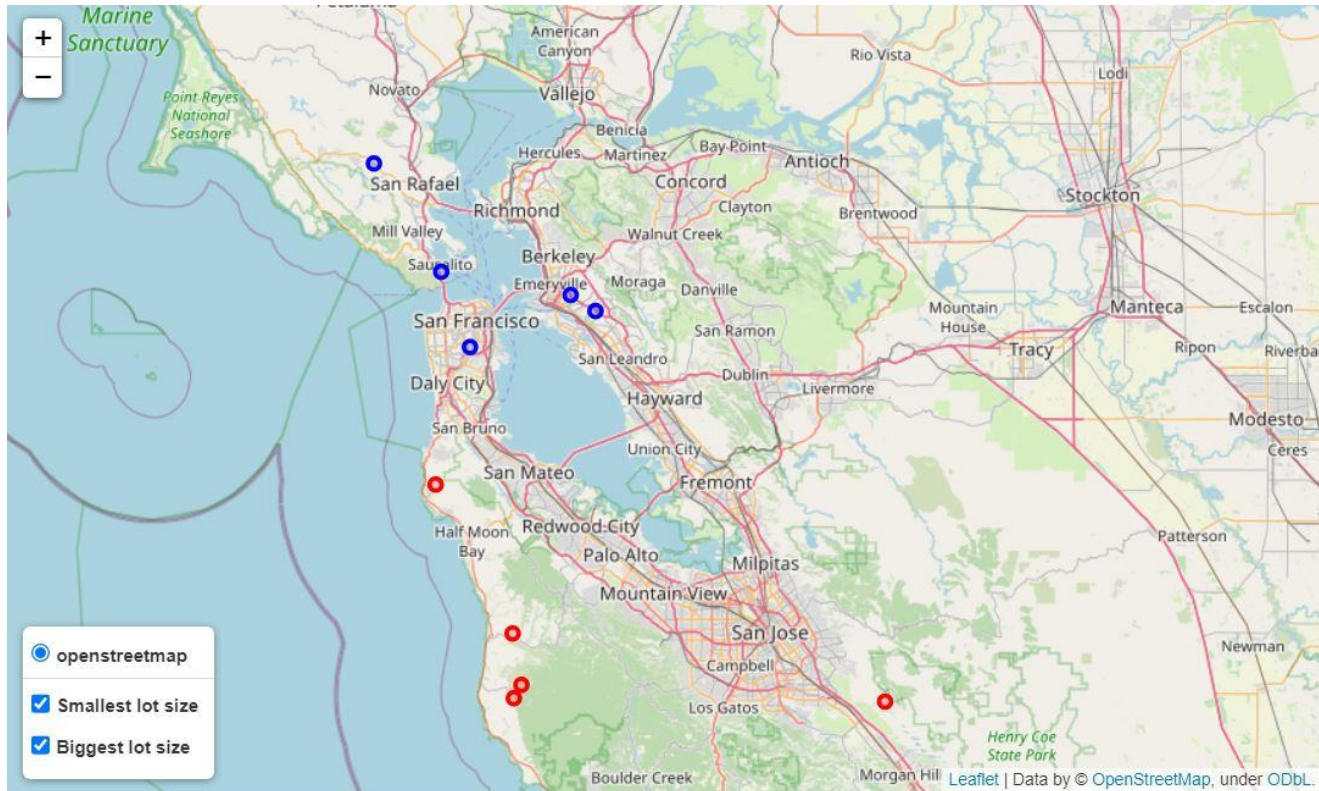


Predicting Bay Area House Prices

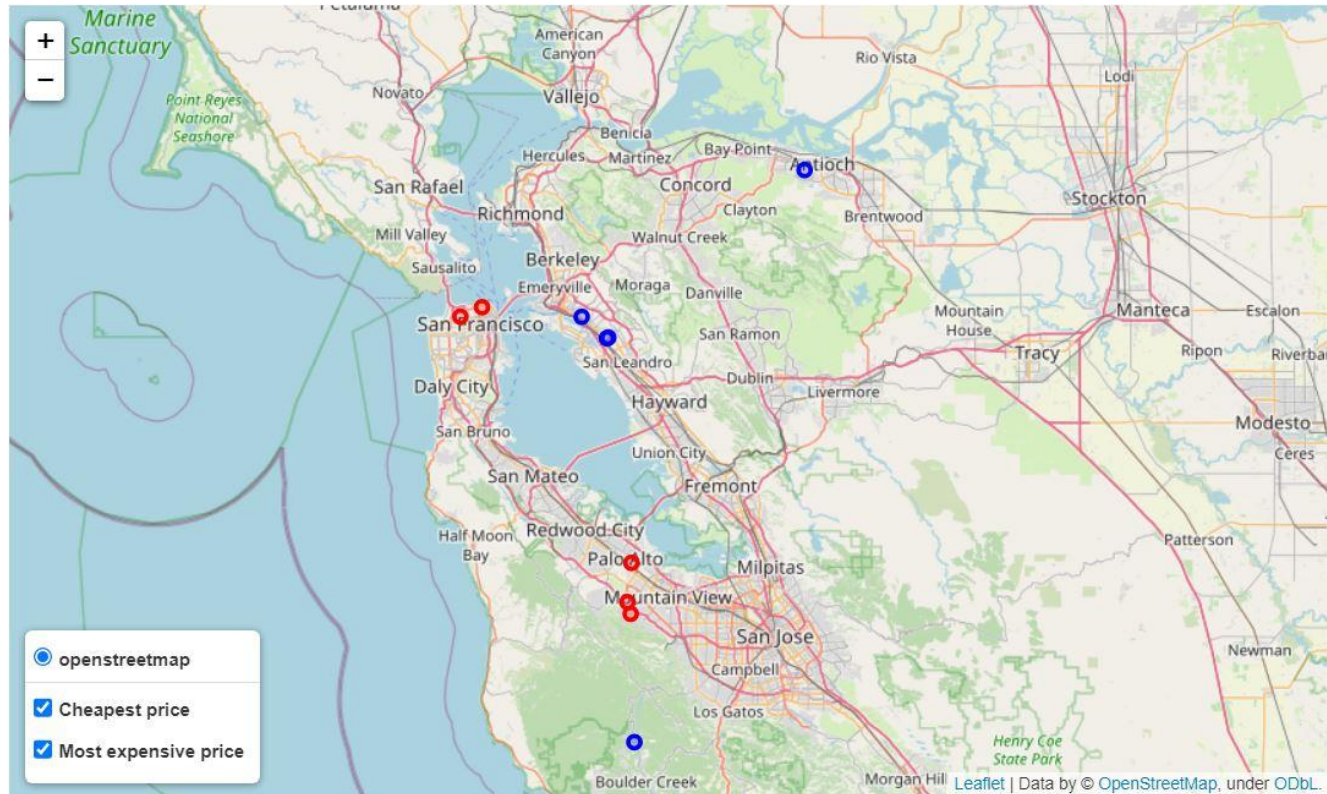
By: Rohan Athalye and Paul Chon



Where are the 5 houses with the biggest lot sizes and 5 houses with the smallest lot sizes located?

