# REAL ESTATE PRICING PREDICTIONS

Elio Aybar, Cristal Garcia, Sunny Li, and Matt Norgren

**GOAL:** PREDICT MEDIAN HOUSE PRICES GIVEN CONVENTIONAL ATTRIBUTES AND PRESENCE OF A TARGET AND/OR STARBUCKS.

**VALUE:** TESTING NON-CONVENTIONAL FEATURES TO SEE IF WORTH EXPLORING IN THE FUTURE.

**HYPOTHESIS:** DIFFERENT TYPES OF NEIGHBORHOODS MAY ATTRACT SPECIFIC BUSINESSES.

# PROBLEM STATEMENT

# PROJECT DEFINITION OF VARIABLES

- Zillow Real Estate:    4.3M X 75

- Starbucks location    25K  X 7
- Target Locations:      1.8K X 47

- USA Public Schools data:        47K X 33
- Food Desert & Economic Data: 73K  X 148
- Hospital Ratings:                4.8K X 28

# FEATURES & TARGET VARIABLE

**Zip code (Key)**

**Zillow**
- NumberOfBedrooms
- Dwelling Type

**School**
- LEVEL_:  Grade level
- Class_Teacher_RATIO:

**Food Desert/Economic Data**
- Urban Indicator
- Occupied Unit Housing Count (2010 Census)
- Low Income Indicator
- Median Family Income
- Food Desert Indicator

**Hospital**
- Hospital Rating
- Emergency Services Flag

**Starbucks Locations**

**Target Locations**

**Pricing**

# LIMITATIONS AND DELIMITATIONS OF DATA

**01**

Hospital: hospital deserts (distance between hospital and home)

**02**

Public School Data: Creating Student-to-Teacher Ratio

**03**

Food Desert Data: converting CensusTract number to Zipcode

**04**

Starbucks and Target

# LIMITATIONS AND DELIMITATIONS OF DATA PART 2

Granularity

Time                                      (1996 – 2017)

Features                                  Sales Data
                                          Zillow Estimates

Complexity of combination                 19K Unique Zip codes

Aggregation of Features                   Number of Bedrooms
                                          Dwelling Type

# IS IT ALL HERE?

| | |
|---|---|
| ✕ | Scarcity of price |
| 📄 | Duplicative feature types |
| 🔀 | Scaling & Aggregation |
| ✓ | Consolidated variable creation    Function creation    Factorization |
| 🔥 | Melting and Reshaping |

```python
p_df['MedianListingPricePerSqft_1Bedroom'] = zip_
p_df['MedianListingPricePerSqft_2Bedroom'] = zip_
p_df['MedianListingPricePerSqft_3Bedroom'] = zip_
p_df['MedianListingPricePerSqft_4Bedroom'] = zip_
p_df['MedianListingPricePerSqft_5BedroomOrMore']

p_df['MedianListingPrice_1Bedroom'] = zip_df['Med
p_df['MedianListingPrice_2Bedroom'] = zip_df['Med
p_df['MedianListingPrice_3Bedroom'] = zip_df['Med
df['MedianListingPrice_4Bedroom'] = zip_df['Med
df['MedianListingPrice_5BedroomOrMore'] = zip_

df['MedianRentalPricePerSqft_1Bedroom'] = zip_
f['MedianRentalPricePerSqft_2Bedroom'] = zip_
['MedianRentalPricePerSqft_3Bedroom'] = zip_
'MedianRentalPricePerSqft_4Bedroom'] = zip_
'MedianRentalPricePerSqft_5BedroomOrMore'] =

ed            1Bedroom'] = zip_df['Medi
             2Bedroom'] = zip_df['Medi
anRentalPrice_3Bedroom'] = zip_df['Medi
RentalPrice_4Bedroom'] = zip_df['Medi
RentalPrice_5BedroomOrMore'] = zip_df

ingPrice_AllHomes'] = zip_df['Med
gPrice_CondoCoop'] = zip_df['Me
ice_DuplexTriplex'] = zip_df
e_SingleFamilyResidence']
```

