# Yelp Expenditure Level & Neighborhood Affluency
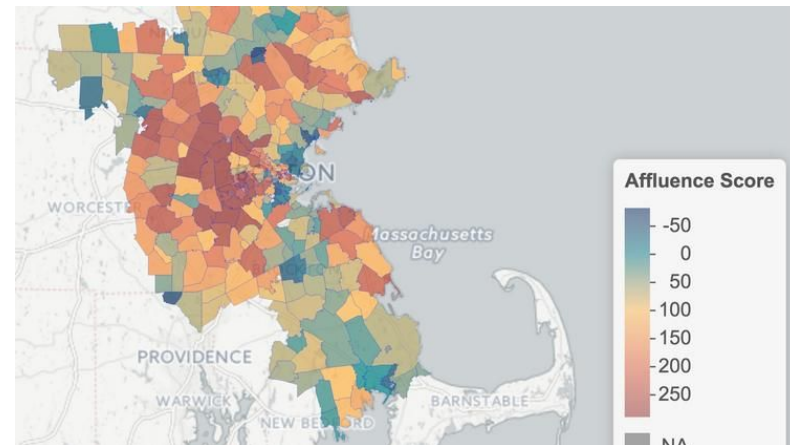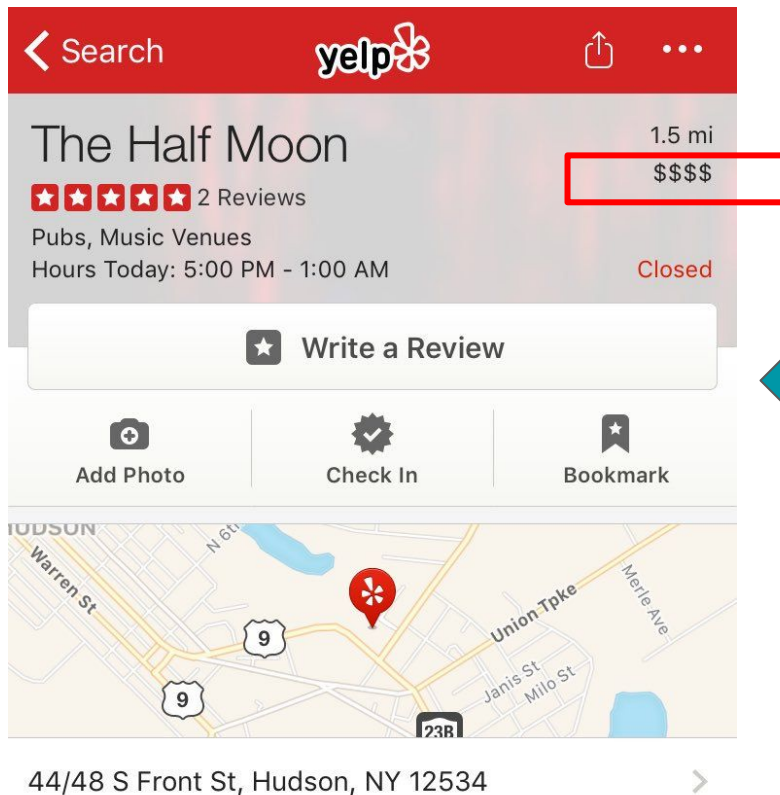
*Presented by <u>Daniel Stern</u>, <u>Rex