

Predicting Home Values

By: Caitlyn Carney

03/23/2021

1
Executive
summary
slide

2
Overview

3
Preparing
Data

4
Exploring
Data

5
Modeling

6
Conclusion

7
Appendix

Big Idea

The number of bedrooms, bathrooms and amount of square footage affects home value prices in southern California.

Goals

Create a regression model to predict home value in southern California as accurately as possible.

The Process

- Gather and clean Zillow data
- Explore data and finding what it means.
- Evaluating and modeling the data.

Appraisal value in southern California is driven by number of bedrooms, bathrooms, as well as square footage.

Square feet, bedrooms, and bathrooms do have an affect the appraisal value of homes in southern California.

Key Findings

Takeaways

Prepare Data

After gathering the data from the Zillow database in the Codeup Sequel server I needed to clean it up like crazy!

How unseemly!

parcelid	id	airconditioningtypeid	architecturalstyletypeid	basementsqft
14634203	2026522	1.0		NaN
11721753	616260		NaN	NaN
11289917	2061546	1.0		NaN
11637029	2554497	1.0		NaN
11705026	1834372		NaN	NaN

So many blank values!

airconditioningtypeid	26358
architecturalstyletypeid	38481
basementsqft	38555
	...
taxdelinquencyyear	37314
censustractandblock	144

I needed something easier to handle with columns not missing half of their values.

Prepare Data

I ended up going from 38,582 rows to 32,542 rows by cleaning.

For more detail please see prepare.py on GitHub

- Dropped columns with less than 35 thousand non null values.
- Dropped unhelpful columns.
 - Renamed columns.
- Handled outliers in some columns.
 - Drop left over null values

bathrooms	bedrooms	square_feet	fips	full_baths
2.0	3.0	1125.0	6059.0	2.0
2.0	3.0	1316.0	6037.0	2.0
2.0	3.0	1458.0	6037.0	2.0
2.0	3.0	1766.0	6037.0	2.0
1.0	2.0	1421.0	6037.0	1.0

Now thats what I call clean!

Prepare Data

After this, I focused my data onto these features, split my data, and scaled the split data.

For further detail please see `prepare.py` on GitHub

Decision time! What columns to focus on!



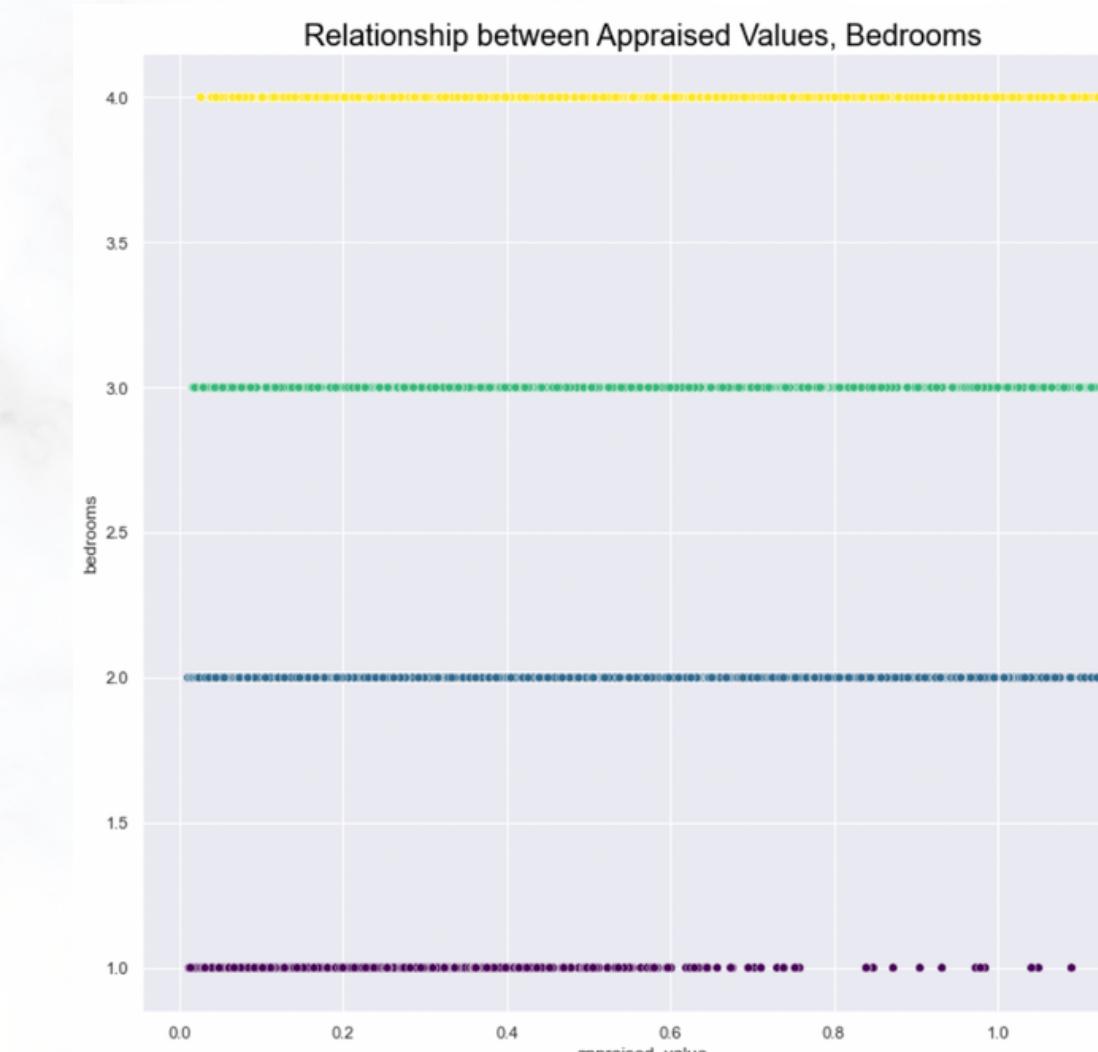
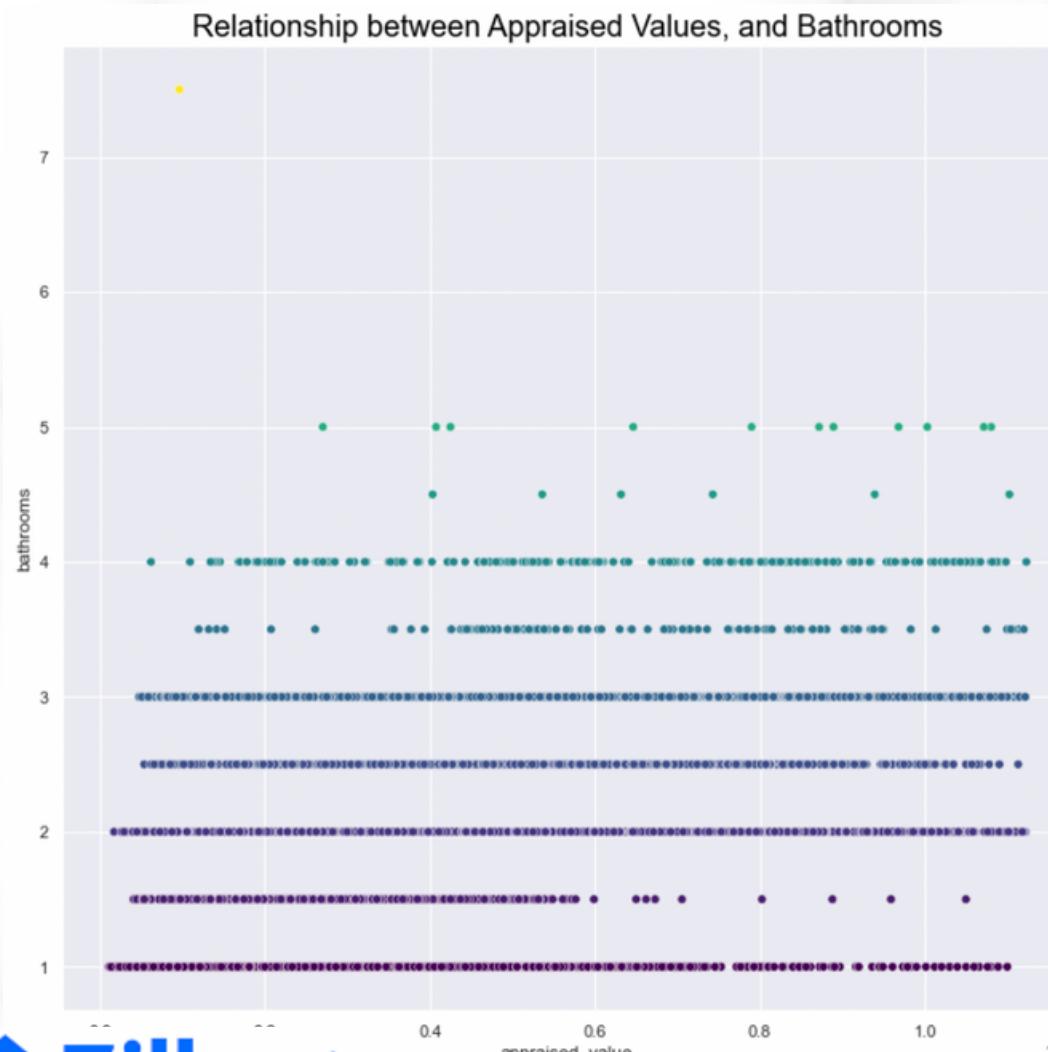
I will be focusing on square feet, bedrooms, and bathrooms today.

