

Survival of Gentrification / Depreciation in Restaurants

F. M. Marsh and J. A. Clithero

September 27, 2016

1 Abstract

2 Introduction

Yelp Dataset Papers:

[Alghunaim, 2015, Byers et al., 2012, Cawkwell et al., 2015, Chepurna and Makrehchi, 2015, Feng and Qian, 2013, Gutierrez, 2014, Hajas et al., 2014, Hu et al., 2014, Liu et al., 2015, Mashhadi et al., 2012, Quattrone et al., 2015]

Zillow Dataset Papers:

Has previously been combined with the yelp dataset Bonnar et al..

Our goal is to use the correlation of two time-series:

1. The monthly median rent, as tracked by Zillow Rental Data.
2. The median restaurant review rating (stars) for each restaurant in a neighborhood.

Zillow rental data can be used to detect appreciating, and depreciating neighborhoods.

As rents rise in a given neighborhood, which types of businesses fare / worse better in the reviews? As rents fall in a given neighborhood, which types of business fare / worse better in the reviews?

We hope to present concrete suggestions to restaurant owners to improve the survivability of their businesses in times of strong appreciation / depreciation in the housing market.

3 Data

3.1 Yelp Academic Dataset

The Yelp Academic Dataset contains five files:

1) `yelp_academic_dataset.business.json`

2) `yelp_academic_dataset.review.json`

3) `yelp_academic_dataset_user.json`

4) `yelp_academic_dataset_checkin.json`

The Yelp Dataset Business file includes data on 77,445 businesses in the metro areas of Las Vegas, NV, Phoenix, AZ, Charlotte, NC Pittsburgh, PA Champaign, IL, Kitchener, Canada, Montreal, Canada, Edinburgh, Scotland, Karlsruhe, Germany.

3.2 Zillow Public Dataset

The Zillow Public Dataset (hereafter Zillow dataset) contains many files.

Zillow divides homes into geographic “neighborhoods” with well defined boundaries. The Zillow Home Value Index (ZHVI) is Zillow’s best estimate of median home price in a neighborhood. ZHVI is reported on a monthly basis for 6,958 neighborhoods across the US.

Median rental price for studio, one, two, three, four and five or more bedroom apartments are reported for a smaller set of about 300 neighborhoods.

The Zillow neighborhood boundaries <http://www.zillow.com/howto/api/neighborhood-boundaries.htm>

Each Zillow neighborhood has geographic boundaries, defined in an associated ESRI arcGIS shape file. Boundary information <http://www.zillowgroup.com/news/7000-neighborhood-boundary-Released-in-2008>

ZHVI methodology. <http://www.zillow.com/research/zhvi-methodology-6032/>
[?]

4 Methods

4.1 Combination of Datasets

Each Yelp business is tagged with a geographic (latitude, longitude) coordinate. In this section, we describe how we sort each Yelp business into its appropriate Zillow neighborhood.

