + pr\_4w$

- Price and Review counts standardized

# Grid Search Results

- K Means Clustering
  - *inits = ["k-means", "random"],*
  - *n_init = range(10,20)*
  - *n_clusters = range(4,10)*

- Agglomerative Clustering

- Hierarchical Clustering
  - *linkage_method = ['complete', 'centroid'..]*
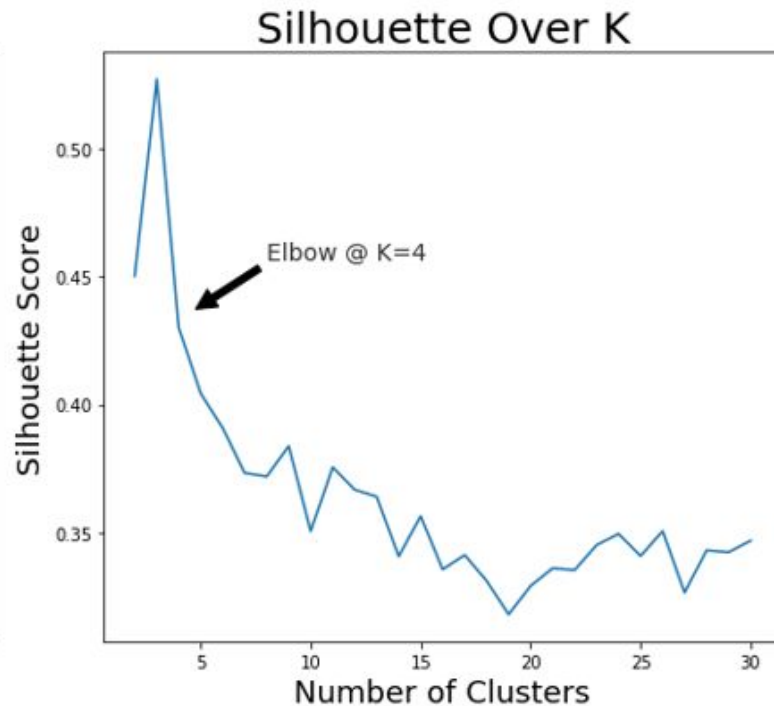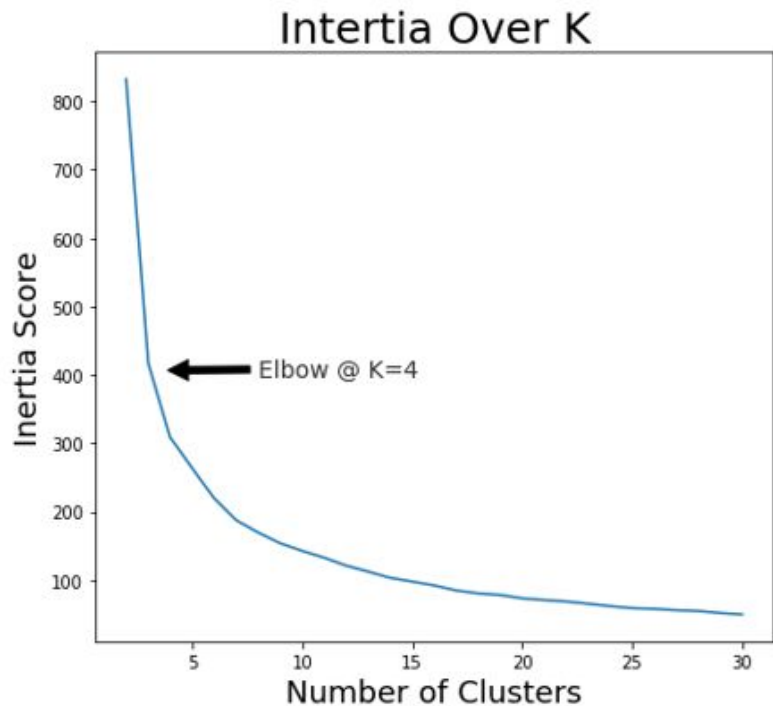  - *affinity = ['euclidean', 'l1', 'manhattan',...]*