Explore Data

Now that the data is clean and easy to read, I needed to take a deep dive into the data and find relationships, abnormalities, etc.

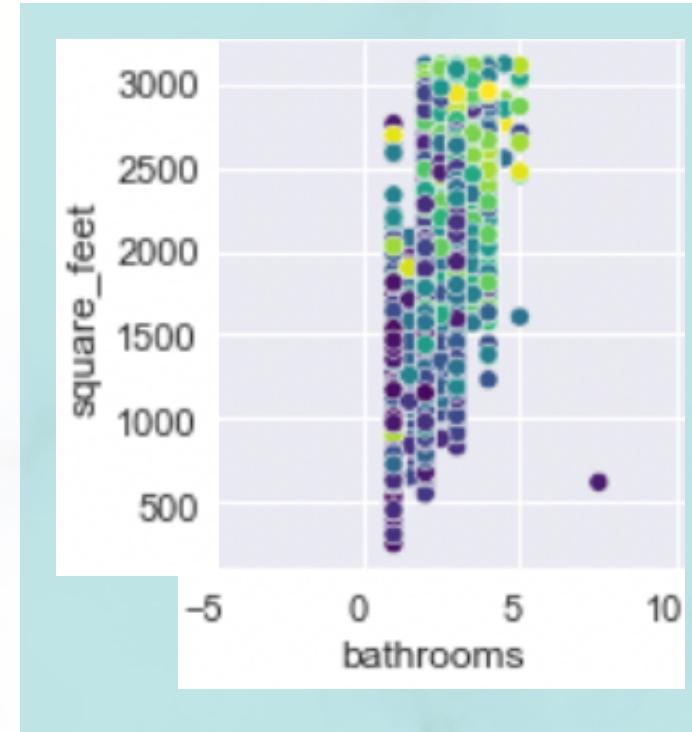
Using Scatterplots to see the Relationships to Appraised Value

- There seems to be a noticeable difference in the amount of homes with a low appraised value and low square footage vs. homes with a higher appraised value's square footage
- Homes with 1 bedroom seem to drop off when appraised values get higher while the 2+ bedroom houses seem to continue rising with the appraised value

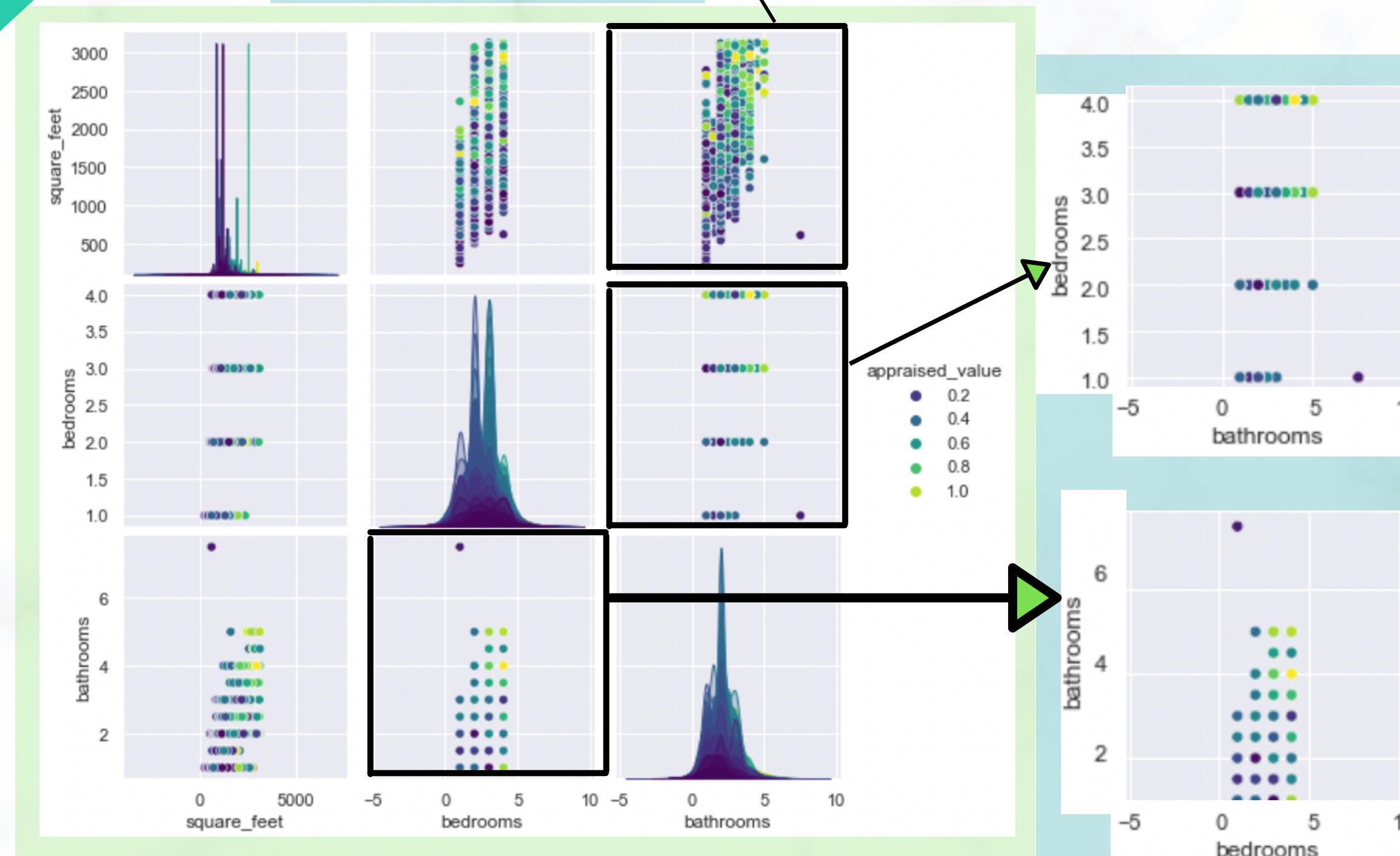


Explore Data

Here I wanted to see how appraised value was affected by the relationship between other feature relationships.



