



Predicting House Prices with Regression Techniques

A Kaggle Project

Tyler Wilbers
NYC Data Science Academy

Outline

- EDA
 - Missingness
 - Imputations
 - Other interesting EDA
 - Neighborhoods vs. Price Exploration
- Preprocessing
 - Skewness and Outliers
 - Box Cox Transformations
 - Standardize Predictors
 - Correlation
 - Multicollinearity
 - Selection
- Modeling
 - Tree methods
 - Linear
 - What we did wrong
 - What we did right
- Future Improvement
 - Kaggle Score
 - Future improvement

Exploratory Data Analysis (EDA)

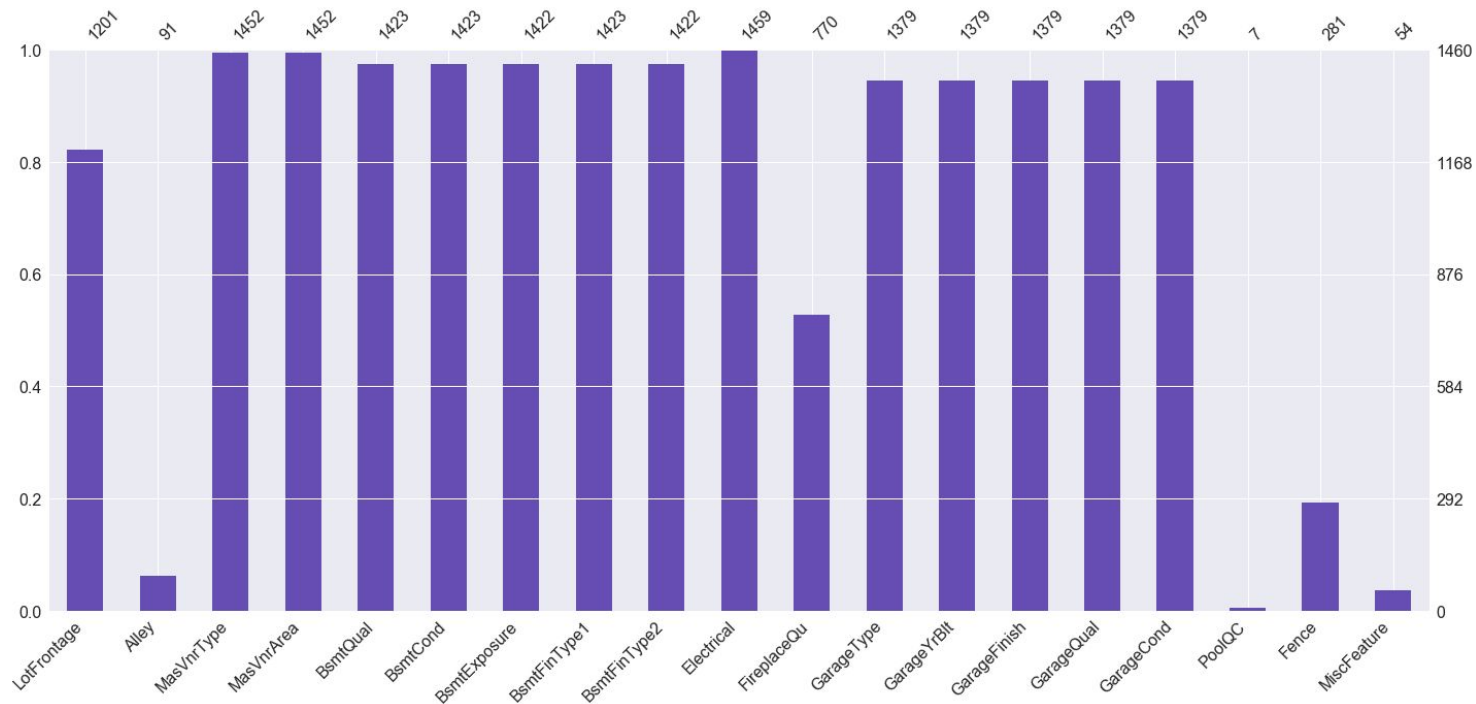


Findings

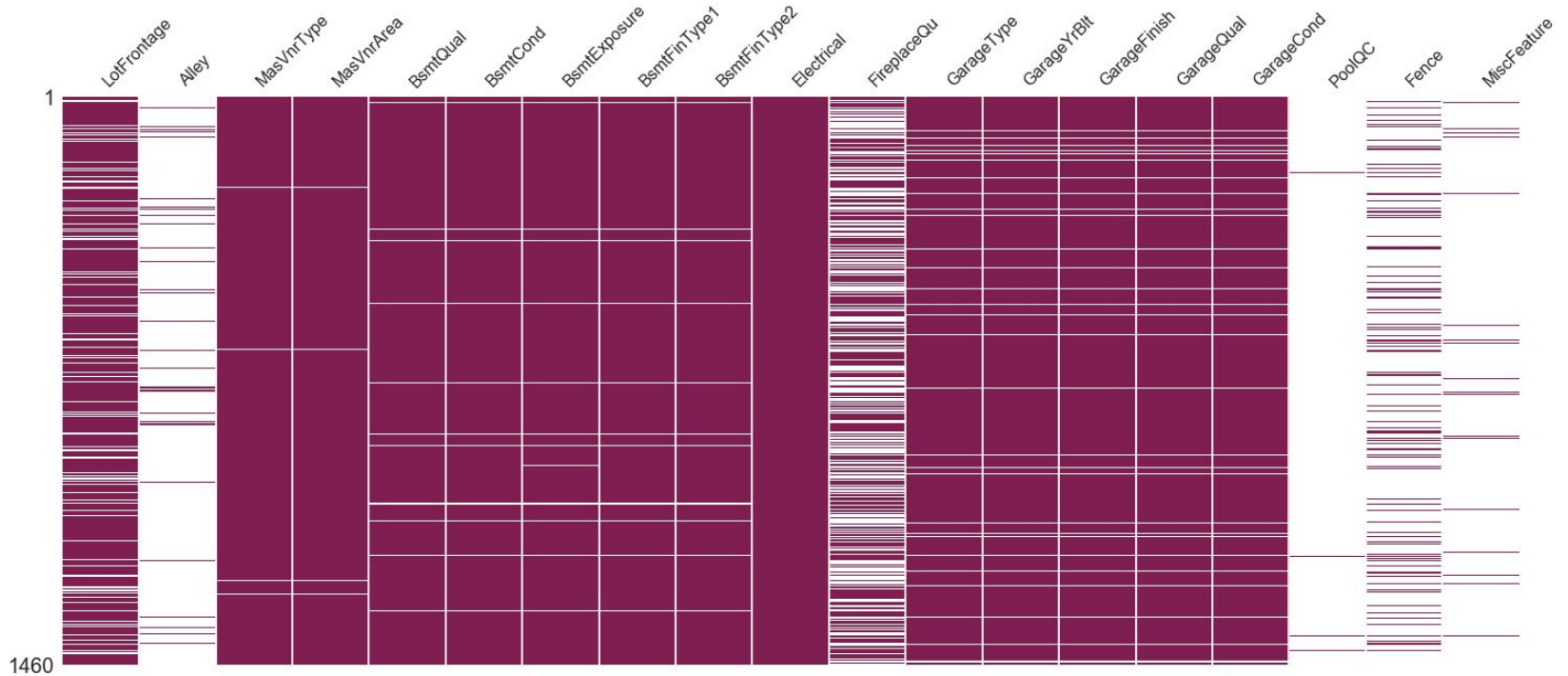


- 80 variables(Exclude: ID)
 - Categorical → 63
 - Continuous → 17
- Counts of all NAs : 6965
- Number of Variables that contains missing values : 19

