

Replication Presentation

Applied Stats II

Due: April 11, 2022

Incorporating Multiple Linear Regression in Predicting the House Prices Using a Big Real Estate Dataset with 80 Independent Variables

The study that was selected for this replication assignment did not have an available R file, so I wrote the code based on the methodology in the paper. This document will go through the steps I took and which parts of the study I altered.

1. Preprocessing

Dealing with NA's

Firstly, I looked at the percent of NA's, the following code showed that 5.8% of the dataset were NA's.

```
1 # calculating the product of dimensions of dataframe
2 totalcells = prod(dim(data))
3 # calculating the number of cells with na
4 missingcells = sum(is.na(data))
5 # calculating percentage of missing values
6 percentage = (missingcells * 100)/(totalcells)
7 print (percentage)
```

The first issue I noticed with the study was that they explained that they replaced all the NA's with the mean of the column, however it seems as though they were treating all the variables as numerical. To get rid of all the NA's I firstly removed the columns with over 1000 NA's, like the study and then in the numerical columns, replaced the NA's with the mean and in the categorical columns, the NA with the mode. The following code was used:

```
1 missing_values <- lapply(data,function(x) { length(which(is.na(x))) })
2 missing_values
3 #$Pool.QC
4 #[1] 2917
5 #$Misc.Feature
6 #[1] 2824
7 #$Fence
8 #[1] 2358
```

```

9  $$Fireplace.Qu
10 #[1] 1422
11 $$Alley
12 #[1] 2732
13
14 #Remove these features from the dataset
15 drop <- c(" Pool.QC", " Misc.Feature", " Fence", " Fireplace.Qu", " Alley")
16 data = data[,!(names(data) %in% drop)]
17
18 #####
19 #to replace NA's in data with means of the column.
20 #####
21
22 #The Mean for Numeric Data
23 for(i in 1:ncol(data)){
24   data[is.na(data[,i]), i] <- mean(data[,i], na.rm = TRUE)
25 }
26
27 #Below is for Categorical data
28
29 calc_mode <- function(x){
30   distinct_values <- unique(x)
31   distinct_tabulate <- tabulate(match(x, distinct_values))
32   distinct_values[which.max(distinct_tabulate)]
33 }
34
35 #Need to put in the mode for the categorical data
36
37 data$Garage.Qual[is.na(data$Garage.Qual)]<-calc_mode(data$Garage.Qual)
38 data$Garage.Cond[is.na(data$Garage.Cond)]<-calc_mode(data$Garage.Cond)
39 data$Garage.Finish[is.na(data$Garage.Finish)]<-calc_mode(data$Garage.Finish)
40 data$Garage.Type[is.na(data$Garage.Type)]<-calc_mode(data$Garage.Type)
41 data$BsmtFin.Type.2[is.na(data$BsmtFin.Type.2)]<-calc_mode(data$BsmtFin.Type
.2)
42 data$BsmtFin.Type.1[is.na(data$BsmtFin.Type.1)]<-calc_mode(data$BsmtFin.Type
.1)
43 data$Bsmt.Exposure[is.na(data$Bsmt.Exposure)]<-calc_mode(data$Bsmt.Exposure)
44 data$Bsmt.Cond[is.na(data$Bsmt.Cond)]<-calc_mode(data$Bsmt.Cond)
45 data$Bsmt.Qual[is.na(data$Bsmt.Qual)]<-calc_mode(data$Bsmt.Qual)
46
47 #Now there are no NA's in the data
48 sum(is.na(data))

```

Now that there were no NA's in the data I conducted exploratory data analysis, as was done in the study.