# DATA EXPLORATION

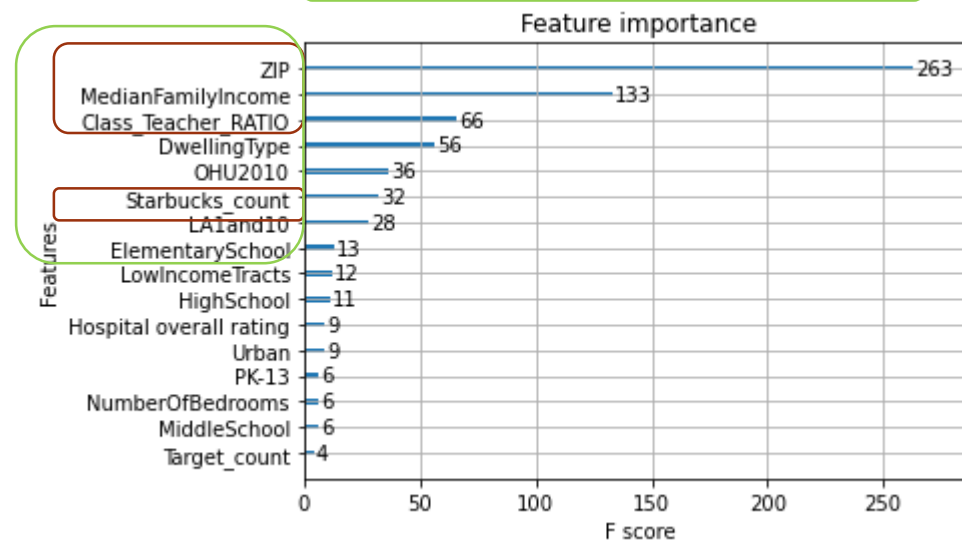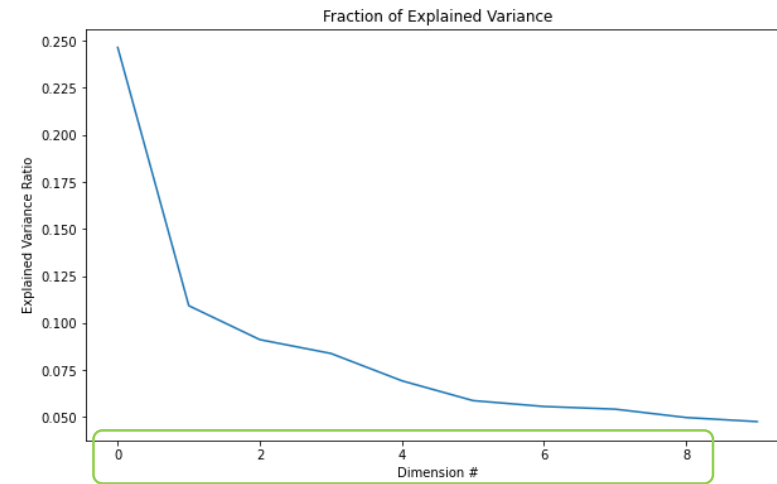