

Read Data into R:

```
Ames <- read.csv("AmesHousing.csv")
```

Next I do a few preliminary checks to see the specifics of this data set:

Look at variable Names:

```
> names(Ames)
[1] "MS.Zoning"      "Lot.Frontage"    "Lot.Area"       "Street"        "Alley"         "Lot.Shape"
[7] "Land.Contour"   "Lot.Config"     "Land.Slope"     "Condition.1"   "Condition.2"   "Bldg.Type"
[13] "House.Style"    "Overall.Qual"  "Overall.Cond"  "Year.Built"    "Year.Remod.Add" "Roof.Style"
[19] "Roof.Matl"      "Exterior.1st"  "Exterior.2nd"  "Mas.Vnr.Type"  "Mas.Vnr.Area"  "Exter.Qual"
[25] "Exter.Cond"     "Foundation"    "Bsmt.Qual"    "Bsmt.Cond"    "Bsmt.Exposure" "BsmtFin.Type.1"
[31] "BsmtFin.SF.1"   "BsmtFin.Type.2" "BsmtFin.SF.2"  "Bsmt.Unf.SF"   "Total.Bsmt.SF" "Heating"
[37] "Heating.QC"     "Central.Air"   "Electrical"    "X1st.Flr.SF"  "X2nd.Flr.SF"  "Low.Qual.Fin.SF"
[43] "Gr.Liv.Area"   "Bsmt.Full.Bath" "Bsmt.Half.Bath" "Full.Bath"    "Half.Bath"    "Bedroom.AbvGr"
[49] "Kitchen.AbvGr"  "Kitchen.Qual"  "TotRms.AbvGrd" "Functional"   "Fireplaces"   "Fireplace.Qu"
[55] "Garage.Type"    "Garage.Yr.Blt"  "Garage.Finish" "Garage.Cars"   "Garage.Area"   "Garage.Qual"
[61] "GarageCond"     "Paved.Drive"   "Wood.Deck.SF"  "Open.Porch.SF" "Enclosed.Porch" "X3Ssn.Porch"
[67] "Screen.Porch"   "Pool.Area"    "Fence"        "Misc.Val"     "Yr.Sold"      "Sale.Type"
[73] "Sale.Condition" "SalePrice"
```

Looking at the Variable Names. We can get rid of Garage Yr Blt due to multicollinearity. Meaning this variable is redundant since we also have the Yr The House