2. Descriptive Visualisations

Firstly examining the distribution of the outcome variable "SalePrice"

```
1 ggplot(data = data, aes(x=SalePrice)) + geom_histogram(color="blue", fill= "
  blue", bins = 70)
```

This produced the following plot:

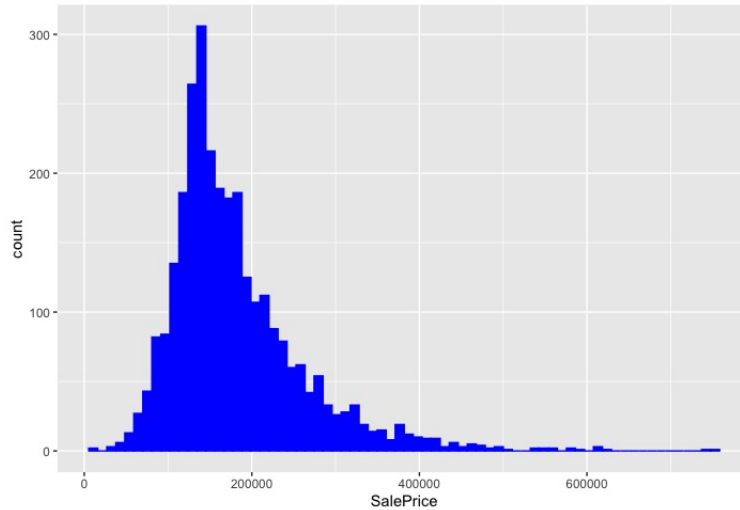


Figure 1: Histogram of the Outcome Variable: SalePrice

The skew to the left is visible and so I followed the steps in the study and log transformed the variable.

```
1 data$SalePrice <- log(data$SalePrice)
2 ggplot(data = data, aes(x=SalePrice)) + geom_histogram(color="blue", fill= "
  blue", bins = 70)
```

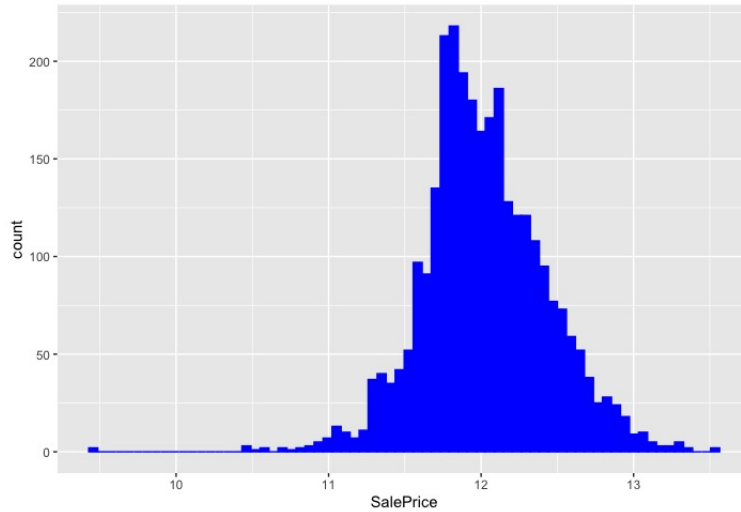


Figure 2: Histogram of the Outcome Variable: Log(SalePrice)

The next visualisation is for the Neighborhood category. I recreated the boxplot from the study below:

```
1 data %>%
2   mutate(Neighborhood = fct_reorder(Neighborhood, SalePrice)) %>%
3   ggplot( aes(x=Neighborhood, y=SalePrice, fill=Neighborhood)) +
4   geom_boxplot() +
5   xlab("") +
6   theme_bw()
```

