# HOUSING PRICE – DATA ANALYSIS AND PREDICTION

ABSTRACT

This paper outlines the data analysis and machine learning methodology to analyze and predict housing price from medium-sized data sample.

By: Zach Nguyen

# Contents

# Executive summary

Despite having been around for a long time, Real Estate Price is one of the most unpredictable commodities in the economy. Many buyers and sellers are often unsure of the true value of their properties or those they want to buy. Traditional market research methods such as Similar home or Precedent price comparison introduce significant personal bias and subjectivity. Thus, financial mistakes are myriad in the bidding and negotiation process. This report aims to describe the process of making use of the big data available recently and cutting-edge data science and machine learning methodologies to predict the price of future Real Estate. As another point of reference for property price that is easy to use and interpret, this model can provide Real Estate participants with an advantage in negotiation. To build the model, a sample dataset was taken from Kaggle's "House Prices: Advanced Regression Techniques" and the R Programming language was use as tool of analysis.

*Section 1: Data preparation* is a minor description of the importing of the sample data and various R packages a practitioner may need when programming and reproducing the methodology. *Section 2: Data cleaning, wrangling and Imputation* describes the data preprocessing phases. Data in real life rarely conforms to ideal expectations of the data scientist. Thus, various techniques to visualize, clean and fill in missing data are conducted. *Section 3: Exploratory Data Analysis* use visualizations to familiarize the analyst with the structure, state of, and correlation between the dependent variable (Sale Price of the property), in which we hope to predict, and the independent variables, from which we base our predictions. This section also introduces a random forest implementation to grasp important variables that the analyst may need to scrutinize.

*Section 4: Feature engineering/Extraction* helps convert the survey data into human-sensible data and extract new features (predictors) based on domain understanding. These features reflect real-life concerns of a typical homebuyer looking for a piece of real estate, as surveyed by the analyst. Binary, binning and feature interaction were the techniques used to achieve feature engineering. In *Section 5: Encoding the data to be used in Machine Learning*, the analyst converts the data into machine-sensible data (using dummy variables), reduce any data bias that may prevent the machine from learning efficiently (with normalization) and remove outliers and non-zero-variance data that would introduce unnecessary noise to the model. Finally, *Section 6: Model Selection and Model Training* explains the rationale in selecting a regression model and its modification, coupled with implementation and results obtained by the model.

Overall, you would expect the model to perform with 11.4% mean error in real life. This means that it is overwhelmingly confident that the real value of a $100,000 property will lie between $90,000 and $110,000. The model's regressive nature also introduce interpretability in the form of co-efficient and may help to explain the reason behind the price of certain properties. *Section 7: Conclusion and Risk assessment* will finalize the report and recommend future improvements to the model.

# Section 1: Data preparation

This section describes the preliminary data preparation for the rest of the analysis. The data sample was imported from Kaggle.com to build the model. It represents data collected on residential homes in Ames, Iowa. Thus, it is highly advisable for an analyst to find and conduct similar techniques, but on dataset regarding to homes in proximate locations.

The following R packages were used to help with the analysis:

```
16  # Loading libraries
17
18  library(dplyr)           ## Used for Data manipulation, data wrangling
19  library(tidyr)           ## Used for Data manipulation, data wrangling
20  library(VIM)             ## Used to visualization of missing variables
21  library(mice)            ## Used for imputation of missing variables
22  library(psych)           ## Used to get Descriptive statistics
23  library(ggplot2)         ## Used to plot graphs
24  library(gridExtra)       ## Used to arrange visualization
25  library(corrplot)        ## Used for drawing correlation plots
26  library(caret)           ## Used for regression
27  library(randomForest)    ## Used to conduct random forest
28
```

Data was also imported, and non-predictor variables such as ID removed:

```
 5  # Importing the data
 6  train <- read.csv("train.csv", stringsAsFactors = FALSE)
 7  test <- read.csv("test.csv", stringsAsFactors = FALSE)
 8
 9  test_id <- test$Id
10  test$Id <- NULL
11  train$Id <- NULL
12
13  test$SalePrice <- NA
14  full <- rbind(train,test)
```

# Section 2: Data cleaning, wrangling and Imputation

This section depicts the steps to clean and fill in missing data from the dataset. Upon closer look at the dataset, we recognize that some of the data is missing, some are illogical/incorrect, and some are too few and far in between to matter.

We find and clean some illogical data points:

```
50
51   # Find any year value above 2010 and replace with suitable values
52   full %>%
53     select(GarageYrBlt, YrSold, YearBuilt, YearRemodAdd) %>%
54     filter_all(any_vars(.>2010))
55   full[which(full$GarageYrBlt  > 2010), "GarageYrBlt"] <- 2007
56
```

```
56
57   # Find illogical year value (Sold before Built, Remodeled before Built, Garage Built before Built)
58   full %>%
59     select(GarageYrBlt, YrSold, YearBuilt, YearRemodAdd) %>%
60     filter(YrSold < YearBuilt | YearRemodAdd < YearBuilt | GarageYrBlt < YearBuilt)
61   full[which(full$YearBuilt > full$YrSold), "YrSold"] <-  2008
62   full[which(full$YearBuilt > full$YearRemodAdd), "YearRemodAdd"] <- 2002
63
```

```
63
64   # Find properties with missing garage year built, set them to year built
65   which(!is.na(full$GarageType) & is.na(full$GarageYrBlt))
66   full[c(2127,2577), c("GarageType", "GarageYrBlt", "YearBuilt")]
67   full[c(2127,2577), "GarageYrBlt"] <- full[c(2127,2577), "YearBuilt"]
68
```

We find sparse categories and assign them to the nearest neighbor:

```
68
69   # Find Predictors catagories with too few data point and assign to nearest neighbors
70   which(full$MSSubClass == 150)
71   full[2819, "MSSubClass"] <- 160
72
73   which(full$TotRmsAbvGrd == 13 | full$TotRmsAbvGrd == 14 | full$TotRmsAbvGrd == 15)
74   full[c(636,1903,2550), "TotRmsAbvGrd"] <- 12
75
76   which(full$Fireplaces == 4)
77   full[2711, "Fireplaces"] <- 3
78
79   which(full$GarageCars == 5)
80   full[1829, "GarageCars"] <- 4
81
```
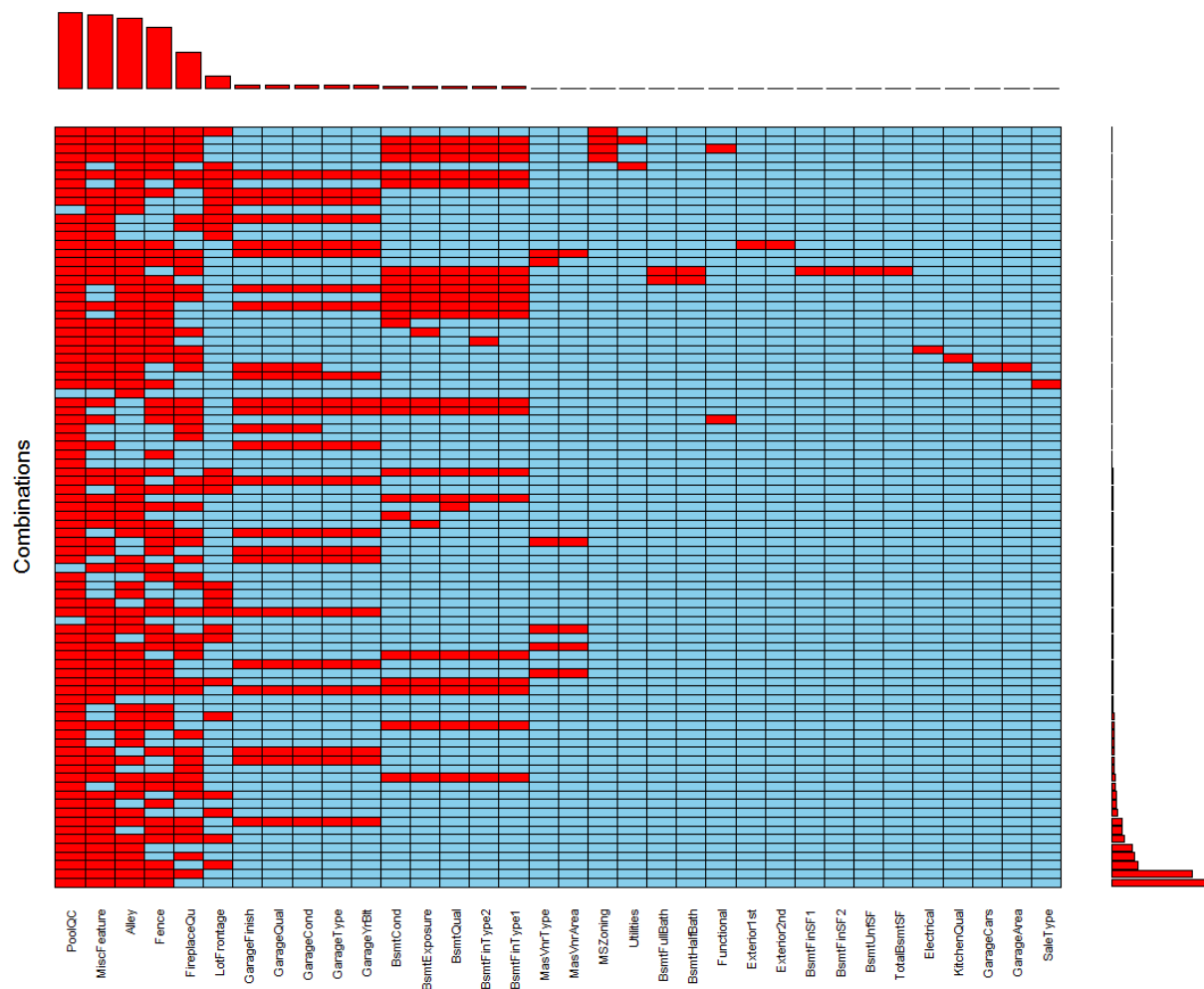
In order to tackle missing data, we first visualize the structure of the missing data:

```
42   # Visualize the missingness of the data
43
44   var_with_missing <- colnames(full[colSums(is.na(full)) >0])
45   var_with_missing <- var_with_missing[-35]
46   aggr(full[,var_with_missing],
47       combined = TRUE, bars = TRUE, sortVars = TRUE, cex.axis = 0.69)
```

We also recognize that while some missing data are Missing Completely at Random (MCAR), some missing data such as PoolQC (the variable with the most number data point missing) does not appear random and has a meaning according to the data dictionary provided. This mean those are Missing Not at Random (MNAR). To reasonably impute the missing data, we will manually fill in variables with MNAR data and use the MICE package to impute the rest (MCAR).

