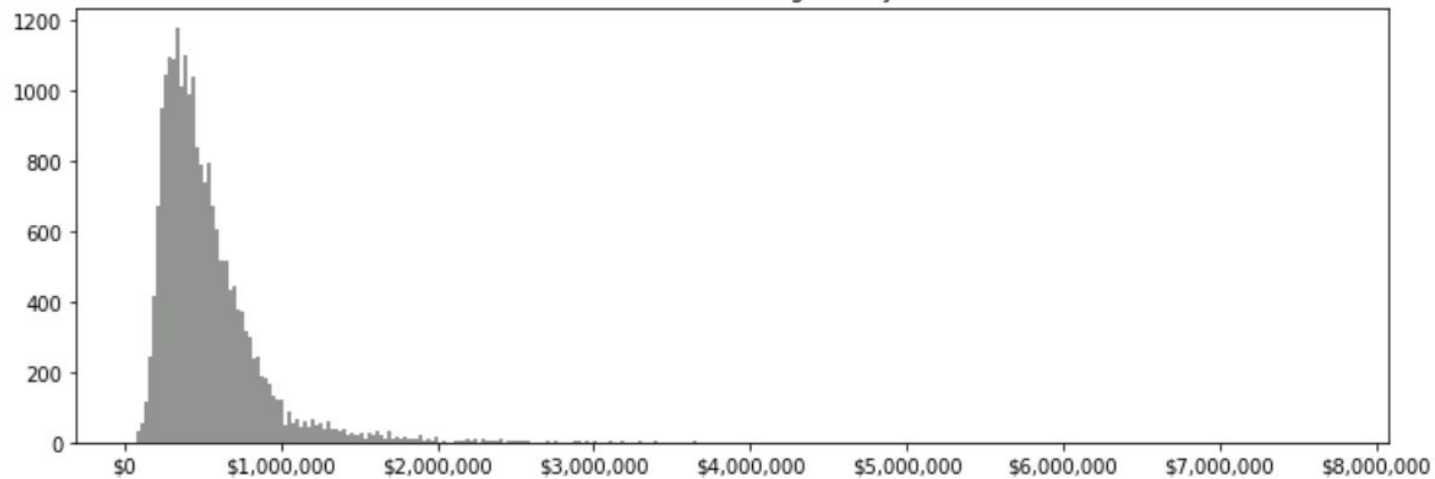# Predicting Home Prices

## King County, Washington

- What are the greatest influencers in the pricing of homes?
- Explore data visually
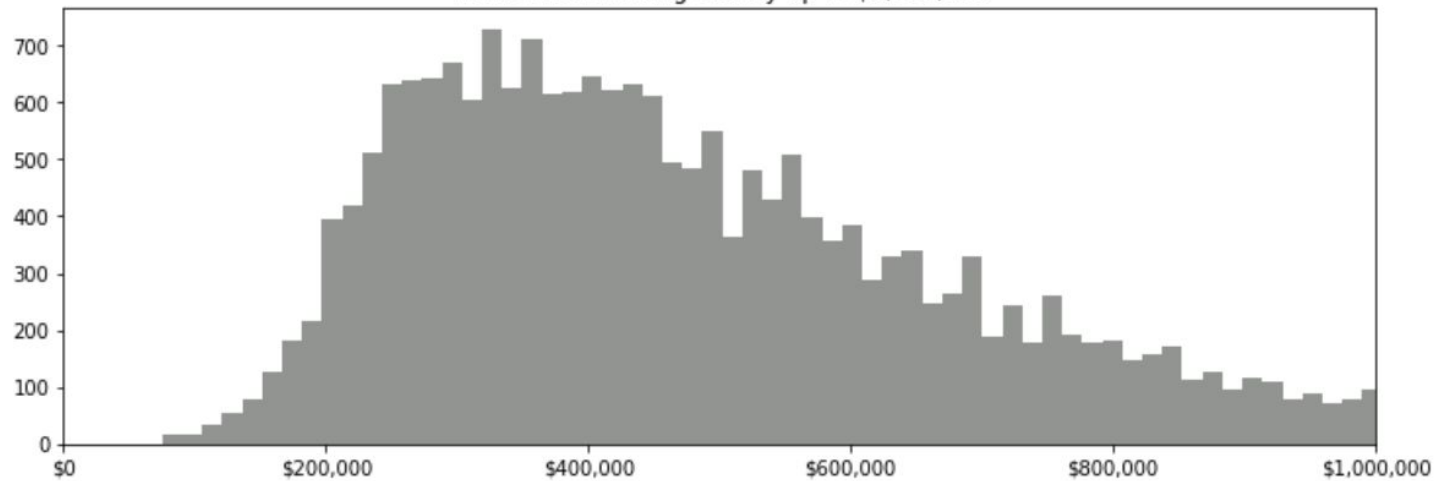- Create linear regression model to predict home prices

# The Dataset

The dataset was made available on Kaggle.com by user harlfoxem.  The data is clean and ready for analysis.