# **Clustering Algorithms**

- Interpreting the grid results
  - Hierarchy – high silhouette score, low clusters counts
  - Agglomerative – good cluster counts, low silhouette score
  - K Means – balanced results

| Data_frame | model | inertia | silhouette | Numb_clusters | Cluster_counts |
|---|---|---|---|---|---|
| All_features | hierarchy | 0.000000 | 0.532652 | 4 | [270, 4, 3, 1] |
| All_features | Agglomerative | 0.000000 | 0.383370 | 4 | [100, 81, 68, 29] |
| scaled_pr_mult_wtot | kmeans | 308.753798 | 0.430094 | 4 | [82, 80, 71, 45] |

# 📌 K Means Best K



Intertia Over K — Inertia Score vs Number of Clusters, with arrow labeled "Elbow @ K=4"

Silhouette Over K — Silhouette Score vs Number of Clusters, with arrow labeled "Elbow @ K=4"

# 📌 Final Model

## K Means Clustering

| Features<br>*variable names: s for Standarized, w for weighted, tot for total* | | | | |
|---|---|---|---|---|
| pr_1s | pr_2ws | pr_3ws | pr_4ws | pr_totws |
| Parameters | | | | |
| n_clusters=4 | algorithm="auto" | init='random' | random_state=42 | |

# Price Distribution: $, $$

- ◉ $ (red)
  - ○ Limited activity across all ranges ($, $$, $$$, $$$$)

- ◉ $$ (orange)
  - ○ Moderate activity in $ and $$
  - ○ Limited activity in $$$ and $$$$

13

Brighton Beach Ny: $$$ (top left, green histogram) — Yelp Price Ratings vs. Number of Businesses in Location

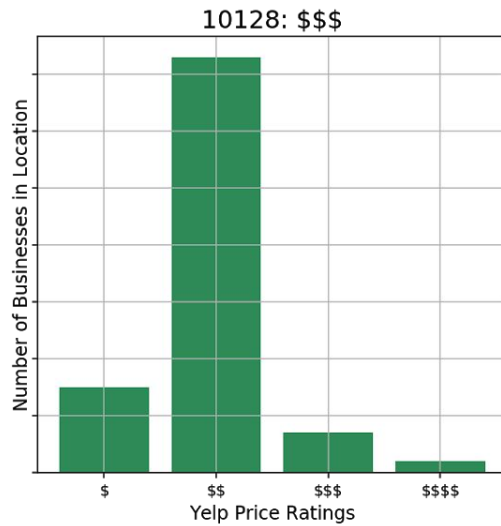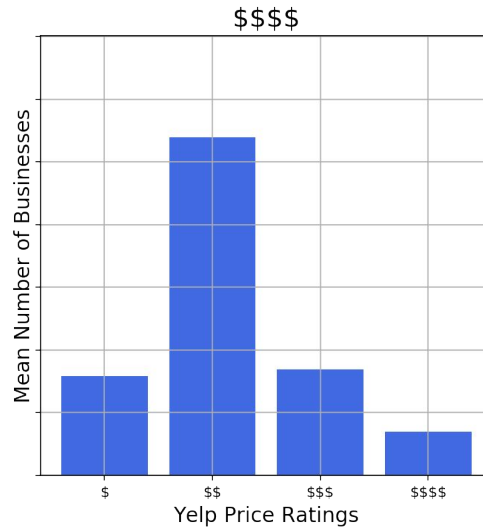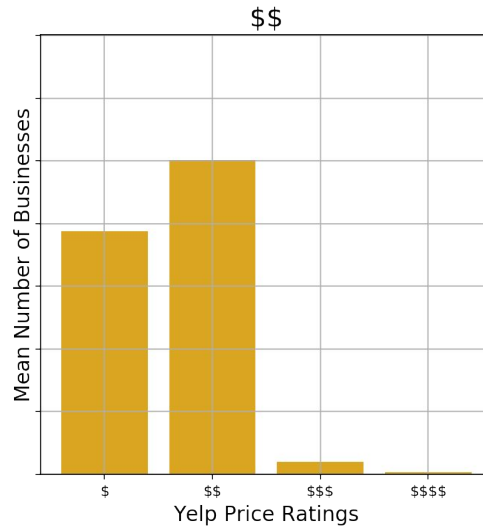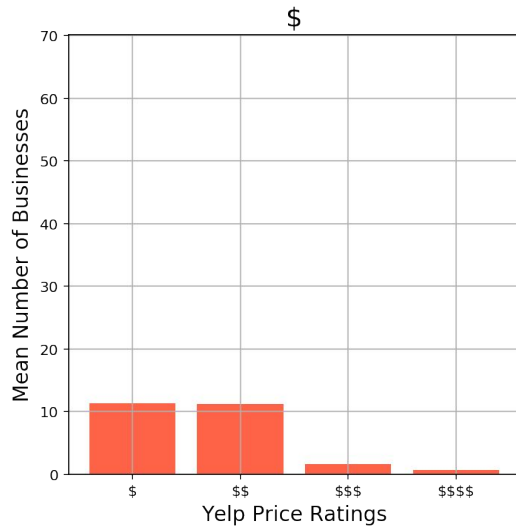10128: $$$ (top right, green histogram) — Yelp Price Ratings vs. Number of Businesses in Location