**Missing completely at random**
- Using regression can help determine suitable value
- Multiple Imputation by Chained Equation (MICE)

**Missing not at random**
- Require understanding of missing structure
- Manual imputation required

We fill in the data Missing not at random:

```
84
85   # Group MNAR vars, where NA indicate the subject is not there
86   var_MNAR <- c("Alley", "BsmtQual", "BsmtCond","BsmtExposure","BsmtFinType1",
87                     "BsmtFinType2","FireplaceQu","GarageType", "GarageYrBlt",
88                     "GarageFinish","GarageQual","GarageCond","PoolQC","Fence",
89                     "MiscFeature")
90
91   # Create a function to fill in missing data
92 - NA_to_NotAvailable <- function (data, predictor) {
93                     levels(data[,predictor]) <- c(levels(data[,predictor]), "Not Available")
94                     data[,predictor][is.na(data[,predictor])] <- "Not Available"
95                     return(data[,predictor])
96                     }
97
98 - for (i in 1:length(var_MNAR)) {
99                     full.MNAR[,var_MNAR[i]] <- NA_to_NotAvailable(full, var_MNAR[i])
100                    }
101
```
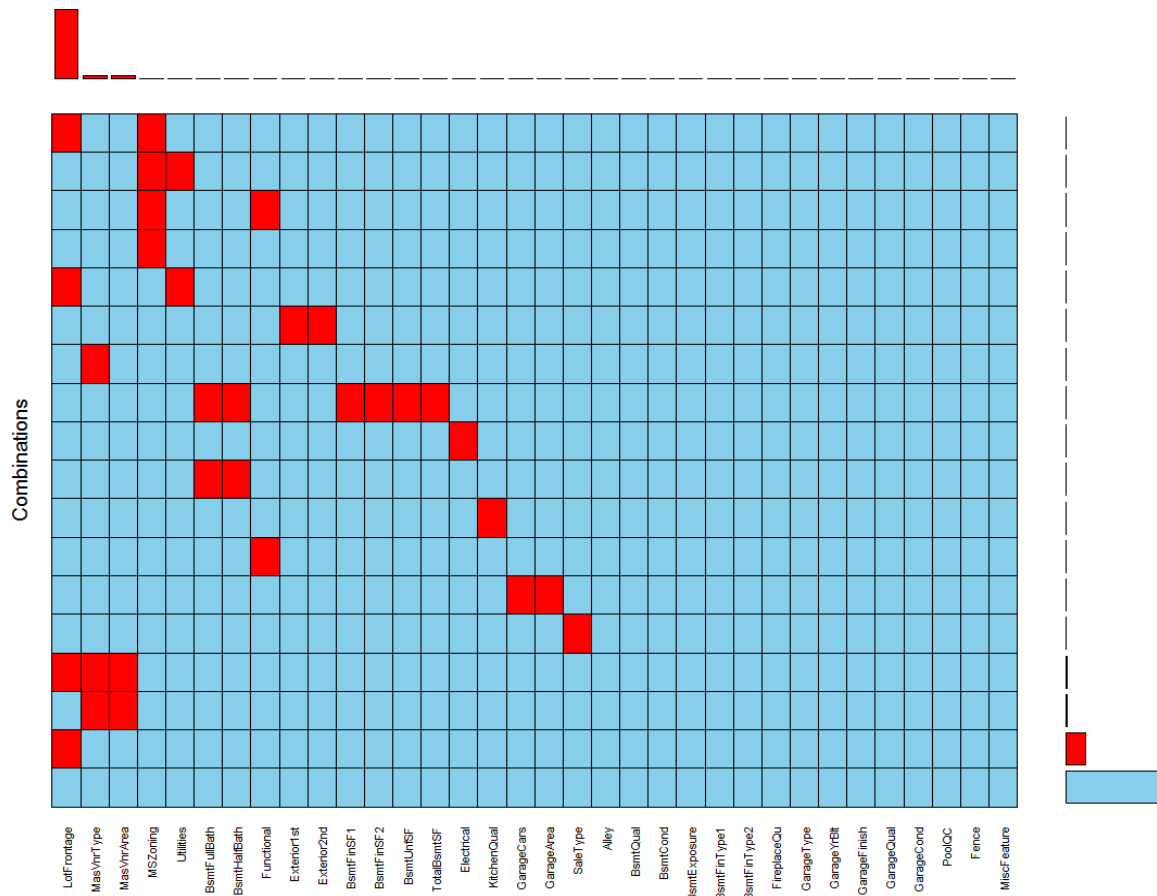
And visualize the data after the imputation:

```
101
102   # Visualize data after MNAR variables imputation
103   aggr(full.MNAR[,var_with_missing],
104         combined = TRUE, bars = TRUE, sortVars = TRUE, cex.axis = 0.69)
105
```

It seems that quite a lot of the missing data has been dealt with. Now we use M.I.C.E to impute the remaining missing data points.

```
112
113   # Group the remaining MAR variables
114   var_MAR <- c("MasVnrType", "MSZoning", "Utilities", "Functional",
115                 "Exterior1st", "Exterior2nd", "Electrical",
116                 "KitchenQual", "SaleType")
117
118   # Convert these into factors for mice imputation
119   full.MNAR <- full.MNAR %>%
120         mutate_at(var_MAR, funs(factor(.)))
121
122   # We use cart method since imputation is fairly linear
123   imputed.data <- mice(full.MNAR[,!(names(full.MNAR) %in% "SalePrice")], m=1, method='cart')
124   full.complete <- mice::complete(imputed.data, action = 1, include = FALSE)
125
```
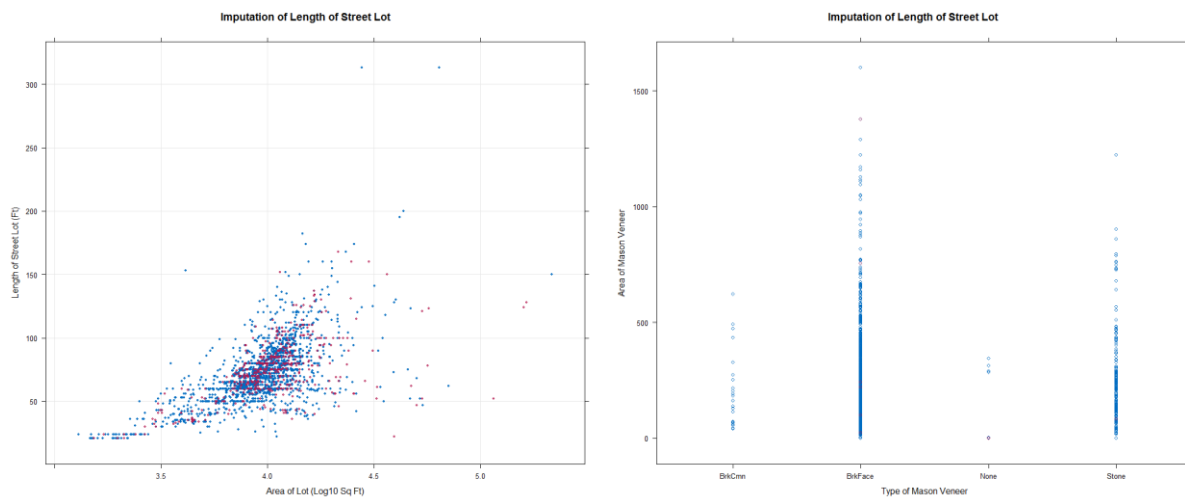
We visualize our imputation to make sure that the values are befitting:

```
125
126   # Visualize Lot Front and Mass Veneer which has the most imputation.
127   # We also use log transformation to linearize graph and aid visualization.
128   xyplot(imputed.data,
129         LotFrontage ~ log10(LotArea), pch = 20, grid = TRUE,
130         xlab = "Area of Lot (Log10 Sq Ft)",
131         ylab = "Length of Street Lot (Ft)",
132         main = "Imputation of Length of Street Lot")
133
134   xyplot(imputed.data,
135         MasVnrArea ~ MasVnrType,
136         xlab = "Type of Mason Veneer",
137         ylab = "Area of Mason Veneer",
138         main = "Imputation of Length of Street Lot")
139
```



The values are quite close, and that's what we want.

```
139
140   # Finalize our data set after imputation
141   train.complete <- full.complete[1:1460,]
142   train.complete$SalePrice <- train$SalePrice
143   test.complete <- full.complete[1460:2919,]
144
```
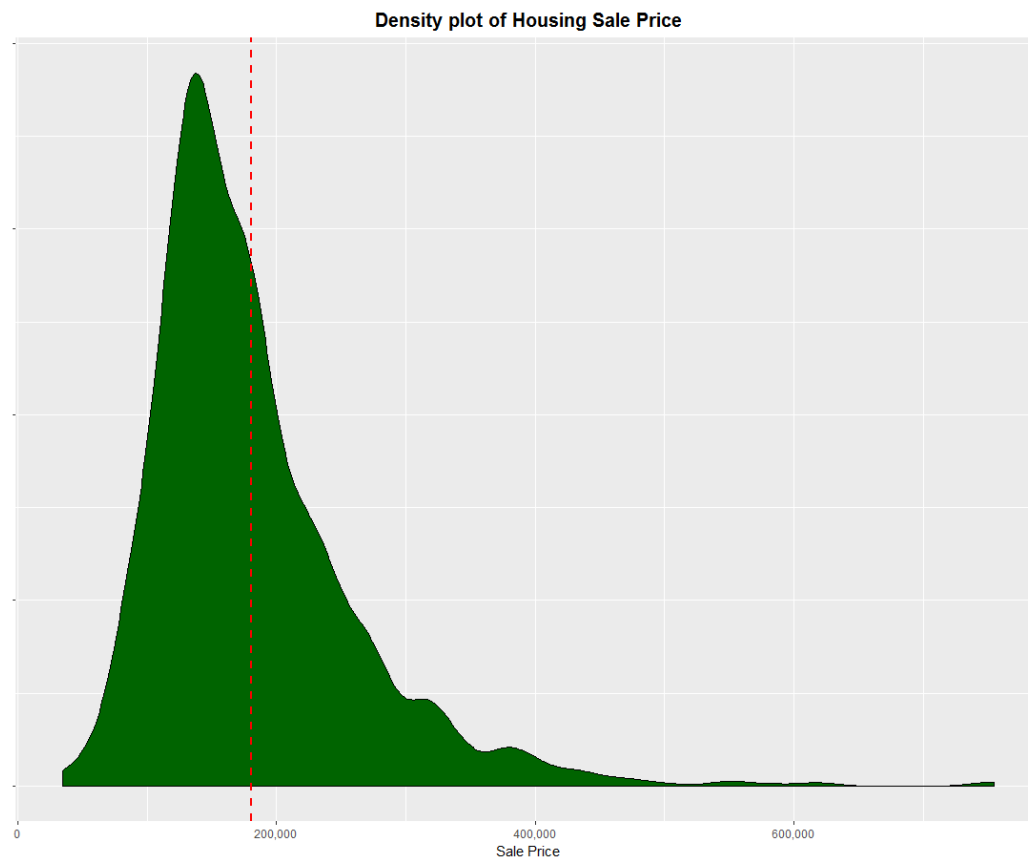
# Section 3: Exploratory Data Analysis

