

PREDICTING HOUSING PRICES WITH PYTHON



Ben Moss

Data Science Immersive Program

General Assembly

January, 2024

INTRODUCTION

- Two Datasets; train.csv and test.csv
- Objective: Effectively utilize Data Munging, EDA, and modeling techniques to predict the house sales prices from an unknown dataset.
- Identify key variables that play the biggest role in predicting housing prices.
- Provide recommendations to those considering their first time home purchase.

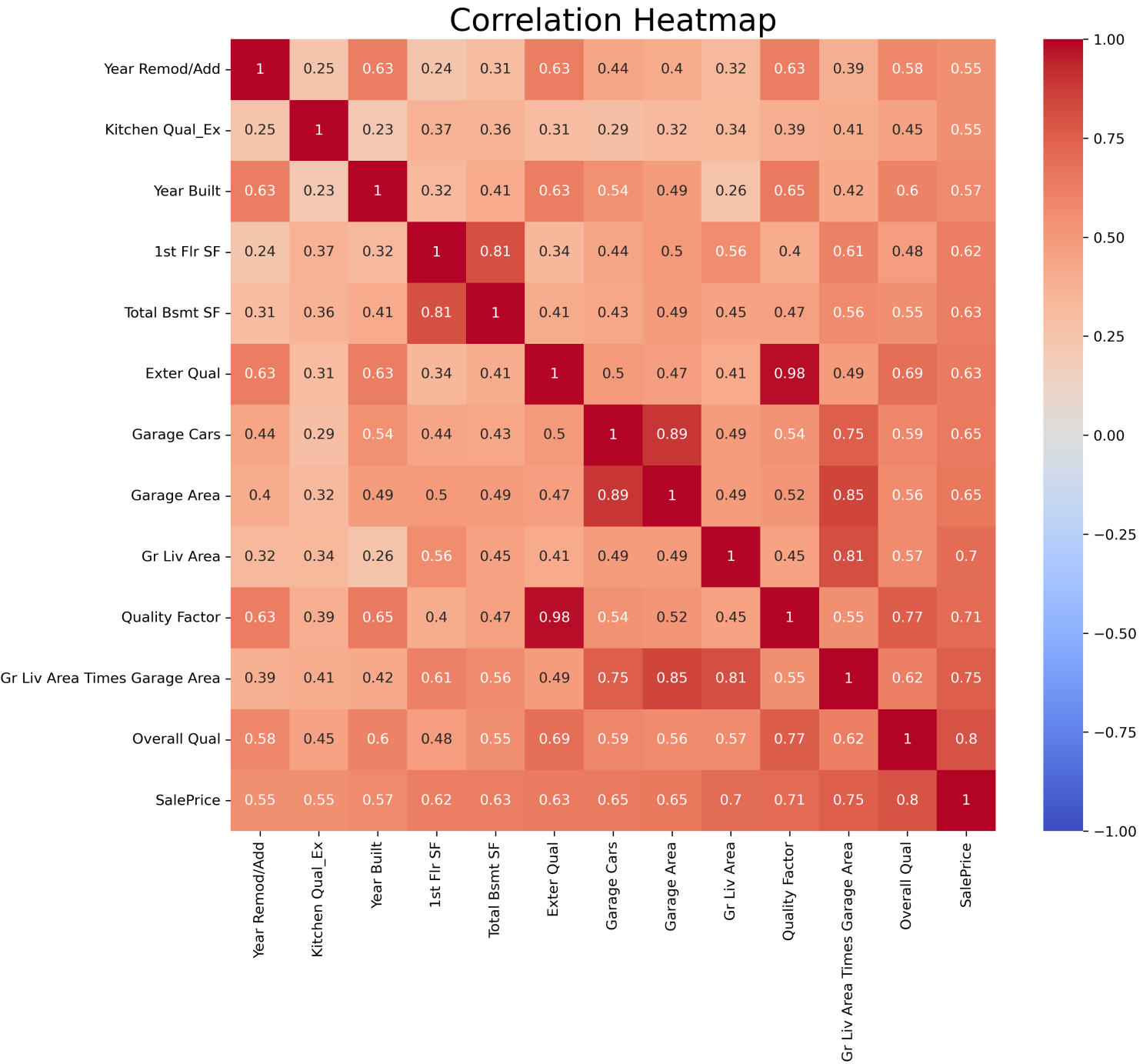


MODELING PIPELINE

- Data Munge; filling NaNs with median values.
- EDA: Heatmaps, histograms, pairplots to discover:
 - Which numerical variables were correlated closely with sales price?
 - Which numerical variables were normally distributed?
 - Which numerical variables had a linear correlation with minimal heteroscedastacity?
- Created Dummy variables and Feature variables that based on commonplace parameters that affected sale price (exterior condition, year build and remodeled etc.)