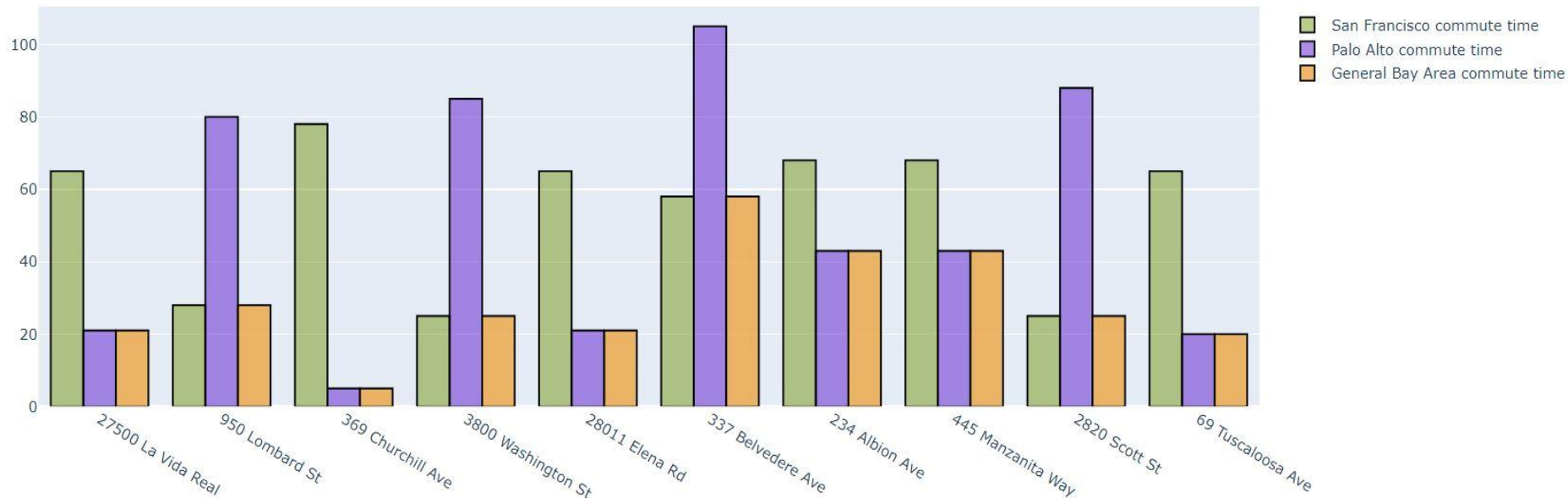


Where are the 5 most expensive and 5 cheapest houses located?



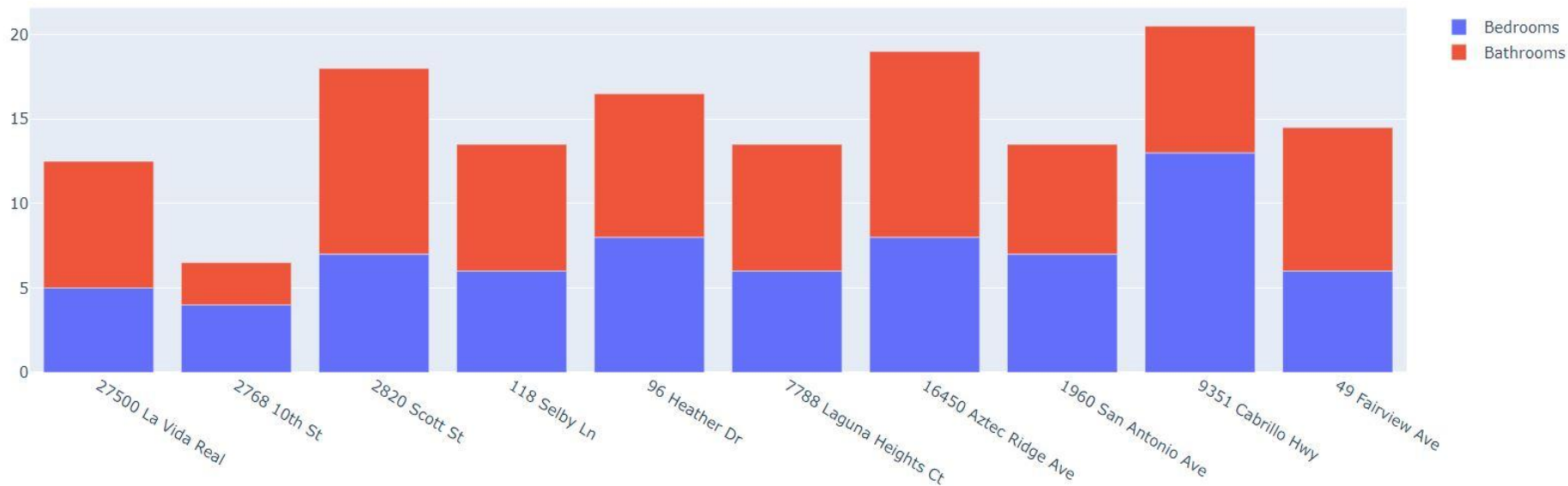
How do commute times by car at 8 AM to San Francisco, Palo Alto, and the general Bay Area compare for the 10 most expensive houses?