Look at the structure of the data set: I normally do this when first looking at a new data set just so I can see how many observations I have and number of variables. I also like to see what variable type each variable is. As we can see from the print out, we have many categorical variables as well as continuous variables.

```
> # Structure of Data:
> str(Ames)
'data.frame': 2952 obs. of 74 variables:
 $ MS.Zoning : Factor w/ 5 levels "C (all)", "FV", ...: 4 3 4 4 4 4 4 4 4 ...
 $ Lot.Frontage : int 141 88 81 93 74 78 41 43 39 60 ...
 $ Lot.Area : int 31770 11622 14267 11168 13839 9762 4926 5085 5389 7506 ...
 $ Street : Factor w/ 2 levels "Grvl", "Pave": 2 2 2 2 2 2 2 2 2 2 ...
 $ Alley : Factor w/ 3 levels "Grvl", "None", ...: 2 2 2 2 2 2 2 2 2 2 ...
 $ Lot.Shape : Factor w/ 4 levels "IR1", "IR2", "IR3", ...: 1 4 1 4 1 1 4 1 4 ...
 $ Land.Contour : Factor w/ 4 levels "Bnk", "HLW5", "Low", ...: 4 4 4 4 4 4 4 4 4 ...
 $ Lot.Config : Factor w/ 5 levels "Corner", "ColSue", ...: 1 5 1 5 5 5 5 5 5 ...
 $ Land.Slope : Factor w/ 3 levels "Gtl", "Mod", "Sev": 1 1 1 1 1 1 1 1 1 ...
 $ Condition.1 : Factor w/ 7 levels "Artery", "Feed", ...: 3 2 3 3 3 3 3 3 3 ...
 $ Condition.2 : Factor w/ 3 levels "Neg", "Norm", "Pos": 2 2 2 2 2 2 2 2 2 ...
 $ Bldg.Type : Factor w/ 5 levels "1Fam", "2fCon", ...: 1 1 1 1 1 5 5 5 1 ...
 $ House.Age : Factor w/ 7 levels "10-20", "21-30", ...: 3 3 3 3 3 5 5 3 3 5 ...
 $ Overall.Qual : Factor w/ 6 levels "5", "6", "7", "8", "9", ...: 5 6 5 5 6 5 5 5 ...
 $ Overall.Cond : int 5 6 5 5 6 5 5 5 ...
 $ Year.Built : int 1968 1961 1958 1964 1997 1998 2001 1992 1995 1999 ...
 $ Year.Remod.Add : int 1968 1961 1958 1964 1998 1998 2001 1992 1996 1999 ...
 $ Roof.Style : Factor w/ 5 levels "Flat", "Gable", ...: 2 2 4 4 2 2 2 2 2 ...
 $ Roof.Matl : Factor w/ 3 levels "CompShg", "Other", ...: 1 1 1 1 1 1 1 1 1 ...
 $ Exterior.1st : Factor w/ 9 levels "AsbShng", "BrkFace", ...: 2 8 9 2 8 8 3 4 3 8 ...
 $ Exterior.2nd : Factor w/ 10 levels "AsbShng", "BrkFace", ...: 9 10 2 9 9 3 4 3 9 ...
 $ Bsmt.Qual : Factor w/ 5 levels "Ex", "Gd", "Av", "Po", "Fa": 4 3 2 3 2 3 3 3 3 ...
 $ Bsmt.Exposure : Factor w/ 3 levels "N", "No", "O": 112 0 108 0 9 20 0 8 0 ...
 $ Exter.Qual : Factor w/ 4 levels "Ex", "Fa", "Gd", ...: 4 4 4 3 4 4 3 3 3 4 ...
 $ Exter.Cond : Factor w/ 5 levels "Ex", "Fa", "Gd", ...: 5 5 5 5 5 5 5 5 5 ...