- It seems that homes with more bathrooms and larger amount of square feet tend to be appraised at a higher amount.
- The yellow which indicated the higher appraisal rates, are more heavily indicated the higher the bathroom number and amount of square feet gets.



We can see a gradual increase in appraisal price in the bedrooms and bathrooms relationships.

- When both features are high we move into higher appraisal values

This is indicated by the green in the top right.

- The lower amount of bedrooms and bathrooms, the lower the appraisal value gets

This is indicated by the dark blue on the bottom left.

Explore Data

Take a deep dive into the bathroom feature. I ran a correlation test for this features.

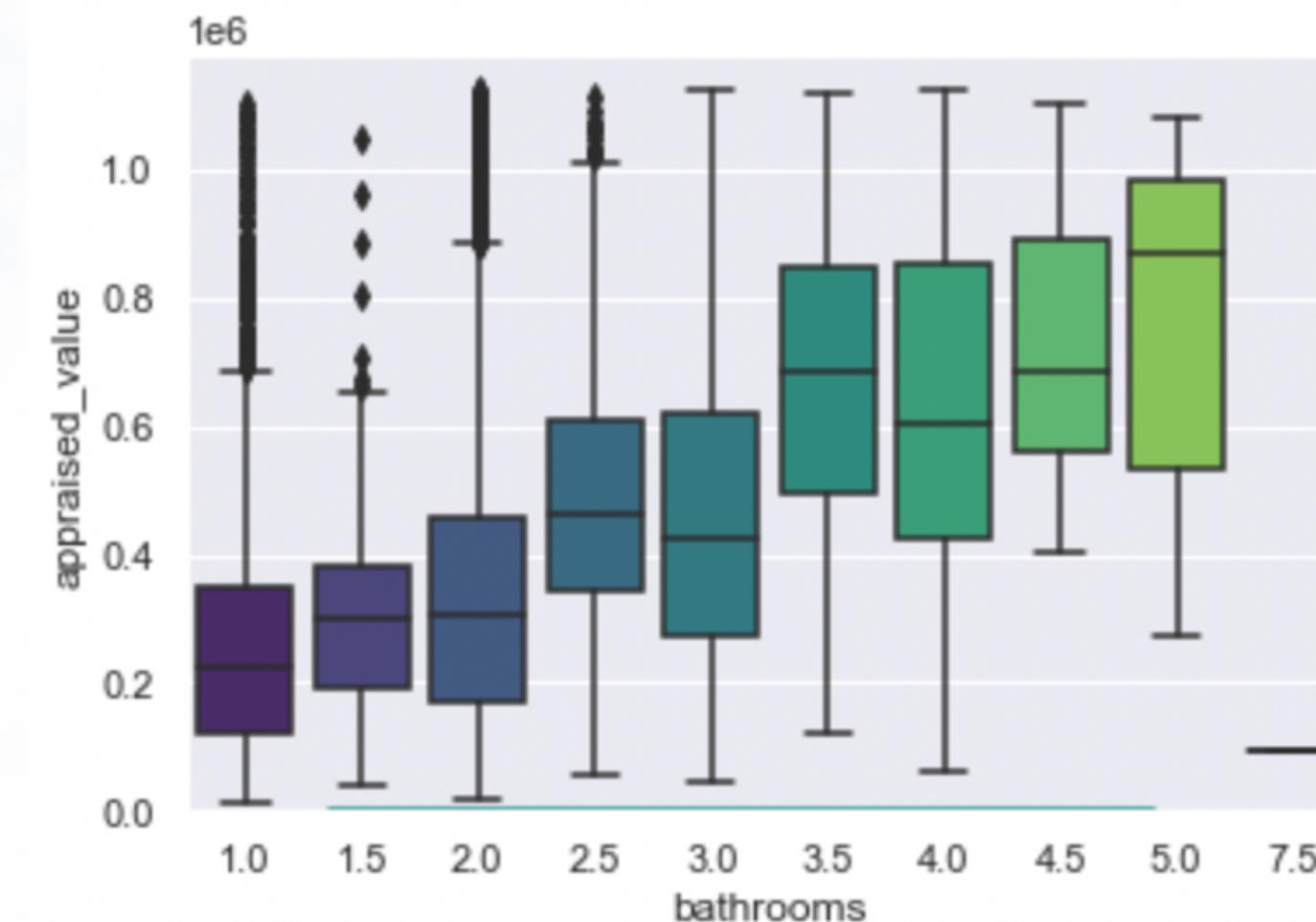
For further detail please see explore.py on GitHub

Null Hypothesis:

"There is no correlation between number of bathrooms and appraised value."

Alternative Hypothesis

"There is a correlation between number of bathrooms and appraised value."



The correlation between Bathrooms and the Appraised value
is: 0.356

The P value between Bathrooms and Appraised Value is: 0.0

I reject the null hypothesis and now move forward with our alternative hypothesis

Explore Data

Take a deep dive into the bedroom feature. By running a correlation test.

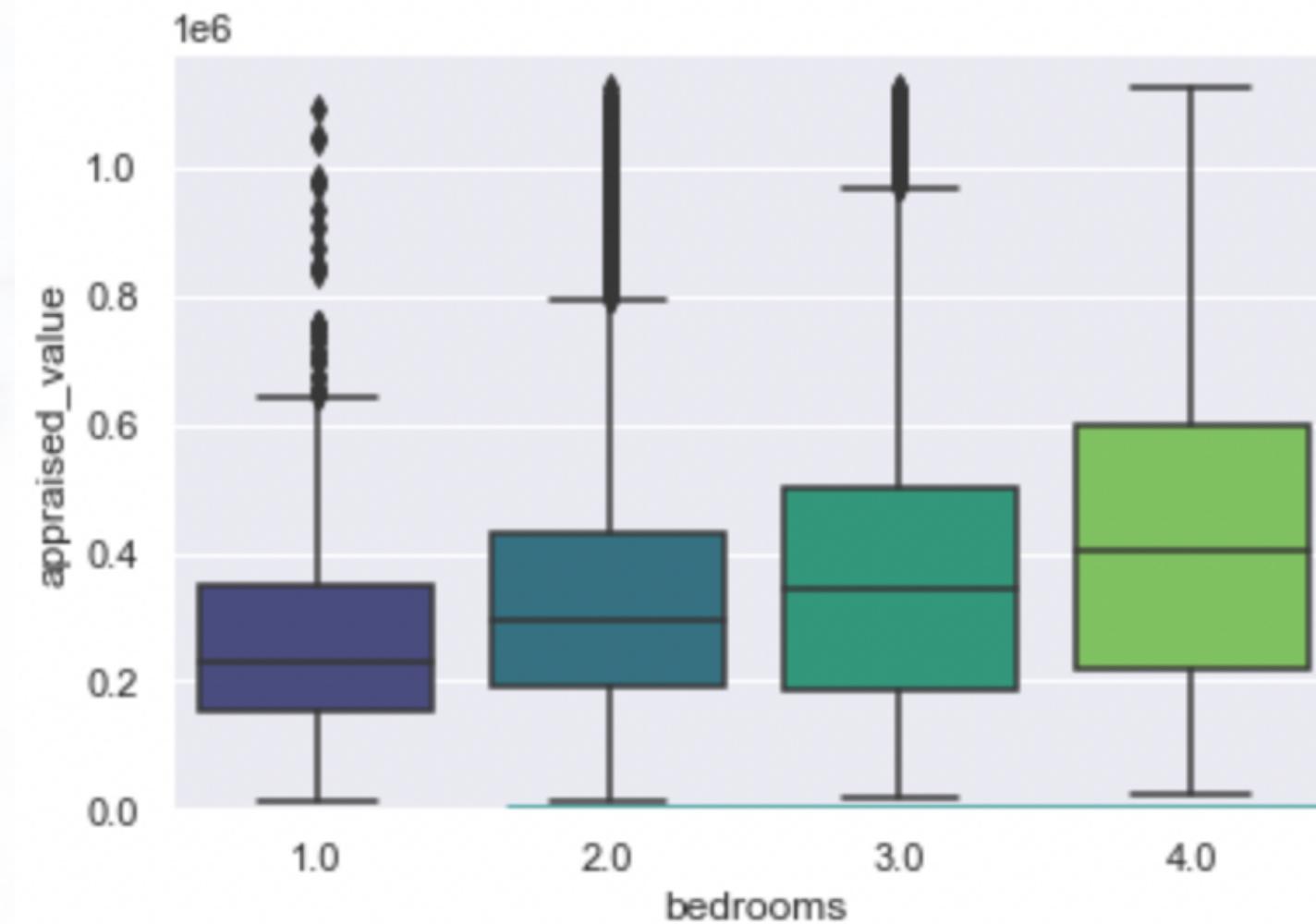
For further detail please see `explore.py` on GitHub