Amount of Missingness of each column



Relationship of missingness b/w each variable



Imputations

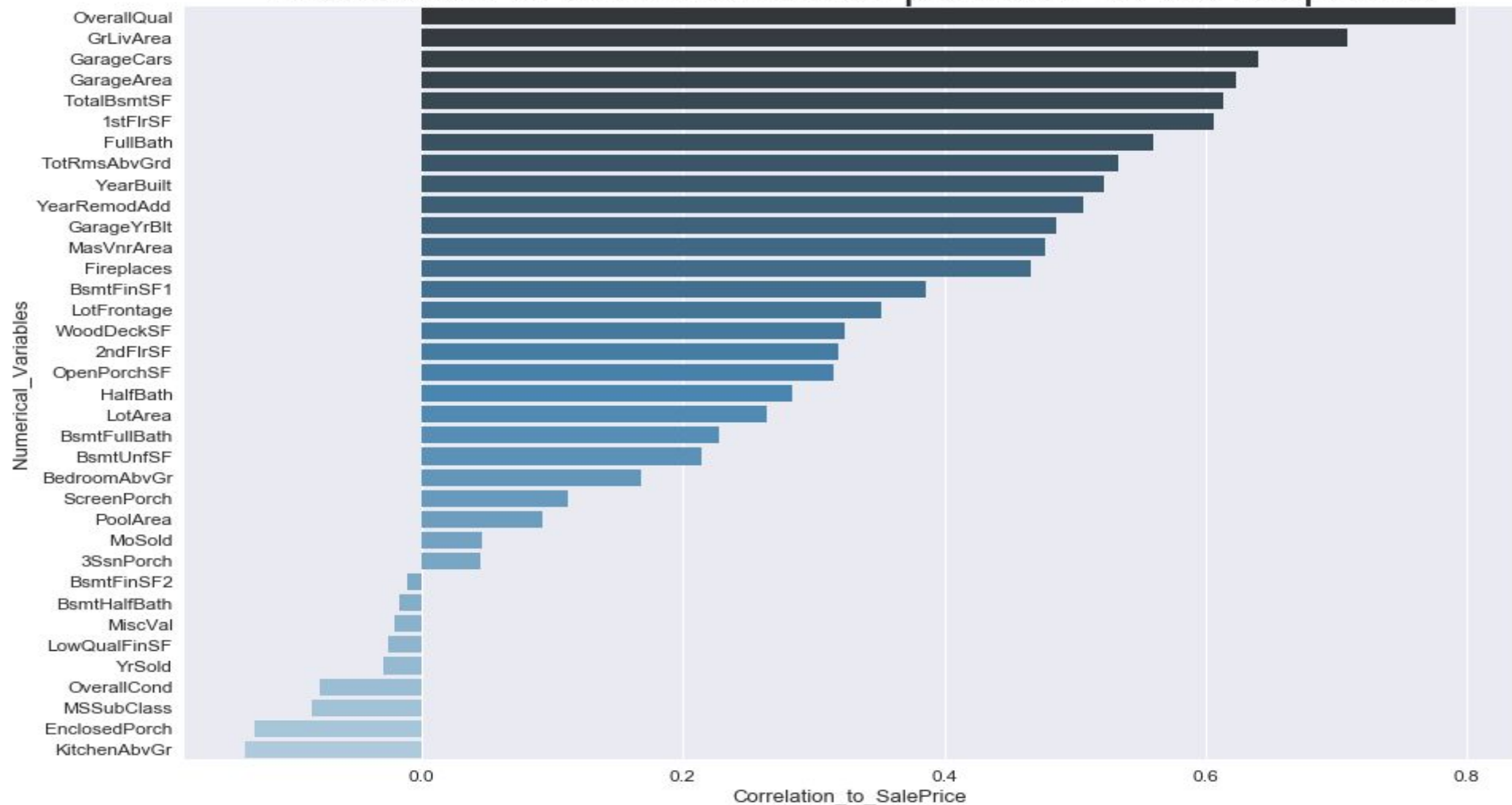
- We took a two part approach:
 - For continuous variable, Group by = neighborhood and impute mean.
 - For categorical variable, we just **dummified**;

MiscFeature	
0	NaN
1	NaN
2	NaN
3	NaN
4	NaN
5	Shed

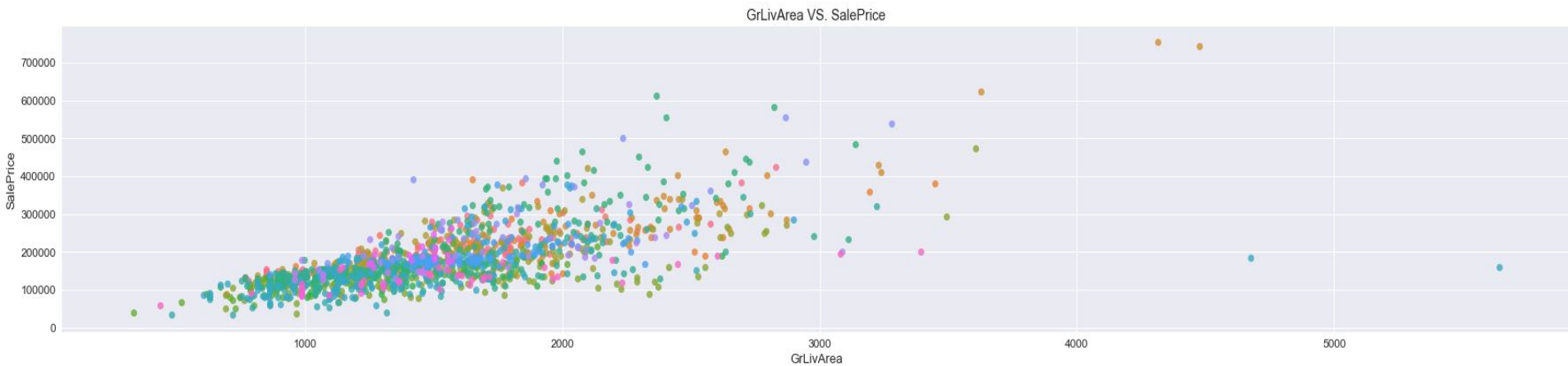


	Gar2	Othr	Shed	TenC
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	1	0

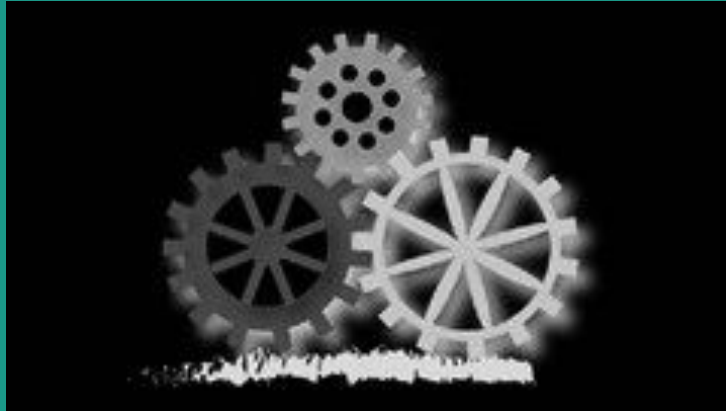
Correlation of each numerical predictor to the response



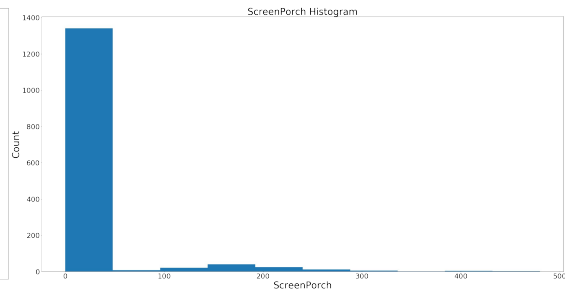
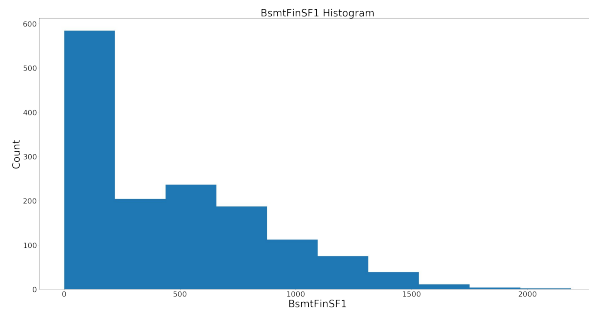
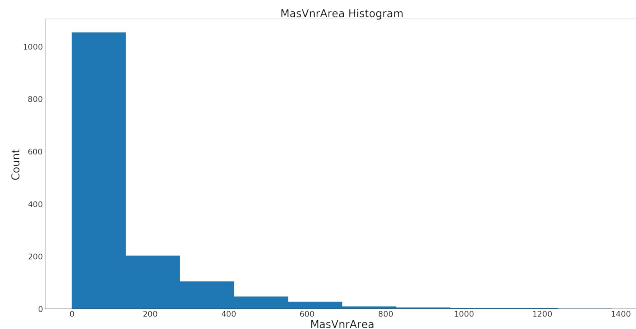
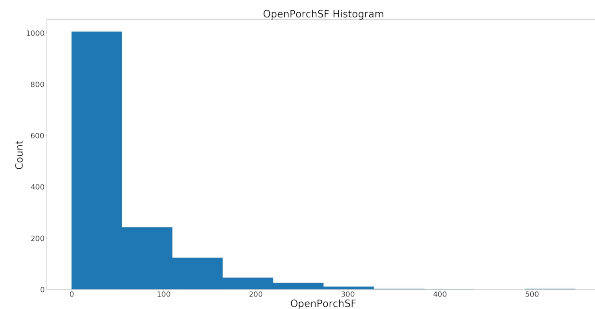
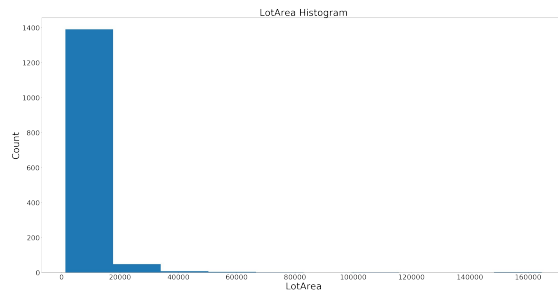
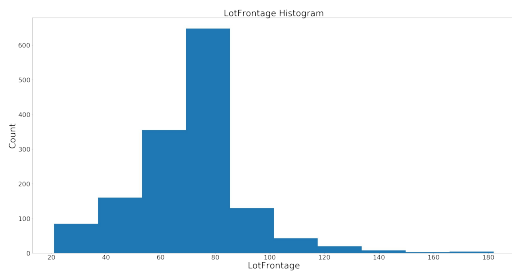
Further Exploration - Group by Neighborhood



Data Preprocessing



Skewness



Skewness (cont.)