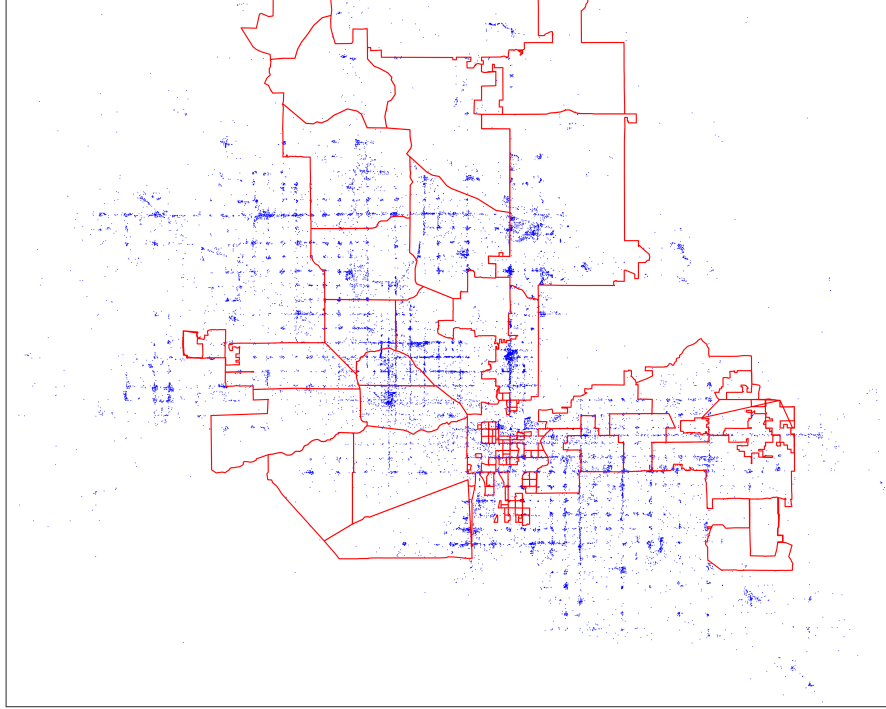


Figure 1: Yelp Businesses (points in blue) and Zillow neighborhood boundaries (lines in red) for the Phoenix, AZ metro area. In §4.1, we describe how we sort each Yelp Business into its appropriate Zillow neighborhood.

To perform this sorting, we employ a two-step approach. In the first step, we test every Yelp business for inclusion in the bounding box of every Zillow neighborhood. In the second step, we test every Yelp business for polygon inclusion in the neighborhoods which bounding boxes it lies within. We use this two-step approach because the first step can rule out all but two or three of the

We introduce the concept of the bounding box which we will define as the smallest range of latitudes and longitudes that include the whole neighborhood polygon. We test each Yelp business for inclusion in the set of 6,958 bounding boxes. In Fig. 2, we see a randomly selected Yelp business, displayed as a red point. We see that this business is included in the bounding boxes of two Zillow neighborhoods.

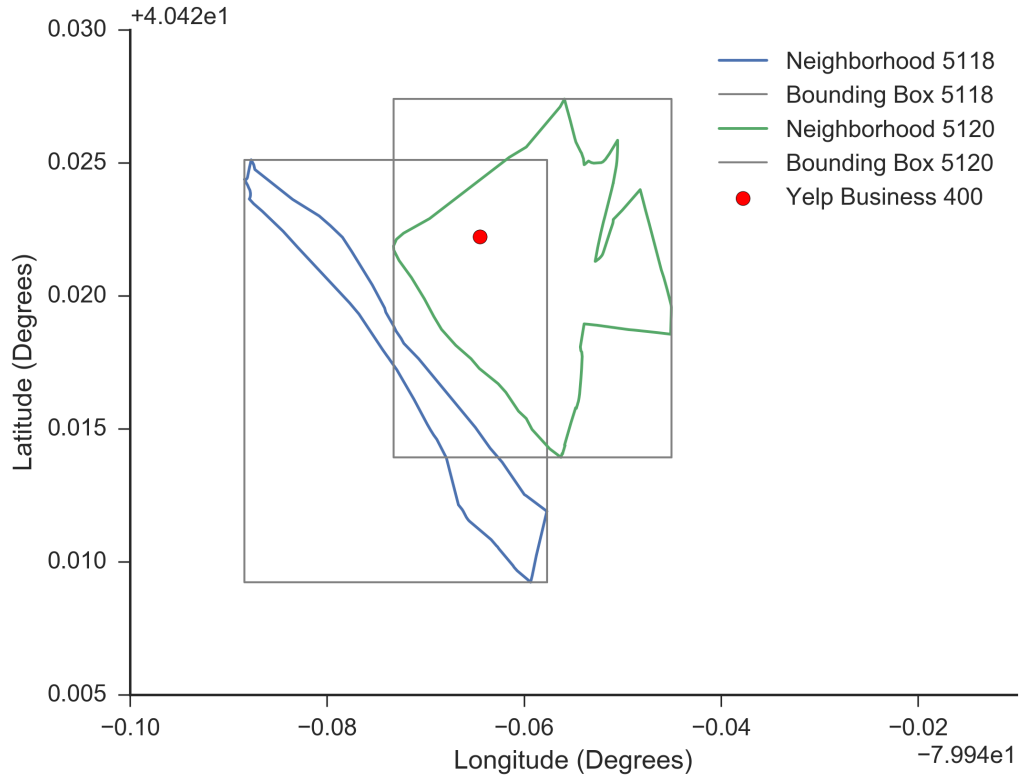


Figure 2: Example of neighborhood and neighborhood bounding box inclusion method. Yelp business 400 (the red point) is included in the bounding boxes of two Zillow neighborhoods. It is only included in one neighborhood polygon, however.