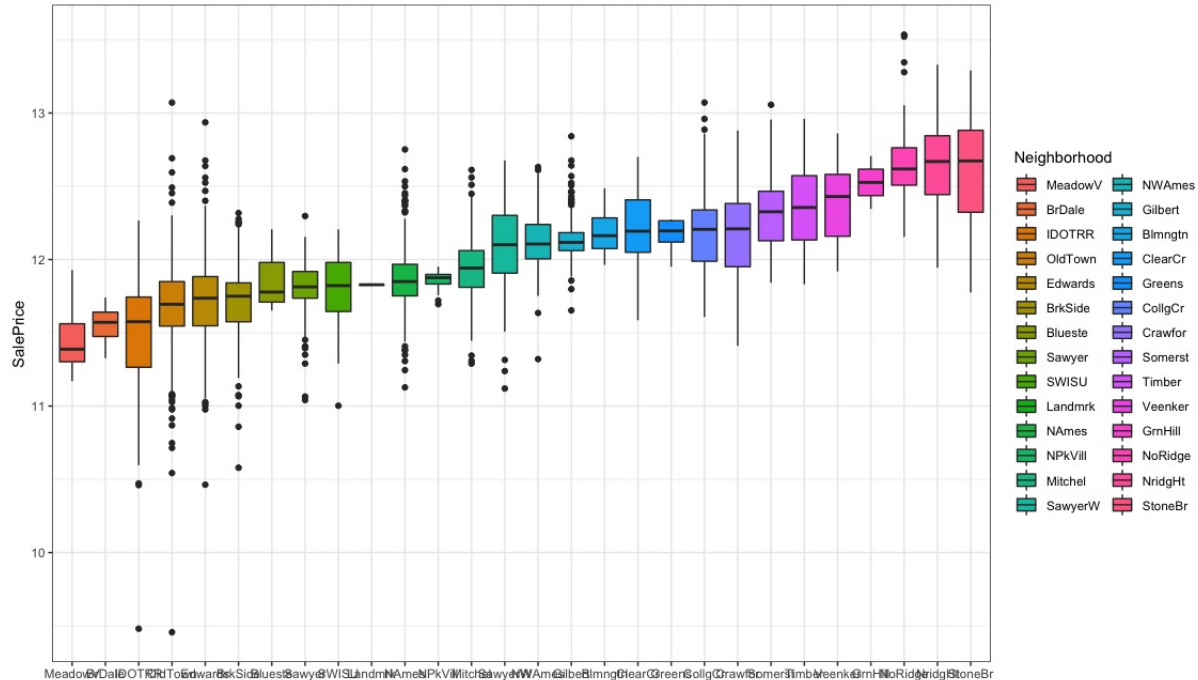


Figure 3: Boxplot of the Variable: Neighborhood

Next the correlations were looked at. The study emphasized the use of understanding if variables were highly correlated, the author stated:

"Multicollinearity assumption: linear regression model assumes that there is little or no multicollinearity in the data. Multicollinearity occurs when the independent variables are highly correlated with each other. We have checked the correlation between all predictors. Some predictors were found to be highly correlated with each other and with the response variable SalePrice. We have identified these highly correlated predictors and removed the predictors with more than 50% correlation from the dataset in Step 10. So, this assumption holds."

Therefore it was very important to replicate what they did. The following code was used to produced a correlation plot, showing correlations that were above 0.5. (As there were 80 variables in the dataset, I set this limit so that it could be legible)

```
1 library(corrplot)
2 corr_simple <- function(data=data, sig=0.5){A
3   df_cor <- data %>% mutate_if(is.character, as.factor)
4   df_cor <- df_cor %>% mutate_if(is.factor, as.numeric)
5   corr <- cor(df_cor)
6   corr[lower.tri(corr, diag=TRUE)] <- NA
7   corr[corr == 1] <- NA
8   corr <- as.data.frame(as.table(corr))
9   corr <- na.omit(corr)
10  corr <- subset(corr, abs(Freq) > sig)
```

```

11 corr <- corr[order(-abs(corr$Freq)),]
12 print(corr)
13 mtx_corr <- reshape2::acast(corr, Var1~Var2, value.var="Freq")
14
15 corplot(mtx_corr, is.corr=FALSE, tl.col="black", na.label=" ")
16 }
17 corr_simple(data)

```

This produced the following plot, which confirms the high correlations which were highlighted in the study:

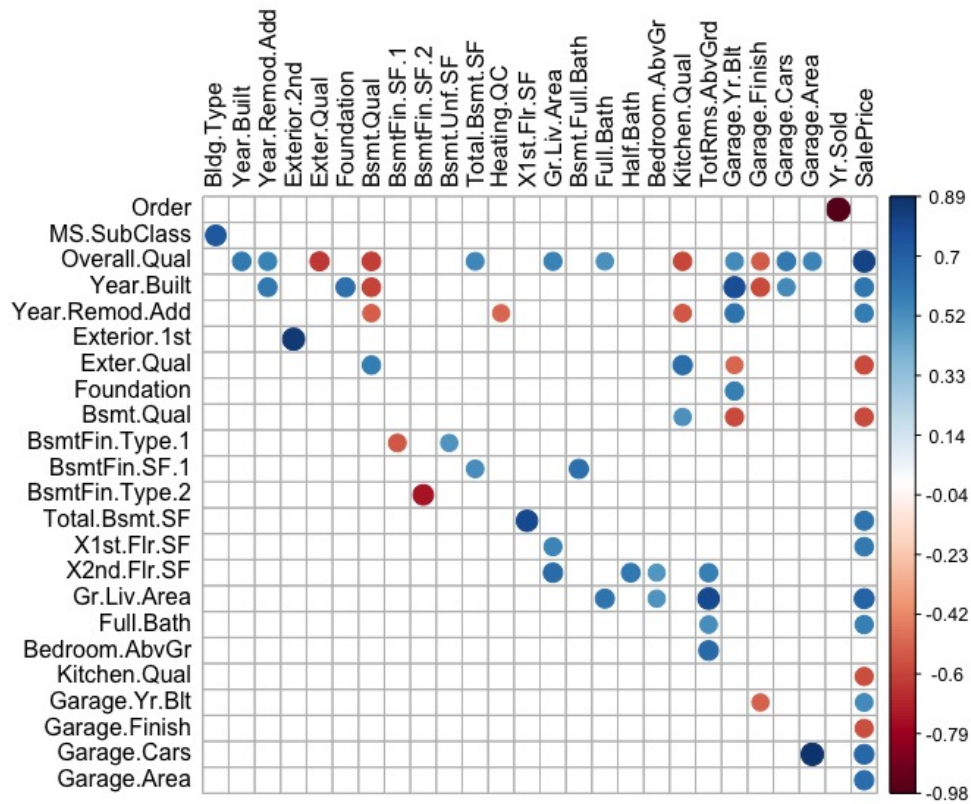


Figure 4: Plot of Correlations over 0.5

Next the features that had near zero variance were removed:

```

1 library(caret)
2 names(data)[nearZeroVar(data)]

```

Checking Linearity

```

1 ggplot(data, aes(x=Overall.Qual, y=SalePrice)) +
2   geom_point()+
3   geom_smooth(method=lm)
4
5 ggplot(data, aes(x=Total.Bsmt.SF, y=SalePrice)) +
6   geom_point()+
7   geom_smooth(method=lm)

```

The following plots show that the data fits the linearity assumption

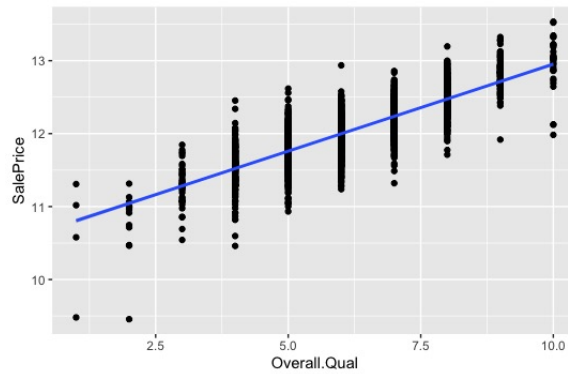


Figure 5: Linearity Overall Quality

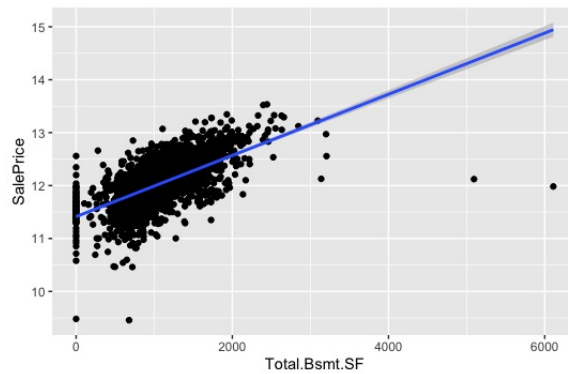


Figure 6: Linearity BSMT SF

3. Random Forest for Feature Selection

By looking through the study, it seemed as though, despite the previous preprocessing steps, the random forest was ran on the original data. I therefore ran the random forest on the original data to explore whether this was the case.

```
1 library(randomForest)
2
3 set.seed(100)
4 train <- sample(nrow(data), 0.5*nrow(data), replace = FALSE)
5 TrainSet <- data[train,]
6 ValidSet <- data[-train,]
7 summary(TrainSet)
8 summary(ValidSet)
9
10 RF_model <- randomForest(SalePrice ~ ., data = TrainSet, importance = TRUE)
11 RF_model
12
13 importance(RF_model)
14 varImpPlot(RF_model)
```