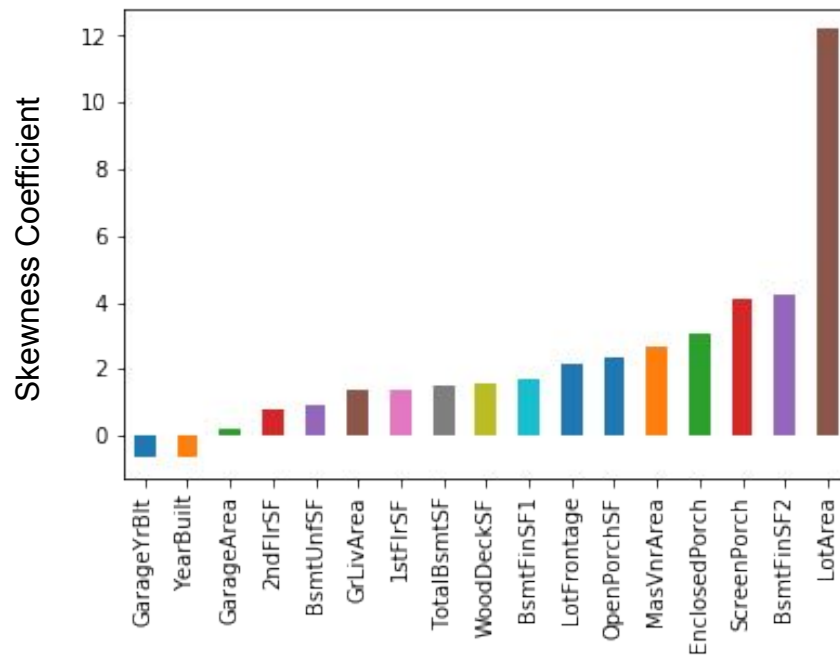
- We used `scipy.stats.skew` which calculates the coefficient of skewness:

$$\frac{\mu_3}{\mu_2^{3/2}}$$

Where μ_i is the central moment.

- Negative skew usually indicates that the tail is on the left side of the distribution, and positive skew indicates that the tail is on the right.
- For normally distributed data, the skewness should be about 0.

Skewness of continuous variables



Box cox transformation



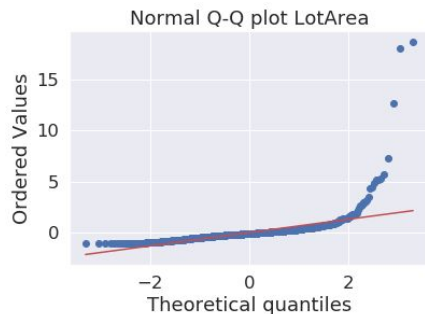
- We decided to perform a one parameter box cox transformation to the skewed variable that have a skewness s such that $2 < s < -2$

$$y_i^{(\lambda)} = \begin{cases} \frac{y_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0, \\ \ln y_i & \text{if } \lambda = 0, \end{cases}$$

- Making this transformation with `scipy.stats.boxcox1p` helped us achieve symmetry, normality, or independence of the error terms.
 - It will also help us stabilize the variance of the distributions and improve the validity of association measures (e.g. correlation).

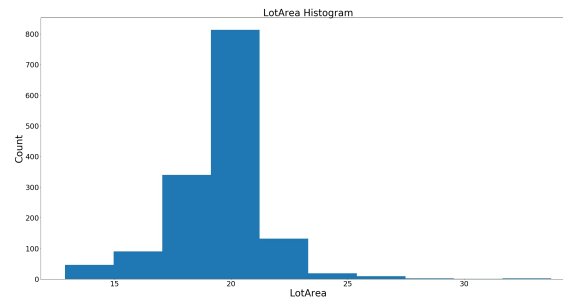
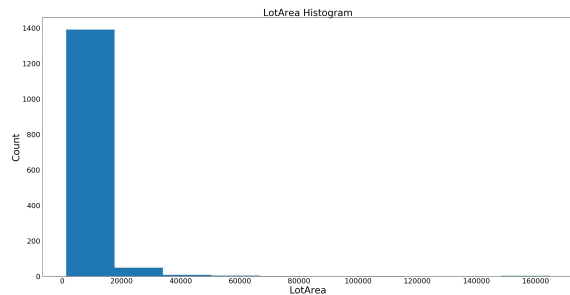
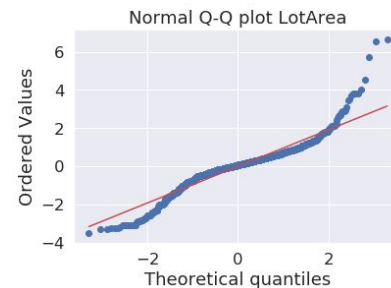
Example - Lot Area

$\text{skew}(\text{LotArea}) = 10.94$



Box-Cox

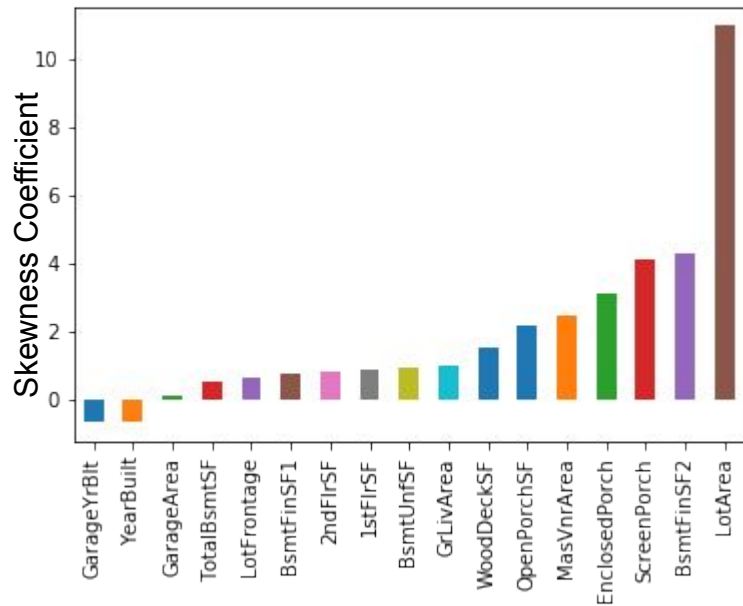
$\text{skew}(\text{LotArea}) = 0.42$



Skewness (After box cox)