 $ Foundation : Factor w/ 5 levels "BrkTil", "Clock", ...: 2 2 2 2 4 4 4 4 4 ...
 $ Bsmt.Qual : Factor w/ 5 levels "Ex", "Fa", "Gd", ...: 5 5 5 5 3 3 3 3 5 ...
 $ Bsmt.Cond : Factor w/ 4 levels "Fa", "Gd", "None", ...: 2 4 4 4 4 4 4 4 ...
 $ Bsmt.Exposure : Factor w/ 5 levels "Av", "Gd", "Mn", ...: 2 4 4 4 4 3 4 4 4 ...
 $ BsmtFin.Type.1 : Factor w/ 7 levels "ALQ", "BLQ", "GLQ", ...: 2 0 1 1 3 3 3 1 3 7 ...
 $ BsmtFin.SF.1 : int 638 450 1000 1000 716 682 618 261 1180 0 ...
 $ BsmtFin.Type.2 : Factor w/ 7 levels "ALQ", "BLQ", "GLQ", ...: 7 4 7 7 7 7 7 7 7 7 ...
 $ BsmtUnf.SF.2 : int 144 0 0 0 0 0 0 0 ...
 $ Total.Bsmt.SF : int 1088 882 2329 2118 928 926 1338 1286 1595 994 ...
 $ Heating : Factor w/ 2 levels "Gas", "Other": 1 1 1 1 1 1 1 1 1 ...
 $ Heating.QC : Factor w/ 4 levels "Ex", "Fa", "Gd", ...: 2 4 4 1 3 3 1 1 1 3 ...
 $ Central.Air : Factor w/ 2 levels "2", "1" ...
 $ Electrical : Factor w/ 4 levels "FuseA", "FuseF", ...: 4 4 4 4 4 4 4 4 4 4 ...
 $ X1st.Flr.SF : int 1656 896 1329 2118 928 926 1338 1280 1616 1628 ...
 $ X2nd.Flr.SF : int 0 0 0 0 781 978 0 0 0 776 ...
 $ Low.Qual.Fin.SF: int 0 0 0 0 0 0 0 0 0 0 ...
```

Number of NA Values: Next, I like to see how many NA's values we have in our dataset. One way I like to do this is by printing a summary of the dataset:

As we can see from the summary output, there are many variables that have at least one NA value. Some being Exterior.2nd, Mas.Vnr.Area, Bsmt.Qual, etc.


```
> complete.cases(ames)
#> [1] TRUE TRUE
#> [22] TRUE TRUE
#> [43] TRUE TRUE
#> [64] TRUE TRUE
#> [85] TRUE TRUE
#> [106] TRUE TRUE
#> [127] TRUE TRUE
#> [148] TRUE TRUE
#> [169] TRUE TRUE
#> [190] TRUE TRUE
#> [211] TRUE TRUE
#> [232] TRUE TRUE
#> [253] TRUE TRUE
#> [274] TRUE TRUE
#> [295] TRUE TRUE
#> [316] TRUE TRUE
#> [337] TRUE TRUE
#> [358] TRUE TRUE
#> [379] TRUE TRUE
#> [400] TRUE TRUE
#> [421] TRUE TRUE
#> [442] TRUE TRUE
#> [463] TRUE TRUE
#> [484] TRUE TRUE
#> [505] TRUE TRUE
#> [526] TRUE TRUE
#> [547] TRUE TRUE
#> [568] TRUE TRUE
#> [589] TRUE TRUE
#> [610] TRUE TRUE
#> [631] TRUE TRUE
#> [652] TRUE TRUE
#> [673] TRUE TRUE
#> [694] TRUE TRUE
#> [715] TRUE TRUE
#> [736] TRUE TRUE
#> [757] TRUE TRUE
#> [778] TRUE TRUE
#> [799] TRUE TRUE
#> [820] TRUE TRUE
#> [841] TRUE TRUE
#> [862] TRUE TRUE
#> [883] TRUE TRUE
#> [904] TRUE TRUE
#> [925] TRUE TRUE
#> [946] TRUE TRUE
```