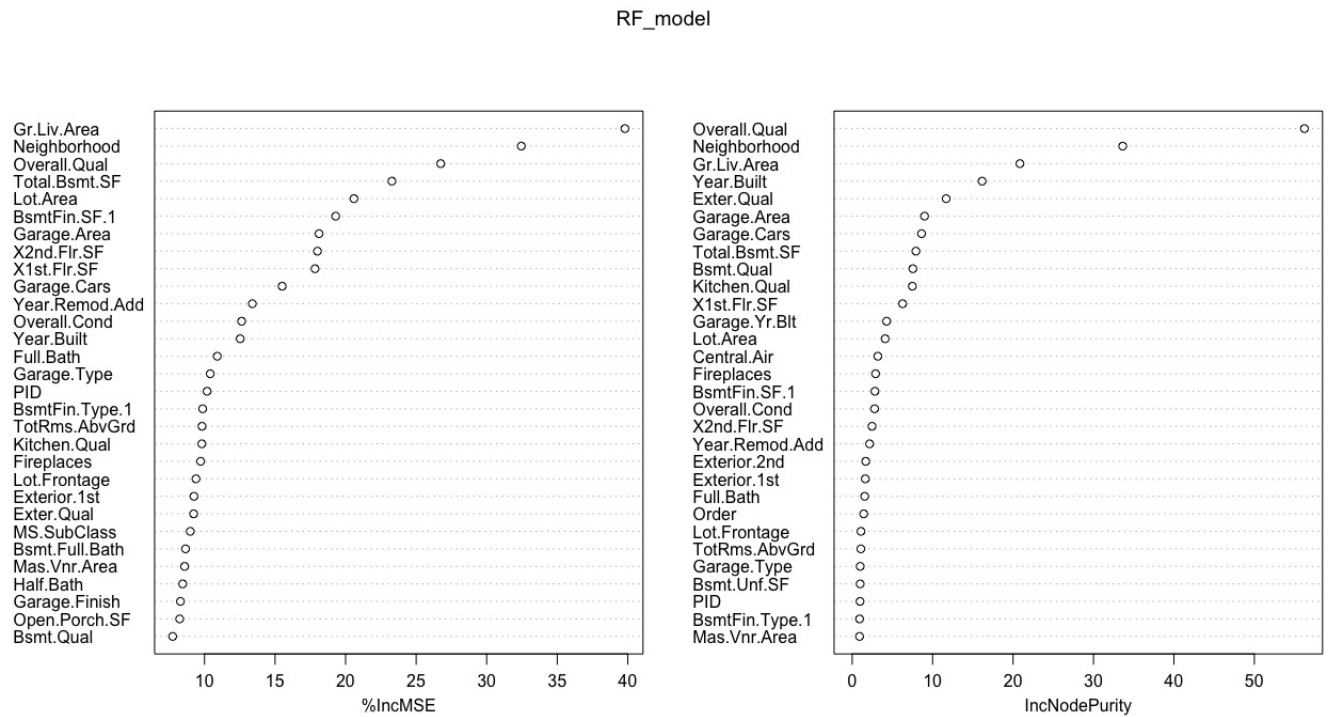


Figure 7: Random Forest on the Original Data

The above plot shows that the most important features that were selected for the regression were from the random forest ran on the original data, not on the data with out highly correlated variables or near zero variance.