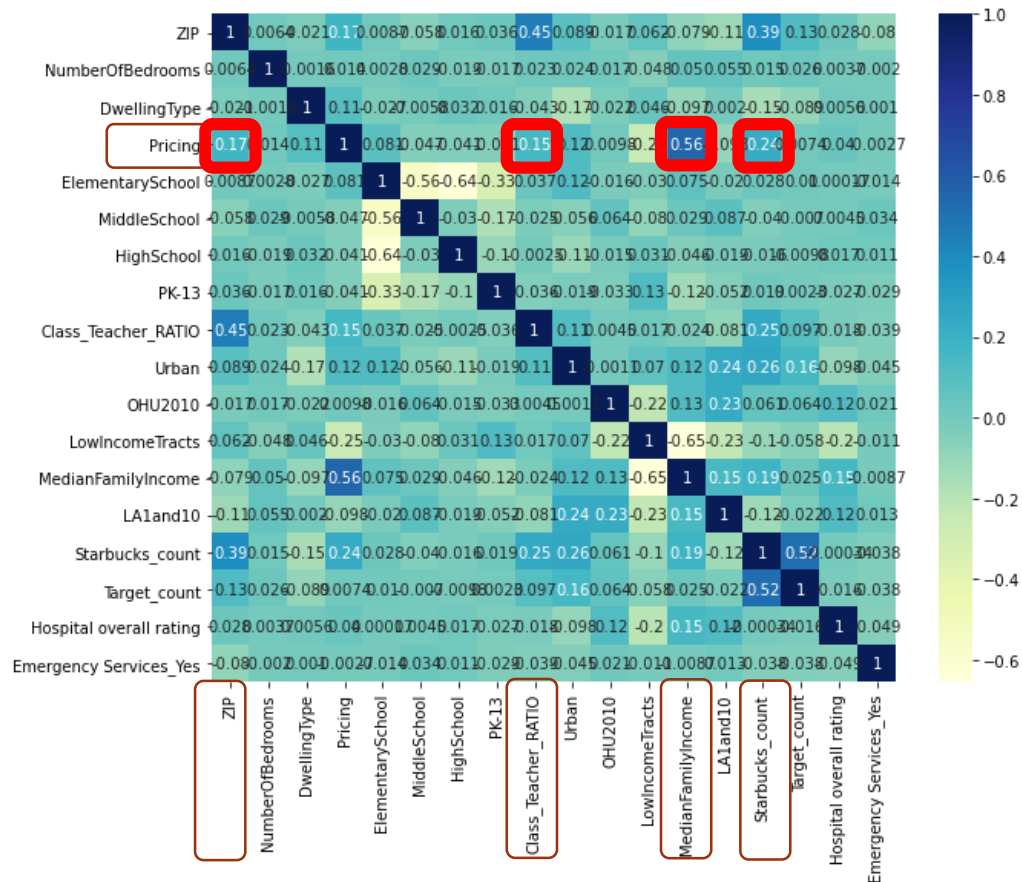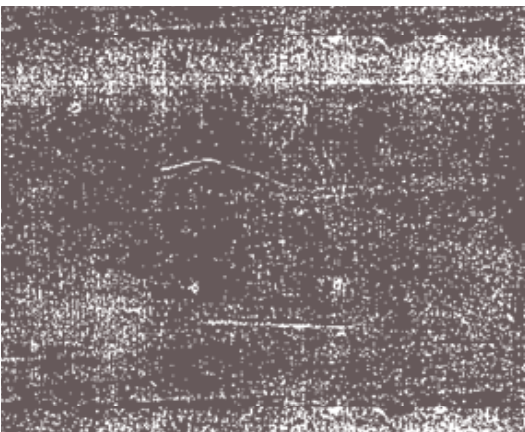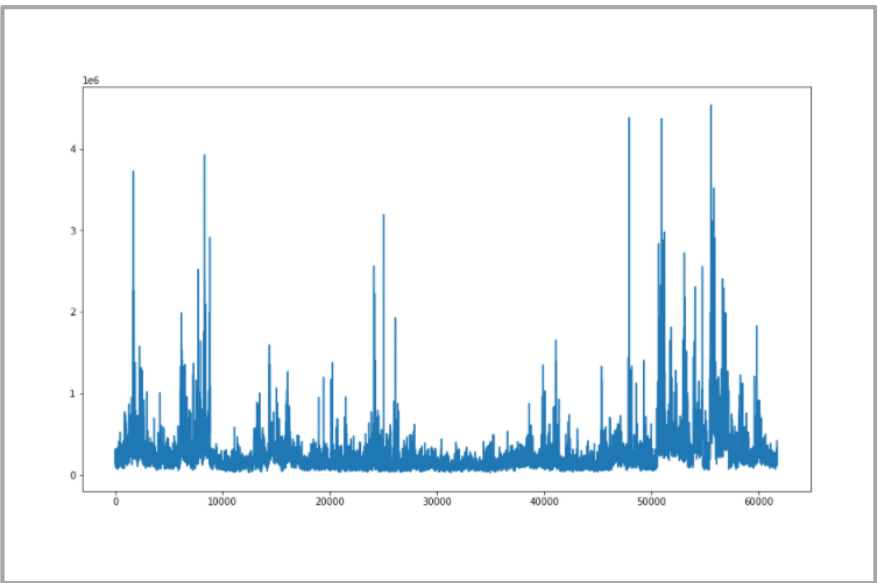Featurization:

- 90%-time spent feature engineering (5.5 weeks)

- Individually worked on datasets separately:
  - Identified and removed features with high collinearity
  - Captured summary statistics

- Consolidated features and visualized relationship

# HEATMAP, PCA & FEATURE IMPORTANCE

# TARGET VARIABLE EXPLORATION

Balance data using SMOTE

```python
#pip install --upgrade scikit-learn
```
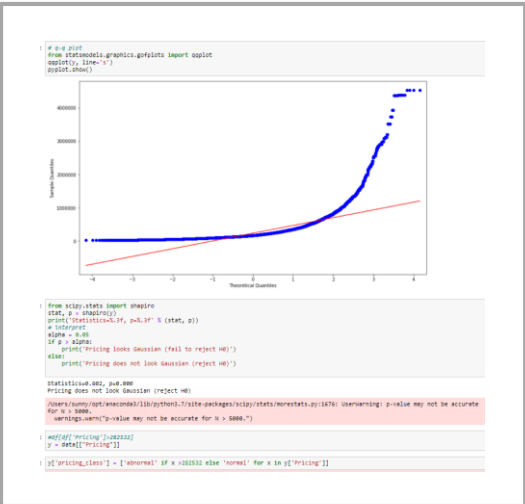
```python
#pip install imblearn
from imblearn.over_sampling import SMOTE
sm = SMOTE(sampling_strategy='auto')

X_sm, y_sm = sm.fit_sample(X, y['pricing_class'])
X_sm.shape, y_sm.shape
```

```
((92576, 18), (92576,))
```

```python
from collections import Counter
print('Original dataset shape {}'.format(Counter(y_sm)))
print('Resampled dataset shape {}'.format(Counter(y_sm)))
```

```
Original dataset shape Counter({'normal': 46288, 'abnormal': 46288})
Resampled dataset shape Counter({'normal': 46288, 'abnormal': 46288})
```

```python
# q-q plot
from statsmodels.graphics.gofplots import qqplot
qqplot(y, line='s')
pyplot.show()
```

```python
from scipy.stats import shapiro
stat, p = shapiro(y)
print('Statistics=%.3f, p=%.3f' % (stat, p))
# interpret
alpha = 0.05
if p > alpha:
    print('Pricing looks Gaussian (fail to reject H0)')
else:
    print('Pricing does not look Gaussian (reject H0)')
```

```
Statistics=0.082, p=0.000
Pricing does not look Gaussian (reject H0)
/Users/sunny/opt/anaconda3/lib/python3.7/site-packages/scipy/stats/morestats.py:1676: UserWarning: p-value may not be accurate for N > 5000.
  warnings.warn("p-value may not be accurate for N > 5000.")
```

```python
y=df['Pricing']>202532
y = data['Pricing']
```

```python
y['pricing_class'] = ['abnormal' if x >202532 else 'normal' for x in y['Pricing']]
```

# METHODOLOGY

- Models
  - Linear Regression
  - Decision Tree
  - Random Forest
  - Support Vector Regression
  - Gradient Boosting
  - KNN for Regression
  - XGBoost

# 1
## Validation Methodology

- Compare Desired Errors with Training and Cross Validation errors

# 2
## Descriptions

- Choose top 3 models

# 3
## Model Selection Justification

- Best performing RMSE

# FRAMEWORK

# RESULTS

# MODELING RESULTS

- Linear Regression
  - Desired Error: 279471.46
  - Training Error: 0.78
  - Cross Validation Error: 0.78

- Support Vector Regression (SVR)
  - RMSE score: 0.65

- Decision Tree
  - Hyperparameter Tuning: max depth = 4, min_samples_leaf = 0.1
  - Desired Error: 210381.81
  - Training Error: 194304.60
  - Cross Validation Error: 195936.04