Flatiron: $$$$ (bottom left, blue histogram) — Yelp Price Ratings vs. Number of Businesses in Location

Rector Park: $$$$ (bottom right, blue histogram) — Yelp Price Ratings vs. Number of Businesses in Location
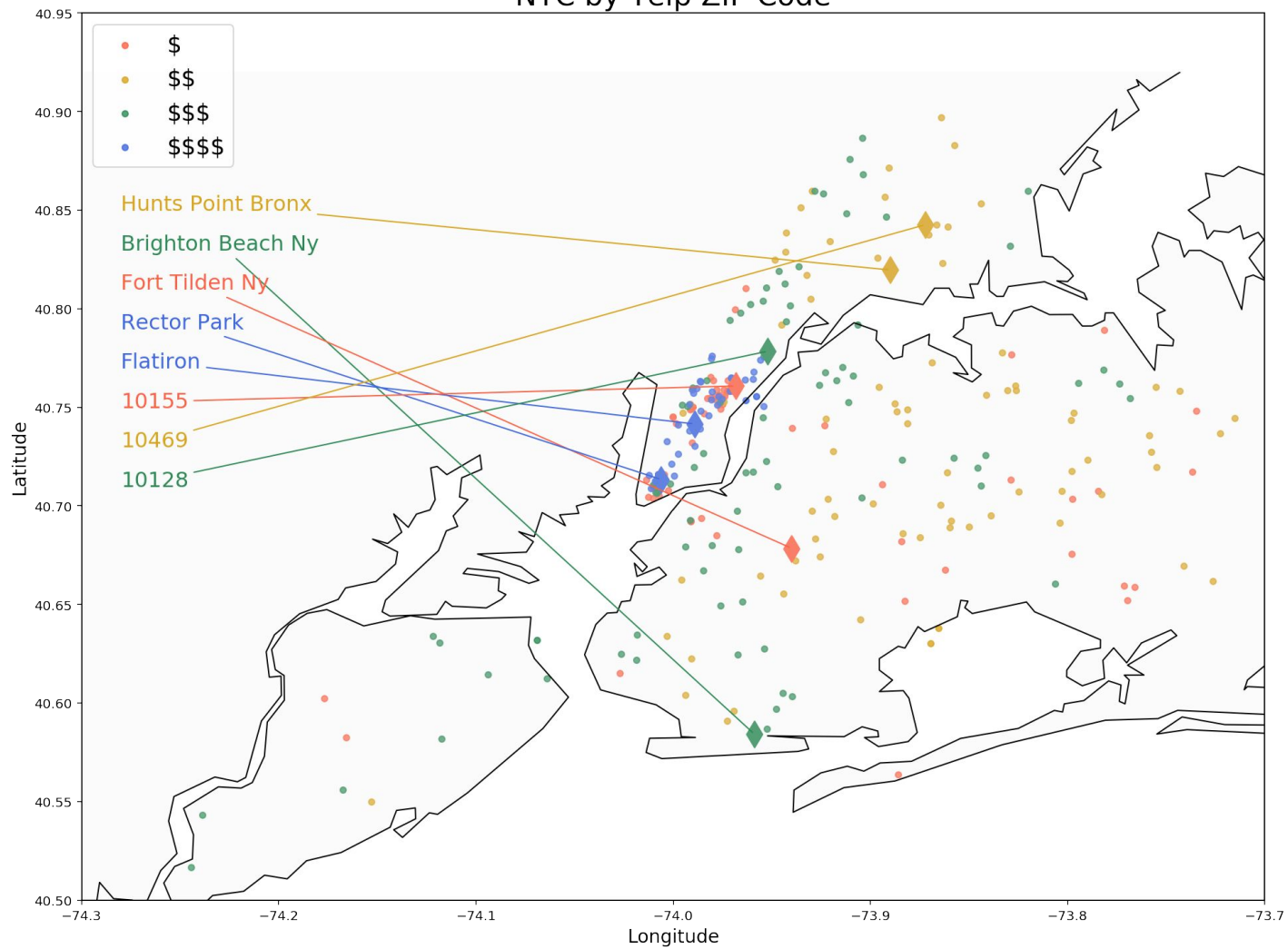
# **Price Distribution: $$$, $$$$**