Commute Times to San Francisco, Palo Alto, and the General Bay Area by car at 8 AM for the 10 most expensive houses

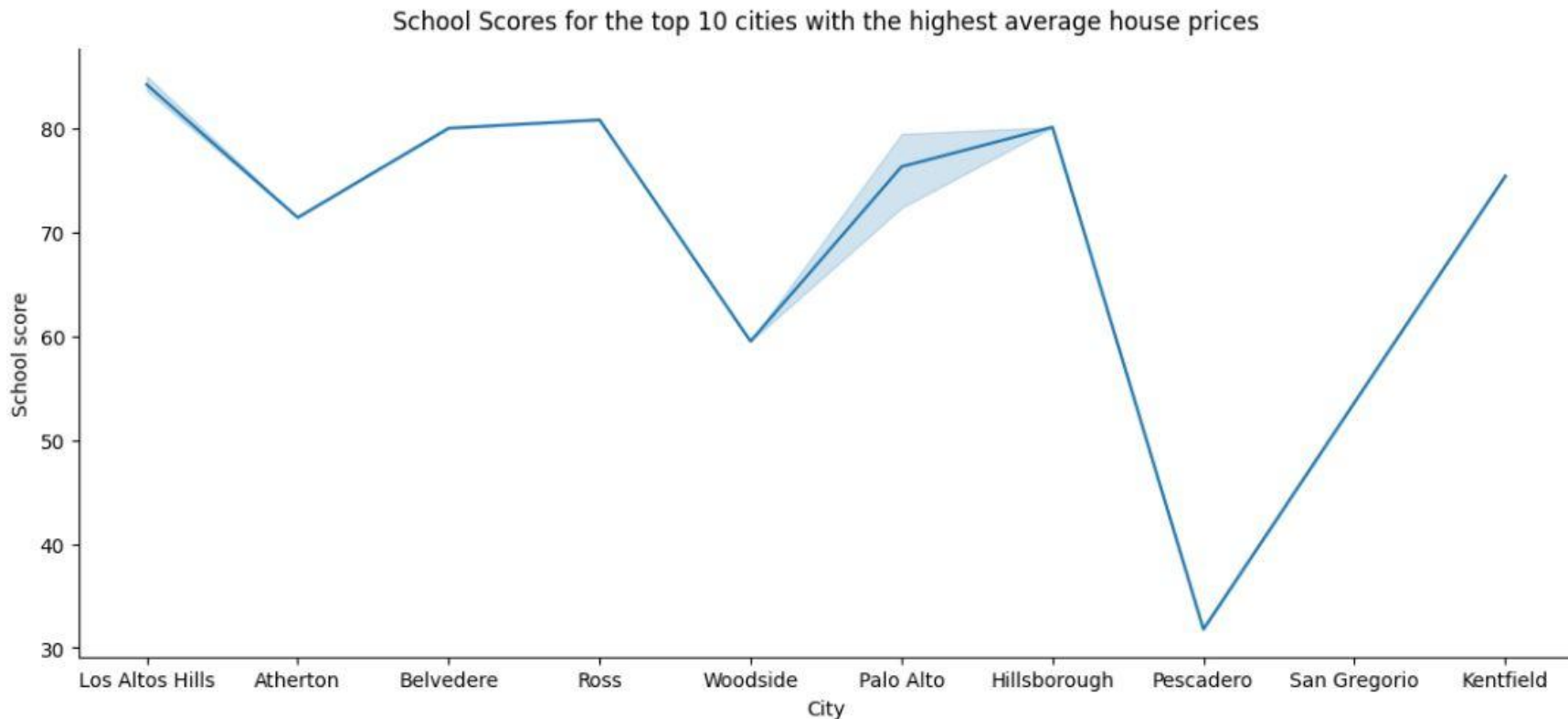


How do the number of bedrooms and bathrooms compare for the top 10 houses with the biggest house sizes?

