

# Housing Prices Analysis

Mandey Brown, Megan Dunnahoo, Emma Hamilton

## Motivation

Housing prices play a central role in the U.S. economy. According to a *Congressional Research Service* article, *Introduction to U.S. Economy: Housing Market*, “at the individual level, roughly 65% of occupied housing units are owner occupied, homes are a substantial source of household wealth in the United States. . . housing accounts for a significant portion of all economic activity, and changes in the housing market can have broader effects on the economy.” Buying a house is considered the most utilized and profitable investment for most of the population. The housing market is also incorporated into gross domestic product (GDP), which is considered the primary measure of economic activity for a country. Also, according to the article, *Introduction to U.S. Economy: Housing Market*, “as of 2020, spending on housing services was about \$2.8 trillion, accounting for 13.3% of GDP. Taken together, spending within the housing market accounted for 17.5% of GDP in 2020.”

In addition to the majority of the population benefiting from predictions of housing prices, many professions and industries would benefit as well. Home appraisers, mortgage lenders, insurers, and tax assessors would be able to more accurately assess the value of a home. Housing price predictions would also prove invaluable for home builders. For this project, the co-owner of Deluxe Homes LLC was interviewed in order to gain more industry insight. The co-owner, Stu Sprecher emphasized the need for flexible pricing predictions that would enable home builders to maintain a profit margin while ordering materials and hiring contractors for each home built. The ability to customize house price predictions to a specific home could prove invaluable for him as an industry professional.

## Methodology

Built off the Kaggle competition, *House Prices - Advanced Regression Techniques*, this project utilizes housing prices compiled by Dean De Cock in 2011, which describes the sale of individual residential property in Ames, Iowa from 2006 to 2010. The data set includes 79 explanatory variables (23 nominal, 23 ordinal, 14 discrete, and 20 continuous) all involved in evaluating home values. In order to build predictive models for housing price, an exploratory analysis was conducted followed by preprocessing. The project focuses on advanced regression models including a Decision Tree, Random Forest, Bagging, LASSO, and out of curiosity, Boosting.

A Decision Tree model is first explored as the output is easily interpreted and its graphical representation can be straightforwardly related to the predicting housing price. Though the downfalls of Decision Trees are known, such as overfitting. Instead of taking the approach of pre-pruning the Decision Tree using  $\chi^2$  test, which is “an algorithm used to find out the statistical significance between parent and child nodes” (Analytics Vidhya), or post-pruning using error estimation, in order to avoid overfitting and get a better understanding of the significance of the data’s variables, other models were explored, including Random Forest.

A Random Forest model was also fit to the data set as it is able to handle large data sets containing higher dimensionality, which was thought to be an appropriate choice as the Ames housing data set contains 79 variables. A Random Forest model was also chosen as it is able to “reduce correlation between trees by injecting more randomness into the tree-growing process” (Greenwell et al). It was also chosen as a means of identifying which of the 79 variables were significant while predicting house price.

In addition to a Random Forest model, a Bagging model was also fit to the data set as it was thought that a Bagging model’s methods of using collections of training data subsets to train multiple decision trees, of

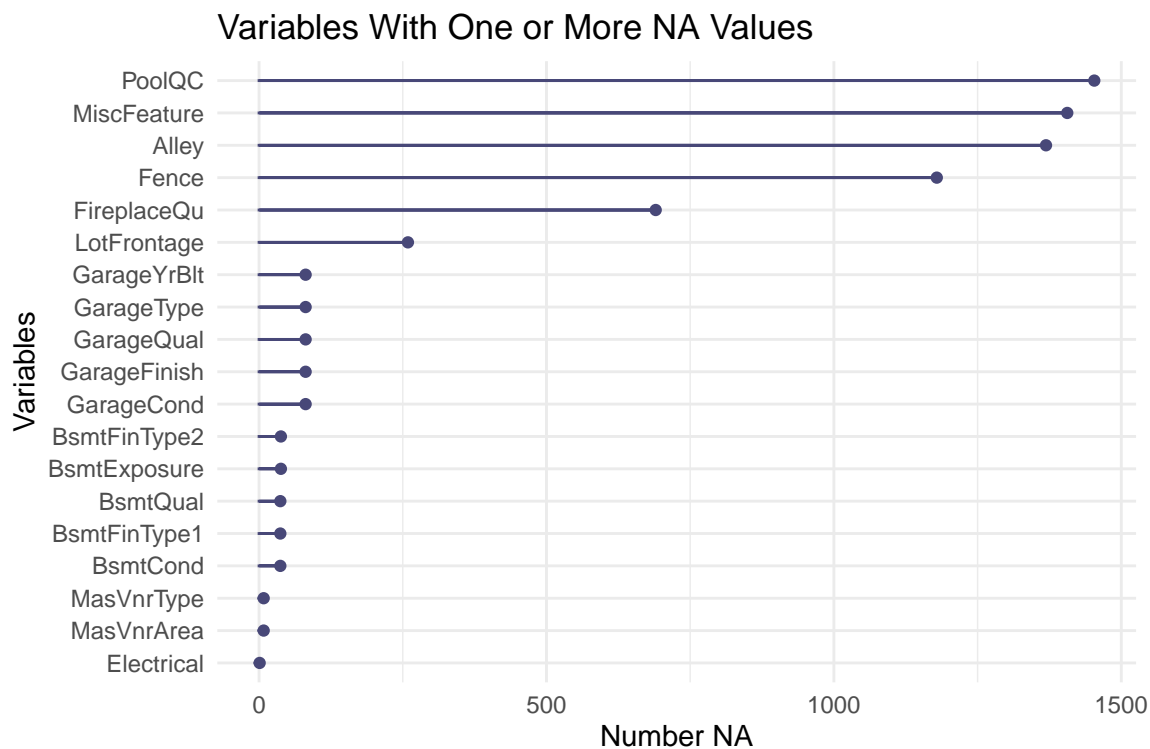
which the average would be used, would not only help avoid overfitting the data, but provide a more robust prediction of housing prices than a single Decision Tree model.