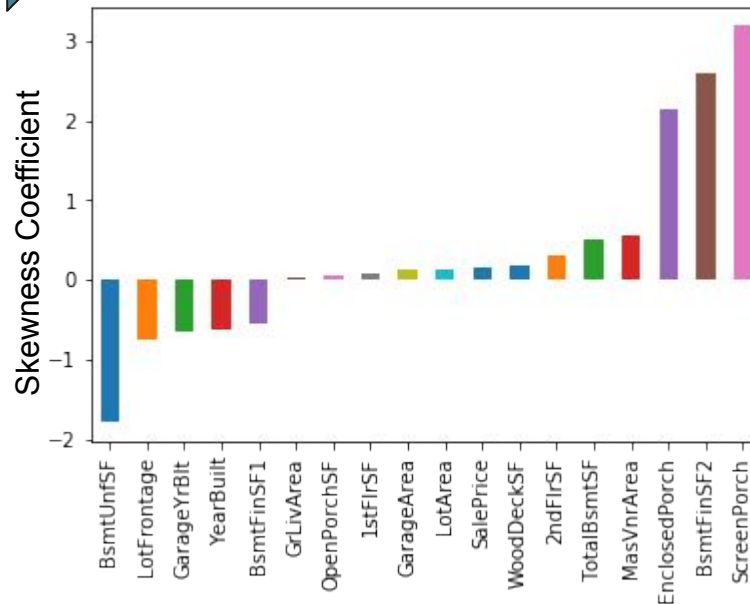


Skewness of continuous variables



Box-Cox

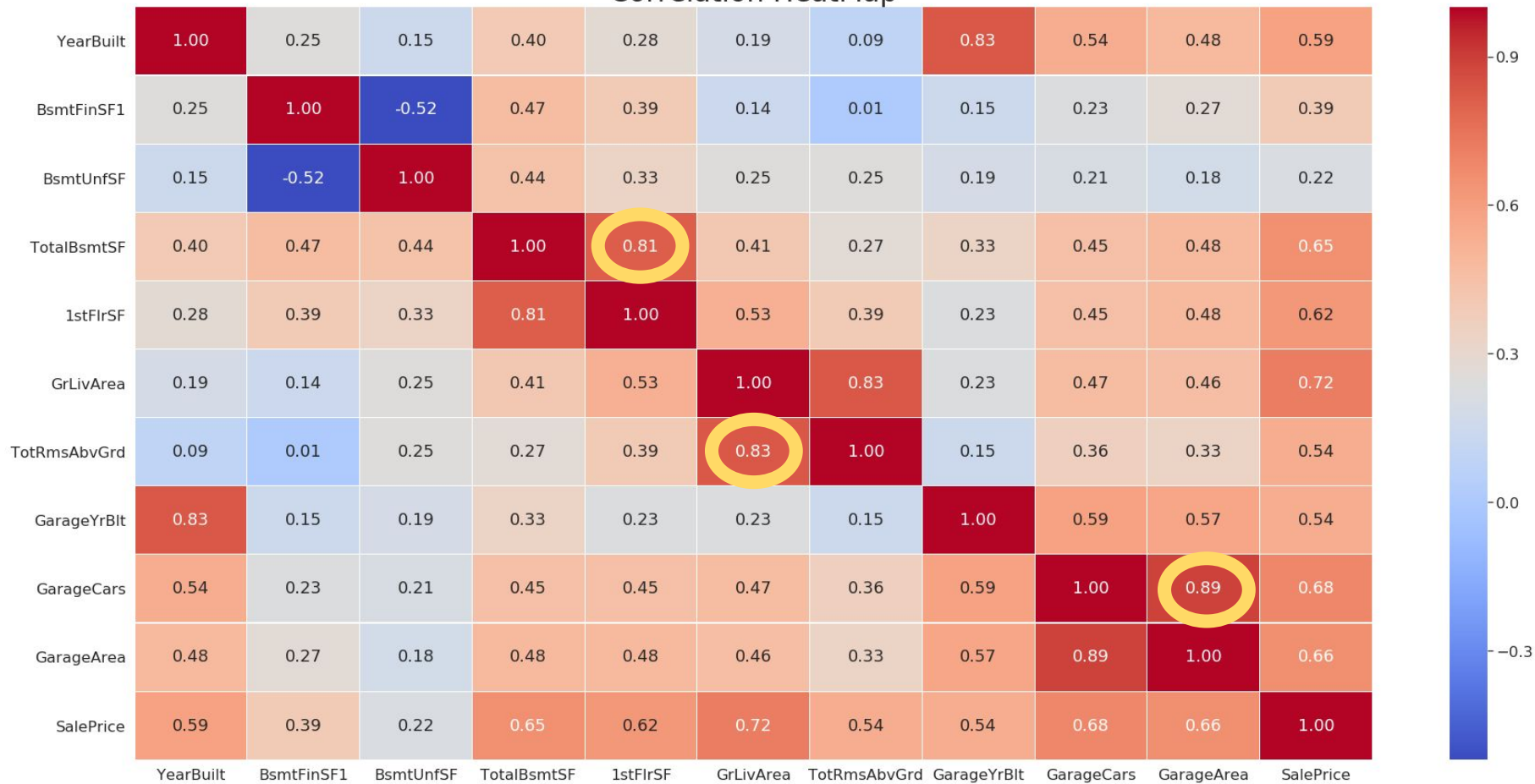
Skewness of continuous variables





Feature Selection - Identifying Multicollinearity

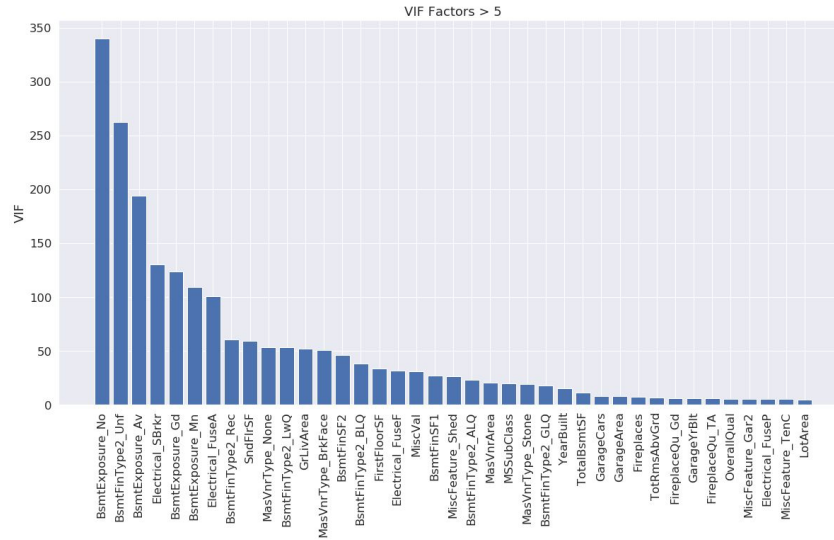
Correlation HeatMap



Variance Inflation Factor

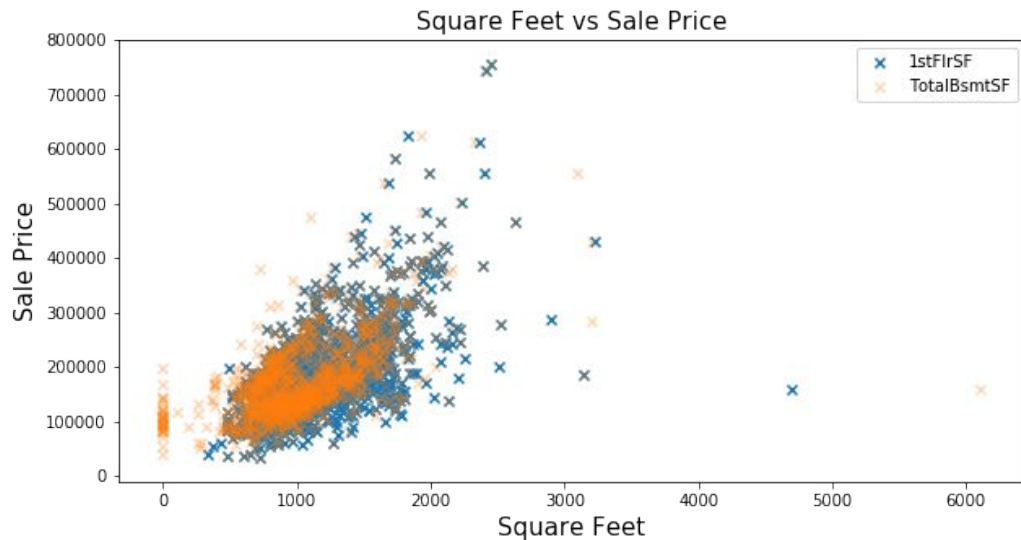
$$VIF_i = \frac{1}{1 - R_i^2}$$

where R_i^2 is the R^2 from a regression of X_i onto all of the other predictors. If R_i^2 is high, then collinearity is present, and so the VIF will be large.



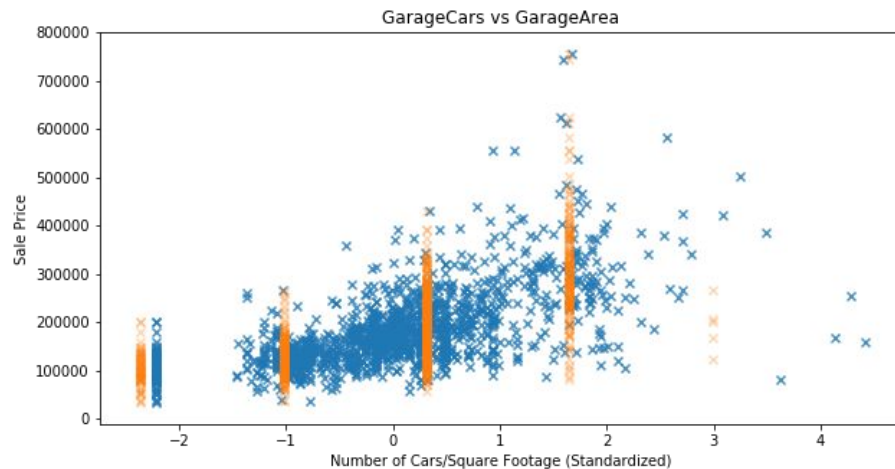
Multicollinearity

- We found that **TotalBsmntSF** is co-linear with **1stFlrSF** ($\rho = 0.82$)
 - $VIF(\text{TotalBsmntSF}) = 11.8922$
 - $VIF(1stFlrSF) = 36.4134$
- Due to the higher VIF we opted to drop **1stFlrSF**.