Null Hypothesis:

"There is no correlation between number of bedrooms and appraised value"

Alternative Hypothesis

"There is a correlation between number of bedrooms and appraised value."



The correlation between Bathrooms and the Appraised value is: 0.177

The P value between Bathrooms and Appraised Value is: 7.71

I reject the null hypothesis and now move forward with our alternative hypothesis

Explore Data

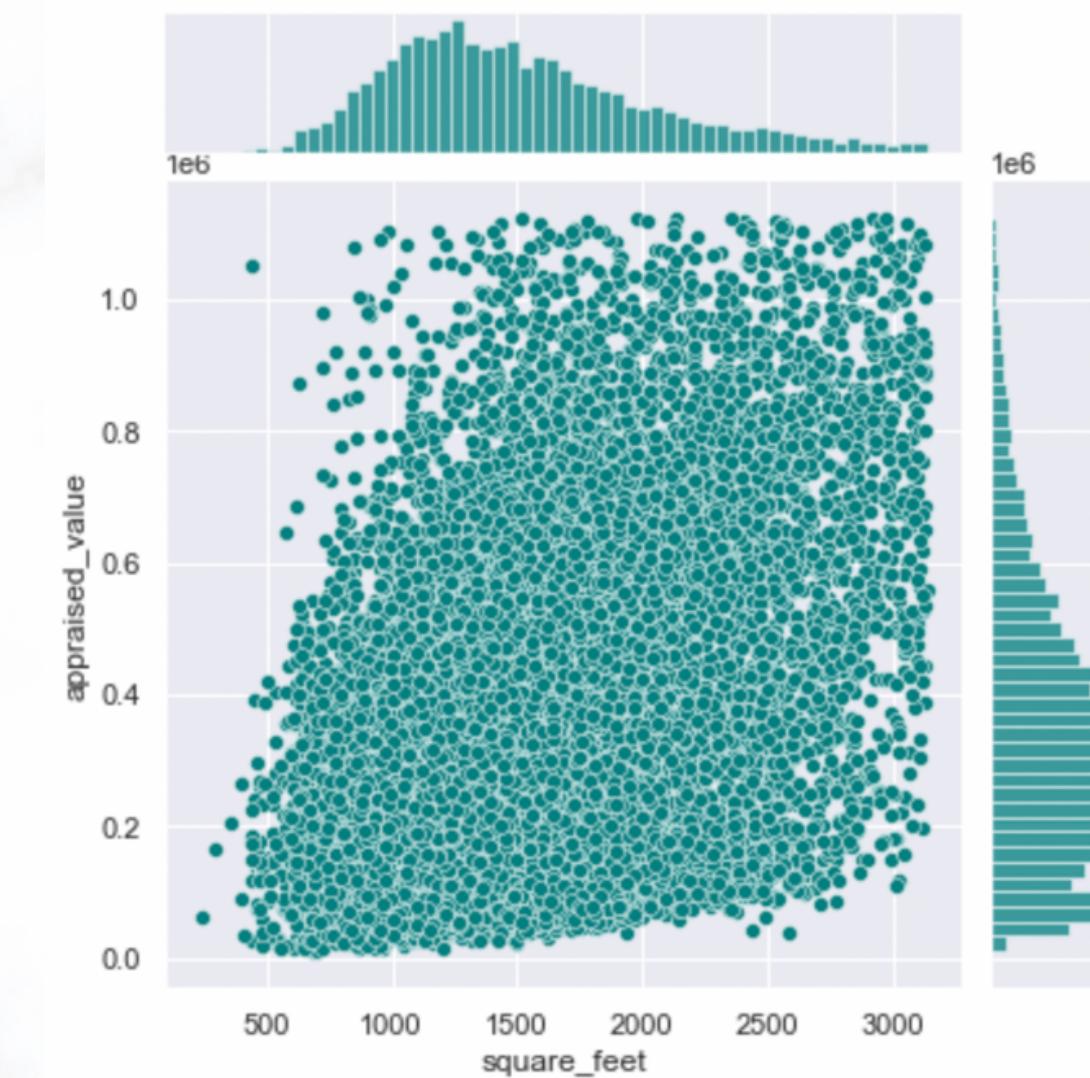
Take a deep dive into the square footage. By running a correlation test.

For further detail please see `explore.py` on GitHub

Null Hypothesis:
"There is no correlation between a homes square footage and appraised value."

Alternative Hypothesis

"There is a correlation between square feet and appraised value."



The correlation between Bathrooms and the Appraised value is: 0.432

The P value between Bathrooms and Appraised Value is: 0.0

I reject the null hypothesis and now move forward with our alternative hypothesis

Model

Now that I have gotten to know the data, it is now time to evaluate and create models!