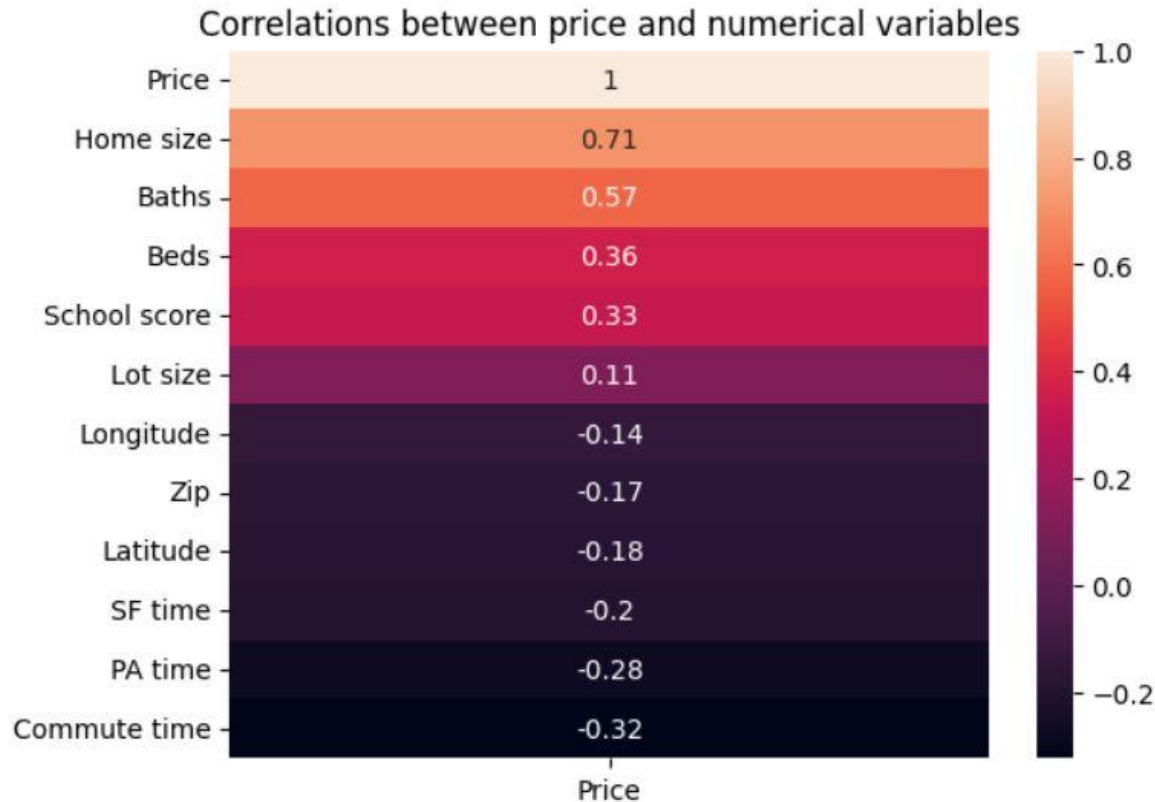
Number of Bedrooms and Bathrooms for the top 10 houses with the biggest house sizes