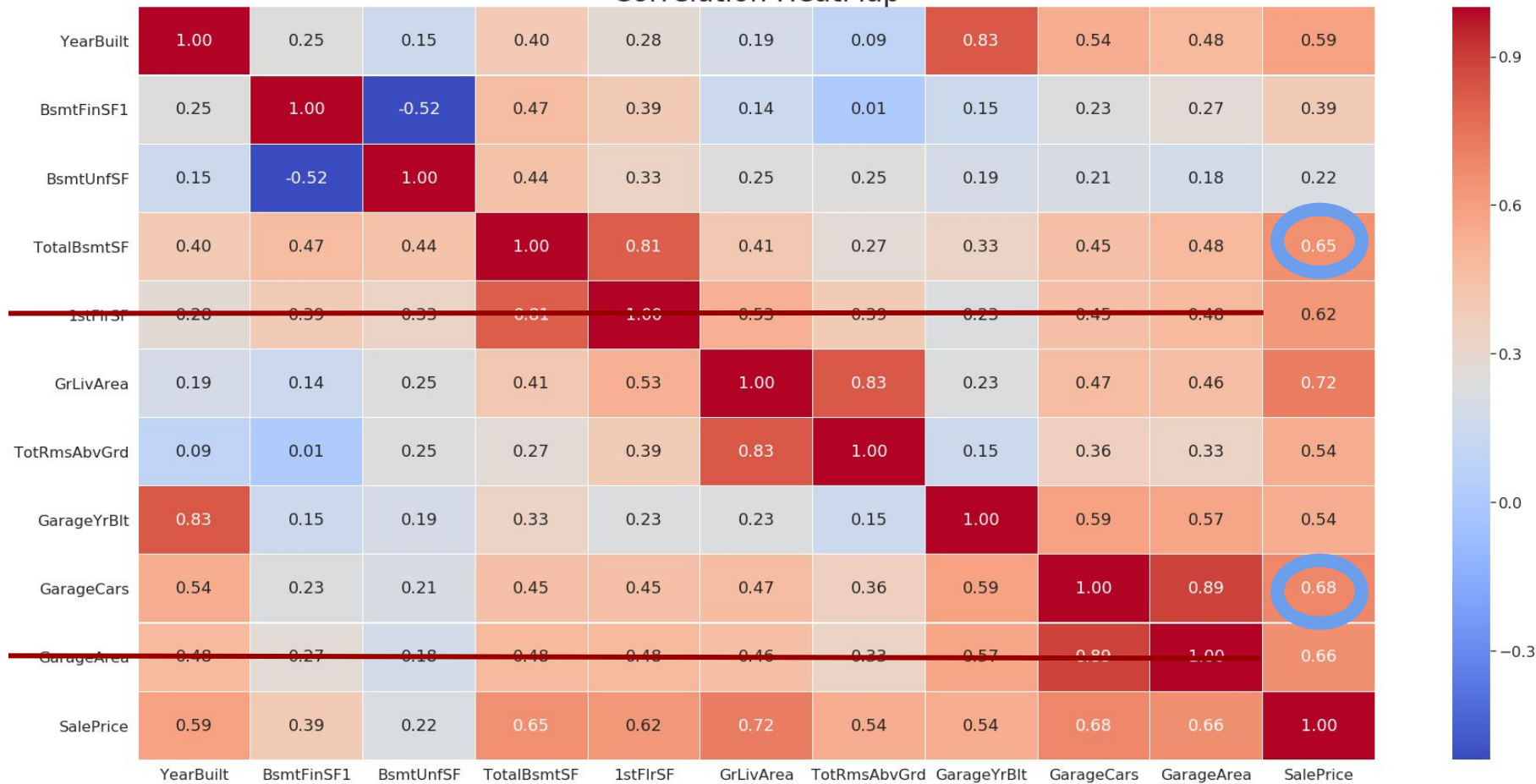


Multicollinearity

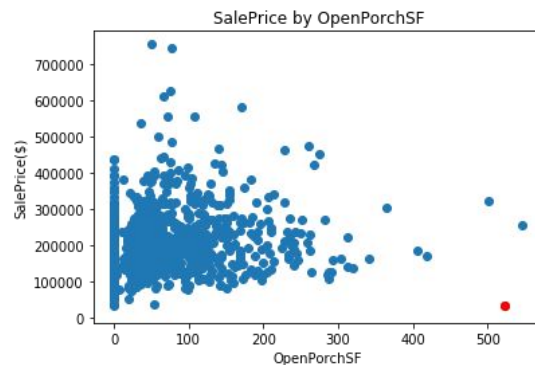
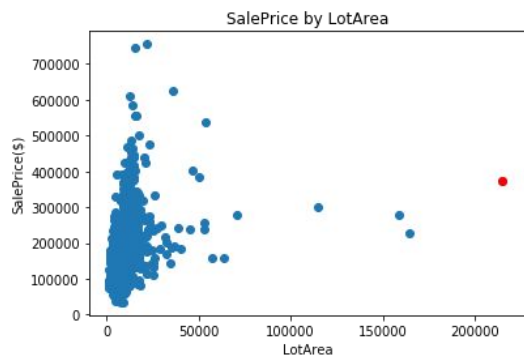
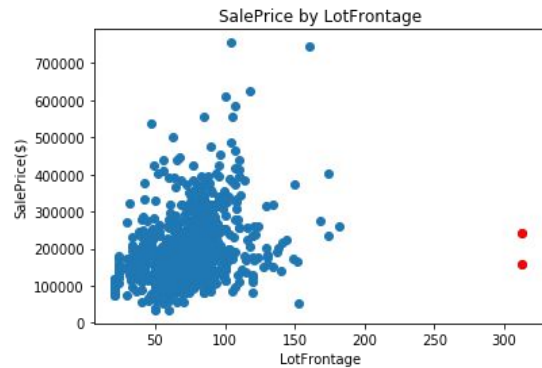
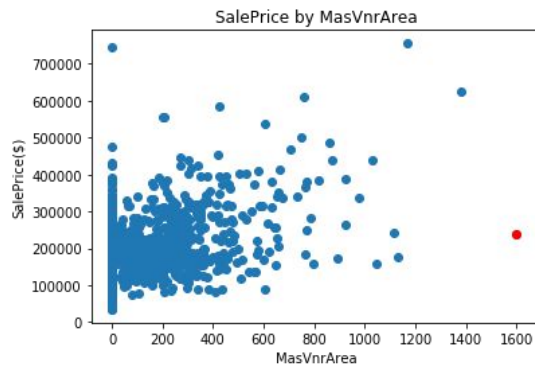
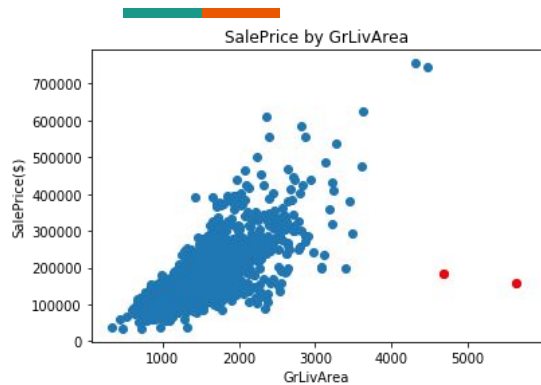
- We found that **GarageArea** is co-linear with **GarageCars** ($\rho = 0.89$).
 - $VIF(\text{GarageCars}) = 8.28583$
 - $VIF(\text{GarageArea}) = 8.18866$
- We opted to drop **GarageArea**.



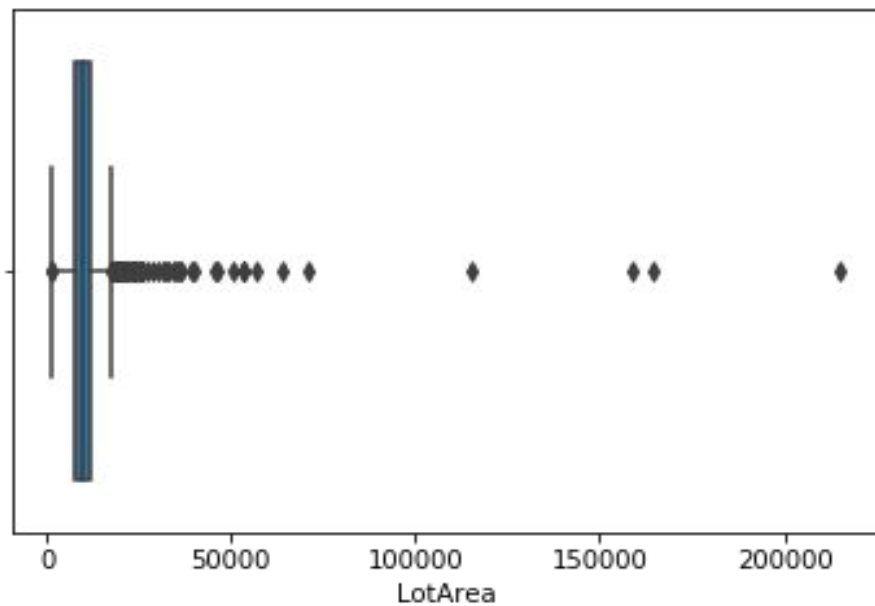
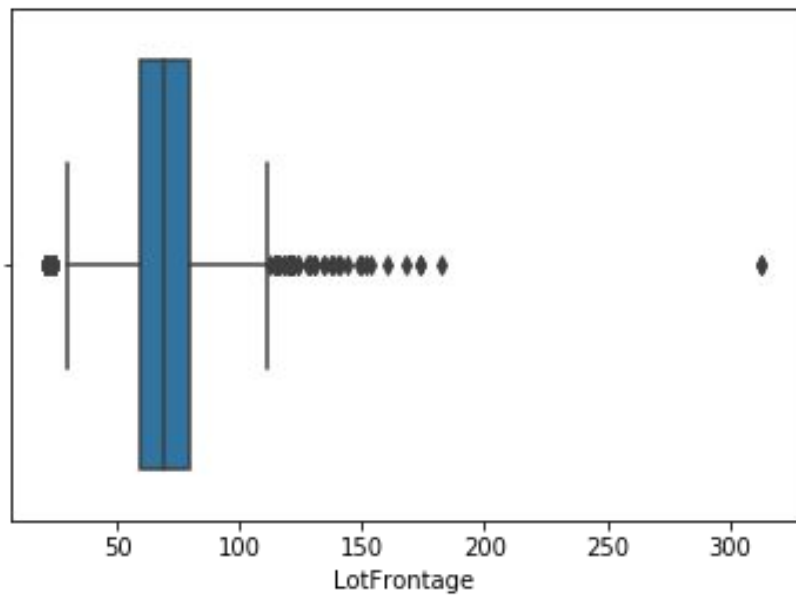
Correlation HeatMap