Chang</u>, <u>Julian Sweet</u> and <u>Jane Liang</u>*
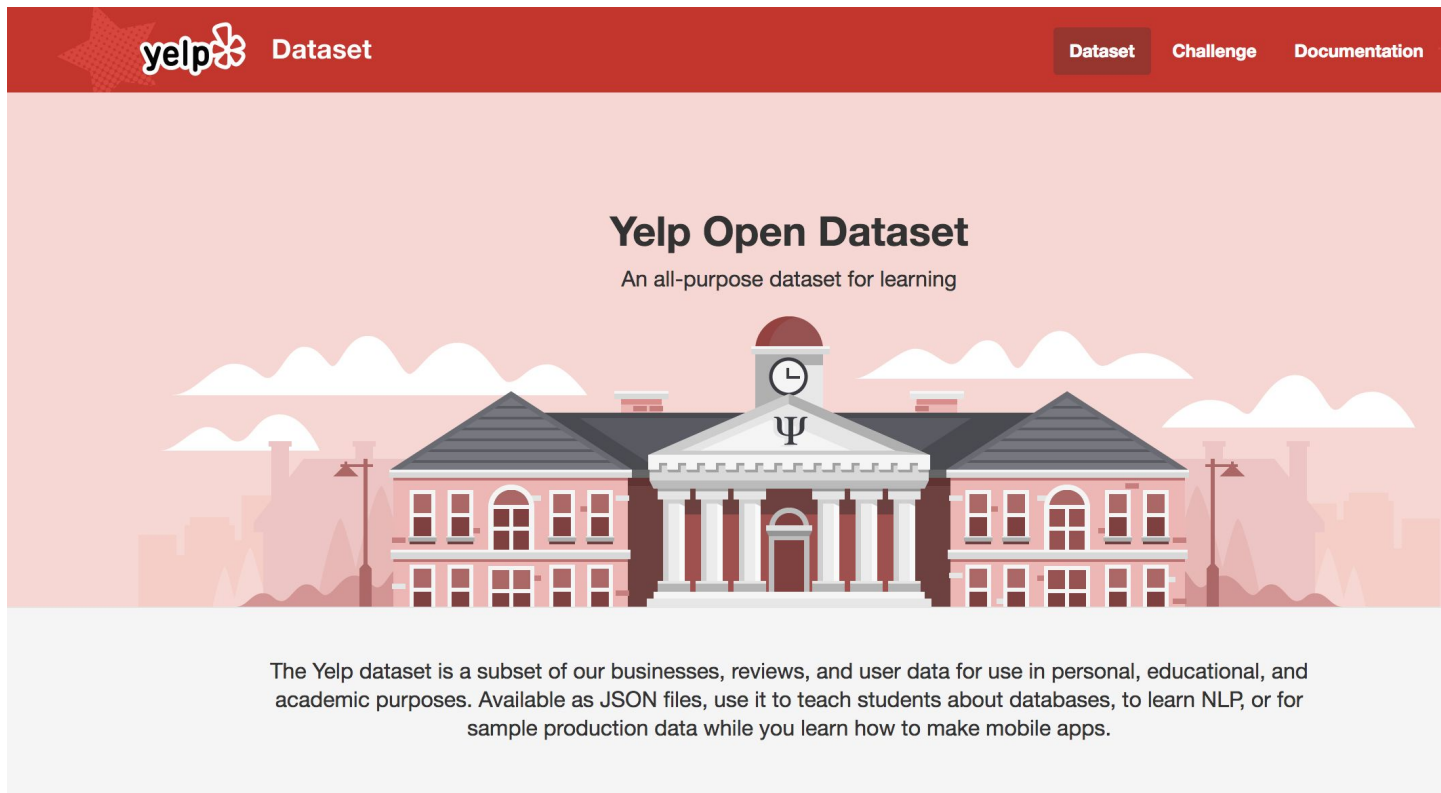
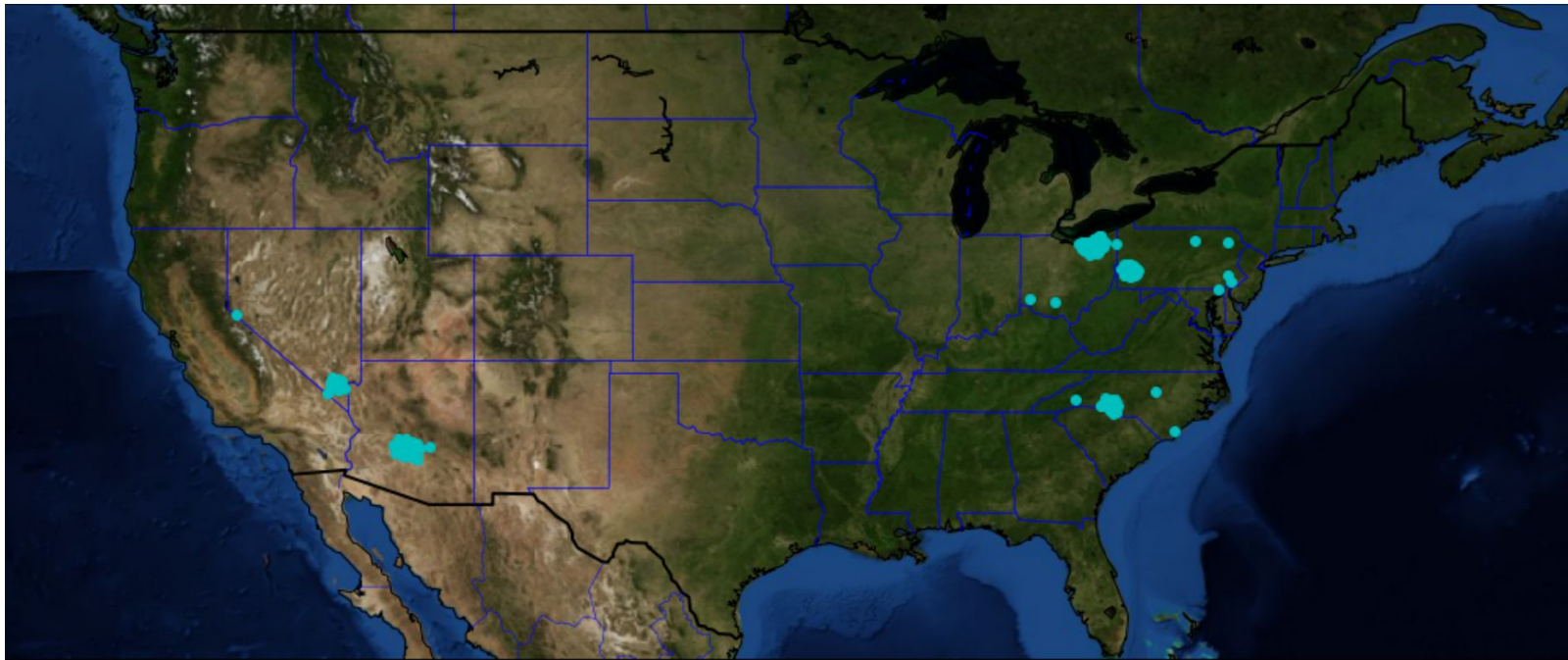# Yelp Open Dataset