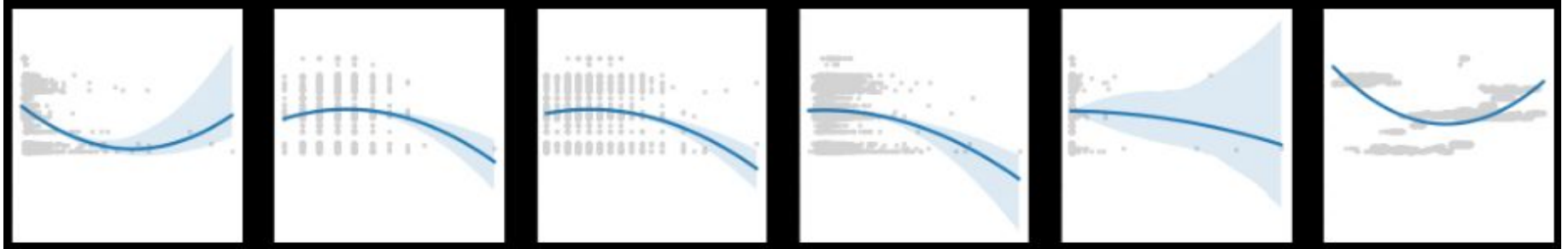
Do cities with expensive house prices also tend to have consistently high school scores?



What are the correlations between price and the numerical variables?



What are the pairwise relationships?



Link to full plot: <https://imgur.com/a/osJwMj5>

Training/Testing Data Strategy

- 80/20 Stratified Shuffle split based on the Baths category

Pipeline:

1. Categorical columns
 - a. Encode
 - b. Impute
2. Numerical columns
 - a. Impute
 - b. Scale
3. PCA (Principal Component Analysis)

Chosen ML Models

1. Linear Regression
 - a. Base test
2. Decision Tree
 - a. Nonlinear
 - b. Comparison
3. Random Forest
 - a. Ensemble learning
4. Support Vector Machine
 - a. Base test, but nonlinear

Performance of Trained Models

- Performance metric: root mean squared error

Model	Mean cross-validation score for training data with all columns	Mean cross-validation score for training data without Address and State columns
Linear Regression	1276886.45	1276827.61
Decision Tree	1535360.50	1503050.50
Random Forest	1181598.45	1144059.25
Support Vector Machine	2141957.29	2141957.08