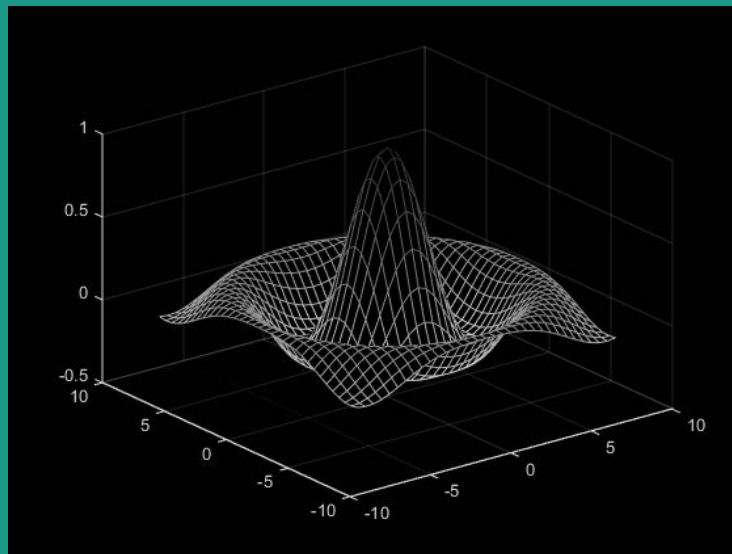
Removing Outliers



Outliers Discovery - Examples



Modeling



Approach

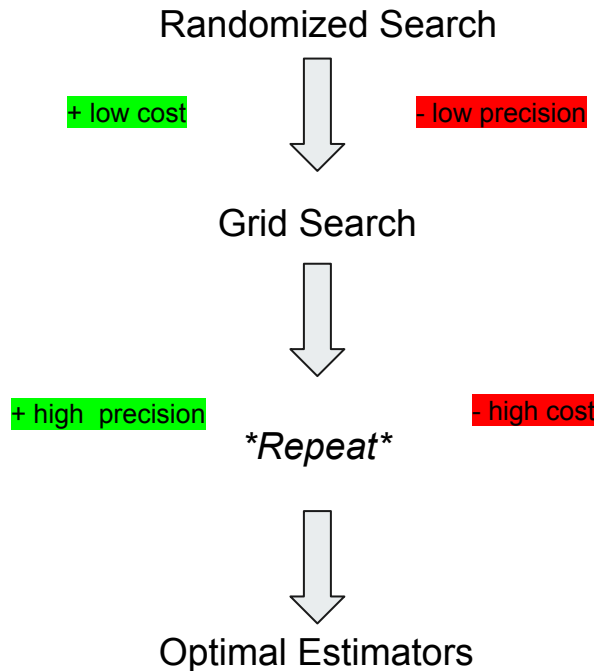


Tree Models:

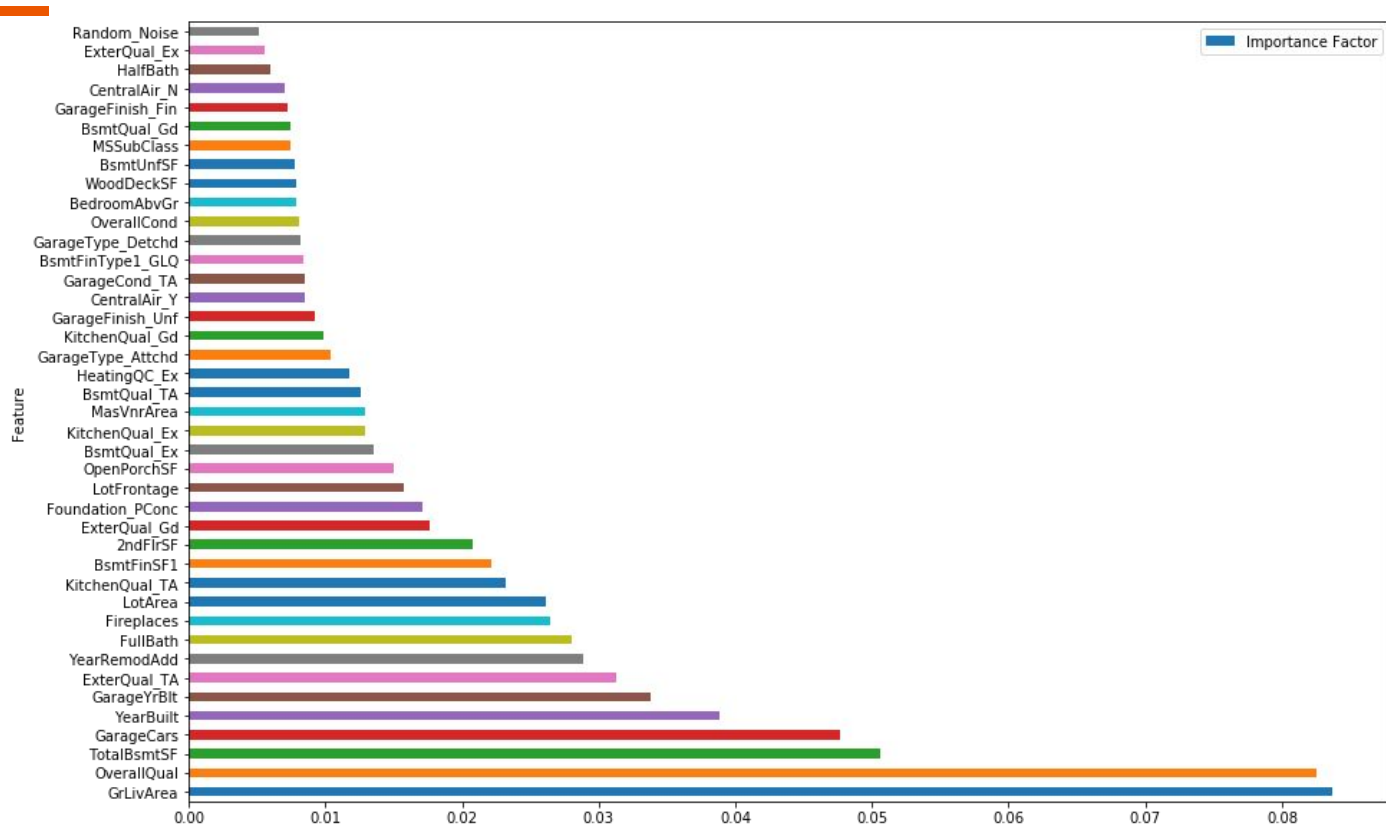
- Random Forest Regressor
- XGBoost

Linear Models:

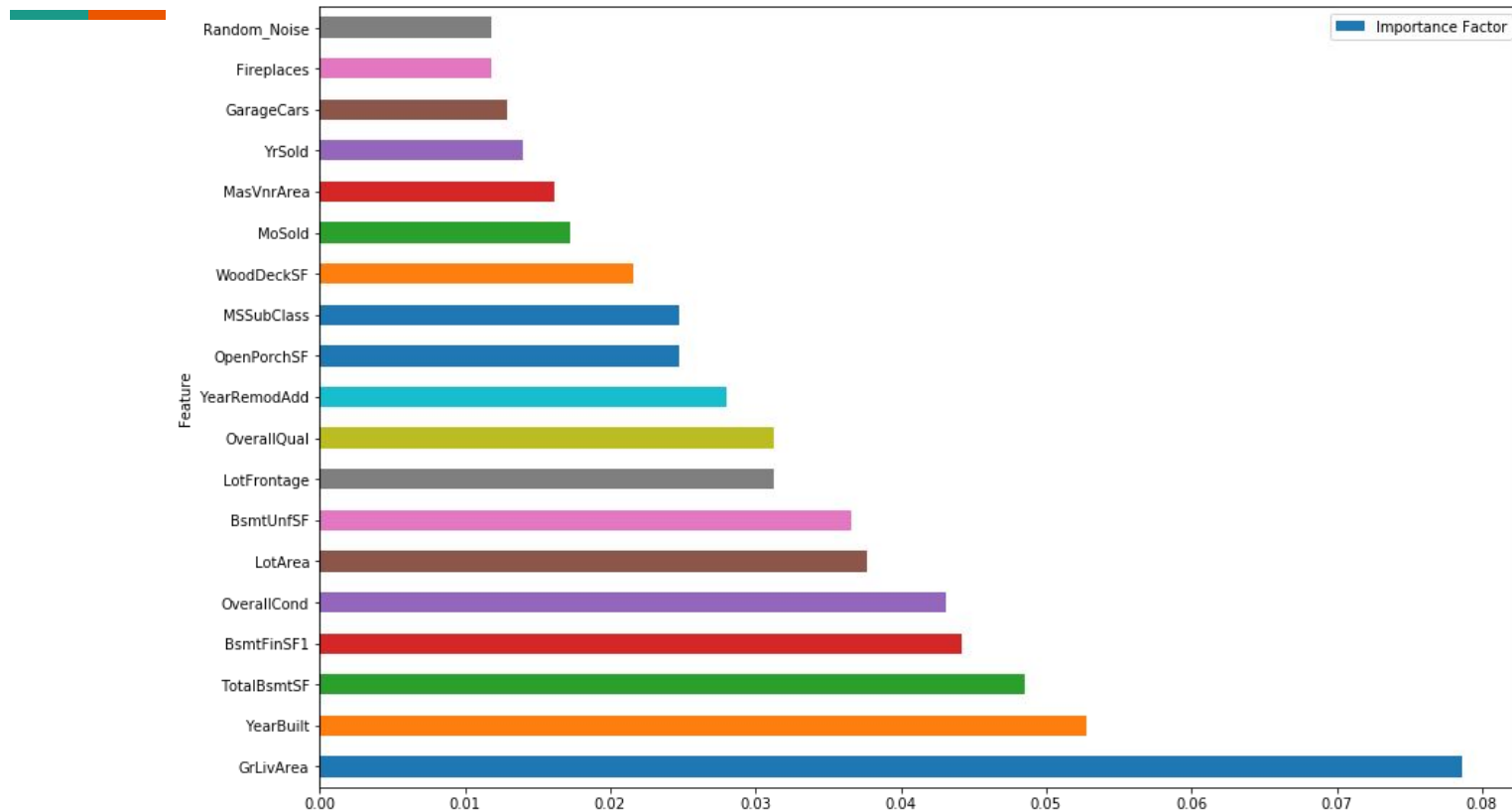
- Lasso
- Ridge
- ElasticNet



Random Forest - Feature Importance



XGBoost - Feature Importance

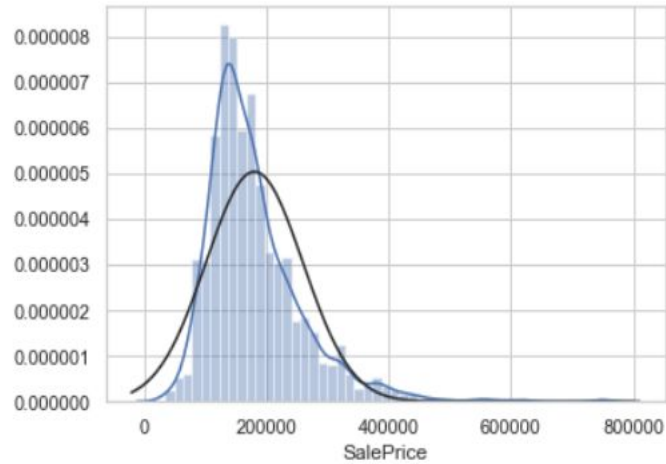


Tree Method Results

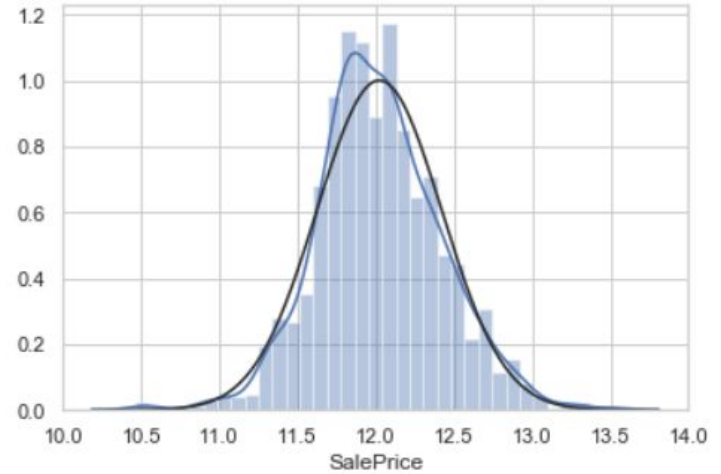
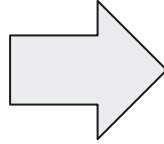