As feature selection is a large component of this project, it was thought that a LASSO model would also prove useful for reducing dimensionality. A LASSO model was thought to offer high prediction accuracy for this data set since the model's method includes shrinking the regression coefficients (some of them to zero), while also reducing variance and minimizing bias.

An AdaBoost Boosting model was also fit as it was thought that a Boosting model may increase predictive accuracy of housing prices via its ability add strength or weight to specific classifiers after taking into account the previous classifier's success. It was thought that a Boosting model would help reduce dimensionality by resulting in significant classifiers being assigned higher weights than less significant classifiers. However, it is believed that a Bagging model will perform better with this data set than the chosen AdaBoost Boosting model, as Boosting does not help avoid over-fitting as a Bagging model does.

## Missing Values

After importing the testing and training data from the Kaggle repository, an initial analysis of missing values was conducted:



It was found that the data for the following variables contained the most Missing Values:

PoolQC > 1250 (>85%) Missing Values  
MiscFeature > 1250 (>85%) Missing Values  
Alley > 1250 (>85%) Missing Values  
Fence > 1000 (>68%) Missing Values  
FireplaceQu > 500 (>34%) Missing Values  
LotFrontage > 250 (>17%) Missing Values

In order to handle the NA/Missing Values, the NA level in the categorical variables were changed to 'None' as these NA values could not be imputed by using the mean, mode, or interpolation of the feature. This was because it was impossible and impractical to impute the value for the quality of a home's pool (PoolQC),

when the home did not come with a pool.

In order to handle the NA/Missing Values of continuous variables, all NA values were converted to zero. This made intuitive sense as it did not make sense to impute any values other than zero for continuous variables such as MasVnrArea or masonry veneer area in square feet if a home did not contain any masonry veneer area.

## Exploratory Data Analysis

In order to better understand the data and gain insight into the relationship between variables, an exploratory data analysis was performed.

### Time variables

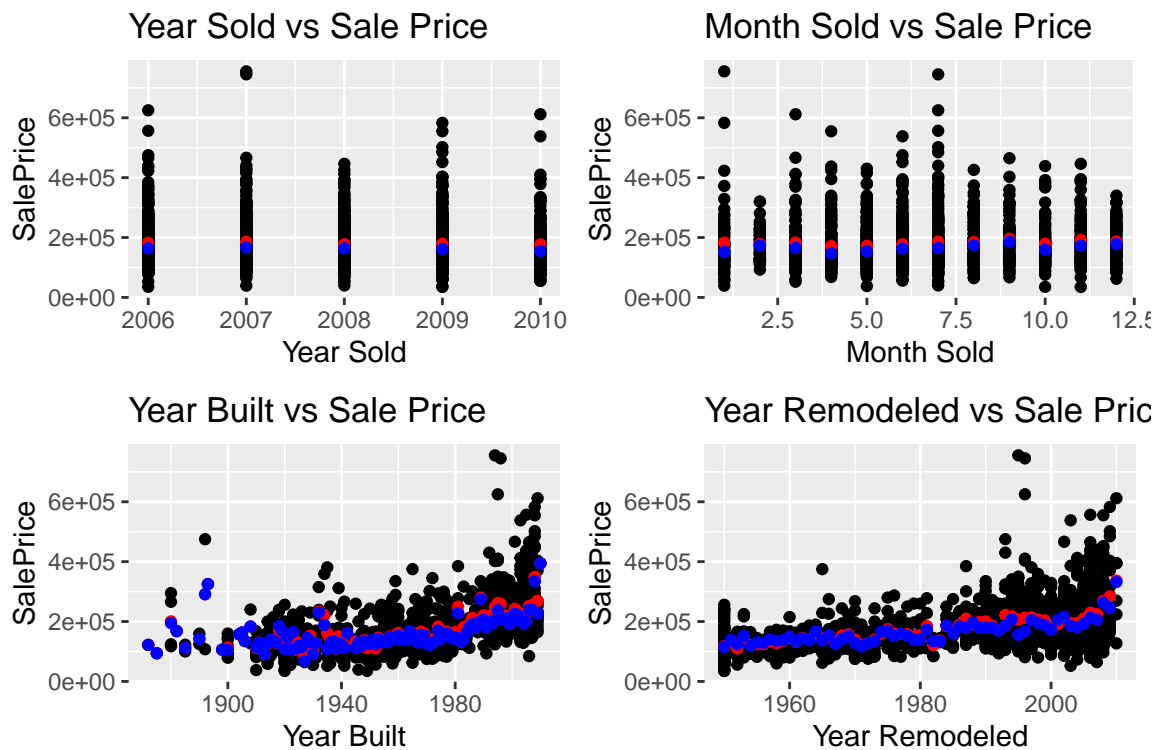
The first series of relationships explored in the Exploratory Data Analysis was the relationships between SalePrice and variables which involve time. These relationships were plotted as:

Sale Price vs. Year Sold

Sale Price vs. Month Sold

Sale Price vs. Year Built

Sale Price vs. Year Remodeled



The rather flat relationship between Sale Price and Year Sold proved interesting as it was thought that this relationship would show to be predominantly positive. It's been hypothesized that the housing market crash of 2007 might have affected this relationship.

The relationship between Sale Price and Month Sold was expected, with the exception of sale prices in January. It was expected that homes would sell for more money in the summer months, but it was not expected that January would prove to have the highest sale price of all months.

The relationship between Sale Price and Year Built was expected with a slightly positive relationship shown.

The relationship between Sale Price and Year Remodeled also was expected with a slightly positive relationship shown.

## Other Variable Subsets

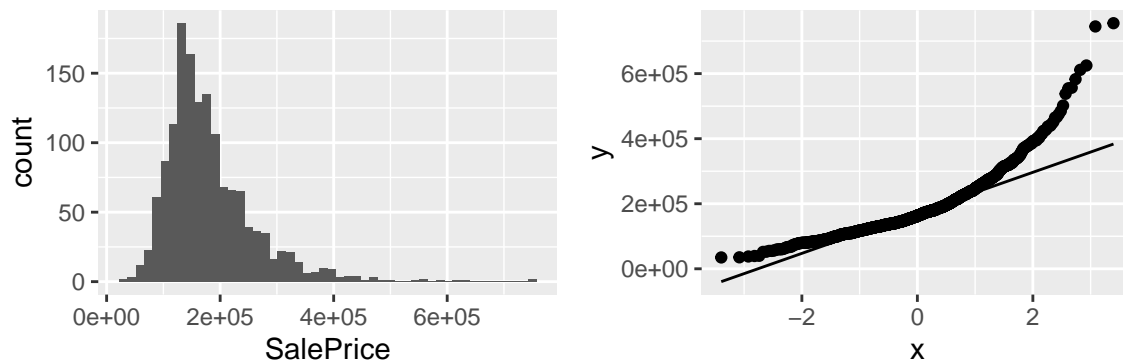
The second part of the Exploratory Data Analysis involved creating ggpairs plots, which shows an overview of relationships between variables. We did this on different subsets of variables indicating size, quality, and condition. We also created a correlation plot on the subset of size indicator variables.



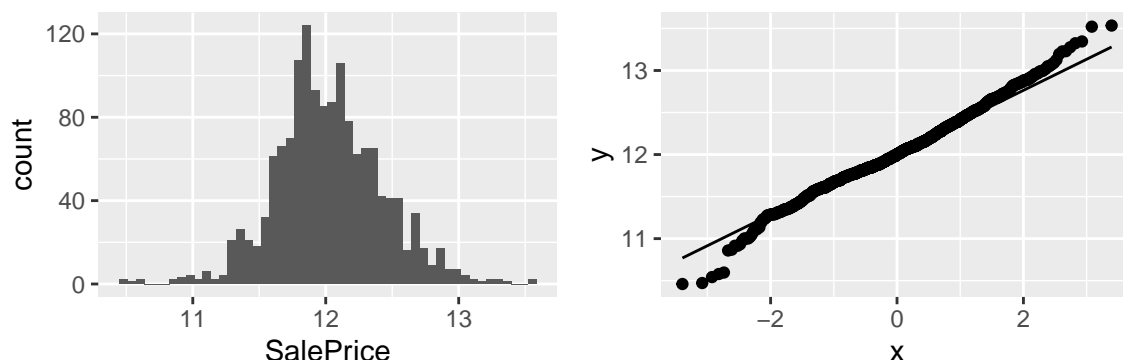
Here, a predominantly positive relationship between Sale Price and the indicator variables was found. However, what stood out most was the skewed distribution of the Sale Price variable across all three scatterplot matrices. Upon further analysis, it was decided that the SalePrice variable should be log transformed.

## Log Transform the Data

The SalePrice variable was log transformed due to having a skewed distribution first discovered in the Exploratory Data Analysis.



Further analysis of the SalePrice variable showed a non-normal right-skewed distribution in the above histogram. This non-normal distribution is also evident in the above (Q-Q) plot where the observations curve off of the line indicating the distribution is skewed.