This section helps the analyst gain intuition on the story the data is telling and aid in feature selection and sanity testing. Having a complete dataset, we can now explore the dataset to gain more understanding and intuition behind the numbers presented. It is, therefore, important to visualize the data structure of the dependent variable, some independent variable that stands out and their correlations.

First, we look at the dependent variable or response variable:

```
154
155    # Visualize the dependent variable SalePrice
156  density_plot_depvar <- ggplot(train.complete, aes(x= SalePrice)) +
157    geom_density(color = "black", fill = "darkgreen") +
158    geom_vline(aes(xintercept = mean(SalePrice)),
159              color = "red", linetype = "dashed", size = 1) +
160    scale_x_continuous(labels = scales::comma) +
161    labs(title = "Density plot of Housing Sale Price", x = "Sale Price", y = "Density") +
162    theme(plot.title = element_text(hjust = 0.5, size = 15, face = "bold"),
163          axis.text.y = element_blank(),
164          axis.title.y = element_blank())
165  density_plot_depvar
166
```

**Density plot of Housing Sale Price**



It seems that the average price of a property is around $175,000. However, many more properties have below average price while the expensive properties "tail off". This makes sense, as a small percentage of houses is
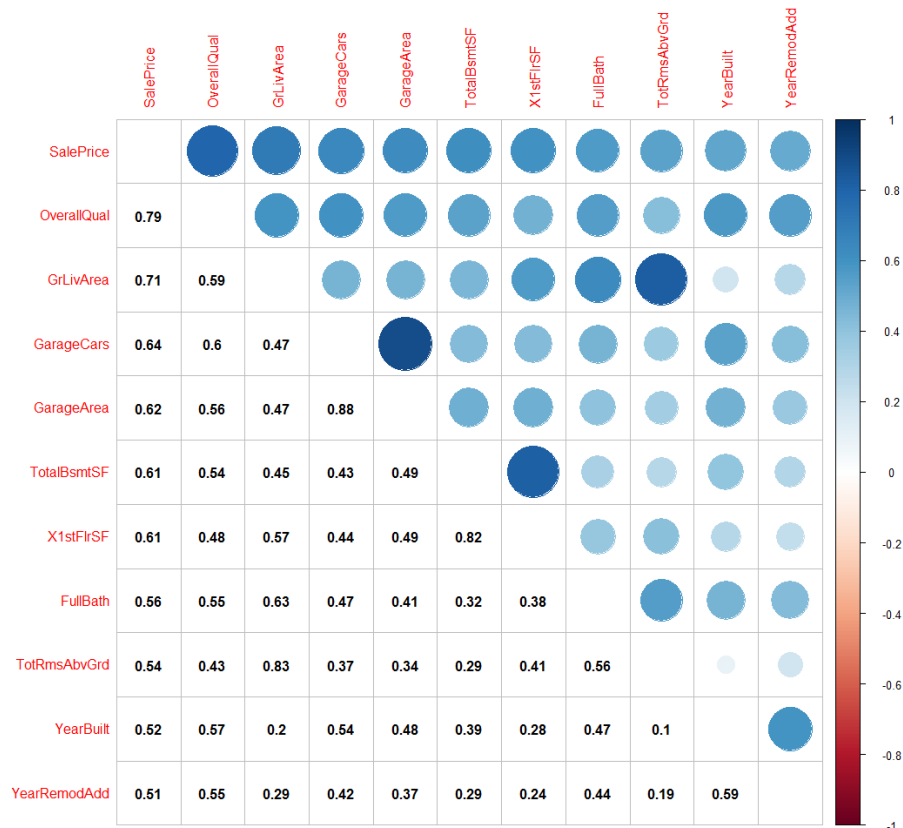
expected to have extraordinary features (mansion, huge land lots, etc. …) that command much higher price. Therefore, the Sale Price variable may be negatively skewed and will be normalized in later parts.

Next, we'll look at correlations of numeric variables to the Sale Price. This may hint us to some numeric variables that contributes significantly to predicting Sale Price and is, therefore, worth the effort.

```
167  # Visualization of numeric variable correlations to SalePrice
168
169  numeric_var <- which(sapply(train.complete, is.numeric))
170  numeric_var_names <- names(numeric_var)
171
172  # Make correlation chart with all numeric variable
173  cor <- cor(train.complete[,numeric_var], use = "pairwise.complete.obs")
174
175  # Sort chart column in descending order
176  cor_sales_vector <- as.matrix(sort(cor[,"SalePrice"], decreasing = TRUE))
177
178  # Find name index of columns, row with high correlation
179  high_cor_var <- names(which(apply(cor_sales_vector,1, function(x) abs(x)> 0.5)))
180
181  # Plot correlation chart
182  cor.data <- as.data.frame(cor)[high_cor_var, high_cor_var] # Subset correlation chart
183  cor.data <- cor.data[order(-cor.data$SalePrice),] # Sort row in descending order
184  corrplot.mixed(as.matrix(cor.data), lower = "number", lower.col = "black", tl.pos = "lt")
185
```

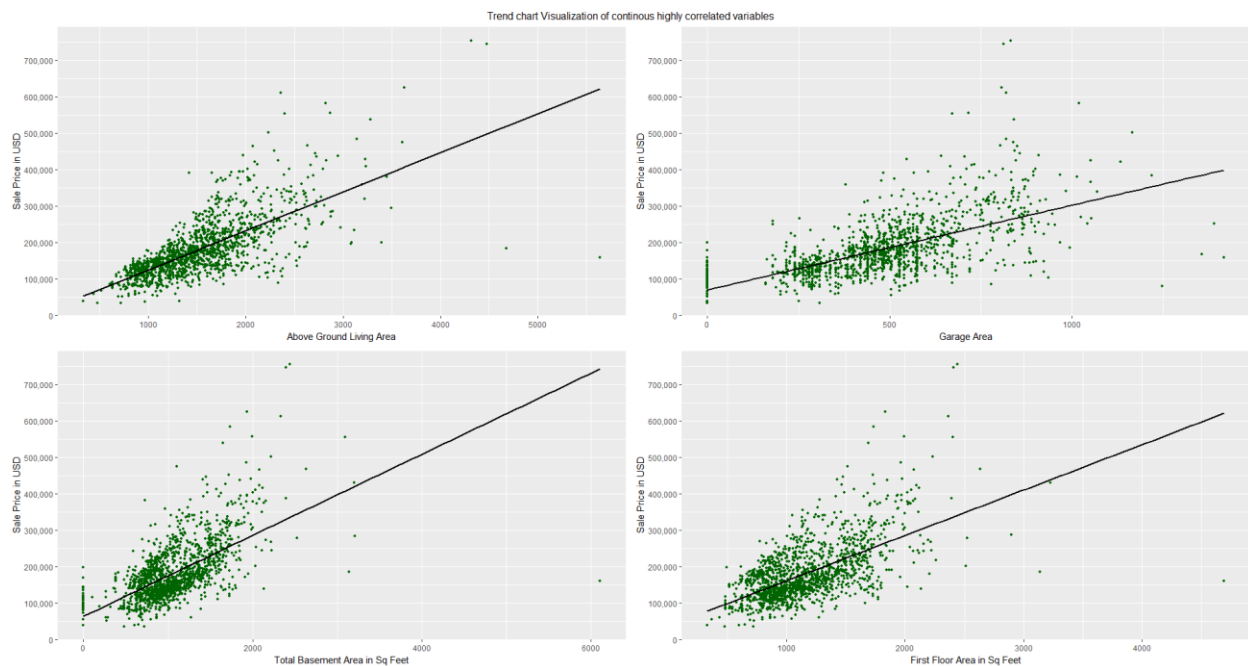| | SalePrice | OverallQual | GrLivArea | GarageCars | GarageArea | TotalBsmtSF | X1stFlrSF | FullBath | TotRmsAbvGrd | YearBuilt | YearRemodAdd |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SalePrice | | | | | | | | | | | |
| OverallQual | 0.79 | | | | | | | | | | |
| GrLivArea | 0.71 | 0.59 | | | | | | | | | |
| GarageCars | 0.64 | 0.6 | 0.47 | | | | | | | | |
| GarageArea | 0.62 | 0.56 | 0.47 | 0.88 | | | | | | | |
| TotalBsmtSF | 0.61 | 0.54 | 0.45 | 0.43 | 0.49 | | | | | | |
| X1stFlrSF | 0.61 | 0.48 | 0.57 | 0.44 | 0.49 | 0.82 | | | | | |
| FullBath | 0.56 | 0.55 | 0.63 | 0.47 | 0.41 | 0.32 | 0.38 | | | | |
| TotRmsAbvGrd | 0.54 | 0.43 | 0.83 | 0.37 | 0.34 | 0.29 | 0.41 | 0.56 | | | |
| YearBuilt | 0.52 | 0.57 | 0.2 | 0.54 | 0.48 | 0.39 | 0.28 | 0.47 | 0.1 | | |
| YearRemodAdd | 0.51 | 0.55 | 0.29 | 0.42 | 0.37 | 0.29 | 0.24 | 0.44 | 0.19 | 0.59 | |