I then ran the random forest on the data which did take into account the preprocessing. Firstly I removed these items from the dataset:

```
1 drop <- c("Overall.Qual",
2          "Gr.Liv.Area",
3          "Garage.Cars",
4          "Garage.Area",
5          "Total.Bsmt.Sf",
6          "X1st.Flr.SF",
7          "Full.Bath",
8          "Tot.Rms.Abv.Gr",
9          "Year.Built",
10         "Year.Remod.Add",
11         "BsmtFin.SF.1",
12         "X1st.Flr.SF",
13         "Garage.Cars",
14         "X2nd.Flr.SF",
15         "Bsmt.Unf.SF",
16         "Central.Air",
17         "Bedroom.AbvGr",
18         "Garage.Finish")
19
20 data = data[,!(names(data) %in% drop)]
```

I then re-ran the Random forest

```
1 set.seed(100)
2 train <- sample(nrow(data), 0.5*nrow(data), replace = FALSE)
3 TrainSet <- data[train,]
4 ValidSet <- data[-train,]
5 summary(TrainSet)
6 summary(ValidSet)
7
8 RF_model2 <- randomForest(SalePrice ~ ., data = TrainSet, importance = TRUE)
9 RF_model2
10
11 importance(RF_model2)
12 varImpPlot(RF_model2)
```

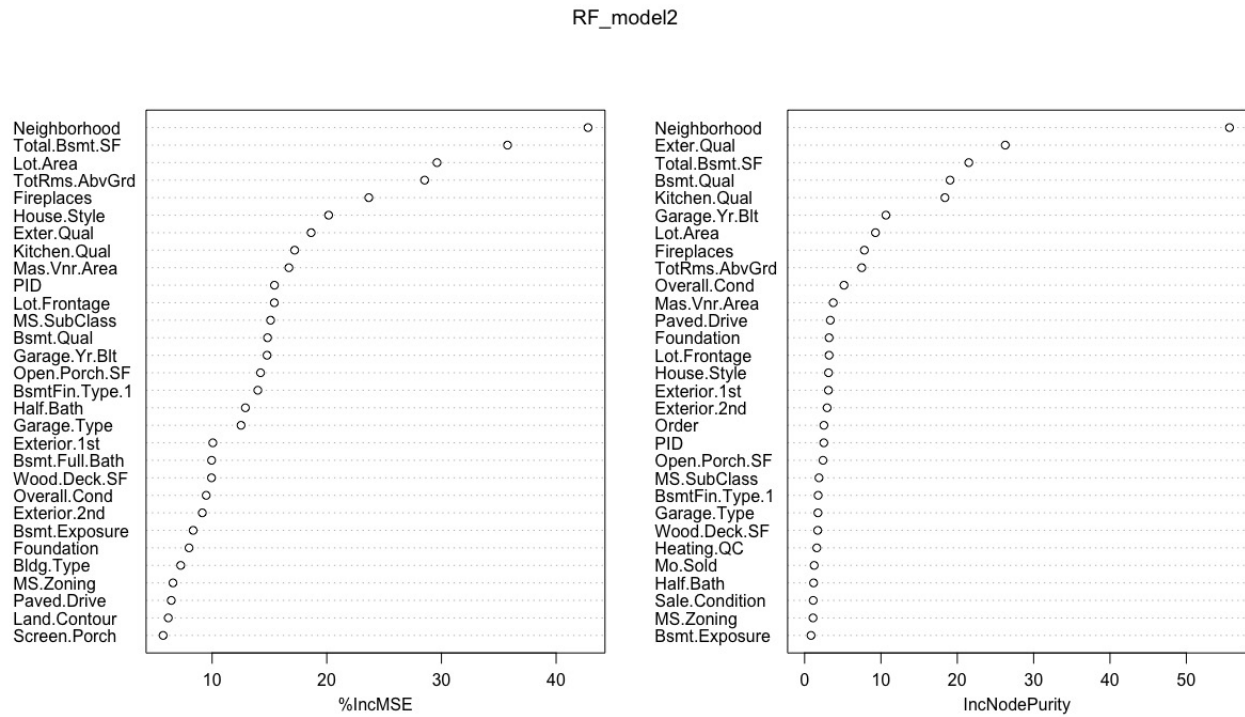


Figure 8: Random Forest on the Processed Data

From here I then ran multiple linear regression on the top 12 most important variable from the random forest.