- Random Forest
  - Hyperparameter Tuning: max depth = 4, min_samples_leaf = 0.1
  - Desired Error: 204785.74
  - Training Error: 201373.10

- Gradient Boosting
  - RMSE: 131471.88186

- KNN for Regression
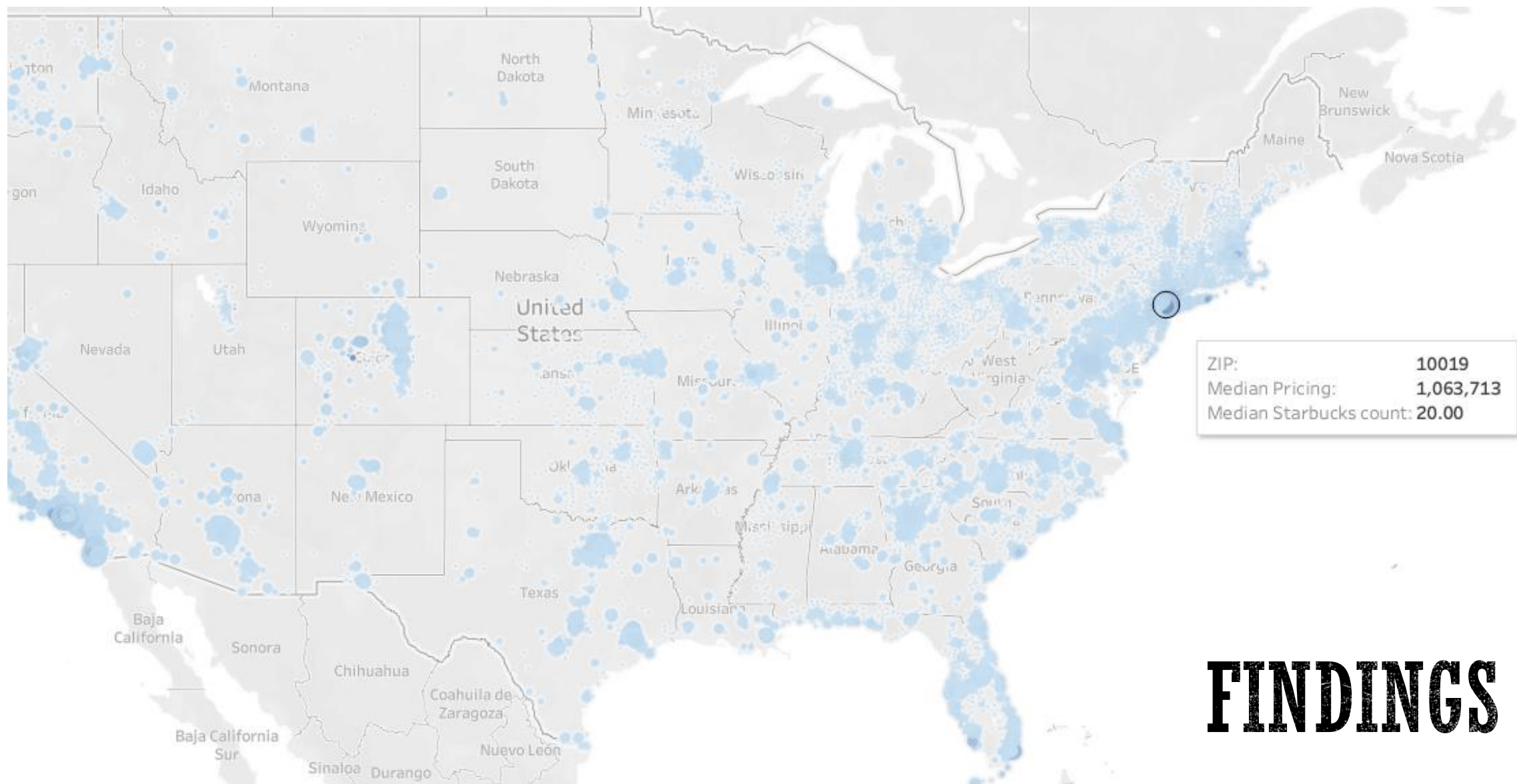  - N_neighbors = 2, RMSE value for k = 2 of 325245.033