The Above print out is from the code `complete.cases()`. What this does is it allows me to see if I have removed all of the missing values after running the `na.omit()` function. If all of the values are set to “TRUE” then that means that the `na.omit()` function worked and now my data set does not have any missing values.

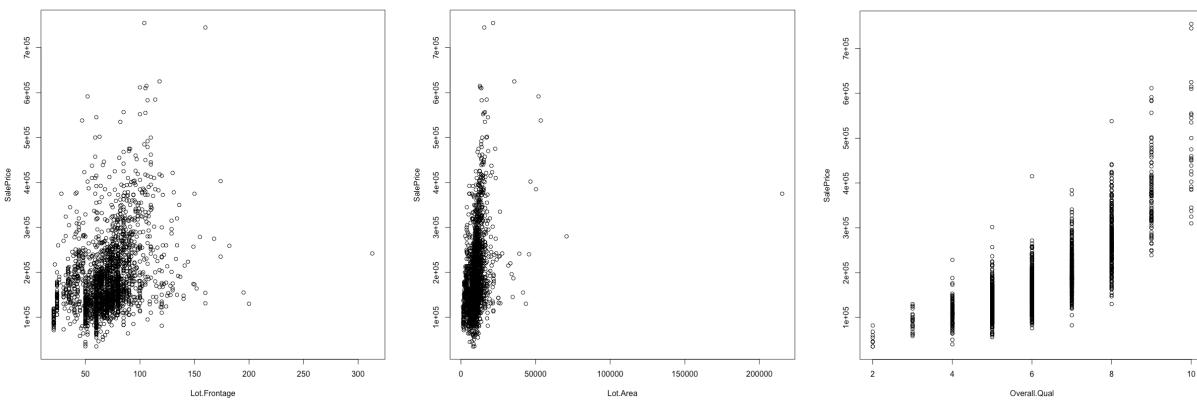
Sale Price vs Independent Variable Plots:

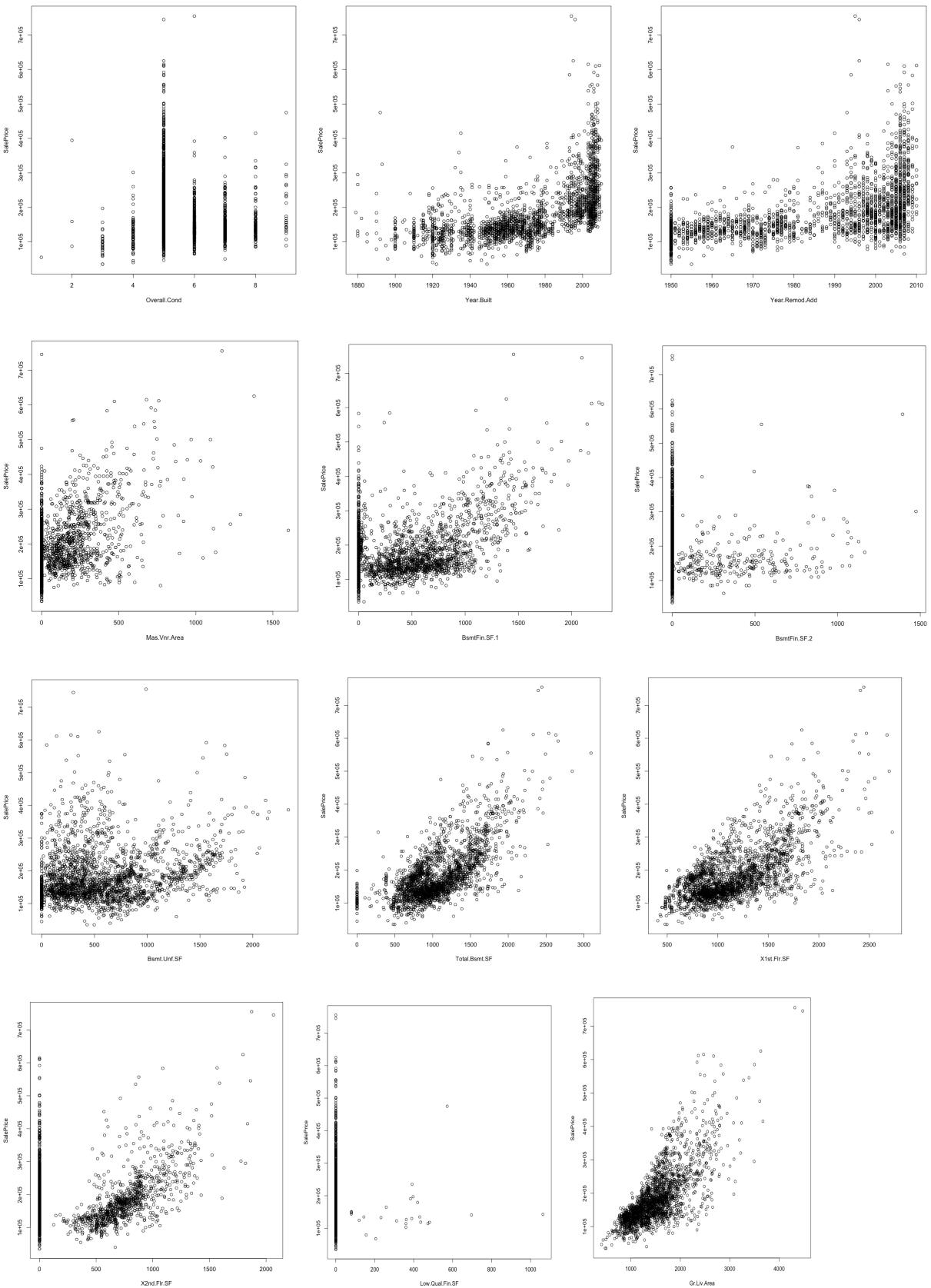
I will now plot SalePrice vs all of the variables in our `num.ames` data frame and see if there is relationship between these two variables. This will give me a good indication which of these variables would be good to keep in our model.