It seems OverallQual score is understandably very highly correlated with SalePrice, along with GrLivArea, 1stFlrSF and TotalBsmtSF. People do prefer bigger houses after all. The YearBuilt and YearRemodAdd variables here also suggest people are looking for newer houses. It is also worth noting that there are highly

correlated variables such as GarageCars and GarrageAreas which can cause multicollinearity and bias the model and must be eliminated in latter parts. For now, let's take a closer look at these variables (continuous and discreet) through visualization:

```
185
186   # Visualization of relationships between predictors highly correlated to SalePrice
187
188 ▾ Visualize_numeric <- function(data, n_predictor, xname) {
189     ggplot(data, aes(x = data[,n_predictor], y = SalePrice)) +
190       geom_point(na.rm = TRUE, col = "darkgreen", size = 1) +
191       geom_smooth(na.rm = TRUE, method = "lm", se=FALSE, color="black", aes(group=1)) +
192       scale_y_continuous(breaks = seq( 0, max(data$SalePrice), by = 100000), labels = scales::comma) +
193       ylab("Sale Price in USD") +
194       xlab(xname)
195   }
196
197   p1 <- Visualize_numeric(train.complete, "GrLivArea", "Above Ground Living Area")
198   p2 <- Visualize_numeric(train.complete, "GarageArea", "Garage Area")
199   p3 <- Visualize_numeric(train.complete, "TotalBsmtSF", "Total Basement Area in Sq Feet")
200   p4 <- Visualize_numeric(train.complete, "X1stFlrSF", "First Floor Area in Sq Feet")
201   grid.arrange(p1, p2, p3, p4, ncol = 2, top = "Trend chart Visualization of continous highly correlated variables")
202
```



Trend chart Visualization of continous highly correlated variables

We can definitely see a trend line in the relationship between these variable and Sale Price, the dependent variable. These variables hold promise as important predictors for Sale Price. However, we also notice some outliers.

```
202
203   # We can also spot the outliers          > which(train$GrLivArea > 4500)
204                                            [1]   524 1299
205   which(train$GrLivArea > 4500)            > which(train$GarageArea > 1250)
206   which(train$GarageArea > 1250)           [1]   582 1191 1299
207   which(train$TotalBsmtSF > 4000)          > which(train$TotalBsmtSF > 4000)
208   which(train$X1stFlrSF > 4000)            [1] 1299
209                                            > which(train$X1stFlrSF > 4000)
                                               [1] 1299
```
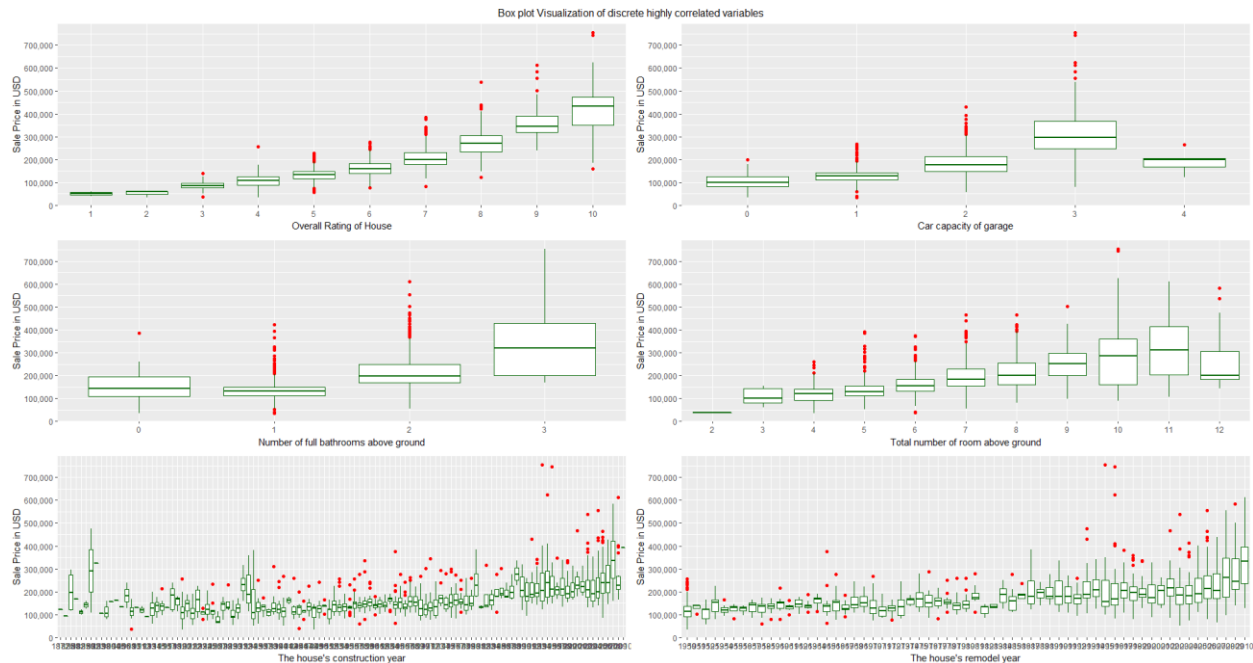
It is necessary to consider removing observation 524 and 1299 at the end. Their low Sale Price despite huge square footage may signify consequence of an imperfect market (e.g.: family, heir transactions, under-the-floor deals …) or failed negotiation. Though we would generally expect some of these noises to get through our data collection, we can create a better fit model by removing them.

For now, we'll visualize discrete category variables from the correlation plot with Boxplots:

```
209
210  # Visualization of discrete numeric variables
211▾ Visualize_category <- function(data, c_predictor, xname) {
212    ggplot(data, aes(x= factor(data[,c_predictor]), y = SalePrice)) +
213      geom_boxplot(na.rm = TRUE, col='darkgreen', outlier.color = "red") +
214      scale_y_continuous(breaks= seq( 0, max(data$SalePrice), by = 100000), labels = scales::comma) +
215      ylab("Sale Price in USD") +
216      xlab(xname)
217  }
218
219  p5 <- Visualize_category(train.complete, "OverallQual", "Overall Rating of House")
220  p6 <- Visualize_category(train.complete, "GarageCars", "Car capacity of garage")
221  p7 <- Visualize_category(train.complete, "FullBath", "Number of full bathrooms above ground")
222  p8 <- Visualize_category(train.complete, "TotRmsAbvGrd", "Total number of room above ground")
223  p9 <- Visualize_category(train.complete, "YearBuilt", "The house's construction year")
224  p10 <- Visualize_category(train.complete, "YearRemodAdd", "The house's remodel year")
225  grid.arrange(p5, p6, p7, p8, p9, p10, ncol = 2, top = "Box plot Visualization of discrete highly correlated variables")
226
```



Box plot Visualization of discrete highly correlated variables

The correlation with SalePrice for these variables can be spotted, especially OverallQual (overall rating). Categories from GarageCars, Full Bath and TotRmsAbvGrd has very discrete boxes (in terms of box height and length) and suggests many areas of differences that could be explored. The Year variables, however, are too spread out (longer boxes) to see a definite trend.

In short, the correlation plot gives us an idea of some classes of important numerical predictors, namely those related to overall quality, garage size and square footage of living spaces. However, we would want to explore important categorial variables as well. Therefore, the plan is to pick up important independent variables through a preliminary random forest implementation.