After log-transforming the SalePrice variable, we now see a normal bell-curve shape in the distribution in the above histogram. The above (Q-Q) plot now shows the observations sticking closely to the line without any curvature away from the line.

## Provide a test set that contains Sale Price

Upon further analysis, it was discovered that an aspect of the Kaggle Competition was that the test data set did not contain a SalePrice column as this is the condition for ranking the effectiveness of the predictive models submitted. Therefore, the training data set was split into new test and training sets in order to analyze the data and fit the models. This allowed us to produce test MSE values, which showed us the accuracy of the fitted models and gave us values to compare the models against each other. We decided on a 70/30 training/test split, which was selected randomly on the training data.

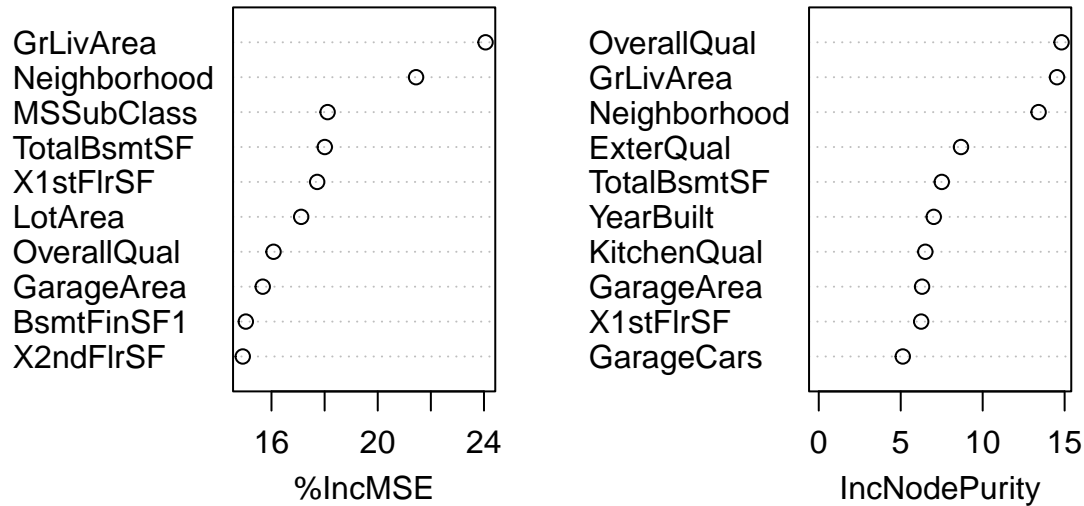
## Decision Tree

A Decision Tree model was first explored as it was thought that the output would be easily interpreted and its graphical representations would be straightforwardly related to predicting Sale Price. In order to fit the Decision Tree, a SalePrice tree was first created. From that tree, the variables used were listed along with omitted variables. The tree was subsequently plotted and the test MSE of the tree was calculated. The test MSE of the Decision Tree was later compared with the test MSE of the other chosen models in order to choose the most accurate model for predicting Sale Price.

```
## [1] "**** Below are variables used in the tree ****"
## [1] OverallQual Neighborhood GrLivArea CentralAir GarageCars
## [1] 81 Levels: <leaf> Id MSSubClass MSZoning LotFrontage LotArea Street ... SaleCondition
## [1] "**** Below are omitted variables ****"
```



## Variables with most Predictive Power

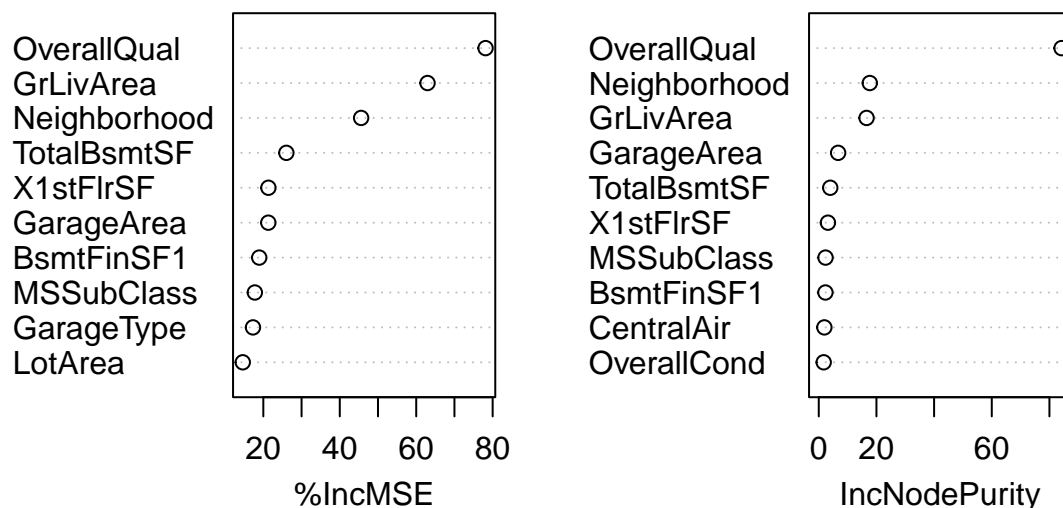


```
## [1] "Random Forest Test MSE = 0.00970699191926917"
```

## Bagging

Next, a Bagging model was fit to the data set as it was thought that a Bagging model's methods of using collections of training data subsets to train multiple decision trees, of which the average would be used, would not only help avoid overfitting the data, but provide a more robust prediction of housing sale prices than a single Decision Tree model. The test MSE of the Bagging model was later compared with the test MSE of the other chosen models in order to choose the most accurate model for predicting Sale Price.