- XGBoost
  - RMSE = 108618

# MODELING RESULTS DIFFERENT FEATURES

- Linear Regression
  - Desired Error: `294214.72`
  - Training Error: `0.84`
  - Cross Validation Error: `0.84`

- Support Vector Regression (SVR)
  - RMSE score: `0.79`

- Decision Tree
  - Hyperparameter Tuning: max depth = 4, min_samples_leaf = 0.1
  - Desired Error: `204245.67`
  - Training Error: `200400.99`
  - Cross Validation Error: `200487.57`

- Random Forest
  - Hyperparameter Tuning: max depth = 4, min_samples_leaf = 0.1
  - Desired Error: `208813.27`
  - Training Error: `205483.11`

- Gradient Boosting
  - RMSE: `173276.42`

- KNN for Regression
  - N_neighbors = 2, RMSE value for k = 2 of `327642.71`

- XGBoost
  - RMSE = `112413`

# XGBOOST

**eXtreme Gradient Boosted trees**

**Most powerful machine learning algorithm up until today**

**Features:**

- Regularized boosting ( prevent overfitting)
- Can handle missing values
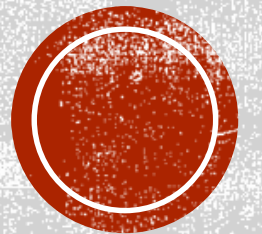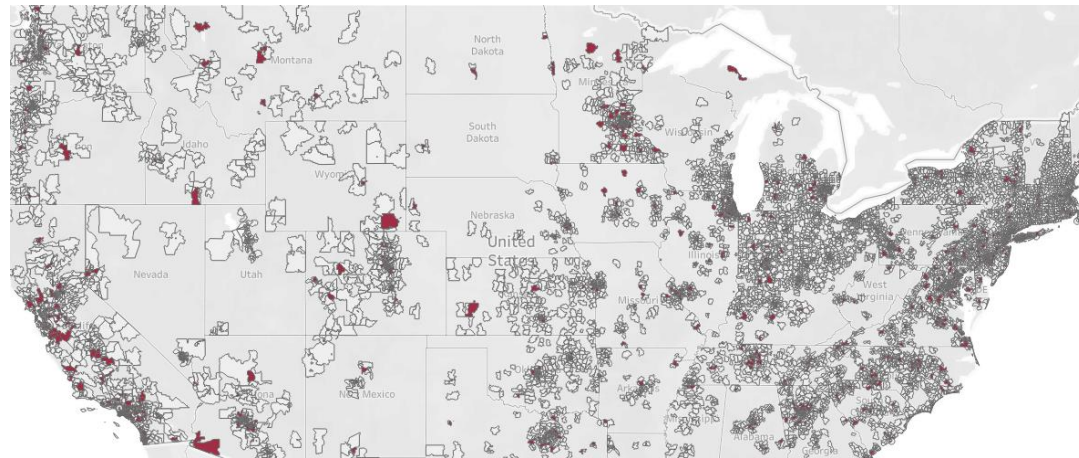- Parallel processing
- Can cross-validation at each iteration

ZIP: 10019
Median Pricing: 1,063,713
Median Starbucks count: 20.00

FINDINGS

CONVENTIONAL FEATURES

Target Density

Starbucks Density

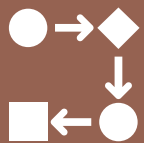Non-conventional features may be economic indicators.

We spent 80% of time on addressing a number of issues around data quality, standards, access.

We have plans for model optimization but didn't have time to achieve it.

CONCLUSIONS

# FUTURE CONSIDERATIONS

- New dataset or extension of features.

- Number of bathrooms, rent cost, Square footage, etc.

- Dollar trees stores, Dunkin donuts, pharmacies, supermarkets, etc

- Complex imputation methods

  - Focus on the hospital specifically, calculate the radii of hospitals, then use radii to weigh the hospital rating for each zip code.

- Robust test

  - Replacing zip code with county/city

- Include time series data to the analysis

# QUESTIONS?

Thanks for Your attention