For further detail please see the evaluate.py and model.py on GitHub

SSE = 911614726222345.5

SSE Baseline = 941263447006771.6

MSE = 28013481845.687

MSE baseline = 28924572767.709

RMSE = 167372.285

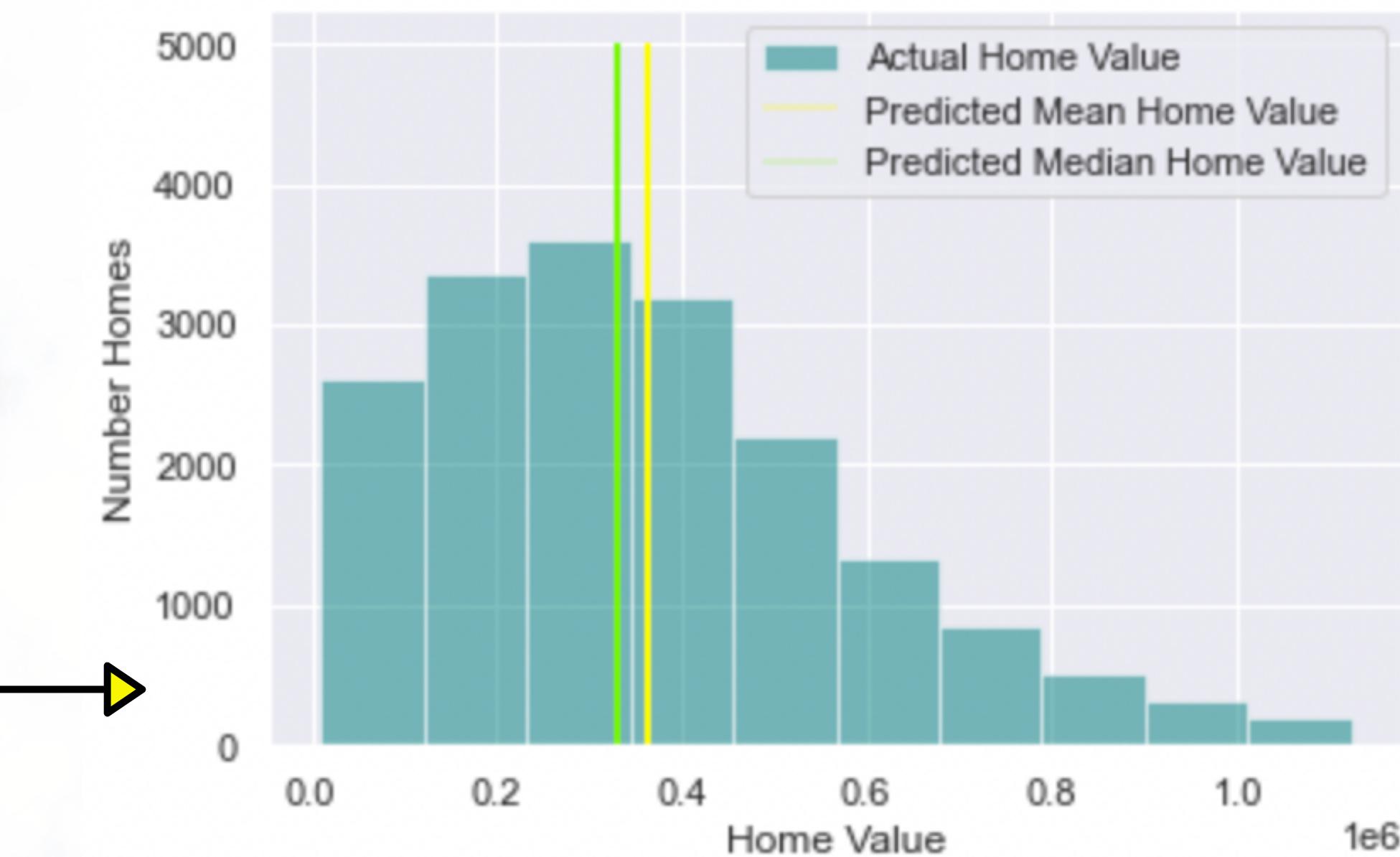
RMSE baseline = 170072.257

My predicted mean was : My predicted median was :

364,215.53

329,107

I will move forward using my predicted median!



Model

Now I take a look at each model and determine the best one!

For further detail please see the evaluate.py and model.py on GitHub

LinearRegression Model
In-Sample: 201813.737
Out-of-Sample: 202696.951

Lasso Model
In-Sample: 20364445665.7664
Out-of-Sample: 20542786732.111

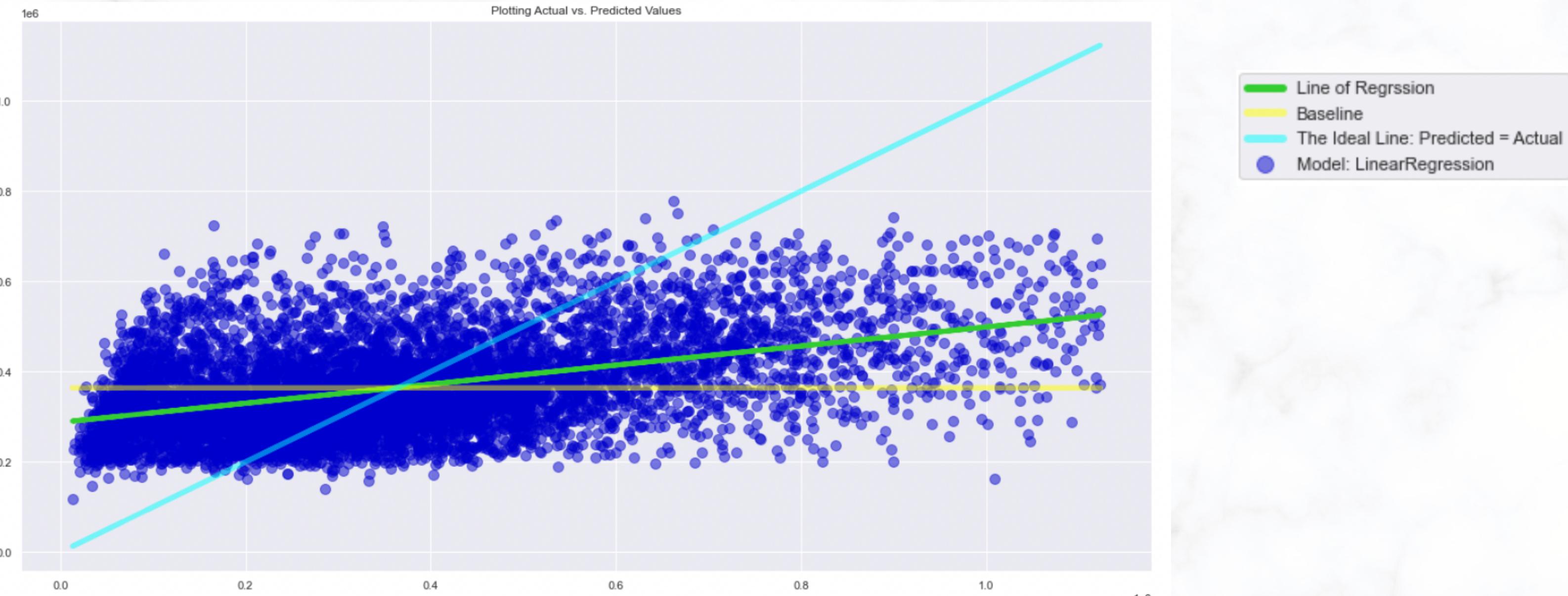
Tweedie Model
In-Sample: 25826248340.196
Out-of-Sample: 26018336332.939