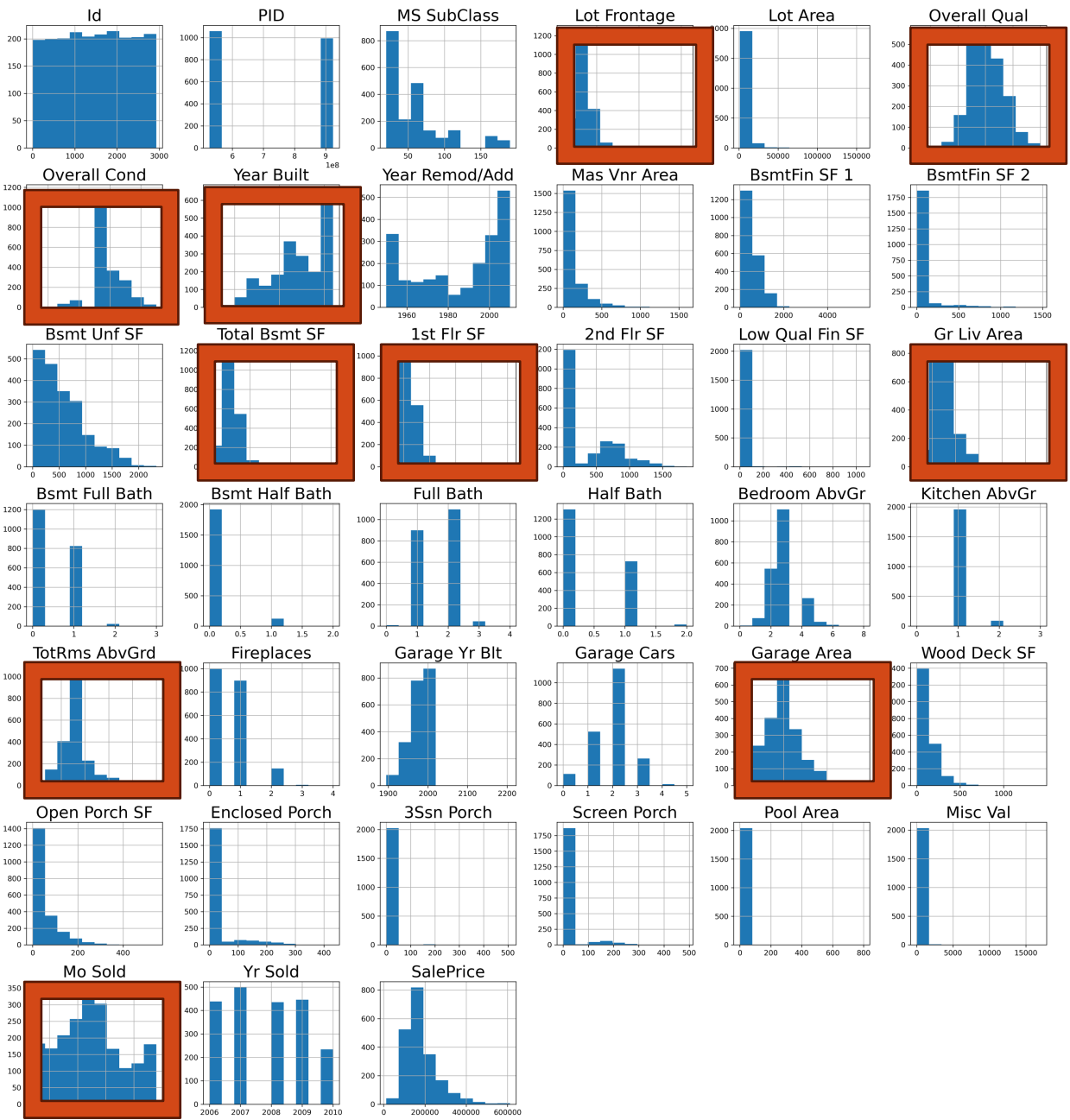


HEATMAP:

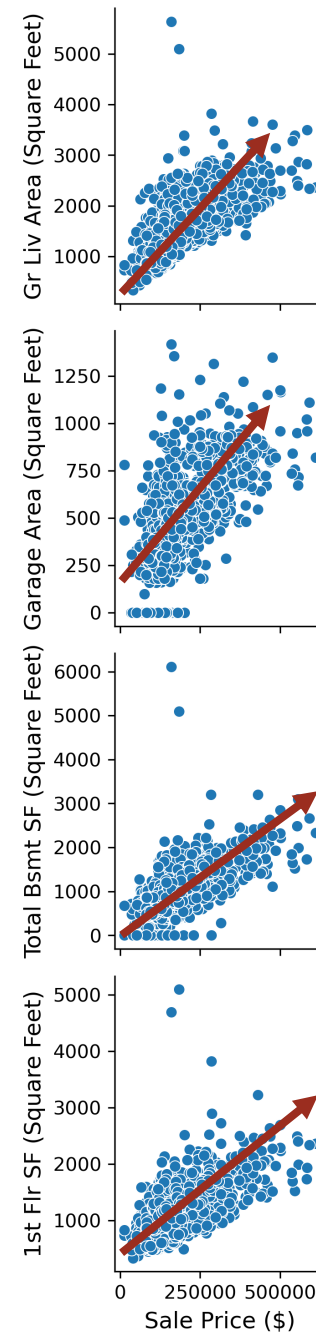


NORMALLY DISTRIBUTED VARIABLES

Training Data Histograms



LINEARITY



TRAIN-TEST-SPLIT

- Small trial and error with training size.
- Feature variables were the best fit correlations, the dummy columns, and feature columns.
- Rescale for uniformity.
- Try Linear, Lasso, Ridge.
- Utilize Polynomial Features for better R^2 and less error.
- Achieved R^2 less than ~2% difference between test and train.
- **RMSE** a little over a tenth of the expected mean (181k)



CHARACTERISTIC (POLYNOMIAL) EQUATION

- Second Order Polynomial with units m^3 as leading coefficient of order-2 term.
- Form : $y = ax^2 + bx + c$
- Why did we choose Ridge?
 - Great with multicollinearity
 - Reducing Complexity
- What if hundreds or thousands of variables?
- Big risk of overfitting
- Force simpler model, often defined as smaller and “more regular” (less varying) coefficients... small Euclidean norm (2D distance formula)



ERROR DISCUSSION

- Extremely high $R^2 \approx 1$, very close between train and test, (less than 0.1% difference) and low RMSE error achieved with suitable choice of parameters.
- For this we used Lasso to achieve a low RMSE and high Kaggle score: 23322.66654



RECOMMENDATIONS

- Biggest predictors of Sales Price:
 - Area-based variables
 - Overall and other quality factors.
 - Neighborhoods also played a role.
- Year built and renovations (remodeling), namely after the year 2000, were also a significant factor.



CONCLUSION

- Utilized training-test-split techniques
- Handled a complex multilinear regression model with a variety of regressions
- Worked with fundamental fitting methods
- Effectively predicted housing prices from an unknown dataset.



SOURCES

- Chance, Beth L. Rossman, Allan J. *Investigating Statistical Concepts, Applications, and Methods (Third Edition)* August 2015 Beth Chance and Allan Rossman San Luis Obispo, California
- Craven, Mark and Page, David *Regression*, Computer Sciences 760 Spring 2018, www.biostat.wisc.edu/~craven/cs760
- Kuiper, Shonda *Introduction to Multiple Regression: How Much Is Your Car Worth?* (2008), Journal of Statistics Education, 16:3, , DOI: 10.1080/10691898.2008.11889579
- Pardoe, Iain *Modeling Home Prices Using Realtor Data*, Lundquist College of Business, University of Oregon, *Journal of Statistics Education* Volume 16, Number 2 (2008), www.amstat.org/publications/jse/v16n2/pardoe.html