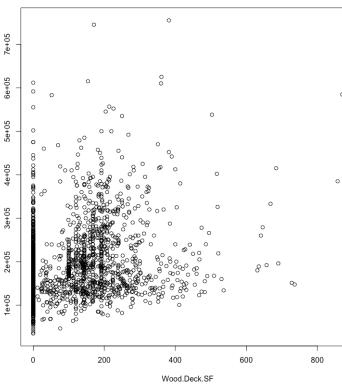
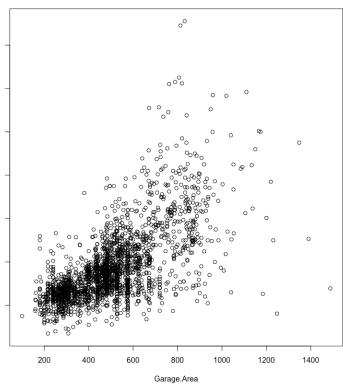
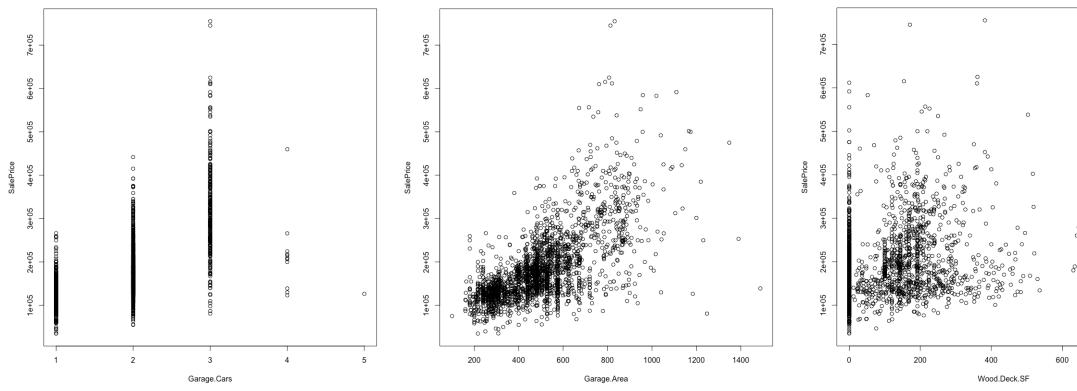
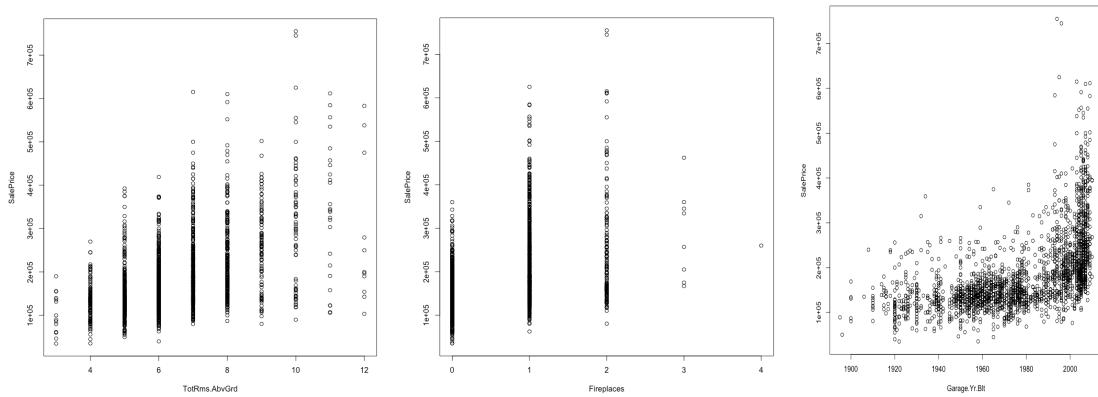
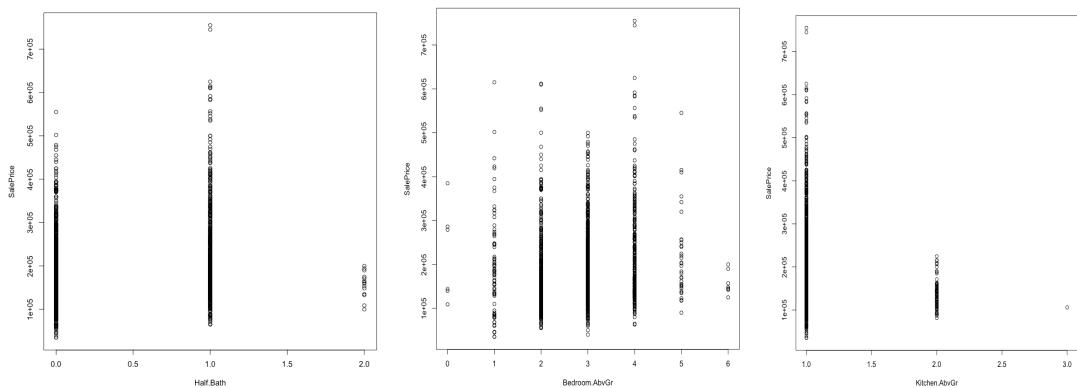
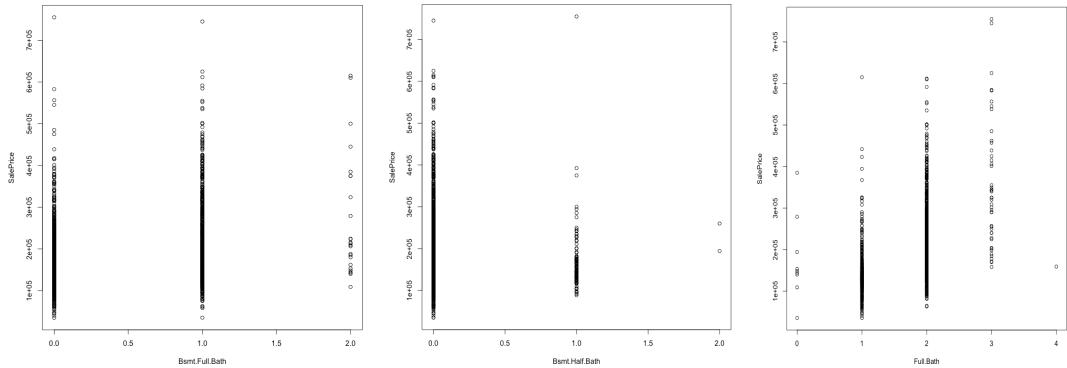
- Looking at the first model is SalePrice vs Lot.Frontage. Doesn’t look like there is a relationship between these two variables as the variance is not consistent. There is a huge cluster at the bottom left of the plot then the variance starts to spread out. However, there could be a non-linear relationship between these two variables. For now, we will keep it in our model.
- Next plot is SalePrice vs Lot.Area. There is no relationship between these variables. Therefore, we will keep this variable out of our model.
- Next plot is SalePrice vs Overall. It looks like there could be a linear relationship between these two variables even though the variance isn’t very consistent. Yet, we will keep this variable in our model for now.
- Next plot is SalePrice vs Overall.Cond, there could be a relationship between these variables. It looks like the plot is almost linear but then the variance shoots up between 4-6 for Overall.Cond. For now will keep this in our model.