https://www.kaggle.com/harlfoxem/housesalesprediction/data

## Home Prices in King County



## Home Prices in King County up to $1,000,000

# **Correlates to price**
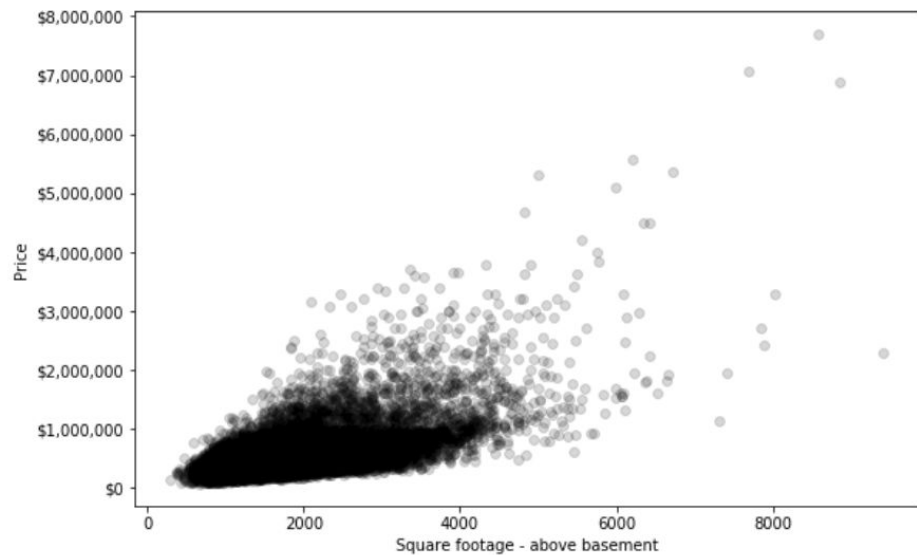
| | |
|---|---|
| sqft_living | 0.702035 |
| grade | 0.667434 |
| sqft_above | 0.605567 |
| sqft_living15 | 0.585379 |
| bathrooms | 0.525138 |
| view | 0.397293 |
| sqft_basement | 0.323816 |
| bedrooms | 0.308350 |
| lat | 0.307003 |
| waterfront | 0.266369 |

# Feature Encoding

To eliminate some of the variability of a feature, we can group each observation into intervals of similarly priced homes. After each observation is in it's bucket, we calculate the average of each bucket, and then assign that average value to a new column in the dataframe.

For example, when I encode the feature 'sqft_living', each home will be grouped into 30 equally populated buckets corresponding to their square footage. I'll calculate the average of each bucket, and it will be assigned to a new feature, which can be used to predict the price.

# Clustering

Kmeans clustering is an unsupervised learning technique used to identify similar groupings in a dataset. We can use clustering to find groups of homes that have similar characteristics, and create models of these smaller subgroups to obtain a more accurate overall model.

# Results of Clustering

# Model using all features

```
        Features    Estimated Coefficients
0       bedrooms            -40546.613566
1      bathrooms             45969.896289
2     sqft_living              117.451295
3        sqft_lot                0.118722
4          floors             5087.980474
5      waterfront            581027.588474
6            view             52480.450568
7       condition             24334.440552
8           grade             93012.019086
9      sqft_above               72.997345
10  sqft_basement               44.453950
11       yr_built            -2702.168375
12   yr_renovated               20.869308
13        zipcode             -586.456083
14            lat            603594.802672
15           long           -209172.985549
16   sqft_living15               19.398367
17     sqft_lot15               -0.390268

Estimated intercept coefficient: 7909720.118896149

Summary Statistics
R-squared value: 0.6962965895543749
Root Mean Squared Error: 193088.65760593285
Mean Absolute Percentage Error (MAPE): 25.838850273124415
```



True prices compared to predicted prices

# Cluster 0 Model

```
                   Features  Estimated Coefficients
0  sqft_living_bucket_price_squared          1.298238e-07
1                       waterfront           5.954120e+05
2                        condition           6.721346e+04
3                            grade           6.211902e+04
4                      sqft_living           9.376187e+01

Estimated intercept coefficient: 135601.22614349728

Summary Statistics
R-squared value: 0.3651722230592157
Root Mean Squared Error: 534305.9167235789
Mean Absolute Percentage Error (MAPE): 18.801434341787772
```