## Variables with most Predictive Power

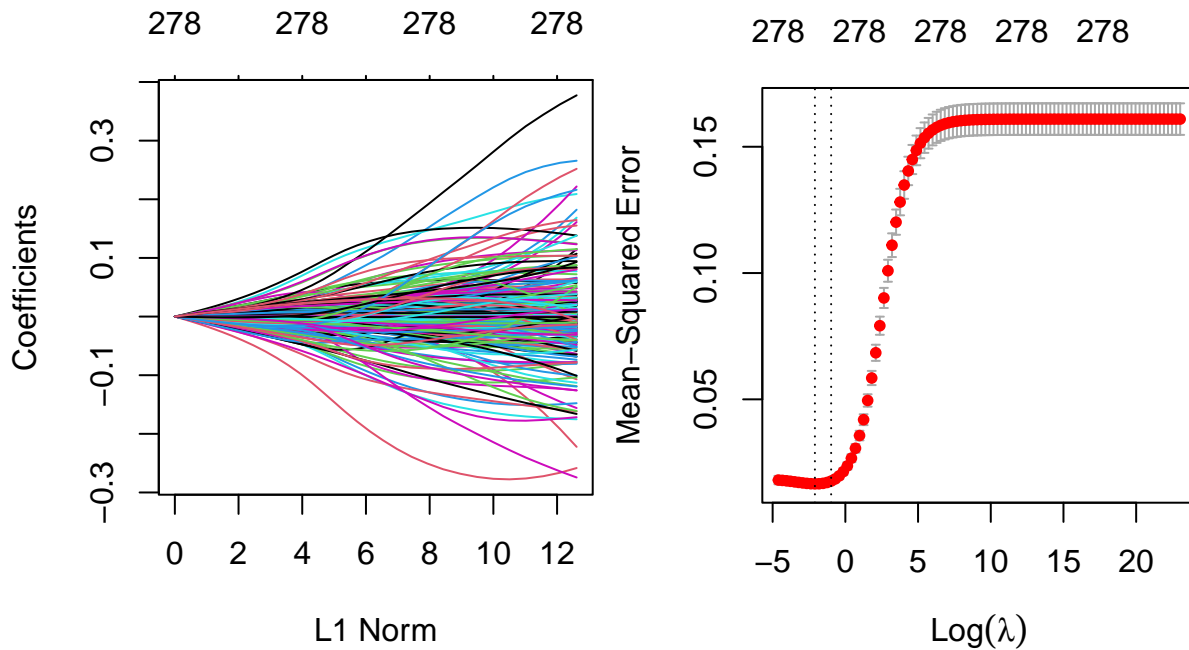


```
## [1] 0.009636063
```

## LASSO

As feature selection was thought to be a large component of the project, a LASSO model was chosen as it offers high prediction accuracy and aids with high dimensionality by shrinking the regression coefficients (some of them to zero).

In order to fit a LASSO model, first, the training and testing data were transformed into matrices and lambda values were added. The coefficients were plotted along with test MSE at different lambda values. The best test MSE of the LASSO model was then calculated. The test MSE of the LASSO model was later compared with the test MSE of the other chosen models in order to choose the most accurate model for predicting Sale Price.