[https://www.yelp.com/dataset](https://www.yelp.com/dataset) (174,000 business attributes)

# Metropolitan Zip Code Clusters

## Five Major Metros

# Distribution of Yelp Prices By Metro Area in Dataset

# Best way to engineer income target feature?

## IRS Income Data

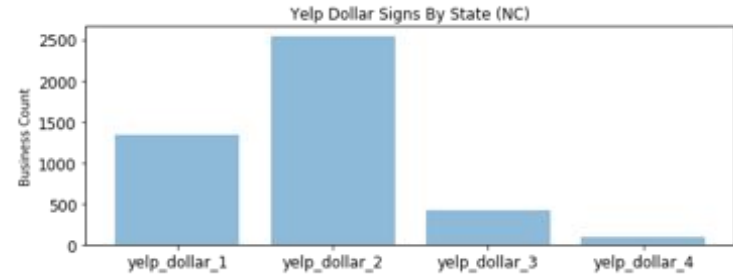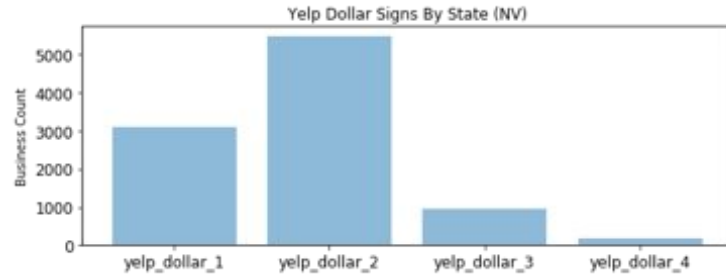