Polynomial Model
In-Sample: 20298379492.33372
Out-of-Sample: 20476129727.74332



RMSE = 167372.285

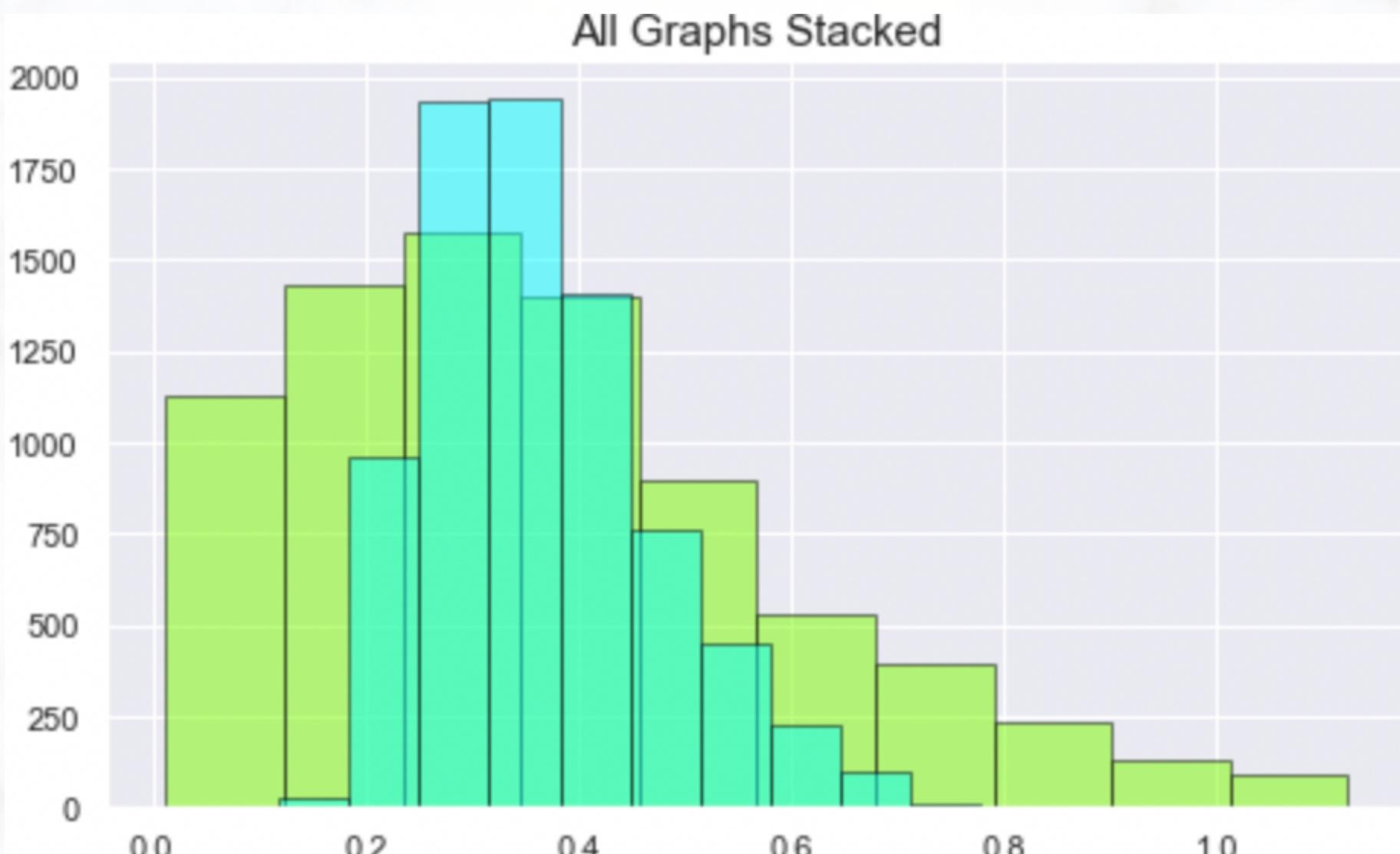
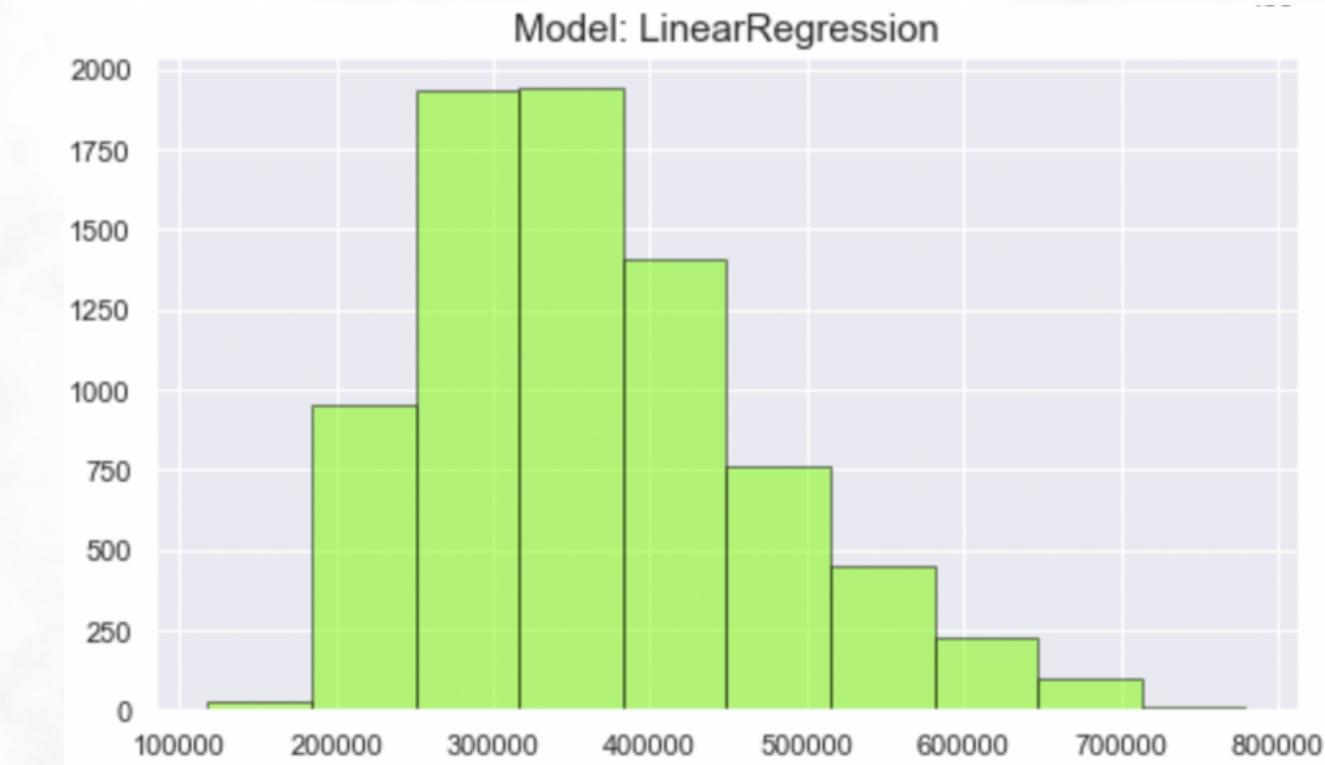
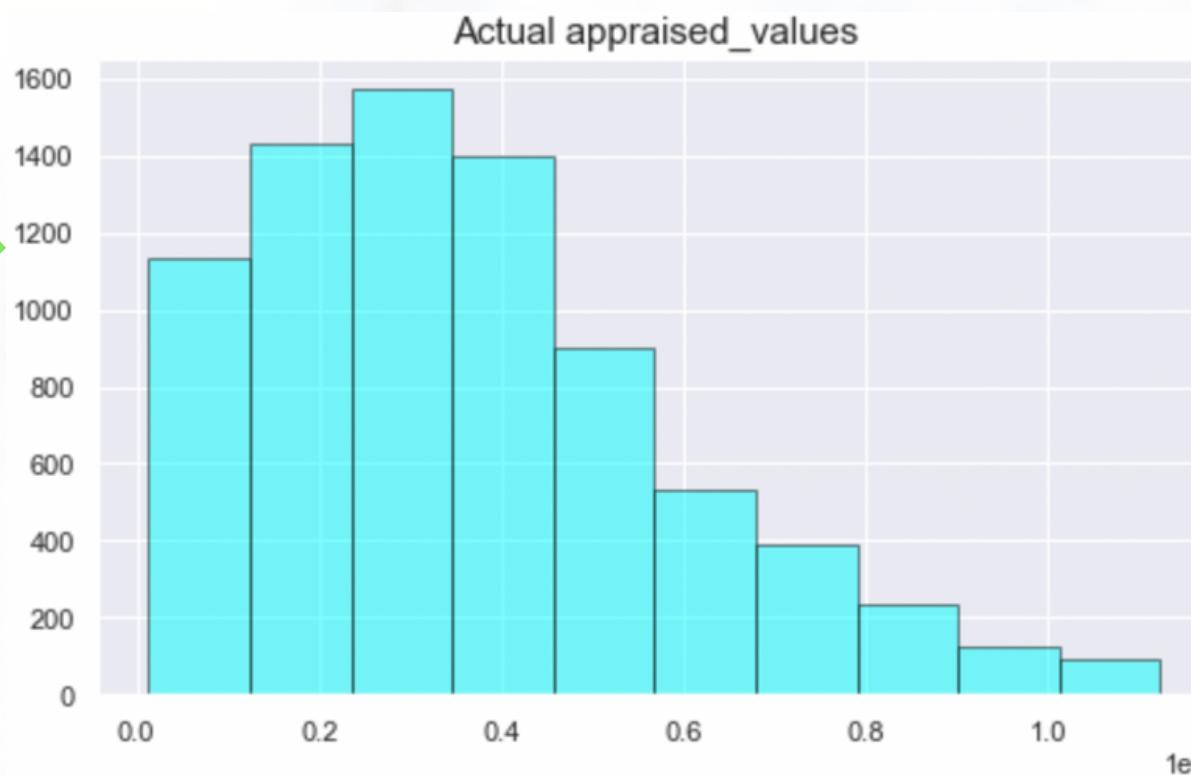
RMSE baseline = 170072.257