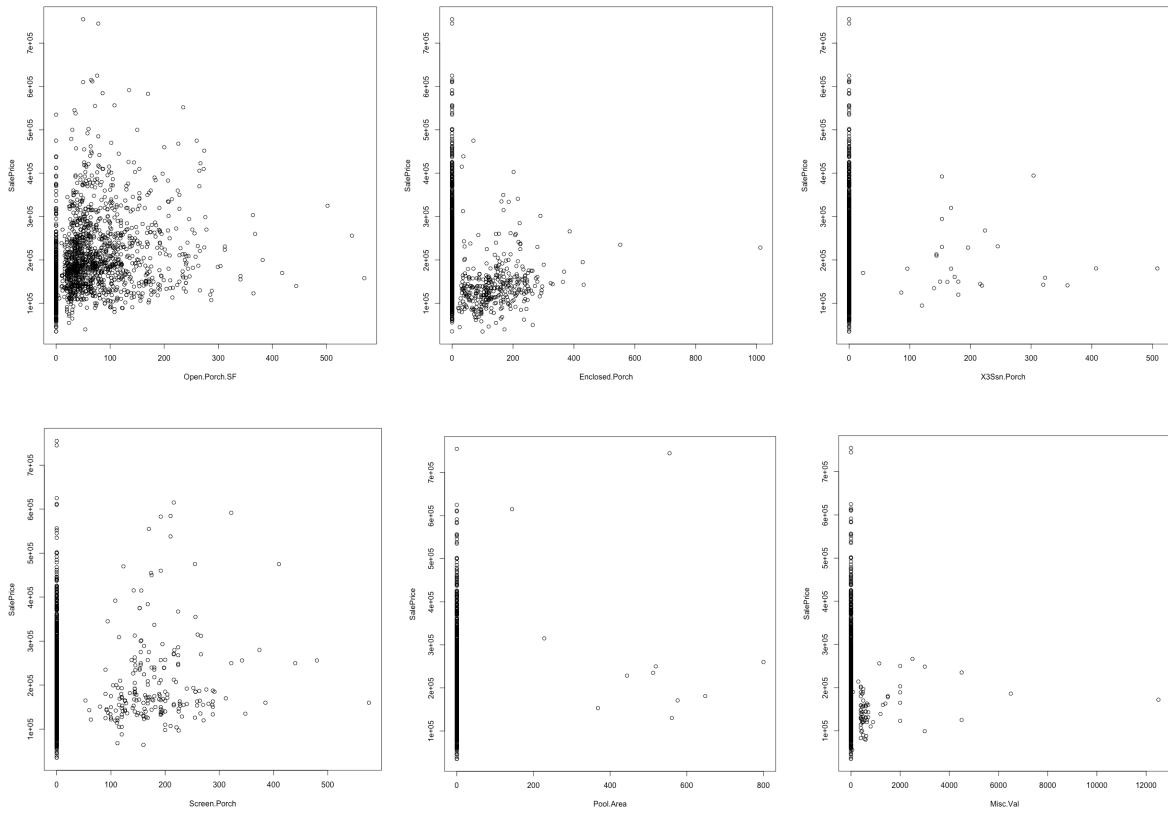
- Next plot is SalePrice vs Year.Built, we can see that there is some sort of positive relationship happening between SalePrice and Year.Built although the variance is not very consistent. We will keep this in our model for now and possibly could take the polynomial of this variable to try and make that variance tighter.
- Next plot is SalePrice vs Years.Remod.Add, we can see that there is a positive relationship between these two variables so we will keep it in our model.
- Next plot is SalePrice vs Mas.Vnr.Area, we can see there there is no relationship between these two variables. Therefore, we will keep this variable out of our model.
- Next plot is SalePrice vs BsmtFin.SF.1. There is no relationship between these two variables. Therefore, we will remove it from our model.
- Next plot is SalePrice vs BsmtFin.SF.2, there is no relationship between these two variables. Therefore, we will remove it from our model.
- Next plot is SalePrice vs Bsmt.Unf.SF, there is no relationship between these two variables. Therefore, we will remove it from our model.
- Next plot is SalePrice vs Total.Bsmt.SF, we can see that there is a positive relationship between these two variables. The variance starts to spread out as square footage increases. We could maybe do a polynomial on this variable to help tighten up the variance. We will keep this in our model for now. It's also important to note that we can see an example of "multicollinearity". Since we have Total Square footage of basement, it's redundant to have the previous two variables BsmFin.SF.1 and BsmFin.SF.2.
- Next plot is SalePrice vs X1st.Fir.SF, we can see here that there is a positive relationship between these two variables but the variance isn't consistent. We will keep this in our model for now.
- Next plot is SalePrice vs X2nd.Flr.SF, There seems to be a lot of outliers for this variable and then a positive relationship between the two. We will keep this in our model for now but will have to address the outliers later.
- Next plot is SalePrice vs Low.Qual.Fin.SF, there is no relationship between these two variables.
- Next plot is SalePrice vs Gr.Live.Area, there is a positive relationship between these two variables. We will keep in model for now.
- Next Plot is SalePrice vs Bsmt.Full.Bath, no relationship here. We will leave this out.