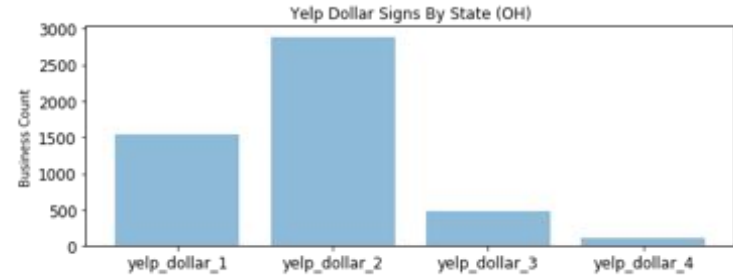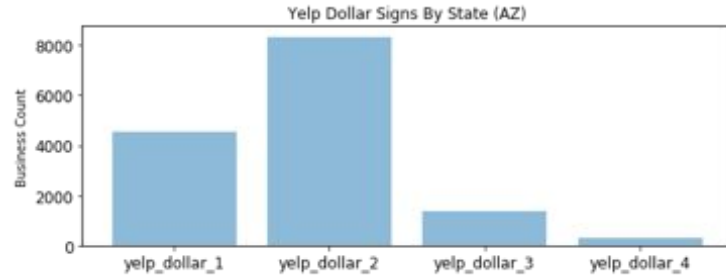### Number of residents in each category by zip code
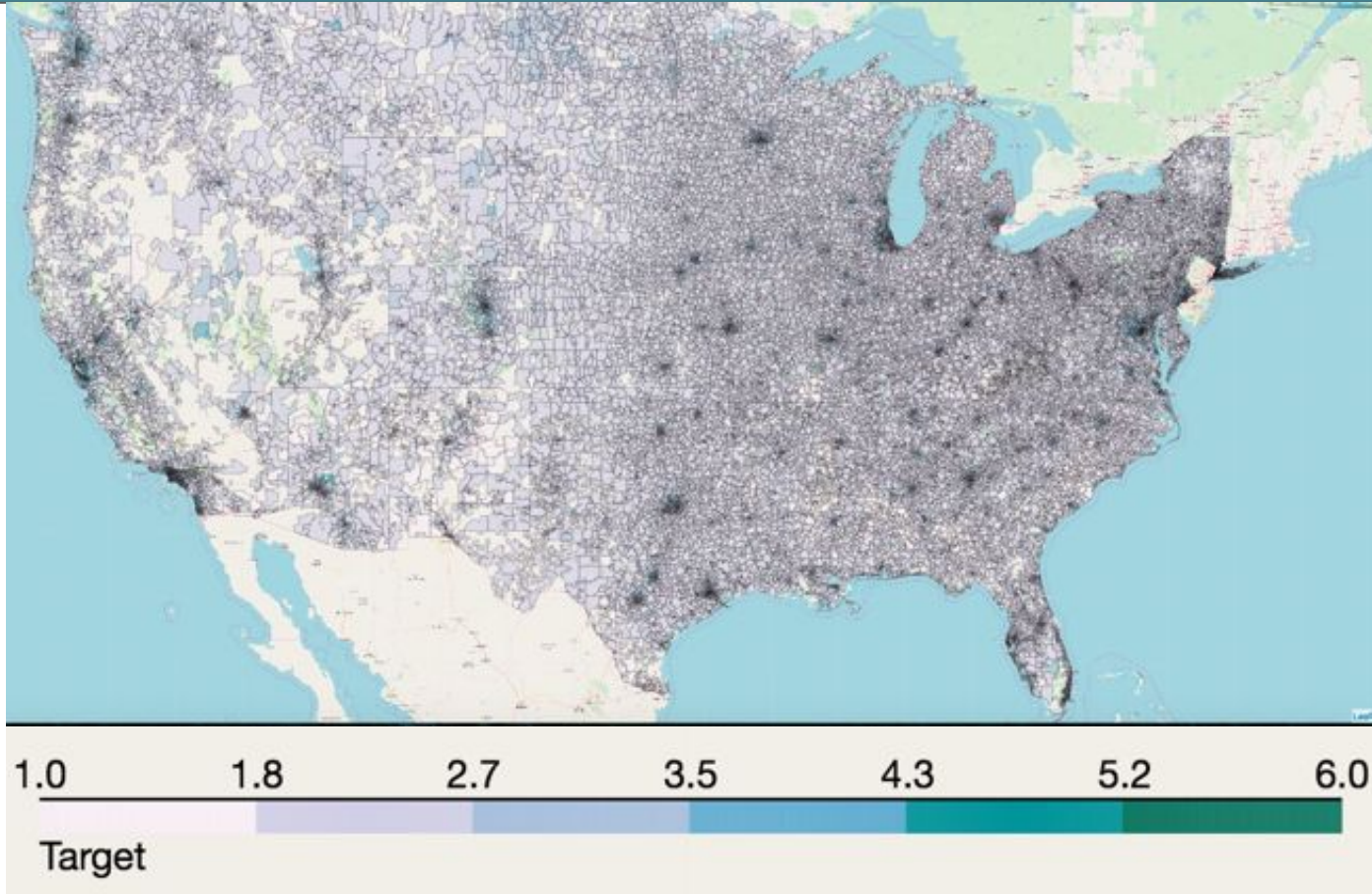
- 1 = $1 under $25,000
- 2 = $25,000 under $50,000
- 3 = $50,000 under $75,000
- 4 = $75,000 under $100,000
- 5 = $100,000 under $200,000
- 6 = $200,000 or more