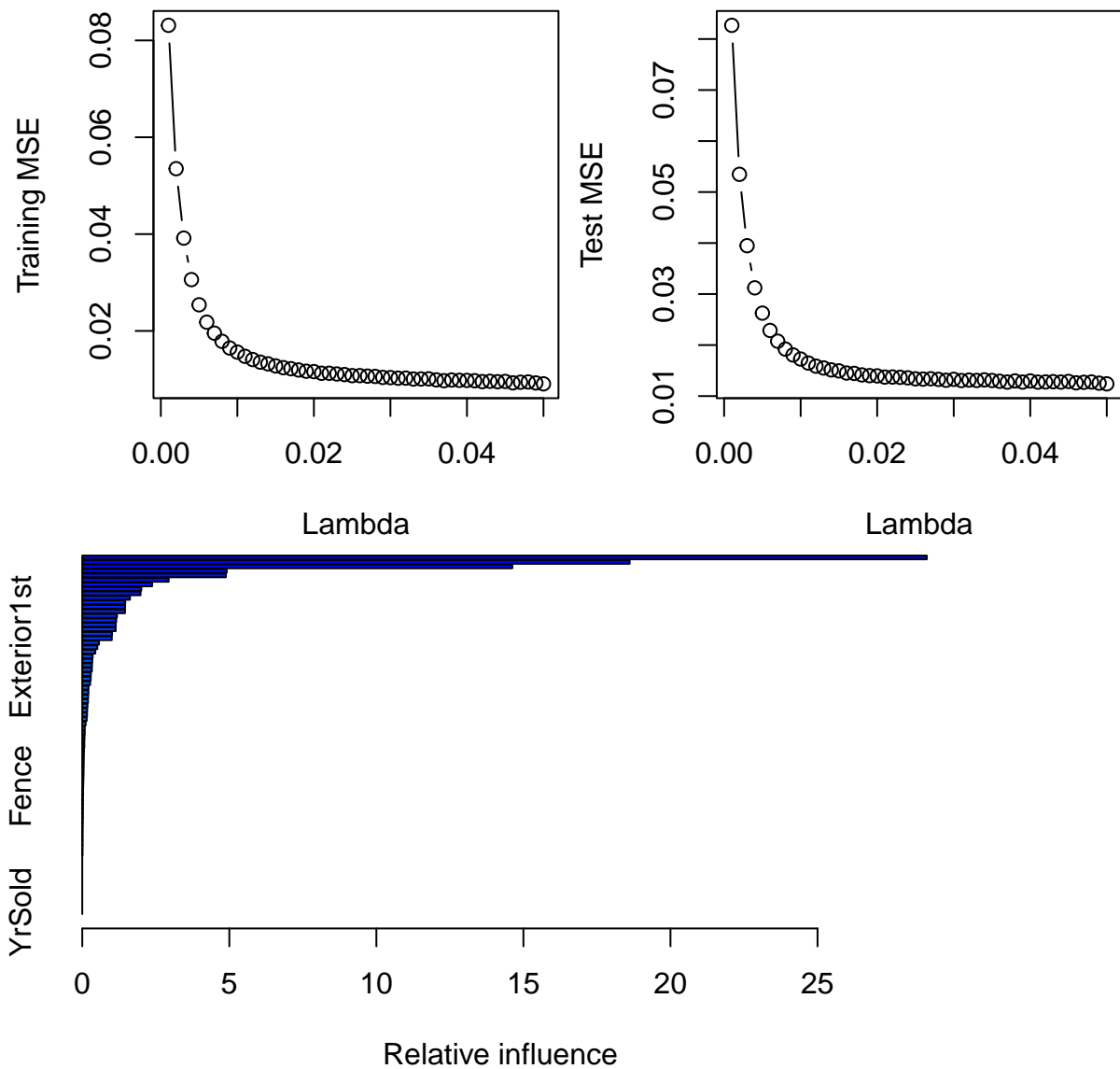


```
## [1] "LASSO Test MSE = 0.123284673944207"
```

## Boosting

An AdaBoost Boosting model was also fit as it was thought that a Boosting model may increase predictive accuracy of housing sale prices via its ability add strength or weight to specific classifiers after taking into account the previous classifier's success. It was also thought that a Boosting model would help reduce dimensionality by resulting in significant classifiers being assigned higher weights than less significant classifiers. The test MSE of the AdaBoost Boosting model was later compared with the test MSE of the other chosen models in order to choose the most accurate model for predicting Sale Price.





##		var	rel.inf
##	OverallQual	OverallQual	28.716246891
##	GrLivArea	GrLivArea	18.611557652
##	Neighborhood	Neighborhood	14.623299752
##	ExterQual	ExterQual	4.911408710
##	TotalBsmtSF	TotalBsmtSF	4.882320561
##	GarageArea	GarageArea	2.936627495
##	OverallCond	OverallCond	2.376467286
##	X1stFlrSF	X1stFlrSF	2.007089543
##	CentralAir	CentralAir	1.980939637
##	MSSubClass	MSSubClass	1.621464759
##	YearRemodAdd	YearRemodAdd	1.461286197
##	GarageType	GarageType	1.456031510
##	BsmtFinSF1	BsmtFinSF1	1.452267045
##	LotArea	LotArea	1.179430782
##	KitchenQual	KitchenQual	1.154701824
##	GarageCars	GarageCars	1.137940832