First, we need to categorize all our independent variables so that R can understand how to feed our data into the machine learning algorithm:

```
228
229    # Specify Factor, ordered factor and numeric values
230  factor_var <- c("MSSubClass", "MSZoning", "Street", "Alley", "LandContour", "LotConfig",
231              "Neighborhood", "Condition1", "Condition2", "BldgType", "RoofStyle", "RoofMatl",
232              "Exterior1st", "Exterior2nd", "MasVnrType", "Foundation", "Heating", "GarageType",
233              "PavedDrive", "MiscFeature", "SaleType", "SaleCondition", "CentralAir", "MSZoning",
234              "HouseStyle", "Fence", "MoSold")
235  ordered_var <- c("LotShape", "Utilities", "LandSlope", "OverallQual",
236              "OverallCond", "ExterQual", "ExterCond", "BsmtQual",
237              "BsmtCond", "BsmtExposure", "BsmtFinType1", "BsmtFinType2", "HeatingQC", "Electrical",
238              "KitchenQual", "Functional", "FireplaceQu", "GarageFinish", "GarageCars",
239              "GarageQual", "GarageCond", "PoolQC", "GarageYrBlt",
240              "BedroomAbvGr", "KitchenAbvGr", "TotRmsAbvGrd", "Fireplaces")
241  numeric_var <- c("LotFrontage", "YearBuilt", "YearRemodAdd", "YrSold", "LotArea", "MasVnrArea",
242              "BsmtFinSF1", "BsmtFinSF2", "BsmtUnfSF", "TotalBsmtSF","X1stFlrSF", "X2ndFlrSF",
243              "LowQualFinSF", "GrLivArea", "GarageArea", "WoodDeckSF", "OpenPorchSF",
244              "EnclosedPorch", "X3SsnPorch", "ScreenPorch", "PoolArea", "MiscVal",  "FullBath",
245              "HalfBath", "BsmtFullBath", "BsmtHalfBath")
246   # Check: length(factor_var) + length(ordered_var) + length(numeric_var) # We have 80 variables
247
```
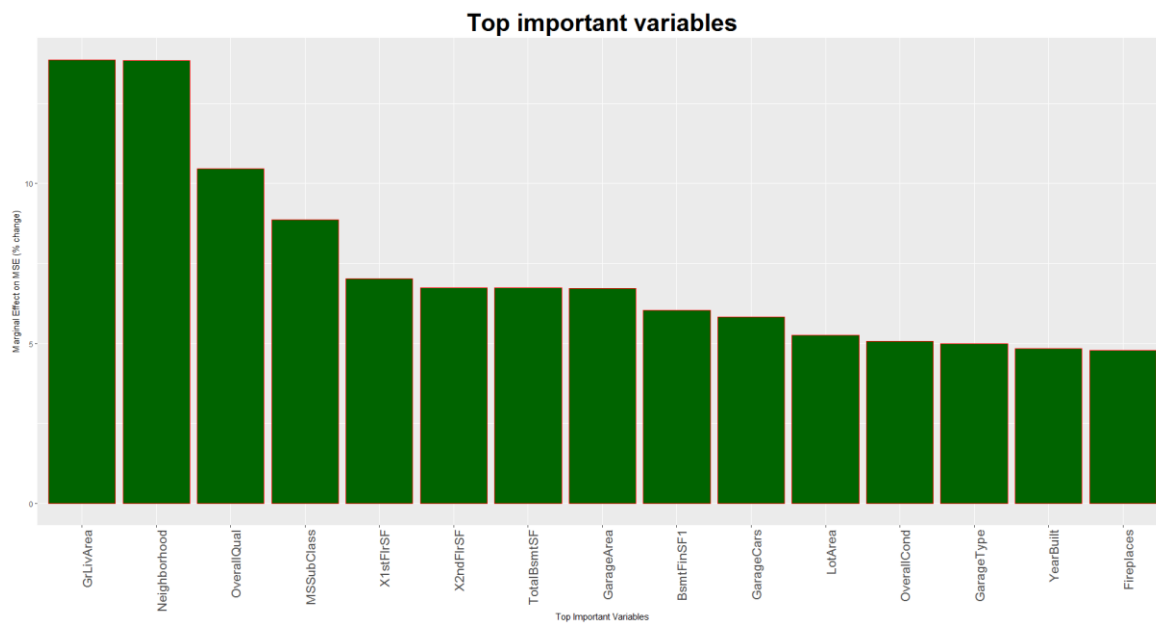
Then we train a random forest and use its feature importance methodology to pick out 20% of the most important predictors it considered. The predictors are chosen based on the sensitivity of the model's MSE to changes in said predictors. The 20% choice is a heuristic used only for data exploration, since 20% of predictors usually net 80% of the predictive power of the Machine Learning model.

```
247
248  # Duplicate the dataset to be used in random forest
249  full.rf <- full.complete %>%
250    mutate_at(factor_var, funs(factor(.))) %>%
251    mutate_at(ordered_var, funs(ordered(.)))
252
253   # str(full.rf)
254
255  # Train random forest and pick out 20% of important variables
256  train.rf <- full.rf[1:1460,]
257  train.rf$SalePrice <- train$SalePrice
258
259  set.seed(123)
260  Rf <- randomForest(x = train.rf[,-80], y=train.rf$SalePrice, ntree=100, importance=TRUE)
261  var.importance <- importance(Rf)
262  var.importance <- data.frame(Variables = row.names(var.importance), MSE = var.importance[,1])
263  var.importance <- var.importance[order(var.importance$MSE, decreasing = TRUE),]
264  var.importance$Variables <- as.factor(var.importance$Variables)
265
266  # Plot important predictor variables
267  ggplot(var.importance[1:(nrow(var.importance)/5),], aes(x = reorder(Variables, desc(abs(MSE))), y = MSE)) +
268    geom_bar(stat = "identity", color = "red", fill = "darkgreen") +
269    labs(title = "Top important variables", x = "Top Important Variables", y = "Marginal Effect on MSE (% change)") +
270    theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 1, size = 16),
271          plot.title = element_text(hjust = 0.5, size = 15, face = "bold"))
272
```

**Top important variables**



Random forest mostly supports our findings from the correlation plot (overall quality, square footage variables, garage and year built). From here, it is clear that most of our efforts should be on these variables. We also discovered important non-numeric categorical variables such as Neighborhood and MSSubClass.
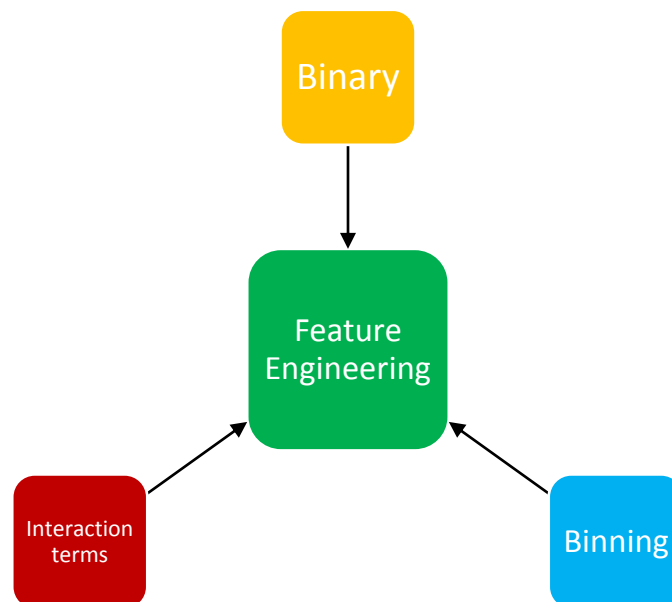
# Section 4: Feature Engineering & Extraction

This section describes the feature engineering process. Feature engineering is extremely important in making accurate predictions. Machine Learning models running on insightful data will obtain much better results than those on raw collected data. Due to importance of domain knowledge to feature engineering, I have done market surveys with a sample of 17 of my friends and acquaintances in Toronto to figure out what are the most common questions one would ask when buying a house. Even though bias (demographic, geographic, personal) maybe introduced into the process, it is reasonable to assume that real-estate seekers in North America have more concerns in common than not.

To begin, let us build a function to visualize extractable predictors to see if the common concerns really match the data:

```
276
277     # Build a function to help visualize available variables to extract new sensible ones
278  full.final <- full.complete
279  Visualize_suspect_var <- function(data, s_predictor, xname) {
280              ggplot(data, aes(x = as.factor(data[,s_predictor]), y = data[,"SalePrice"])) +
281                  geom_bar(stat = "summary", fun.y = "mean", fill = "green") +
282                  geom_bar(stat = "summary", fun.y = "median", fill = "blue", alpha = 0.7) +
283                  scale_y_continuous(labels = scales::comma) +
284                  geom_hline(yintercept= median(data$SalePrice), linetype="dashed", color = "red", size = 1.5) +
285                  geom_hline(yintercept= mean(data$SalePrice), linetype="dashed", color = "orange", size = 1.5) +
286                  geom_label(stat = "count", aes(label = ..count.., y = ..count..)) +
287                  ylab("Sale Price in USD") +
288                  xlab(xname)
289              }
```

Let us also categorize the methods in which we can extract new features from outstanding ones:
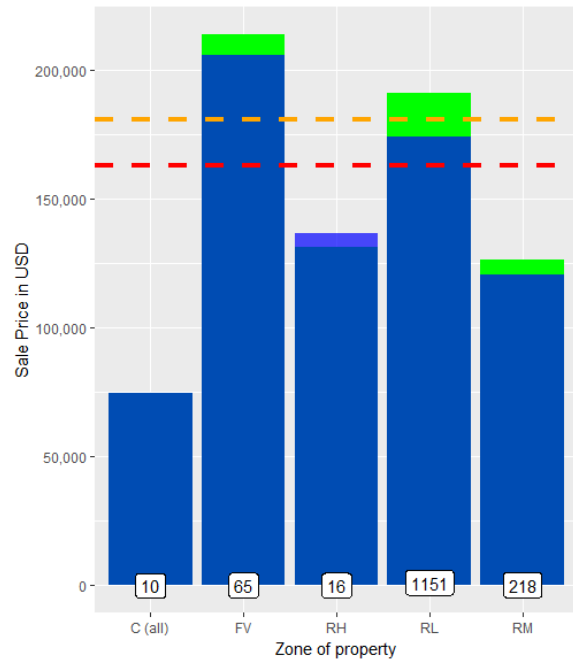