1 to 6 – Categorical, but also ordinal

Weighted average? No

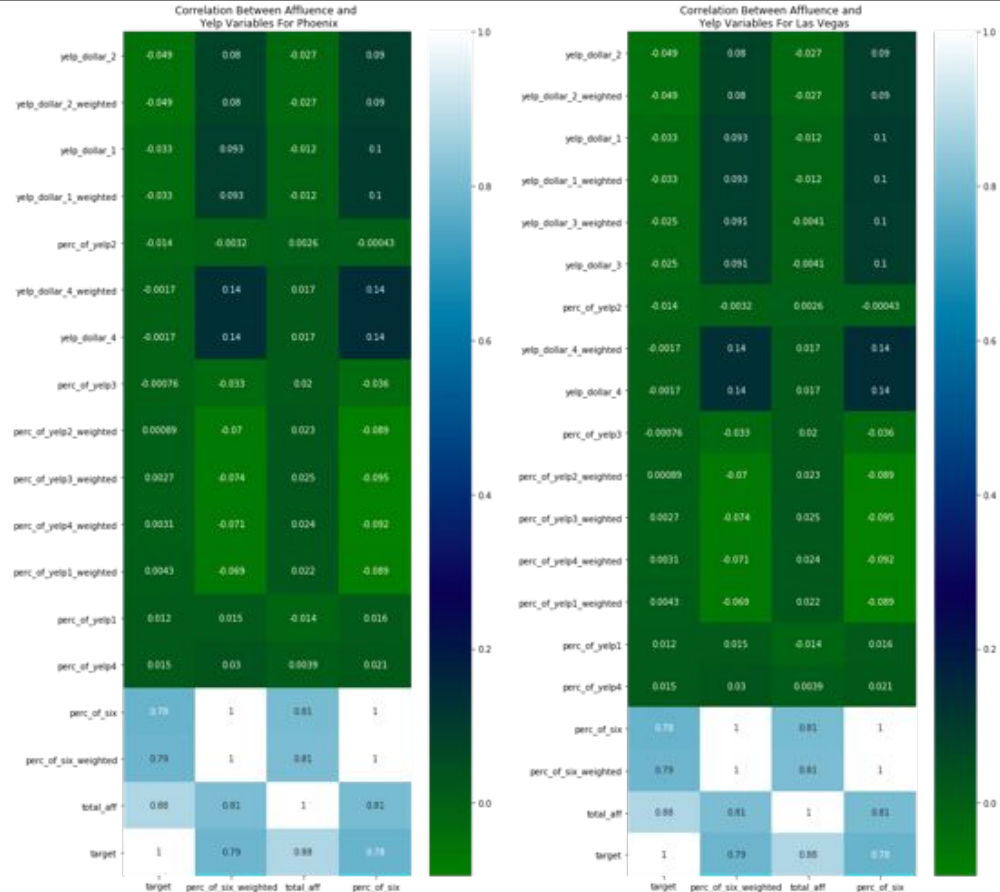Instead, median category #

# Sample Heatmap of Yelp + Affluency Data

- Very low correlation between Yelp & Affluence-related data.

- Correlation seems to be similar for each distinctive metro area.

# Engineered Feature – Metro Distance

- For 5 largest metros, Metro centroid was assigned to city hall for the respective metro. Latitude and longitude were determined by lookup.

- Specialized GeoPy library used to to determine distance in miles from city hall to each zip code.

- This served as an engineered feature to determine if proximity, or lack thereof to center of the metro correlated with affluence.