Model

Here I visualize the actual appraised value vs. predicted appraised value.

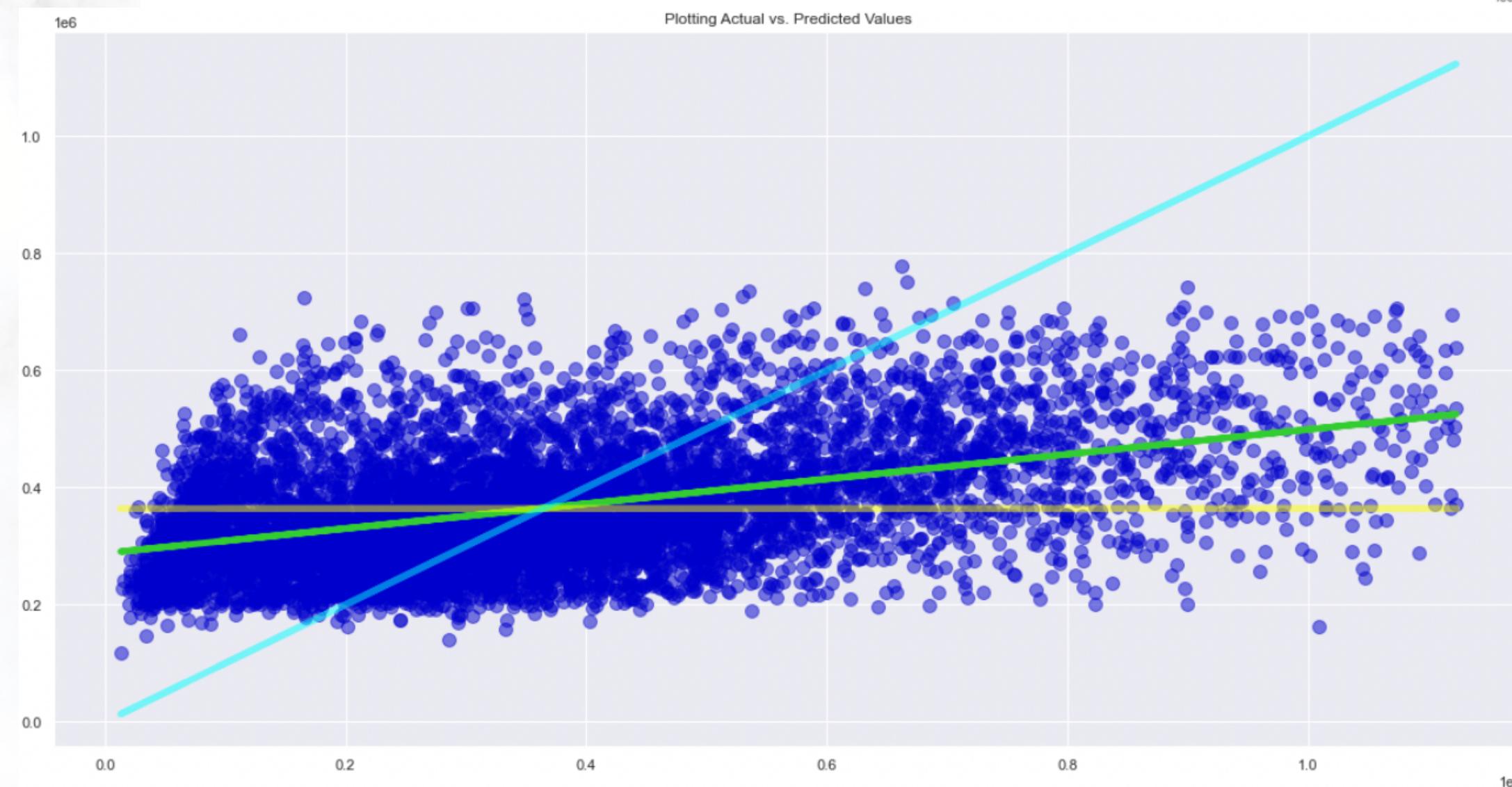
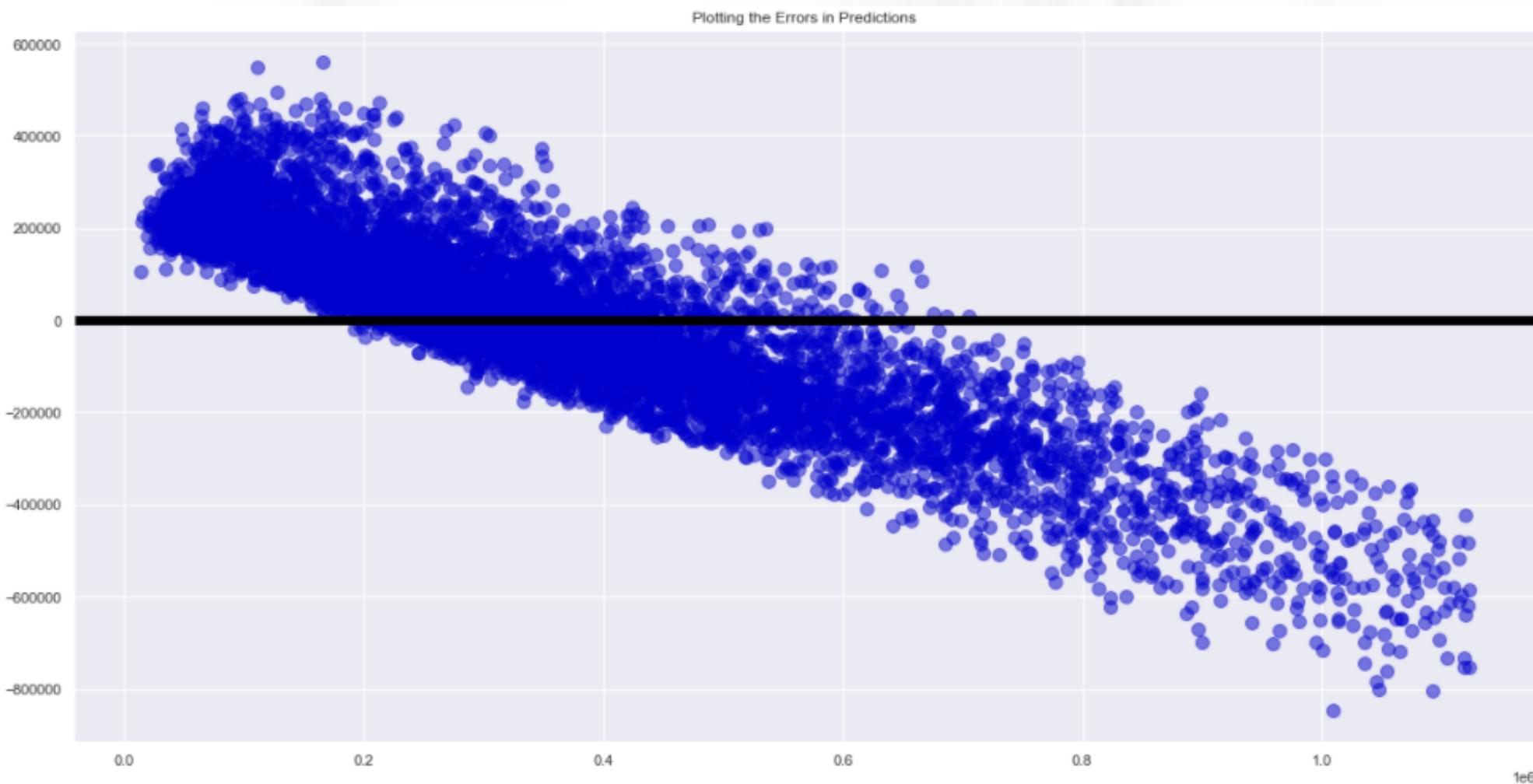
For further detail please see the evaluate.py and model.py on GitHub