Some new features can be extracted using the survey's recurring questions and the visualization of the data to recognize distinct binary categories. These features describe whether the property is densely populated, new, bought in peak season, has deformed shapes, has detached garage, has paved drive, has been remodelled and bought under a mortgage. These all seem legitimate concerns that would be affect the price of a property. Thus, we devise the categories of the respective variables into reasonable binaries.

```
292
293   # Is the location densely populated?
294   Visualize_suspect_var(train.complete, "MSZoning", "Zone of property")
295   full.final$IsDensePop <- ifelse(full.final$MSZoning %in% c("FV", "RL"),"Yes", "No")
296   full.final$IsDensePop <- as.factor(full.final$IsDensePop)
297
```
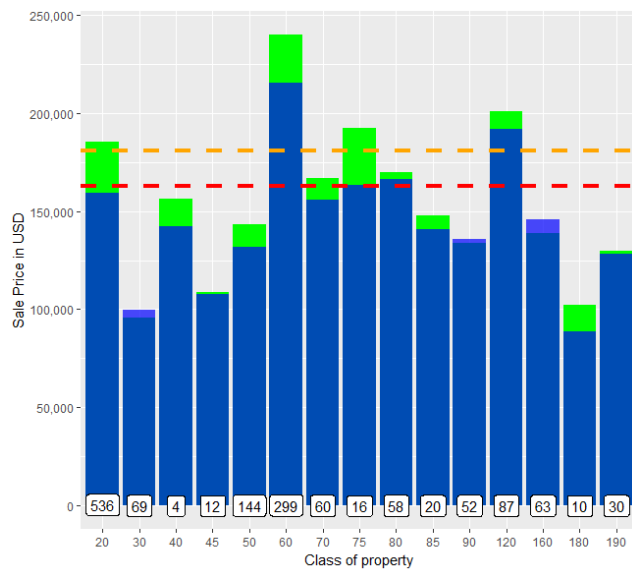


```
297
298   # Is the property newer than most?
299   Visualize_suspect_var(train.complete, "MSSubClass", "Class of property")
300   full.final$IsNewer <- ifelse(full.final$MSSubClass %in% c("20","60", "120","160"), "Yes", "No")
301   full.final$IsNewer <- as.factor(full.final$IsNewer)
302
```
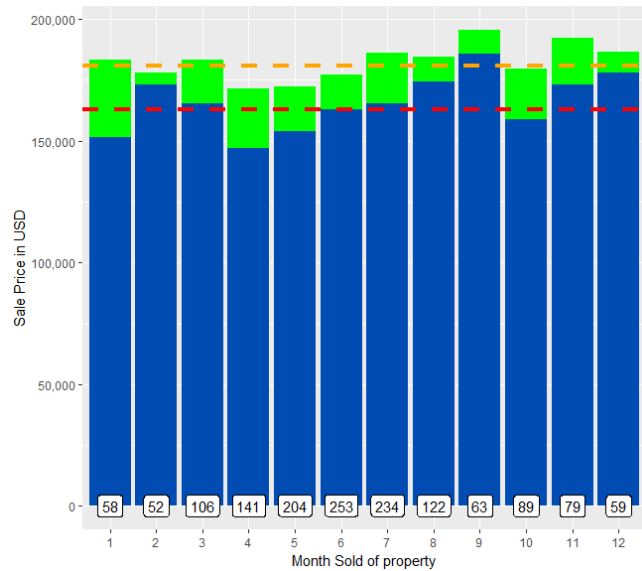
```
302
303   # Was the property bought during peak season?
304   Visualize_suspect_var(train.complete, "MoSold", "Month Sold of property")
305   full.final$IsPeak <- ifelse(full.final$MoSold %in% c("3","4", "5","6", "7", "8"), "Yes", "No")
306   full.final$IsPeak <- as.factor(full.final$IsPeak)
307
```
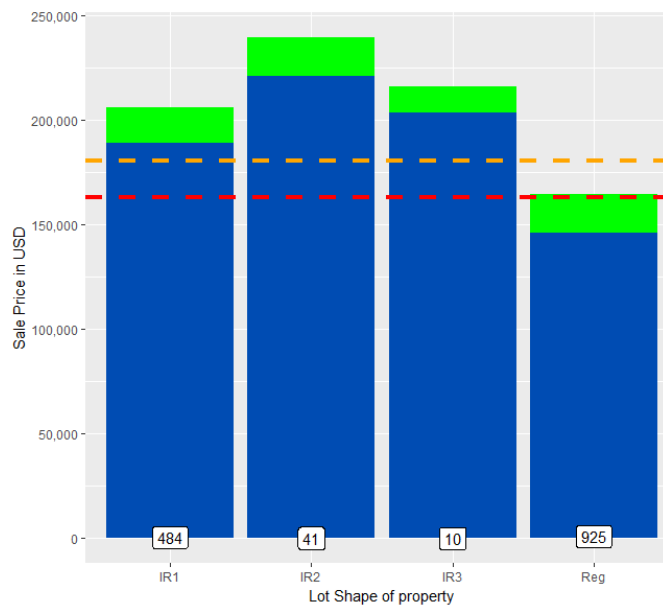


```
312
313   # Is lot shape regular, any deformities?
314   Visualize_suspect_var(train.complete, "LotShape", "Lot Shape of property")
315   full.final$IsRegular <- ifelse(full.final$LotShape %in% c("Reg"), "Yes", "No")
316   full.final$IsRegular <- as.factor(full.final$IsRegular)
317
```
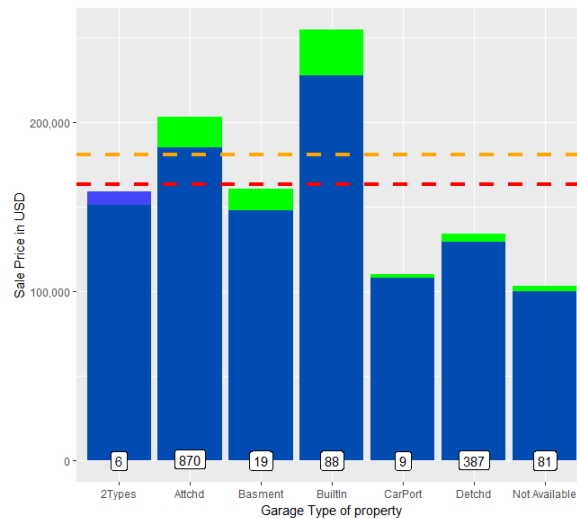
```
321
322  # Is Garage Detached?
323  Visualize_suspect_var(train.complete, "GarageType", "Garage Type of property")
324  full.final$IsDetached <- ifelse(full.final$GarageType %in% c("Detchd", "Not Available", "CarPort"),
325                                  "Yes", "No")
326  full.final$IsDetached <- as.factor(full.final$IsDetached)
327
```
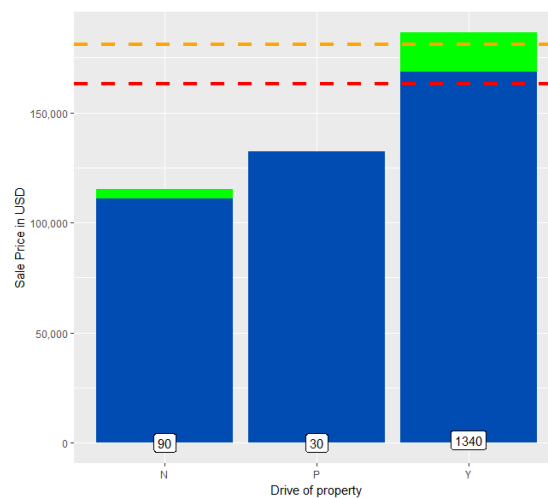


```
327
328  # Is the Drive Paved?
329  Visualize_suspect_var(train.complete, "PavedDrive", "Drive of property")
330  full.final$IsPaved <- ifelse(full.final$PavedDrive %in% c("Y"), "Yes", "No")
331  full.final$IsPaved <- as.factor(full.final$IsPaved)
332
```
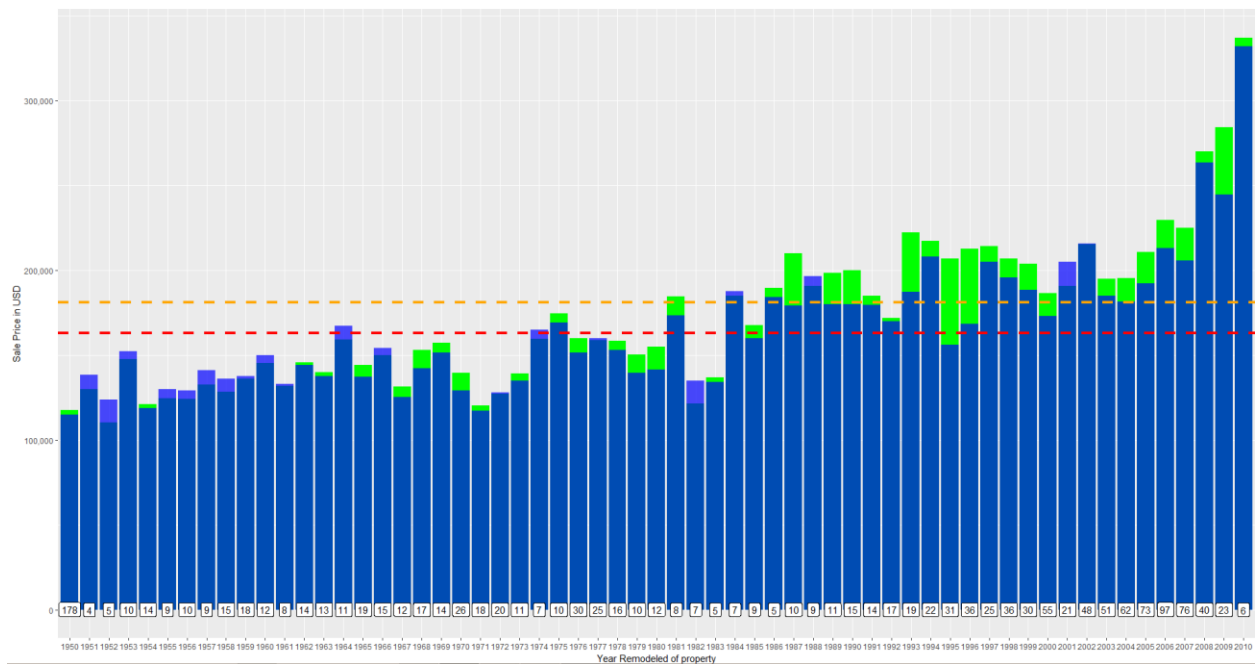


```
332
333  # Has the house been changed/remodeled?
334  Visualize_suspect_var(train.complete, "YearRemodAdd", "Year Remodeled of property")
335  full.final$IsRemod <- ifelse(full.final$YearBuilt != full.final$YearRemodAdd, "Yes", "No")
336  full.final$IsRemod <- as.factor(full.final$IsRemod)
337
```
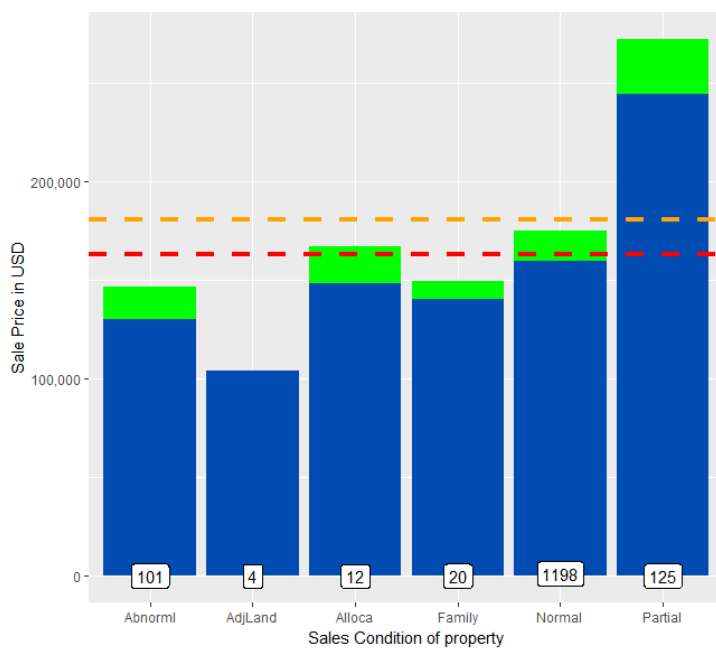
```
343
344    # Was the property bought under loan/mortgage?
345    Visualize_suspect_var(train.complete, "SaleCondition", "Sales Condition of property")
346    full.final$Invested <- ifelse(full.final$SaleCondition %in% c("Partial"), "Yes", "No")
347    full.final$Invested <- as.factor(full.final$Invested)
348
```

```
342
343    # Does the property have a fireplace?
344    full.final$HasFireplace <- ifelse(full.final$FireplaceQu %in% c("Not Available"), "No", "Yes")
345    full.final$HasFireplace <- as.factor(full.final$HasFireplace)
346
347    # Does the property have a basement?
348    full.final$HasBsmt <- ifelse(full.final$BsmtQual %in% c("Not Available"), "No", "Yes")
349    full.final$HasBsmt <- as.factor(full.final$HasBsmt)
350
351    # Is the property new?
352    full.final$IsNew <- ifelse(full.final$YearBuilt == full.final$YrSold, "Yes", "No")
353    full.final$IsNew <- as.factor(full.final$IsNew)
354
355    # Does the property have a second floor?
356    full.final$Has2ndFlr <- ifelse(full.final$X2ndFlrSF != 0, "Yes", "No")
357    full.final$Has2ndFlr <- as.factor(full.final$Has2ndFlr)
358
359    # Does the property have a Garage?
360    full.final$HasGarage <- ifelse(full.final$GarageArea != 0, "Yes", "No")
361    full.final$HasGarage <- as.factor(full.final$HasGarage)
362
363    # Does the property have a WoodDeck?
364    full.final$HasWD <- ifelse(full.final$WoodDeckSF != 0, "Yes", "No")   # Has Wood Deck?
365    full.final$HasWD <- as.factor(full.final$HasWD)
366
367    # Does the property have a Open Porch?
368    full.final$HasOporch <- ifelse(full.final$OpenPorchSF != 0, "Yes", "No")   # Has Open Porch?
369    full.final$HasOporch <- as.factor(full.final$HasOporch)
370
371    # Does the property have a Close Porch?
372    full.final$HasCporch <- ifelse(full.final$EnclosedPorch != 0, "Yes", "No")   # Has Closed Porch?
373    full.final$HasCporch <- as.factor(full.final$HasCporch)
374
375    # Does the property have a Screen Porch?
376    full.final$HasSporch <- ifelse(full.final$ScreenPorch != 0, "Yes", "No")   # Has Screen Porch?
377    full.final$HasSporch <- as.factor(full.final$HasSporch)
378
379    # Does the property feels like its protecting the owners?
380    Visualize_suspect_var(train.complete, "Fence", "Fence of property") # Feels Protected?
381    full.final$Protected <- ifelse(full.final$Fence %in% c("GdPrv", "Not Available"), "Yes", "No")
382    full.final$Protected <- as.factor(full.final$Protected)
383
```
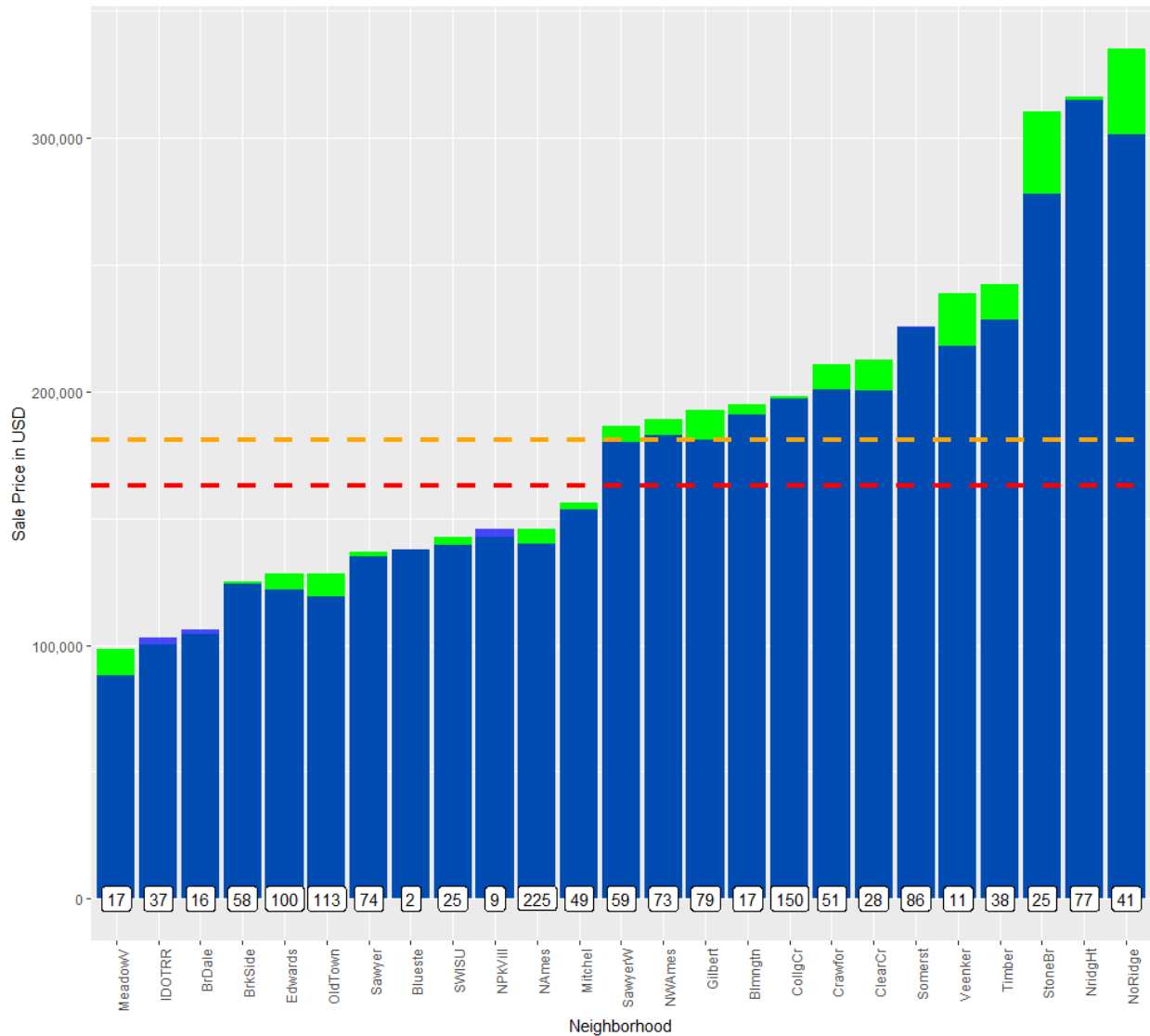
The Neighborhood feature can also be extracted and divided into four distinct categories: Bad, Fair, Good, Excellent. Neighborhood is binned to engineer the NeighborhoodQual feature:

```
387    ggplot(train.complete, aes(x = reorder(Neighborhood, SalePrice, FUN = mean), y = SalePrice)) +
388      geom_bar(stat = "summary", fun.y = "mean", fill = "green") +
389      geom_bar(stat = "summary", fun.y = "median", fill = "blue", alpha = 0.7) +
390      scale_y_continuous(labels = scales::comma) +
391      geom_hline(yintercept= median(train.complete$SalePrice), linetype="dashed", color = "red", size = 1.5) +
392      geom_hline(yintercept= mean(train.complete$SalePrice), linetype="dashed", color = "orange", size = 1.5) +
393      geom_label(stat = "count", aes(label = ..count.., y = ..count..)) +
394      theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 1)) +
395      ylab("Sale Price in USD") +
396      xlab("Neighborhood")
397
```

```
398
399   full.final$NeighborhoodQual <- ifelse(full.final$Neighborhood %in% c("MeadowV"), "Bad",
400                    ifelse(full.final$Neighborhood %in% c("OldTown", "Edwards", "BrkSide", "Sawyer", "Blueste", "SWISU", "NAmes", "NPkVill"),"Fair",
401                        ifelse(full.final$Neighborhood %in% c("Mitchel", "SawyerW", "Gilbert", "NWAmes",
402                                   "Blmngtn", "CollgCr"), "Good", "Excellent")))
403   full.final$NeighborhoodQual <- as.ordered(full.final$NeighborhoodQual) # Categorize neighborhood by cluster visually
404
```