4. Multiple Linear Regression

```
1 data <- data[, c("SalePrice",
2                 "TotRms.AbvGrd",
3                 "Neighborhood",
4                 "Fireplaces",
5                 "MS.SubClass",
6                 "Total.Bsmt.SF",
7                 "Exter.Qual",
8                 "Lot.Area",
9                 "Kitchen.Qual",
10                "House.Style",
11                "Mas.Vnr.Area",
12                "PID",
13                "Lot.Frontage"
14                )]
15 model1 <- lm(SalePrice ~., data = data)
16 summary(model1)
```

The output this created was quite messy and difficult to interpret as the neighborhood variable has 28 categories.

For this reason, I grouped the Neighborhood variable by the residuals in the model using the following code used in a tutorial in Term 1:

```
1 zip_groups <- data %>% # our dplyr code for creating a 3-level grouping for
  zipcode
2 mutate(resid = resid(model1)) %>%
3 group_by(Neighborhood) %>%
4 summarise(median = median(resid),
5           n = n()) %>%
6 arrange(desc(median)) %>%
7 mutate(sum_n = cumsum(n),
8        ZipGroup = ntile(sum_n, 3))
9
10 data <- data %>% # joining the zipcode groups to our original dataset
11 left_join(select(zip_groups, Neighborhood, ZipGroup), by = "Neighborhood")
```

I then replaced the Neighborhood variable with the ZipGroup, and reran the regression.

```

Call:
lm(formula = SalePrice ~ ., data = data)

Residuals:
    Min       1Q   Median       3Q      Max
-2.54405 -0.09018  0.00857  0.10951  0.67071

Coefficients:
              Estimate      Std. Error t value      Pr(>|t|)
(Intercept)  11.46891092023249  0.03957645533157  289.791 < 0.0000000000000002 ***
TotRms.AbvGrd  0.03570455286286  0.00326768904132   10.927 < 0.0000000000000002 ***
ZipGroup      0.06469445964678  0.00541263079103   11.952 < 0.0000000000000002 ***
Fireplaces    0.10674730024137  0.00634259491387   16.830 < 0.0000000000000002 ***
MS.SubClass   -0.00065457591449  0.00010957061106   -5.974  0.000000002596 ***
Total.Bsmt.SF  0.00026160290290  0.00001234663097   21.188 < 0.0000000000000002 ***
Exter.QualGd   0.15884617378735  0.01171254894334   13.562 < 0.0000000000000002 ***
Exter.QualFa  -0.30310047203535  0.03464273914855   -8.749 < 0.0000000000000002 ***
Exter.QualEx   0.24260509636516  0.02731505513048    8.882 < 0.0000000000000002 ***
Lot.Area       0.00000330842855  0.00000051746077    6.394  0.00000000188 ***
Kitchen.QualFa -0.40859111681425  0.03169969633288  -12.889 < 0.0000000000000002 ***
Kitchen.QualGd -0.11672330278547  0.01881595278181   -6.203  0.00000000631 ***
Kitchen.QualPo -0.36086984624860  0.19852550030252   -1.818  0.069205 .
Kitchen.QualTA -0.26032415460133  0.02057452771589  -12.653 < 0.0000000000000002 ***
House.Style1.SUnf -0.09130022787161  0.04701941042979   -1.942  0.052263 .
House.Style1Story -0.04088058371125  0.01388971690133   -2.943  0.003274 **
House.Style2.SFin  0.13519028135094  0.07204577702940    1.876  0.060694 .
House.Style2.SUnf  0.08966753406225  0.04233528409835    2.118  0.034257 *
House.Style2Story  0.09341866289868  0.01439840484588    6.488  0.00000000102 ***
House.StyleSFoyer  0.09629439979117  0.02539399522751    3.792  0.000152 ***
House.StyleSLvl   0.11645061570917  0.02142059471535    5.436  0.00000058879 ***
Mas.Vnr.Area     0.00009623511123  0.00002411772492    3.990  0.000067641897 ***
PID              -0.00000000009098  0.00000000002103   -4.326  0.000015718891 ***
Lot.Frontage     0.00046065923465  0.00021430784633    2.150  0.031675 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1971 on 2906 degrees of freedom
Multiple R-squared:  0.768,    Adjusted R-squared:  0.7662
F-statistic: 418.2 on 23 and 2906 DF,  p-value: < 0.0000000000000022

```

Figure 9: Regression, with Zipped Groups

I then used the `avPlots()` function from the "car" package in R to view some of the variables in the regression.