◉ $$$ (green)
  ○ Highest activity in $$
  ○ Increased activity in $$$ and $$$$

◉ $$$$ (blue)
  ○ Highest activity in $$$ and $$$$
  ○ High activity in $$

14

# Mean for All NYC Zips:
# $, $$, $$$, $$$$

- ◉ $ (red)

- ◉ $$ (orange)

- ◉ $$$ (green)

- ◉ $$$$ (blue)
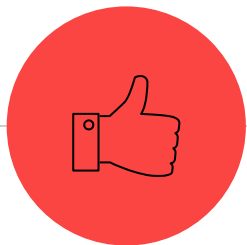
NYC by Yelp ZIP Code

# Recommendations

- Feed other models with cluster results
  - This information can be useful as an economic rating variable in a predictive model

- Pay to use the Yelp Fusion VIP API
  - Will allow for commercial–scale queries

- Beware of ZIP query results
  - Yelp returned out of state results for some NYC–based zipcodes
  - e.g. 10015 returned results in Tucson, AZ

# Next Steps

- Gathering and testing more data
  - Beyond the top 100, best match results

- Expanding class functionality
  - Enable collection of new training data
  - Automate model optimization

- Scaling the model
  - Train on other large metropolitan areas and check consistency of results

# Questions?

**How to find/contact us**

- ◉ Bernard Kurka
  - ○ [linkedin.com/in/bernardkurka](linkedin.com/in/bernardkurka)
  - ○ [bkexcel2014@gmail.com](bkexcel2014@gmail.com)
- ◉ Thomas Ludlow
  - ○ [linkedin.com/in/thomas-w-ludlow-jr-4568a1b](linkedin.com/in/thomas-w-ludlow-jr-4568a1b)
  - ○ [tludlow@gmail.com](tludlow@gmail.com)
- ◉ Brittany Allen
  - ○ [linkedin.com/in/brittadjacent](linkedin.com/in/brittadjacent)
  - ○ [thebrittallen@gmail.com](thebrittallen@gmail.com)