- Next plot is SalePrice vs Bsmt.Half.Bath, no relationship here, we will leave this out of our model.
- Looks like we could have a positive relationship between SalePrice and Full.Bath, but the variance is not consistent. We will keep this in our model for now.
- No relationship between SalePrice and Half.Bath
- No relationship between SalePrice and Bedroom.AbvGr.
- No relationship between SalePrice and Kitchen.AbvGr.
- Looks like there could be a relationship between SalePrice and Total Rooms above ground. We will keep this in our model for now.
- No relationship between SalePrice and Fireplaces.
- Looks like there could be a positive relationship between SalePrice and Year the garage was built, we will keep this in our model for now.
- No relationship between SalePrice and Garage.Cars
- Looks like there is a relationship between SalePrice and the Area of the garage. Variance isn't consistent, we will keep this in our model for now.
- No relationship between SalePrice and Wood.Deck.SF.
- Looks like there could be a negative relationship between SalePrice and Open.Proch.SF, we will keep in model for now.
- No relationship between SalePrice and Enclosed.porch
- No relationship between SalePrice and X3Ssn.Porch.
- No relationship between SalePrice and Screen.Porch.
- No relationship between SalePrice and Pool.Area.





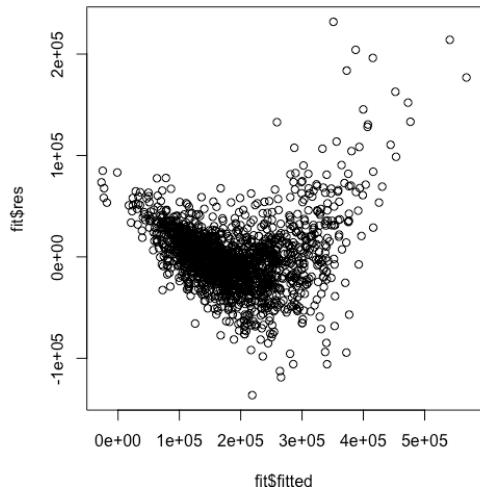




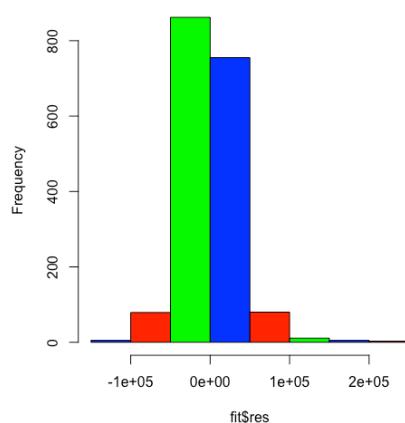
Diagnostic Checks:

I conducted the following checks on my first model after looking at the different plots to see which variables I would want to include. The model is ok but need to do something about TotRms.AbvGrd and Garage.Yr.Blt and Open.Porch.SF as these are not statistically significant.

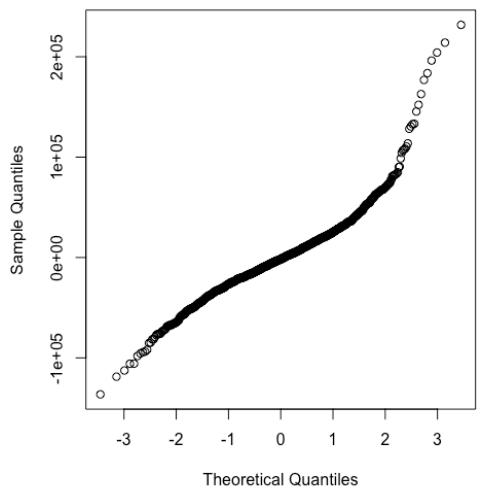
Diagnostic Check Model 1



Diagnostic Check Model Histogram



Normal Q-Q Plot

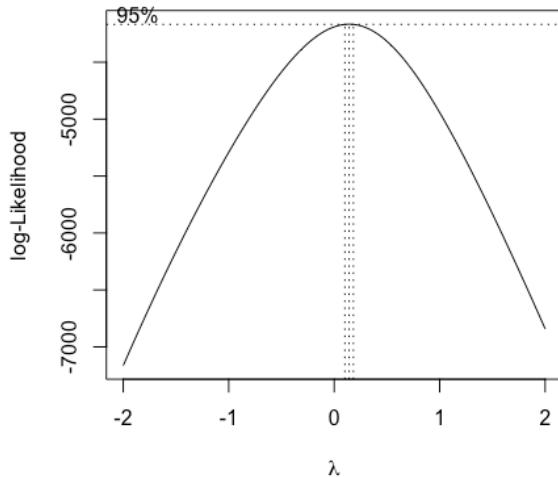


Shapiro-Wilk normality test

```
data: fit$res  
W = 0.93176, p-value < 2.2e-16
```

To address both of these issues, I ran a boxcox to see if I could fix these three variables. After running the boxcox, the result had me take the log of SalePrice. The summary print out was a stronger model but TotRms.AbvGrd and Garage.Yr.Blt still had a P-Value higher than 0.05. After looking at the relationship between these two variables and SalePrice again, I decided to

remove these two variables all together. The Diagnostic checks are much better all around and I am happy with this being my model for my numeric values.