True prices compared to predicted prices

# Cluster 1 Model

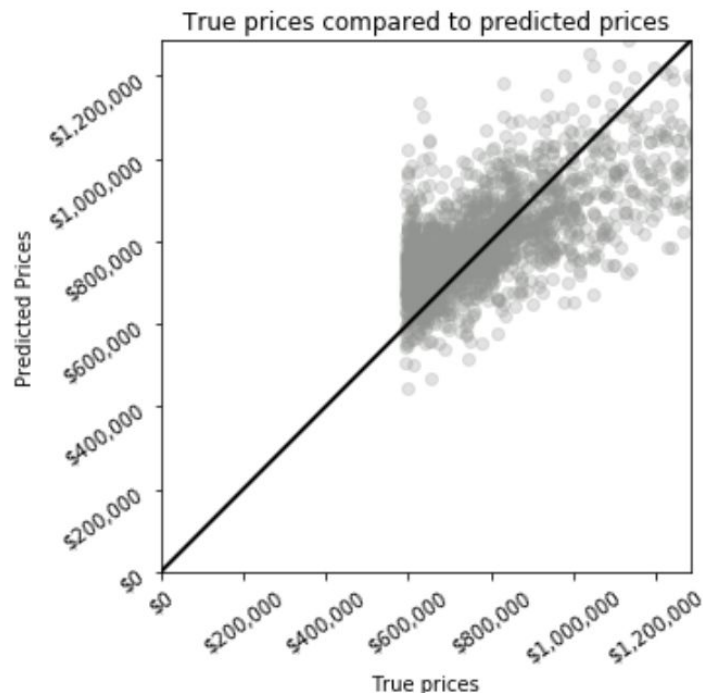|    | Features | Estimated Coefficients |
|----|----------|------------------------|
| 0  | bathrooms_bucket_price_squared | 4.058957e-07 |
| 1  | grade_bucket_price | 2.143519e-01 |
| 2  | bathrooms | -4.196916e+04 |
| 3  | waterfront | 2.167078e+05 |
| 4  | condition | 1.528197e+04 |
| 5  | grade | 4.867775e+04 |
| 6  | yr_built | -1.160438e+03 |
| 7  | zipcode | -6.985400e+02 |
| 8  | lat | 2.144938e+05 |
| 9  | long | -5.242644e+05 |
| 10 | sqft_living | 6.546359e+01 |
| 11 | sqft_living15 | 4.725114e+01 |

Estimated intercept coefficient: -3770611.0657177963

Summary Statistics
R-squared value: 0.44342874856398595
Root Mean Squared Error: 139180.63909621062
Mean Absolute Percentage Error (MAPE): 12.935603453041134



True prices compared to predicted prices

# Cluster 2 Model

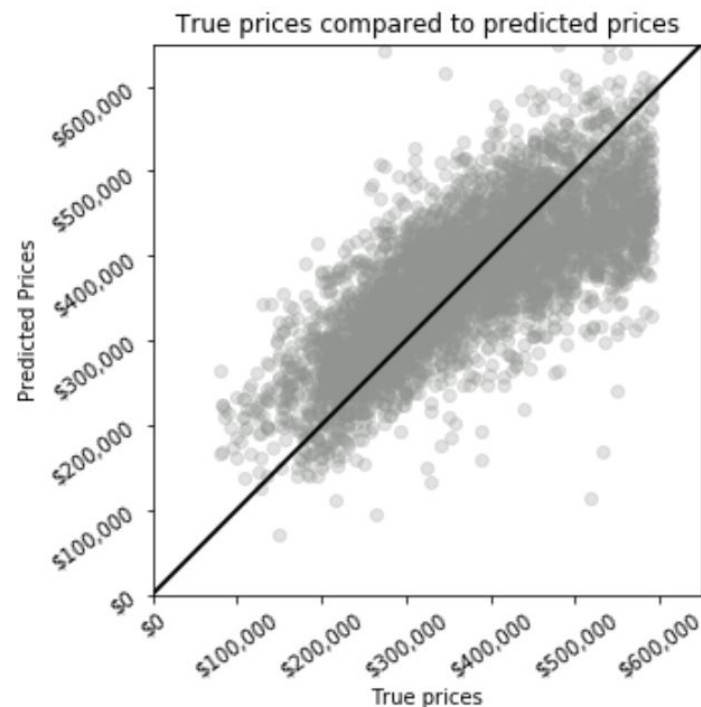|    | Features | Estimated Coefficients |
|----|----------|------------------------|
| 0  | bedrooms | -1.229909e+04 |
| 1  | grade_bucket_price_sqrt | 3.231988e+02 |
| 2  | sqft_living_bucket_price_squared | 7.410774e-07 |
| 3  | bathrooms | 2.437804e+04 |
| 4  | waterfront | 1.295705e+05 |
| 5  | condition | 1.240750e+04 |
| 6  | grade | 3.218932e+04 |
| 7  | yr_built | -8.241704e+02 |
| 8  | zipcode | -4.188423e+01 |
| 9  | lat | 4.057994e+05 |
| 10 | long | 3.945118e+04 |
| 11 | sqft_living | 1.685261e+01 |

Estimated intercept coefficient: -8974715.60852908

Summary Statistics
R-squared value: 0.5664798337761238
Root Mean Squared Error: 75550.38559401885
Mean Absolute Percentage Error (MAPE): 17.925521103105808



True prices compared to predicted prices

# Clustering and Feature Encoding

9.28% improvement

| | MAPE |
|---|---|
| All Features - no clustering | **25.84%** |
| Cluster 0 | 18.80% |
| Cluster 1 | 12.94% |
| Cluster 2 | 17.93% |
| Average of clusters | **16.56%** |

# Conclusion

It should be noted that this model is only trained on home sales in King County, Washington between 2014 and 2015. It's use is limited to the same geographical area and time bounds, and should not be used to make predictions of home prices far into the future.