We then test for point-in-polygon inclusion using an implementation of a ray-casting method in `Python` [?]. For each Yelp business, we only test the Zillow neighborhoods whose bounding boxes it lies within.

4.2 Description of Combined Dataset

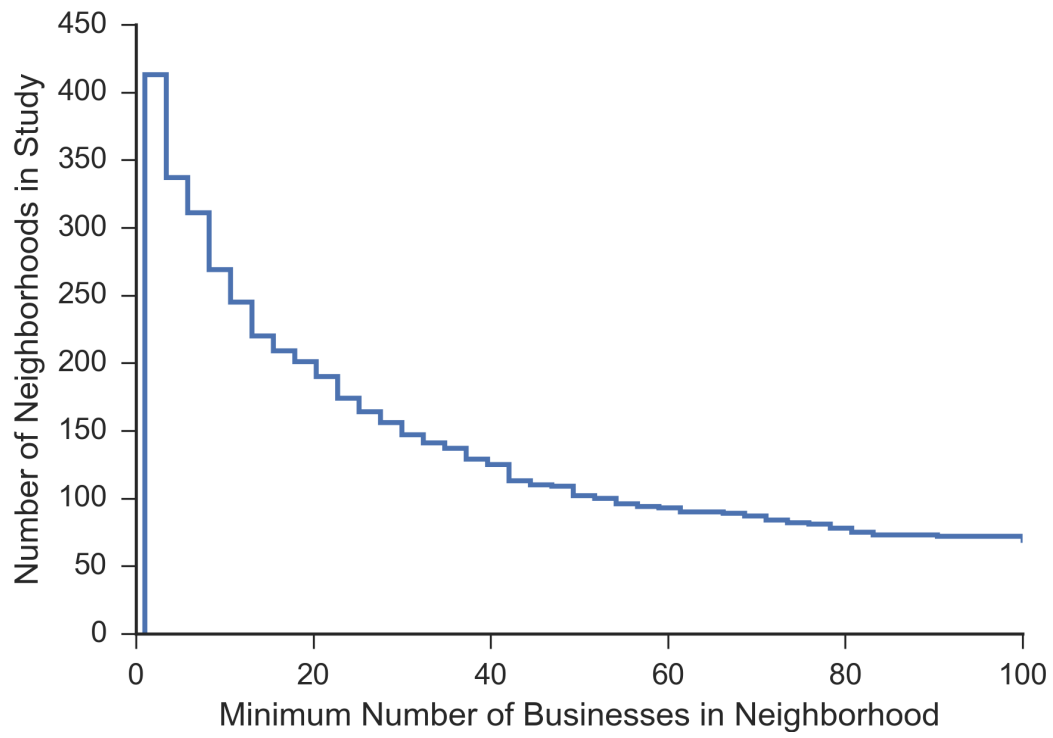


Figure 3:

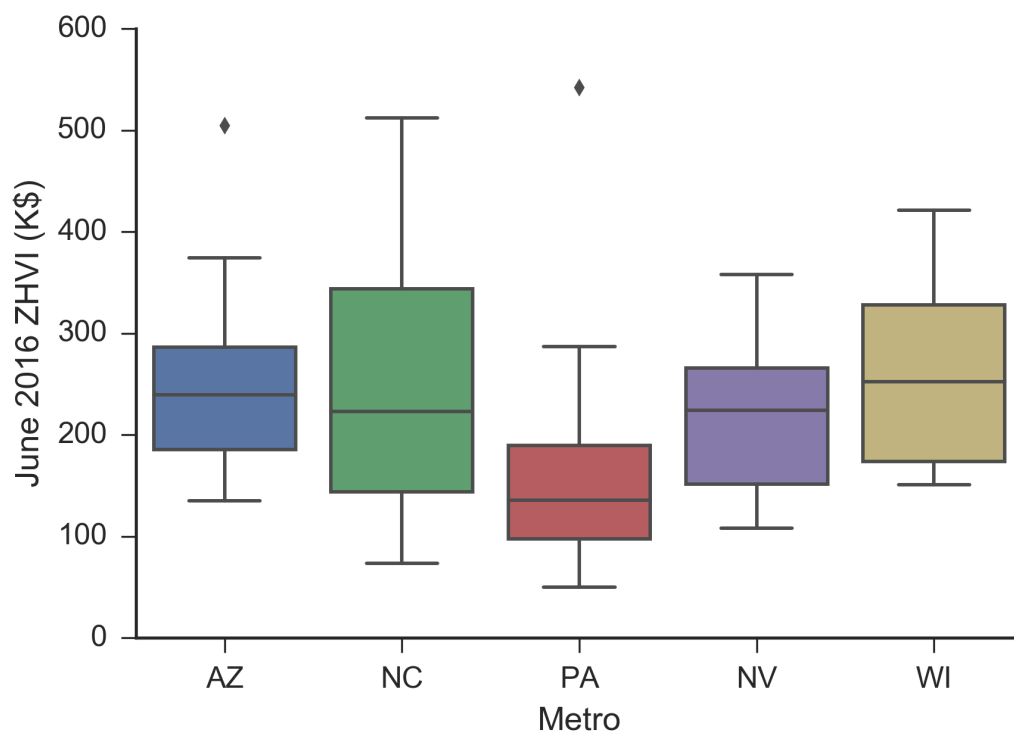


Figure 4:

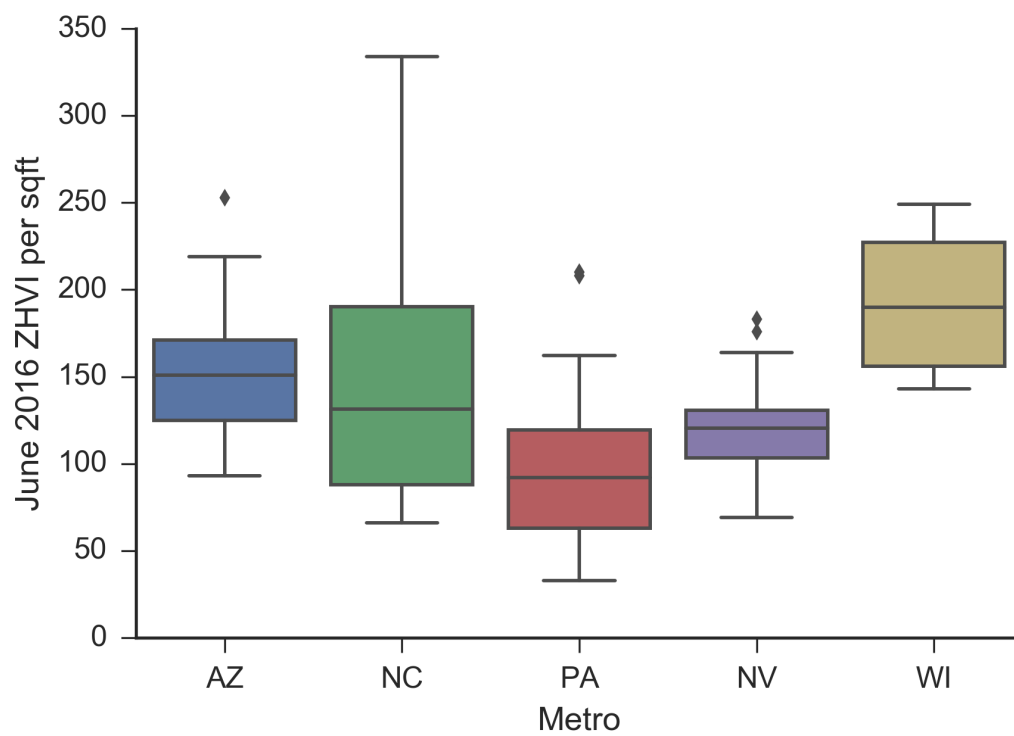


Figure 5:

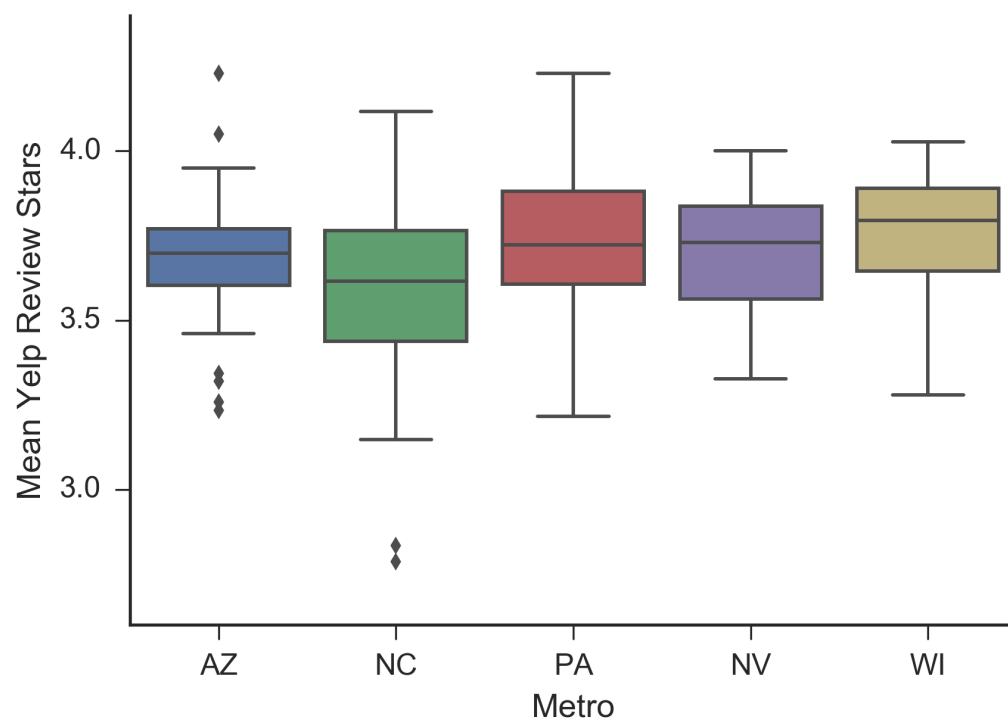


Figure 6:

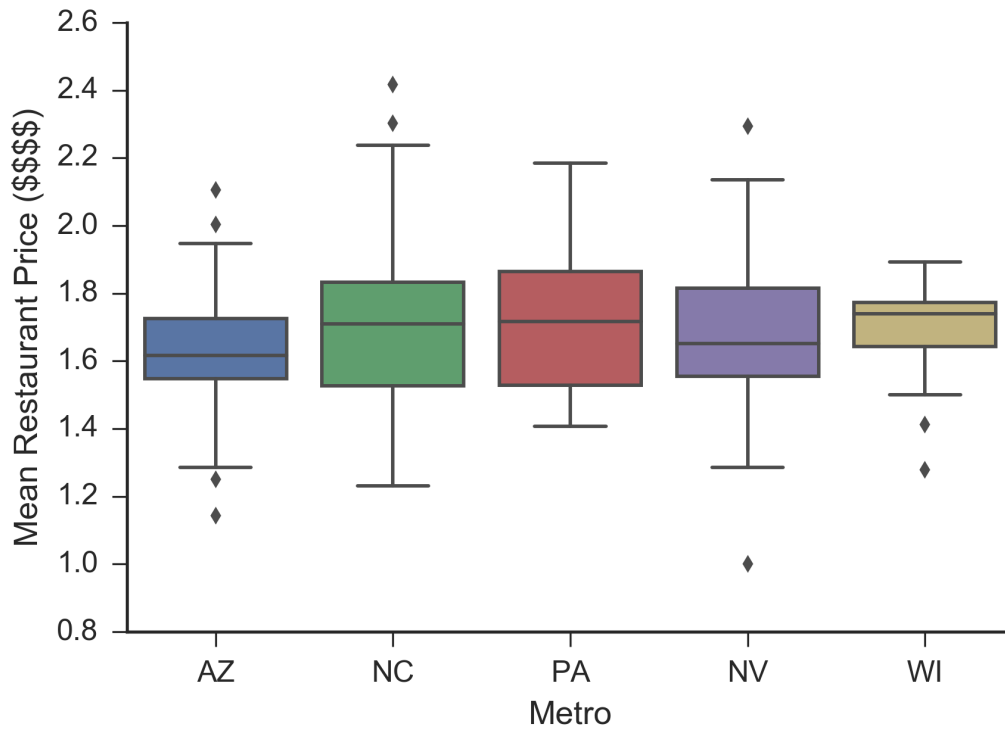


Figure 7:

| Metro | Original | Retained | Rate |
|---------------------|----------|----------|-------|
| Phoenix, AZ | 32,615 | 19,958 | 61.2% |
| Charlotte, NC | 6,162 | 4,853 | 78.7% |
| Las Vegas, NV | 21,233 | 9,694 | 45.7% |
| Pittsburgh, PA | 3,754 | 2,749 | 73.2% |
| Madison, WI | 2,802 | 1,355 | 48.3% |
| Champaign, IL | x | 0 | 0% |
| Kitchener, Canada | x | 0 | 0% |
| Montreal, Canada | x | 0 | 0% |
| Edinburgh, Scotland | x | 0 | 0% |
| Karlsruhe, Germany | x | 0 | 0% |