The following features are also extracted using some of the survey's recurring questions. Contrary to its real value, the property price is often determined by subjective human generalizations. Therefore, we may be able to more accurately predict price if we incorporate some aggregate features such as total bathrooms, total living space and house age.
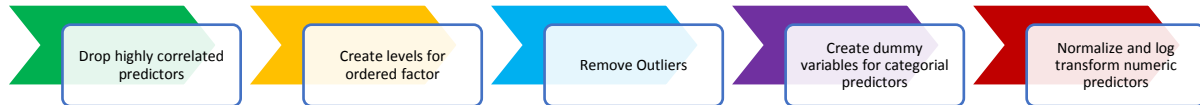
```
406
407   # Aggregating some features to better reflect human decision making mentality
408   full.final$TotBath <- full.final$FullBath + (full.final$HalfBath*0.5) +
409                      full.final$BsmtFullBath + (full.final$BsmtHalfBath*0.5) # Total number of bathrooms
410   full.final$TotSF <- full.final$GrLivArea + full.final$TotalBsmtSF # Total living space
411   full.final$HouseAge <- full.final$YrSold - full.final$YearBuilt # Total house Age
412
```

# Section 5: Data encoding for Machine Learning

This section displays the steps to encode the data for machine learning. Now that we have accomplished a reasonable tailoring of the data and derived the final dataset, we will have to convert the features into their machine-readable forms. These include:

| Drop highly correlated predictors | Create levels for ordered factor | Remove Outliers | Create dummy variables for categorial predictors | Normalize and log transform numeric predictors |

- Dropping highly correlated variables: Correlated independent variables can cause multi-collinearity and overfitting in the model, thereby reducing the interpretability and generalizability power of the model (reducing the precision of a predictor's impact on the independent variable).

```
418
419   # Drop variables with lack of information/ duplicate information
420
421   Visualize_suspect_var(train.complete, "Utilities", "Utilities of property") # Very few observation to matter
422
423   cor.high <- cor
424   cor.high[lower.tri(cor.high ,diag = TRUE)] <- NA
425   cor.high <- subset(na.omit(data.frame(expand.grid(dimnames(cor.high)), value = c(cor.high))), value > .7)
426                         # Subset to find pairs with high correlation
427                         # cor.high
428
429   var.drop <- c("YearRemodAdd", "Utilities", "GarageYrBlt", "GarageArea", "GarageCond", "TotalBsmtSF",
430                 "TotalRmsAbvGrd", "BsmtFinSF1")
431   full.final <- full.final[,!(names(full.final) %in% var.drop)]   # Drop these variables
432
```