## FireplaceQu	FireplaceQu	1.135751444
## SaleCondition	SaleCondition	1.011012295
## MSZoning	MSZoning	1.003169041
## BsmtQual	BsmtQual	0.568865994
## Exterior1st	Exterior1st	0.512101236
## Functional	Functional	0.443052628
## ExterCond	ExterCond	0.351713696
## SaleType	SaleType	0.340458055
## Exterior2nd	Exterior2nd	0.330585212
## GarageYrBlt	GarageYrBlt	0.325577967
## GarageQual	GarageQual	0.292292664
## GarageCond	GarageCond	0.287472438
## Condition1	Condition1	0.268567681
## YearBuilt	YearBuilt	0.212149486
## BsmtUnfSF	BsmtUnfSF	0.209035694
## WoodDeckSF	WoodDeckSF	0.199432529
## BsmtCond	BsmtCond	0.194336819
## BsmtFinType1	BsmtFinType1	0.177307305
## Id	Id	0.166807131
## ScreenPorch	ScreenPorch	0.162221936
## BsmtExposure	BsmtExposure	0.150257243
## OpenPorchSF	OpenPorchSF	0.115028463
## FullBath	FullBath	0.086225436
## HeatingQC	HeatingQC	0.085835574
## BsmtFullBath	BsmtFullBath	0.076250445
## Fireplaces	Fireplaces	0.069755731
## Alley	Alley	0.069570557
## BedroomAbvGr	BedroomAbvGr	0.058094005
## LotFrontage	LotFrontage	0.054694329
## LotShape	LotShape	0.051908421
## X2ndFlrSF	X2ndFlrSF	0.048400463
## LotConfig	LotConfig	0.045552841
## EnclosedPorch	EnclosedPorch	0.043836923
## LowQualFinSF	LowQualFinSF	0.040564019
## HalfBath	HalfBath	0.035403540
## Fence	Fence	0.030943782
## Electrical	Electrical	0.026268402
## RoofStyle	RoofStyle	0.023973876
## PavedDrive	PavedDrive	0.021481764
## LandContour	LandContour	0.021358610
## MiscVal	MiscVal	0.019718920
## HouseStyle	HouseStyle	0.018668038
## TotRmsAbvGrd	TotRmsAbvGrd	0.018349795
## BldgType	BldgType	0.015217350
## Foundation	Foundation	0.013902940
## LandSlope	LandSlope	0.011101216
## MasVnrArea	MasVnrArea	0.010938763
## MoSold	MoSold	0.010919914
## X3SsnPorch	X3SsnPorch	0.010371747
## Heating	Heating	0.007502745
## BsmtFinSF2	BsmtFinSF2	0.006912420
## Street	Street	0.000000000
## Utilities	Utilities	0.000000000
## Condition2	Condition2	0.000000000

```

## RoofMatl      RoofMatl  0.000000000
## MasVnrType    MasVnrType 0.000000000
## BsmtFinType2  BsmtFinType2 0.000000000
## BsmtHalfBath  BsmtHalfBath 0.000000000
## KitchenAbvGr  KitchenAbvGr 0.000000000
## GarageFinish  GarageFinish 0.000000000
## PoolArea      PoolArea  0.000000000
## PoolQC        PoolQC    0.000000000
## MiscFeature    MiscFeature 0.000000000
## YrSold         YrSold    0.000000000

## [1] 0.01239545

```

## Table of Test MSE Values

In order to compare the test MSE values of the chosen models, a table was created. The table shows that the MSE for the Bagging and Random Forest models to be the smallest with Bagging having a slightly smaller MSE. However, the Random Forest model was ultimately selected as it has better interpretability when compared to Bagging, and the difference in MSE is not drastic.

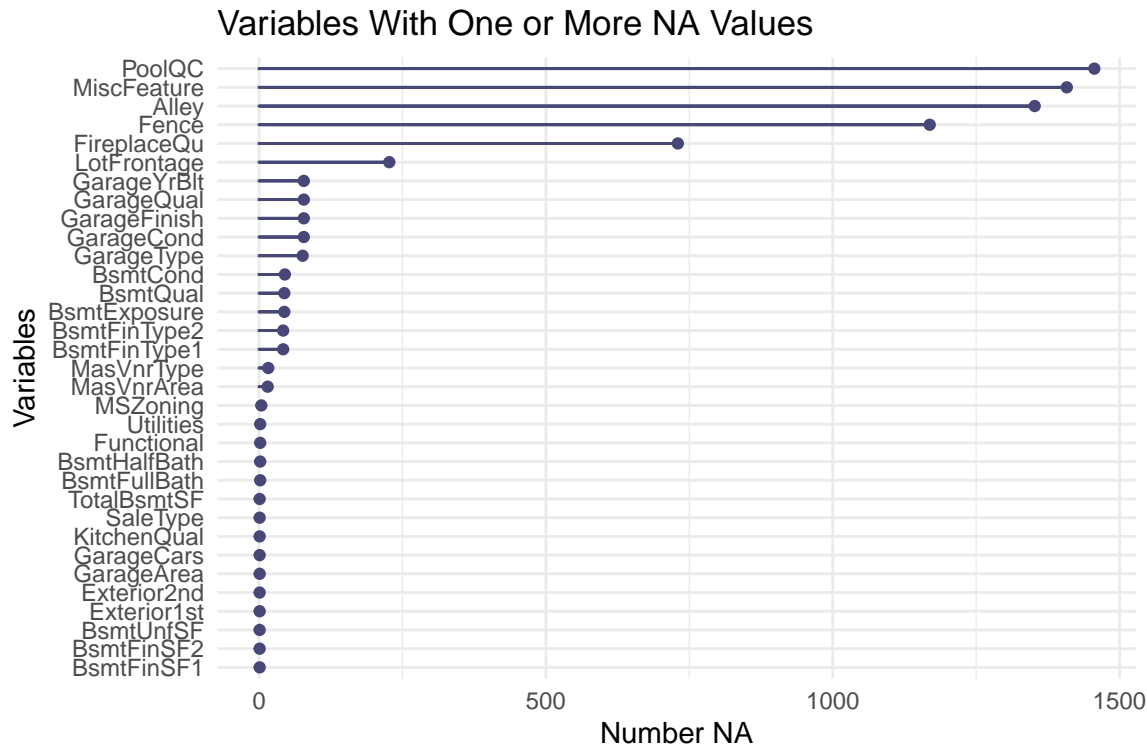
Table 1: Test MSE Values for Different Models