Baseline Accuracy

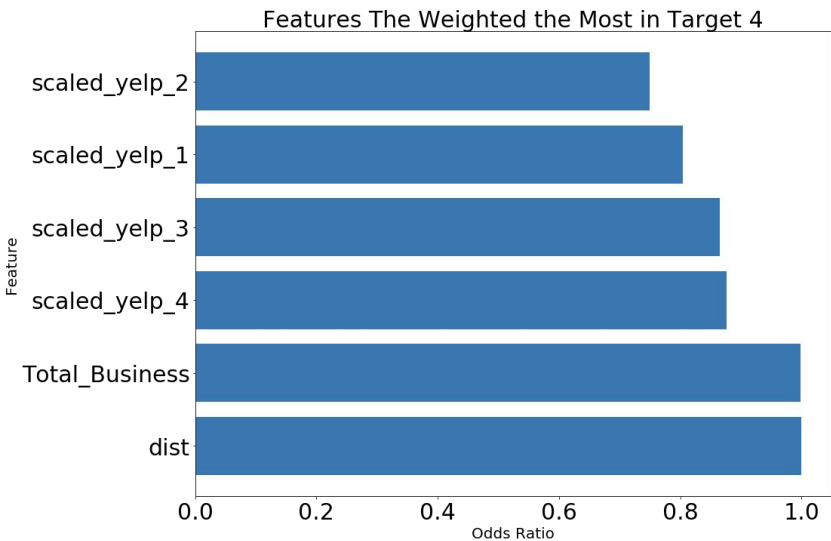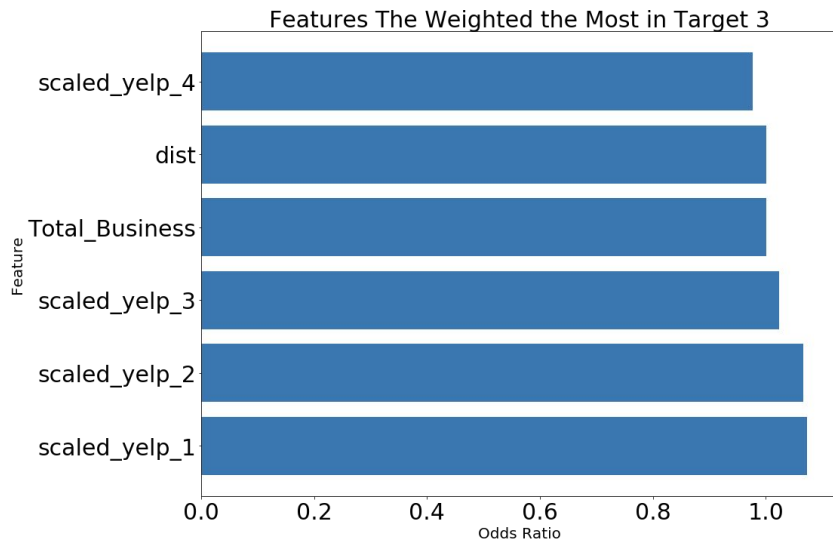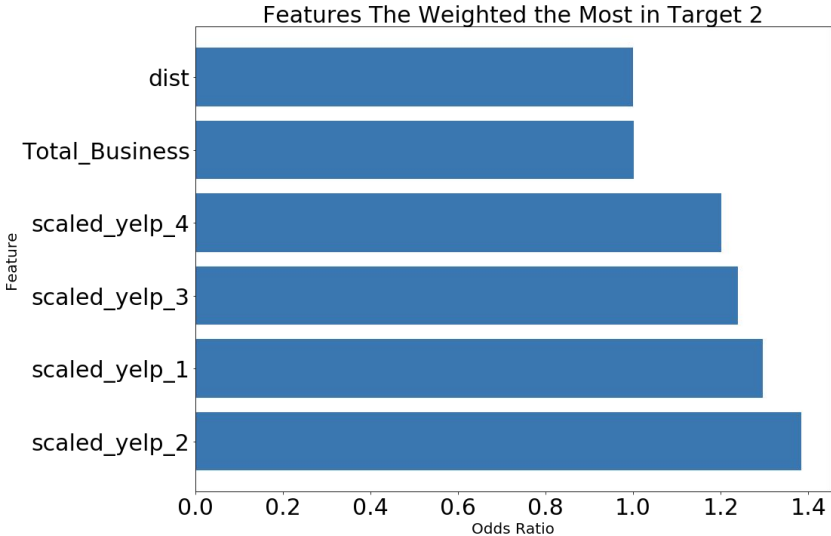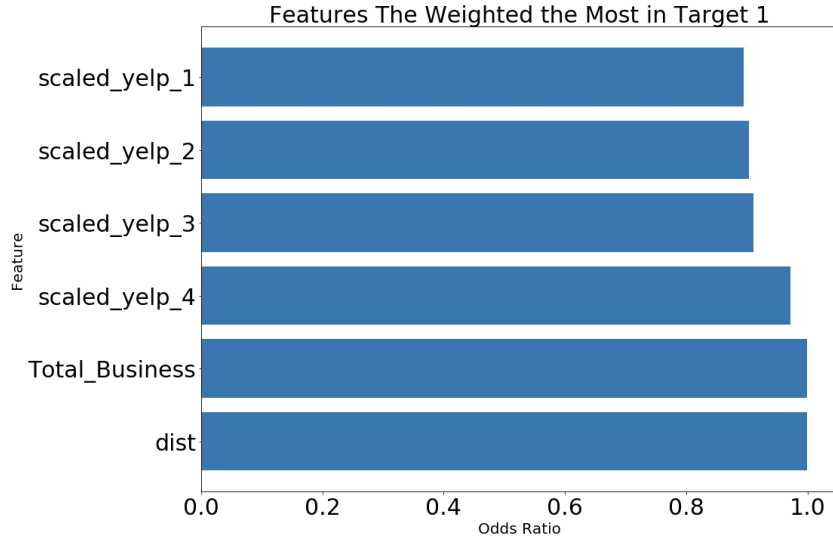**63.19%**

# Machine Learning: Logistic Regression Model

Modeling Accuracy :

Train :  **64.5%**                          Test :  **59.3%**

Intercepts :

| Target 1 | Target 2 | Target 3 | Target 4 |
|:---:|:---:|:---:|:---:|
| **-0.378** | **1.080** | **0.134** | **-0.836** |

Poor Model performance due to :

1. **Limited Scope:** currently not representative of the overall U.S. population.
2. **Limited Features:** aggregate expenditure level and distance from downtown for any given zip code. This is not sufficient when building our model.
3. **Other ML models**.