- Create levels for ordered factors: We create levels for the algorithm to use as dummy variables (with 0 or 1 values for each observation). This is to get around the fact that the regression can only understand numbers and not string categories. Ordered factors will also give the Machine Learning algorithm more information to work with than ordinary factors.

```
432
433   # Convert the remaining into suitable variable type and add correct levels to ordered
434
435   ordered_var <- ordered_var[!ordered_var %in% c("Utilities", "GarageCond", "GarageYrBlt")]
436   full.final <- full.final %>%
437     mutate_at(factor_var, funs(factor(.))) %>%
438     mutate_at(ordered_var, funs(ordered(.)))
439
440   str(full.final[,ordered_var])
441
442   full.final$LotShape <- ordered(full.final$LotShape, levels = c("IR3", "IR2", "IR1", "Reg"))
443   full.final$LandSlope <- ordered(full.final$LandSlope, levels = c("Sev", "Mod", "Gtl"))
444   full.final$ExterQual <- ordered(full.final$ExterQual, levels = c("Po", "Fa", "TA", "Gd", "Ex"))
445   full.final$ExterCond <- ordered(full.final$ExterCond, levels = c("Po", "Fa", "TA", "Gd", "Ex"))
446   full.final$BsmtQual <- ordered(full.final$BsmtQual, levels = c("Not Available", "Fa", "TA", "Gd", "Ex"))
447   full.final$BsmtCond <- ordered(full.final$BsmtCond, levels = c("Not Available","Po", "Fa", "TA", "Gd"))
448   full.final$BsmtExposure <- ordered(full.final$BsmtExposure, levels = c("Not Available","No", "Mn", "Av", "Gd"))
449   full.final$BsmtFinType1 <- ordered(full.final$BsmtFinType1, levels = c("Not Available","Unf", "LwQ", "Rec", "BLQ", "ALQ", "GLQ"))
450   full.final$BsmtFinType2 <- ordered(full.final$BsmtFinType2, levels = c("Not Available", "Unf", "LwQ", "Rec", "BLQ", "ALQ", "GLQ"))
451   full.final$HeatingQC <- ordered(full.final$HeatingQC, levels = c("Po", "Fa", "TA", "Gd", "Ex"))
452   full.final$Electrical <- ordered(full.final$Electrical, levels = c("Mix", "FuseP", "FuseF", "FuseA", "SBrkr"))
453   full.final$KitchenQual <- ordered(full.final$KitchenQual, levels = c("Fa", "TA", "Gd", "Ex"))
454   full.final$Functional <- ordered(full.final$Functional, levels = c("Sal", "Sev", "Maj1", "Maj2", "Mod", "Min2", "Min1", "Typ"))
455   full.final$FireplaceQu <- ordered(full.final$FireplaceQu, levels = c("Not Available", "Po","Fa", "TA", "Gd", "Ex"))
456   full.final$GarageFinish <- ordered(full.final$GarageFinish, levels = c("Not Available", "Unf","RFn", "Fin"))
457   full.final$GarageQual <- ordered(full.final$GarageQual, levels = c("Not Available", "Po","Fa", "TA", "Gd", "Ex"))
458   full.final$PoolQC <- ordered(full.final$PoolQC, levels = c("Not Available","Fa", "Gd", "Ex"))
459   full.final$TotRmsAbvGrd <- ordered(full.final$TotRmsAbvGrd, levels = c("2","3", "4", "5", "6", "7", "8", "9", "10", "11", "12", "13", "14", "15"))
460   full.final$Fireplaces <- ordered(full.final$Fireplaces, levels = c("0","1", "2", "3", "4"))
461
```

- Remove outliers: The outliers identified previously may cause lower the model's accuracy and its ability to generalize outside the test set.

```
462
463   # Removing outliers
464   full.final$SalePrice <- full$SalePrice
465   full.final <- full.final[-c(524,1299),]
466
```

- Create dummy variables: The dummy variables will greatly increase our number of predictors. Therefore, we will also remove some dummy variable with less than 10 positive observations.

```
466
467  ## Prepocess numeric variables
468
469  nums <- unlist(lapply(full.final, is.numeric))
470  full.numeric <- full.final[,nums]
471  full.numeric <- full.numeric[,!(names(full.numeric) %in% "SalePrice")]
472  full.factor <- full.final[,!nums]
473  cat('There are', length(full.numeric), 'numeric variables, and', length(full.factor), 'factor variables')
474
475  # Create dummy variables
476
477  full.dummy <- as.data.frame(model.matrix(~.-1, full.factor))
478  dim(full.dummy)
479
480  # Taking out a few categorical/dummy variables with very little information
481  var.lowinfo <- which(colSums(full.dummy[1:nrow(full.final[!is.na(full.final$SalePrice),]),])<10)
482  colnames(full.dummy[var.lowinfo])
483
484  full.dummy <- full.dummy[,-var.lowinfo] #removing variables
485  dim(full.dummy)
486
```
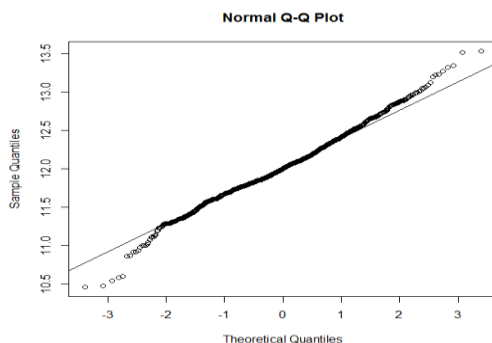
- Normalization & Log transformation: We observe skewness in some independent variable and the response variable. Therefore, we will log transform these variables. We will also normalize the numeric predictors to center and scale them for the Machine Learning algorithm.

```
474
475  # Log transformation of numeric variables
476  describe(full.numeric)
477
478  for(i in 1:ncol(full.numeric)){
479    if (abs(skew(full.numeric[,i]))>0.8){
480      full.numeric[,i] <- log(full.numeric[,i] +1)
481    }
482  }
483
484  # Normalizing the numeric variables
485  num.process <- preProcess(full.numeric, method=c("center", "scale"))
486  print(num.process)
487  full.numeric.scaled <- predict(num.process, full.numeric)
488
```

```
504
505  # Dealing with skewness in response variable
506  skew(full.final$SalePrice)
507  full.encode$SalePrice <-  log(full.final$SalePrice)
508  qqnorm(full.encode$SalePrice)
509  qqline(full.encode$SalePrice)     # Log transform and visualize SalePrice
510
```



Normal Q-Q Plot

Finally, we combine the two subsets to an encoded dataset. We will use this to derive the encoded training set and the encoded test set.

```
500
501  # Combining the predictor variables
502  full.encode <- cbind(full.numeric.scaled, full.dummy)
503
```

```
510
511  # Divide encoded data into train and test set
512  train.encode <-  full.encode[1:1458,]
513  test.encode <- full.encode[1459:2917,]
514
```

# Section 6: Machine Learning Model Selection and Training

This section attempts to train a suitable machine learning algorithm to predict the sale price of future properties. The problem here is a prediction of continuous value for the response variable. Thus, we can begin with conducting a simple regression. However, with only a few predictors showing strong predictive power out of over 200 features, we will want to introduce a regularization parameter to penalize the number and/or effect of features used by our regression. Thus, we will use an elastic net via cross-validation to best tune our regularization parameters (alpha as degree of mix between LASSO and Ridge regression and lambda as the shrinkage parameter of the penalty).

```
518
519  set.seed(123)
520  mod.control <-trainControl(method="cv", number=5)
521  grid <- expand.grid(alpha = 1, lambda = seq(0.001,0.1,by = 0.001))
522
523  regressor <- train(x=train.encode[,!(names(train.encode) %in% "SalePrice")], y=train.encode$SalePrice,
524                method='glmnet', trControl= mod.control, tuneGrid= grid)
525  regressor$bestTune
526
```

```
   alpha lambda
2      1  0.002
```

Our hyperparameter tuning shows best results for LASSO regression with a shrinkage parameter of 0.002. This will also result in a RMSE score of ~ **11.4%:**

```
> min(regressor$results$RMSE)
[1] 0.1138914
```

Let's also check if Lasso Regression did its job by cutting down on many of the variables.

```
529
530  # Check the number of variables being selected for the best model.
531  contrib <- varImp(regressor,scale=F)$importance
532  var.selected <- length(which(contrib$Overall!=0))
533  var.nselected <- length(which(contrib$Overall==0))
534
535  cat('L-Regression uses', var.selected, 'variables in its model, and did not select', var.nselected, 'variables.')
536
```

```
L-Regression uses 110 variables in its model, and did not select 128 variables.
```

It has removed 128 variables and used 110. The regularized regression seems to be using only independent variables with statistical significant. It therefore looks to be a viable model.

# Section 7: Conclusion and Risk assessment

Throughout this report, we have prepared, cleaned and explored the dataset collected. Then, we engineered new features with domain understanding to help our predictive model perform better and chose a suitable machine learning model for the prediction task. Our model is likely to score 11.4% precision with any test set and has been rigorously devised to generalize well in reality.

However, improvements can certainly be investigated. The following are risks and improvements that users of this report can make to improve the model's performance:

- First, the user should care to use a training dataset that accurately reflects the geographic, social and demographic trends relevant to real estate transactions in the setting of the consulting property.
- Second, the user should explore more algorithm options. Although simple linear regression models are usually best when predicting a generally linear problem such as this, results could be enhanced with incorporation of non-linear terms. If accuracy over speed is the user's primary concern, he/she should explore using neural networks, gradient boosting trees or other combination and ensemble of models to derive predictions.
- Lastly, the user should only consider the results of the algorithm as a referred generalized price point. Some rare properties and real-estates can have characteristics that make them extra-ordinary and will most likely not be accounted for in the model.