Table 1: min n = 0

| | median | mean | sum | len |
|-------|-----------|-----------|-------------|-----|
| state | | | | |
| AZ | 31.609536 | 62.769432 | 2636.316160 | 42 |
| NC | 2.776164 | 3.466144 | 225.299334 | 65 |
| NV | 11.553754 | 22.659192 | 997.004433 | 44 |
| PA | 1.192682 | 1.586050 | 41.237288 | 26 |
| WI | 0.057833 | 0.066647 | 0.999710 | 15 |

Table 4: n = 20

The combined Zillow and Yelp Dataset contains

| Metro | median (B/N) | mean (B/N) | Total B | Total N |
|----------------|--------------|------------|---------|---------|
| Phoenix, AZ | 116.5 | 467.6 | 19,640 | 42 |
| Charlotte, NC | 43.0 | 66.1 | 4,296 | 65 |
| Las Vegas, NV | 178.0 | 218.6 | 9,619 | 44 |
| Pittsburgh, PA | 43.5 | 94.8 | 2,464 | 26 |
| Madison, WI | 36.0 | 71.4 | 1,071 | 15 |

Table 2: min n = 20

| City | median (B/N) | mean (B/N) | Total B | Total N |
|-----------------|--------------|------------|---------|---------|
| Charlotte | 43.0 | 66.1 | 4296 | 65 |
| Henderson | 105.0 | 153.0 | 2907 | 19 |
| Las Vegas | 241.5 | 242.0 | 5809 | 24 |
| Madison | 36.0 | 71.4 | 1071 | 15 |
| Mesa | 491.5 | 509.5 | 3057 | 6 |
| North Las Vegas | 903.0 | 903.0 | 903 | 1 |
| Phoenix | 625.0 | 763.30 | 11450 | 15 |
| Pittsburgh | 43.5 | 94.8 | 2464 | 26 |
| Scottsdale | 2100.0 | 1531.7 | 4595 | 3 |
| Tempe | 26.5 | 29.9 | 538 | 18 |

Table 3:

4.3 Determination of Most Common Yelp Categories and Chains

Each Yelp business has user-generated tags, that allow other users to determine what genre the business is. In the full Yelp Dataset, there are 892 distinct categories. In the full Yelp

| category | counts |
|---------------------------|--------|
| Restaurants | 25071 |
| Shopping | 11233 |
| Food | 9250 |
| Beauty & Spas | 6583 |
| Health & Medical | 5121 |
| Nightlife | 5088 |
| Home Services | 4785 |
| Bars | 4328 |
| Automotive | 4208 |
| Local Services | 3468 |
| Active Life | 3103 |
| Fashion | 3078 |
| Event Planning & Services | 2975 |
| Fast Food | 2851 |
| Pizza | 2657 |
| Mexican | 2515 |
| Hotels & Travel | 2495 |
| American (Traditional) | 2416 |
| Sandwiches | 2364 |
| Arts & Entertainment | 2271 |

Table 5: The 20 most common Yelp tags in the Full Dataset ($n = 77,445$).

Dataset ($n = 77,445$)

Each Yelp business has a user-generated name. Many prominent chains have locations in neighborhoods across the country

4.4 Computation of Price Trend

5 Results

Ideas:

1. For each chain, how does chain density (spatial, and fractional) relate to ZHVI.
2. How does “chainy-ness” ($\frac{N_{\text{chain}}}{N_{\text{local}}}$) relate to neighborhood ZHVI / size.
3. For each chain, how does chain review correlate with neighborhood ZHVI / size.

| counts | name |
|--------|------------------------|
| 483 | Starbucks |
| 365 | Subway |
| 345 | McDonald's |
| 200 | Walgreens |
| 180 | Taco Bell |
| 155 | Pizza Hut |
| 147 | Burger King |
| 144 | Wendy's |
| 134 | The UPS Store |
| 120 | Panda Express |
| 119 | Dunkin' Donuts |
| 118 | Bank of America |
| 114 | Great Clips |
| 108 | Wells Fargo Bank |
| 107 | Circle K |
| 97 | Domino's Pizza |
| 95 | Chipotle Mexican Grill |
| 95 | Jimmy John's |
| 93 | KFC |
| 88 | US Post Office |

Table 6: The 20 most common Yelp restaurant names in the Full Dataset ($n = 77,445$).