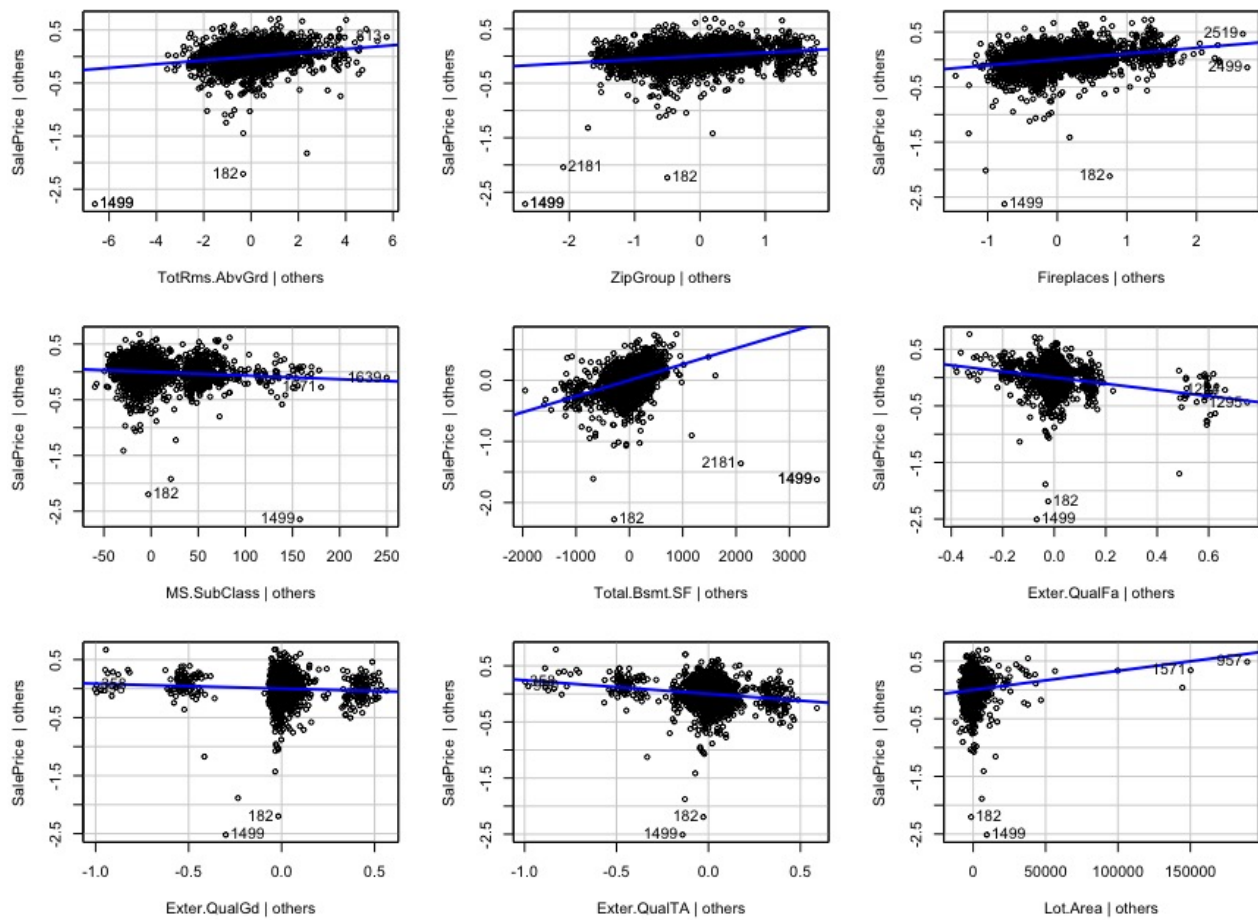


Figure 10: Regression Plot for some of the variables

4. Conclusions

The multivariate linear regression above explains less of the data than was described in the study, however this follows the correct preprocessing steps and is therefore representing a more accurate model.