Model	Train CV Scores (RMSE)
Random Forest	0.1361
XGBoost	0.1258

Distribution of Sale Price with Transformation



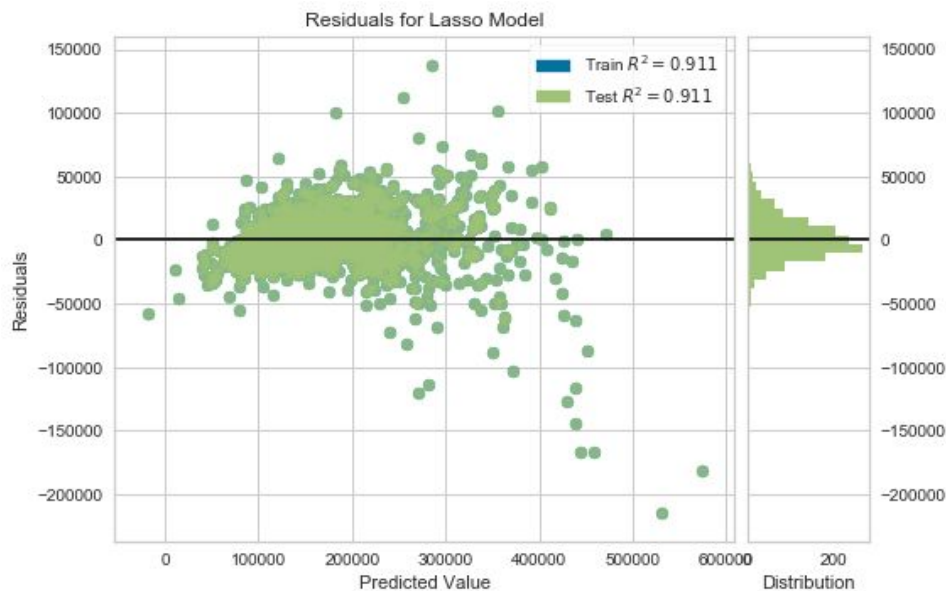
Before transformation



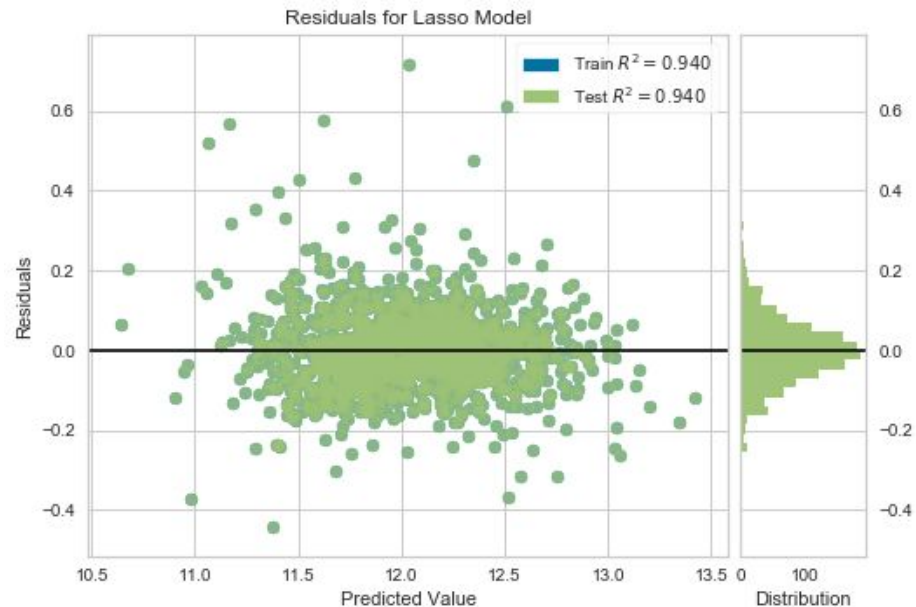
After log transformation

Log Transform for Normally Distributed Residuals

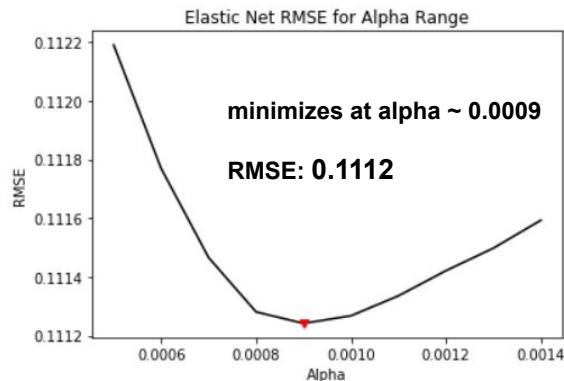
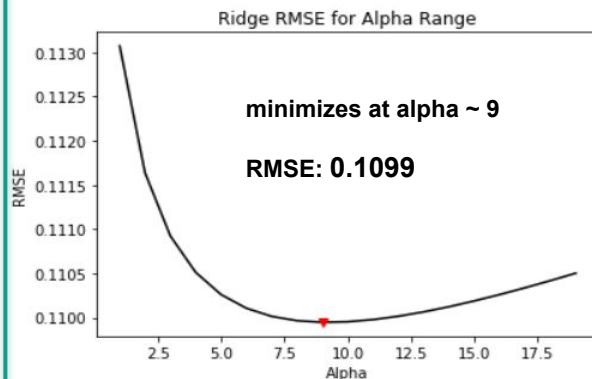
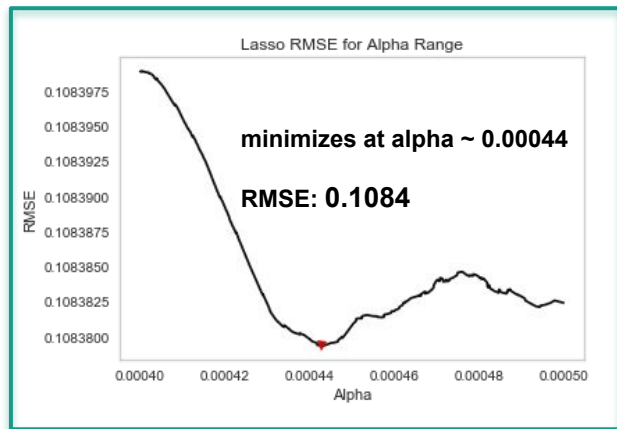
Standard target variable



Log Transform



Regularized Regression: Hyperparameter Testing w/ GridSearch

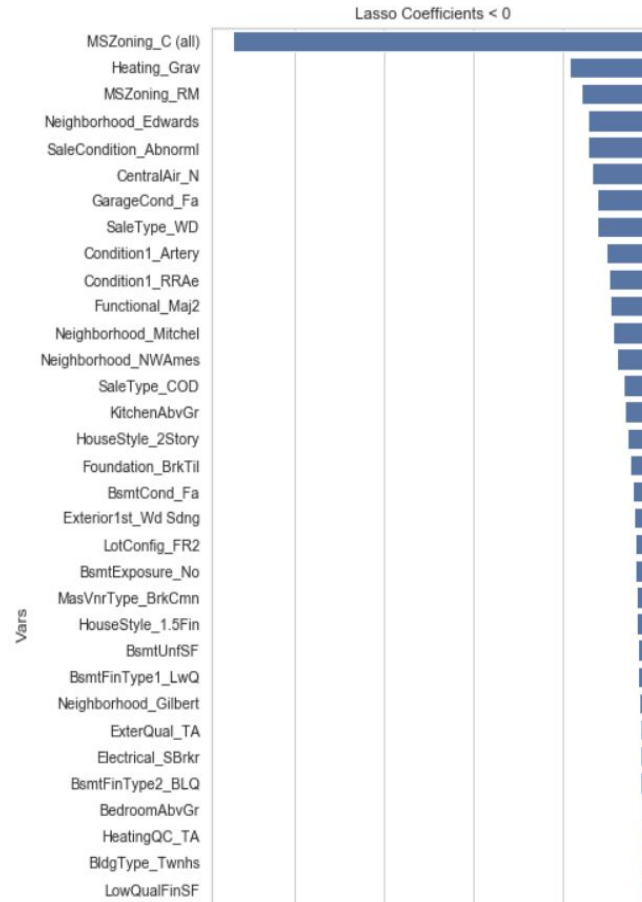
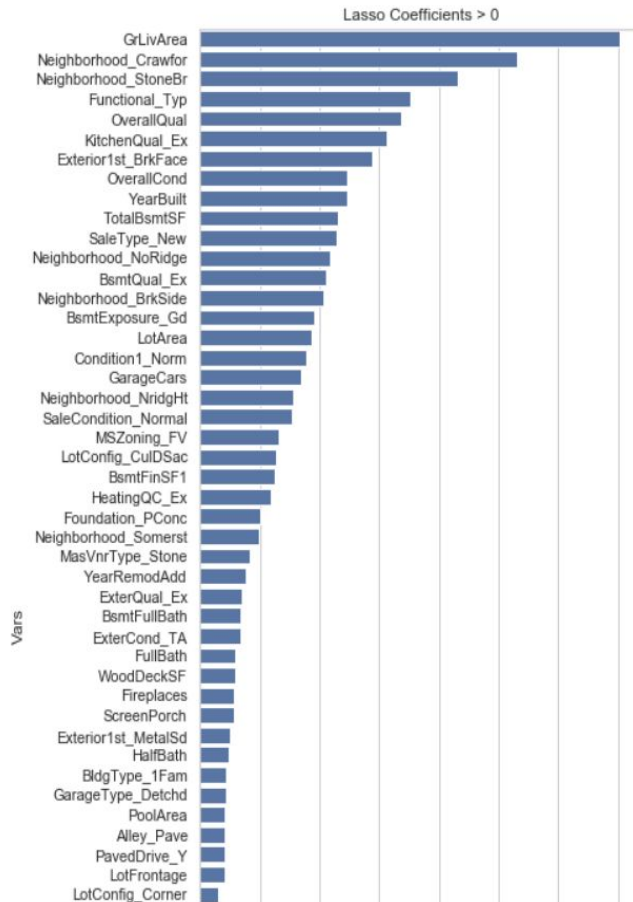


Lasso

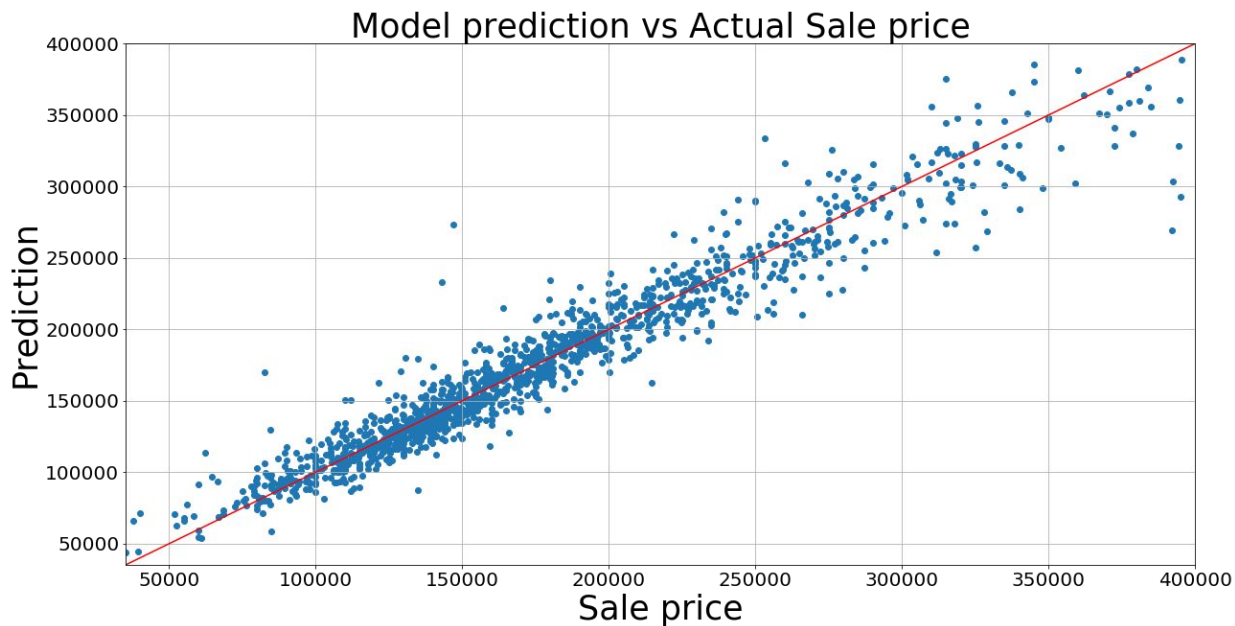
Coefficients

Original # of
Features:
285

of Features
After Lasso:
105



Lasso Predictions on Train



Model Cross Validation Scores



Model	Train CV Scores (RMSE)
Random Forest	0.1361
XGBoost	0.1258
Lasso	0.1084
Ridge	0.1099
Elastic Net	0.1112

Future Improvement



Conclusion

- We scored multiple models, but in the end Lasso scored the highest on Kaggle.
 - $\text{RMSLE} = 0.11544$



House Prices: Advanced Regression Techniques

Predict sales prices and practice feature engineering, RFs, and gradient boosti...

Getting Started · Ongoing · 🗎 tabular data, regression



355/4295

Top 9%



Q & A