5.1 Statics Results

In this section, we compute static results for each neighborhood.

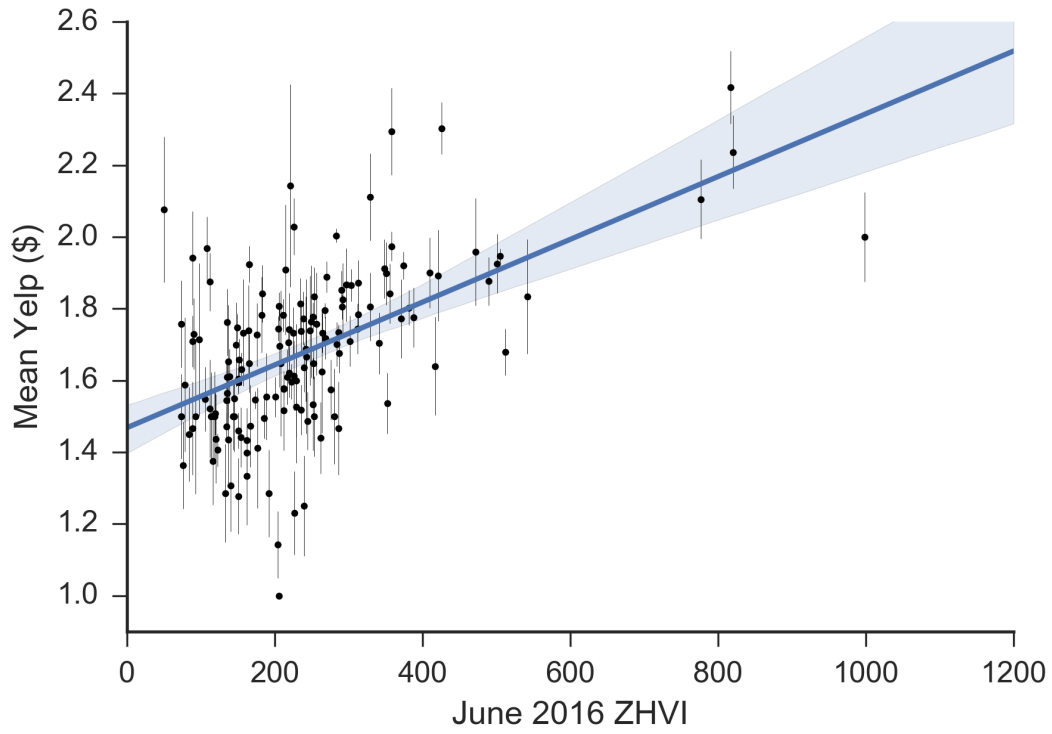


Figure 8:

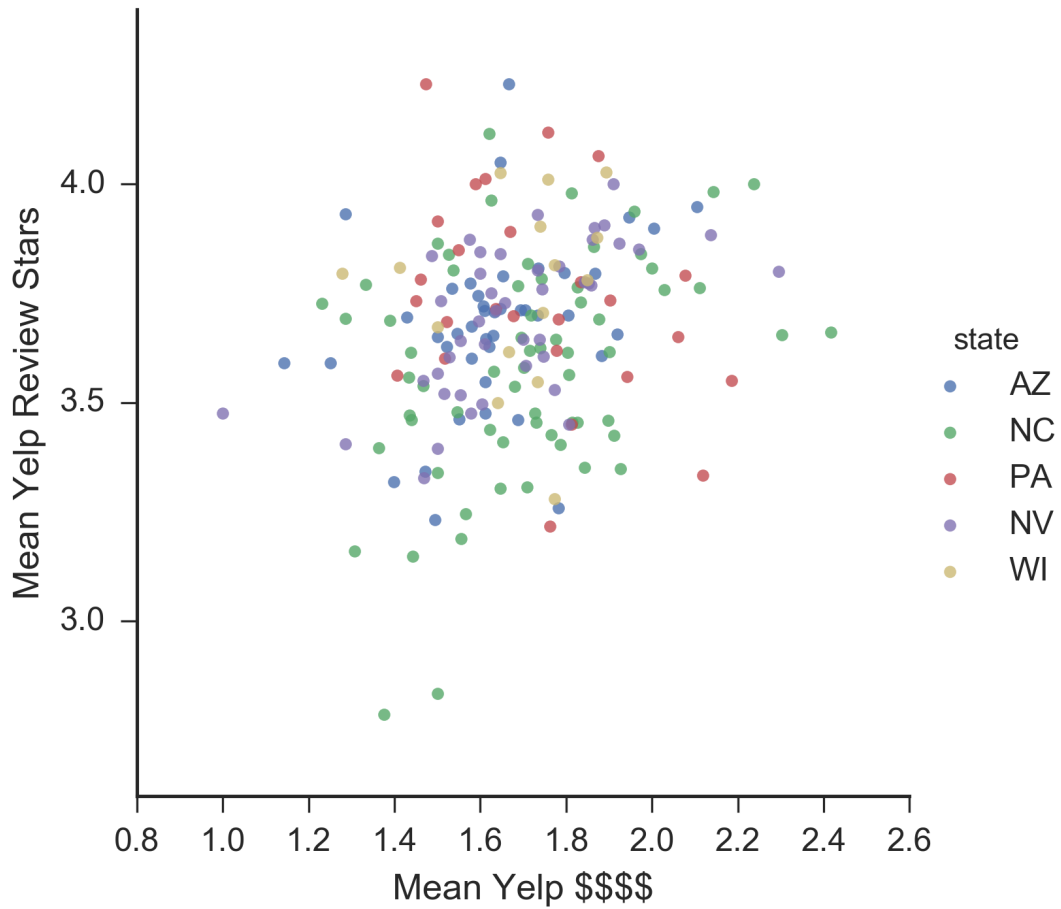


Figure 9:

References

- A. Alghunaim. *A Vector Space Approach for Aspect-Based Sentiment Analysis*. PhD thesis, Massachusetts Institute of Technology, 2015.
- C. Bonnar, F. Cordeiro, and J. Michelman. With a little help from yelp.
- J. W. Byers, M. Mitzenmacher, and G. Zervas. Thegroupon effect on yelp ratings: a root cause analysis. In *Proceedings of the 13th ACM conference on electronic commerce*, pages 248–265. ACM, 2012.

- P. B. Cawkwell, L. Lee, M. Weitzman, and S. E. Sherman. Tracking hookah bars in new york: Utilizing yelp as a powerful public health tool. *JMIR public health and surveillance*, 1(2), 2015.
- I. Chepurna and M. Makrehchi. Exploiting class bias for discovery of topical experts in social media. In *2015 IEEE International Conference on Data Mining Workshop (ICDMW)*, pages 64–71. IEEE, 2015.
- H. Feng and X. Qian. Recommendation via user’s personality and social contextual. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 1521–1524. ACM, 2013.
- L. A. Gutierrez. *Noise reduction in user generated datasets*. PhD thesis, RENSSELAER POLYTECHNIC INSTITUTE, 2014.
- P. Hajas, L. Gutierrez, and M. S. Krishnamoorthy. Analysis of yelp reviews. *arXiv preprint arXiv:1407.1443*, 2014.
- L. Hu, A. Sun, and Y. Liu. Your neighbors affect your ratings: on geographical neighborhood influence to rating prediction. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 345–354. ACM, 2014.
- J. Liu, J. Shang, C. Wang, X. Ren, and J. Han. Mining quality phrases from massive text corpora. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 1729–1744. ACM, 2015.
- A. Mashhadi, G. Quattrone, L. Capra, and P. Mooney. On the accuracy of urban crowd-sourcing for maintaining large-scale geospatial databases. In *Proceedings of the Eighth Annual International Symposium on Wikis and Open Collaboration*, page 15. ACM, 2012.
- G. Quattrone, L. Capra, and P. De Meo. There’s no such thing as the perfect map: Quantifying bias in spatial crowd-sourcing datasets. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*, pages 1021–1032. ACM, 2015.