Model

.....
Here I visualize the errors in my prediction model and actual vs. predicted.

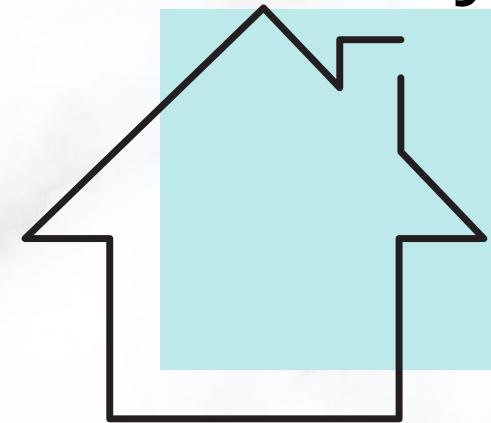
For further detail please see the evaluate.py and model.py on GitHub



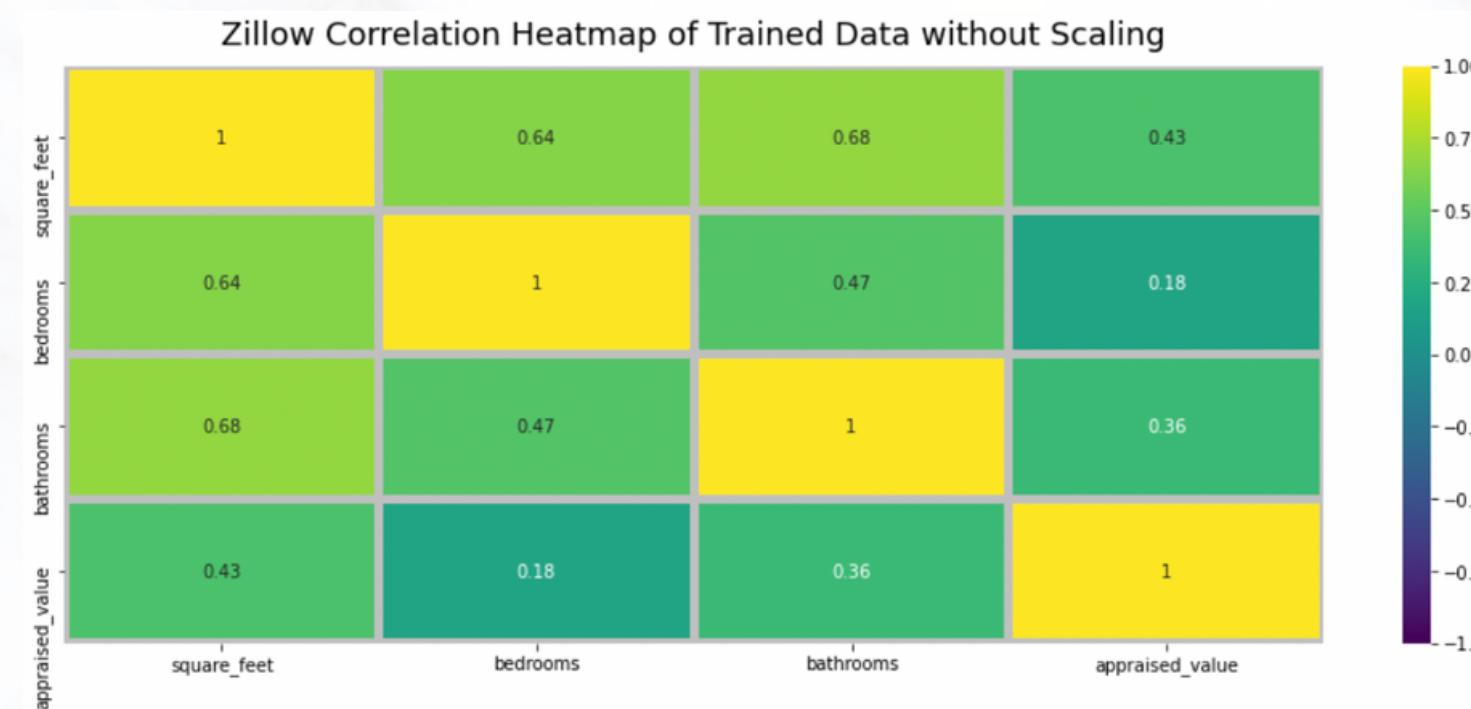
Conclusion

By gathering, cleaning, exploring, and modeling the Zillow data; I learned a lot.

In southern California, a homes appraised price is affected in many ways.



Square feet, number of bedrooms, and the number of bathrooms are just to name the top 3 drivers.



By looking at these top 3 drivers we can see a correlation between them and appraisal values.

Conclusion

By gathering, cleaning, exploring, and modeling the Zillow data; I learned a lot.

Even though these three drivers give a good prediction of appraisal values, we can't discount the possibility that features such as zip code, city, etc. have a role in the pricing as well.

Although these other features may have a smaller correlation that doesn't mean that together don't affect appraisal value in a big way.



Appendix

Follow these links for more information on the creation of the project, steps taken, and myself.

Link to the GitHub Repository

Zillow Home Value Prediction Modeling

- README.md
- acquire.py
- prepare.py
- explore.py
- evaluate.py
- model.py
- project-final.ipynb

Link to the Trello Board

Zillow Predicting Trello Board

Link to the my Linkedin

Caitlyn Carney