Best ML Model Performance using Test Set

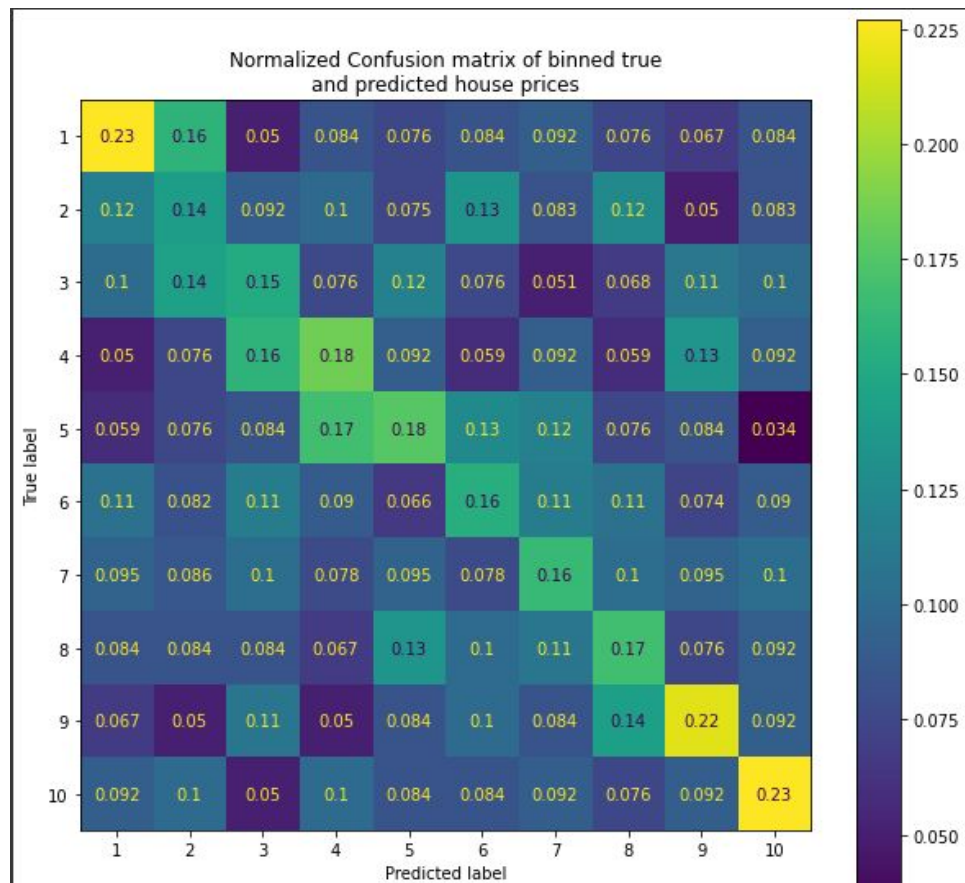
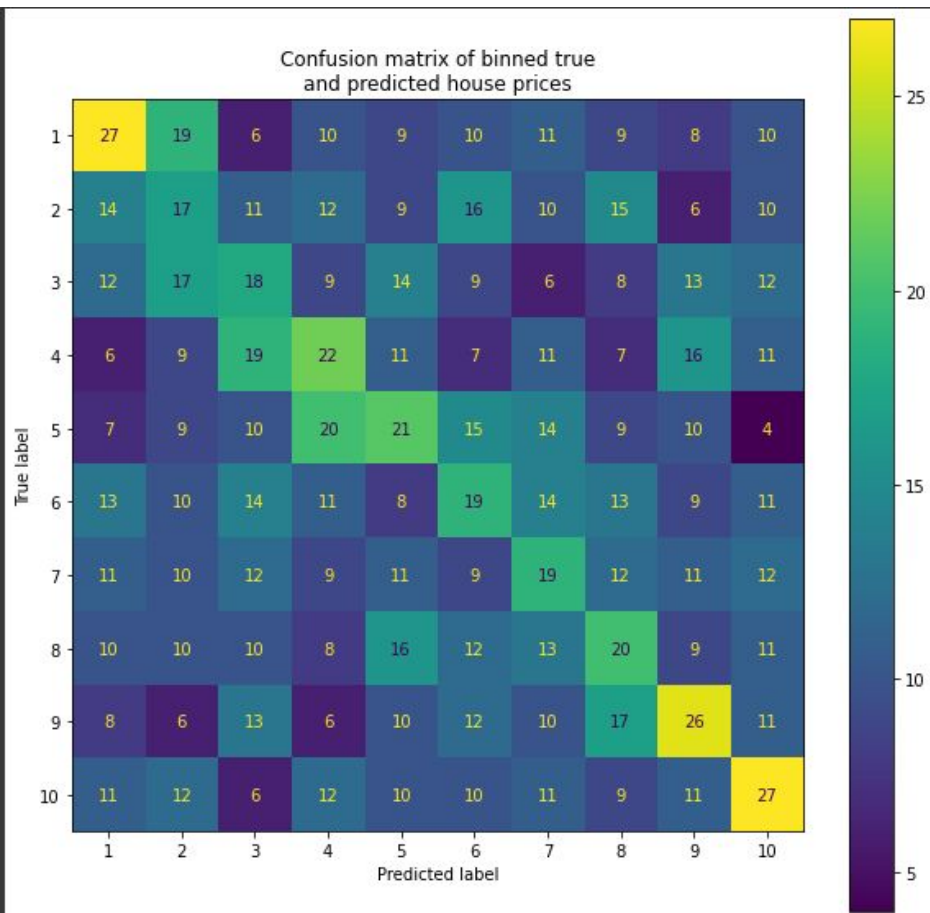
Model: Random Forest

Features:

- max_features=30
- n_estimators=38

RMSE (train): 1161480

RMSE (test): 1237363



Challenges and Thoughts

- Creating the confusion matrices
- Pre-processing
 - Deciding what to do
 - Feature selection
- Future
 - Try more ensemble learning
 - Try to better visualize the model performance