Model	MSE
Decision Tree	0.0394
Random Forest	0.0097
Bagged Forest	0.0096
Boosted Forest	0.0124
LASSO	0.0166

## Predict Sale Price using Random Forest

The final step of the project was to take the test data supplied from Kaggle and use our best model, we decided to use the Random Forest model for the reasons previously explored, and predict the SalePrice using the model and the test data.

The data had to be configured similarly to that of the training data. The NA values needed to be addressed in the same way and in addition to that the levels of the training data needed to be present in the test data before prediction could even be done.



Once the test data was properly configured, the Random Forest model was fit using the entire training set and the prediction of SalePrice using the configured test data was made.

## Results

In this project, five models were created and fit with training data. The top two models, Bagging and Random Forest, resulted in test MSE values of 0.0096 and 0.0097, respectively, on the portion of training data set aside for testing. The Random Forest model was selected to predict Sale Price values on the Kaggle test set. The predicted Sale Price values were submitted to Kaggle and a score of 0.14931 was given. According to the Kaggle competition, this score is calculated by taking the root mean squared error between the logarithm of the predicted value and the logarithm of the observed sales price. This corresponds to a place of 3074 out of 5396 submissions. However, most of the top 1000 submissions had RMSE scores of about 0.12, which isn't very far off from the score obtained using the Random Forest model created in this project. Overall, the selected model was proven to be fairly accurate at predicting Sale Price.

## Outside Exploration

An original hope for the project was to see if this model and data could be applied outside of the original population. We worked with the co-owner of Deluxe Homes LLC, Stu Sprecher who is a licensed contractor to explore what we had been finding and get some industry insights. The first problem that was discovered was that all of the models were showing that predictors such as Neighborhood and OverallQuality were very significant in predicting sale price. This was an issue as for starters, the neighborhoods of Ames, Iowa cannot be easily compared to a neighborhoods in Berthoud, Colorado. Additionally, these quality ratings are very subjective, what Stu may rate an 8, Dean De Cock who compiled this data, might rate as a 6 for example. Stu also pointed out that this model and data fail to consider certain aspects involved in building a home that can impact sale price. An example he kept returning to was the soil quality. Where he has started to build in Severance, Colorado has much poorer soil quality then where he was building in Berthoud, Colorado. As such, it costs a lot for to reinforce the foundation to account for that poor soil quality and as such the price of the home will increase. Additionally, prices of building materials has shifted recently, making production

of a home more costly which also impacts that sale price.

Once it was concluded that fitting the model using data provided from Stu in the Berthoud, Colorado area would not be ideal, we considered finding similar data in Ames but more recent. Through further considerations that would also prove problematic as the housing market has shifted in the last decade and it is hard to determine if the significance of a predictor such as Neighborhood has also shifted. Ultimately, it appears that the Ames, Iowa housing data appears to simply be a snap shot in time. It can help predict sale prices for that area in that given time frame but with so many influences outside of the data it is hard to use this data, and the models that are fit with the data, to predict housing prices for other locations or even dates.

## References

Ames, Iowa: Alternative to the Boston Housing Data as an ... <http://jse.amstat.org/v19n3/decock.pdf>.

“Convert Character to Factor in R: Vector, Data Frame Columns & Variable.” Statistics Globe, 14 June 2021, <https://statisticsglobe.com/convert-character-to-factor-in-r>.

Greenwell, Bradley Boehmke & Brandon. “Hands-on Machine Learning with R.” Chapter 11 Random Forests, 1 Feb. 2020, <https://bradleyboehmke.github.io/HOML/random-forest.html#fn29>.

Holtz, Yan. “Correlation Matrix with GGALLY.” – The R Graph Gallery, [https://www.r-graph-gallery.com/199-correlation-matrix-with-ggally.html#:~:text=The%20ggpairs\(\)%20function%20of,is%20displayed%20on%20the%20right](https://www.r-graph-gallery.com/199-correlation-matrix-with-ggally.html#:~:text=The%20ggpairs()%20function%20of,is%20displayed%20on%20the%20right).

“House Prices - Advanced Regression Techniques.” Kaggle, <https://www.kaggle.com/c/house-prices-advanced-regression-techniques>.

Sprecher, Stu. “Deluxe Homes LLC Housing Prices.” 1 Dec. 2021.

“Tree Based Algorithms: Implementation in Python & R.” Analytics Vidhya, 26 Aug. 2021, <https://www.analyticsvidhya.com/blog/2016/04/tree-based-algorithms-complete-tutorial-scratch-in-python/>.