Diagnostic Checks After BoxCox:

```

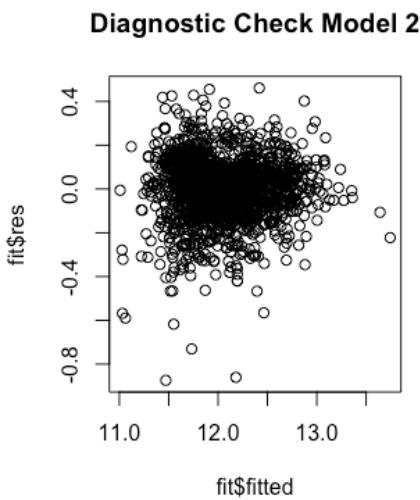
Call:
lm(formula = num.ames$logSalePrice ~ Overall.Qual + Year.Built +
    Year.Remod.Add + BsmtFin_SF_1 + Total.Bsmt_SF + X1st.Flr_SF +
    Gr.Liv.Area + Wood.Deck_SF + Open.Porch_SF, data = num.ames,
subset = train)

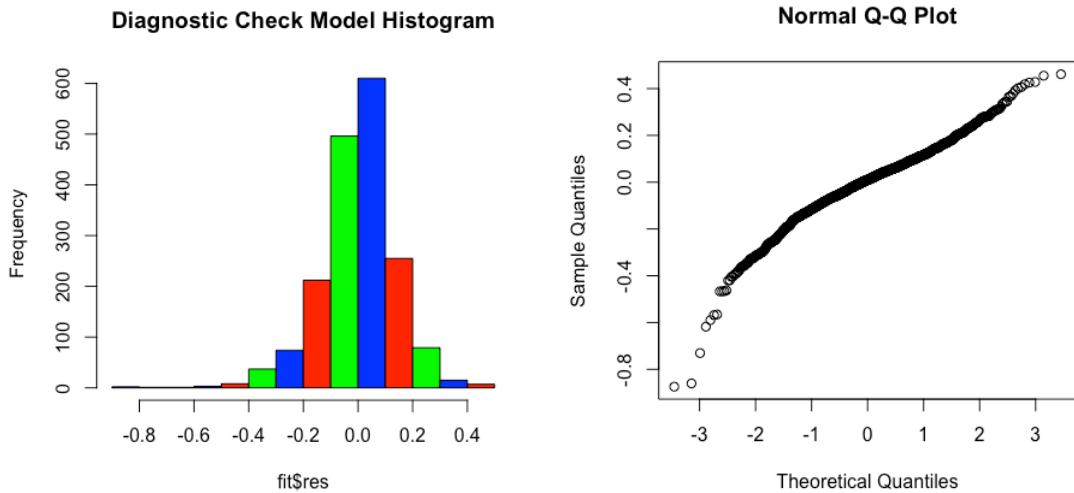
Residuals:
    Min      1Q  Median      3Q     Max 
-0.87421 -0.07030  0.00998  0.07985  0.46140 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 2.423e+00 3.927e-01  6.171 8.35e-10 ***
Overall.Qual 9.480e-02 3.797e-03 24.971 < 2e-16 ***
Year.Built   1.670e-03 1.532e-04 10.901 < 2e-16 ***
Year.Remod.Add 2.550e-03 2.137e-04 11.934 < 2e-16 ***
BsmtFin_SF_1 1.204e-04 8.476e-06 14.208 < 2e-16 ***
Total.Bsmt_SF 1.043e-04 1.446e-05 7.214 7.99e-13 ***
X1st.Flr_SF  6.190e-05 1.566e-05  3.952 8.04e-05 *** 
Gr.Liv.Area   2.826e-04 9.192e-06 30.746 < 2e-16 ***
Wood.Deck_SF 1.120e-04 2.808e-05  3.987 6.95e-05 *** 
Open.Porch_SF 1.733e-04 5.420e-05  3.198 0.00141 ** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1374 on 1790 degrees of freedom
Multiple R-squared:  0.8848,    Adjusted R-squared:  0.8842 
F-statistic:  1527 on 9 and 1790 DF,  p-value: < 2.2e-16

```



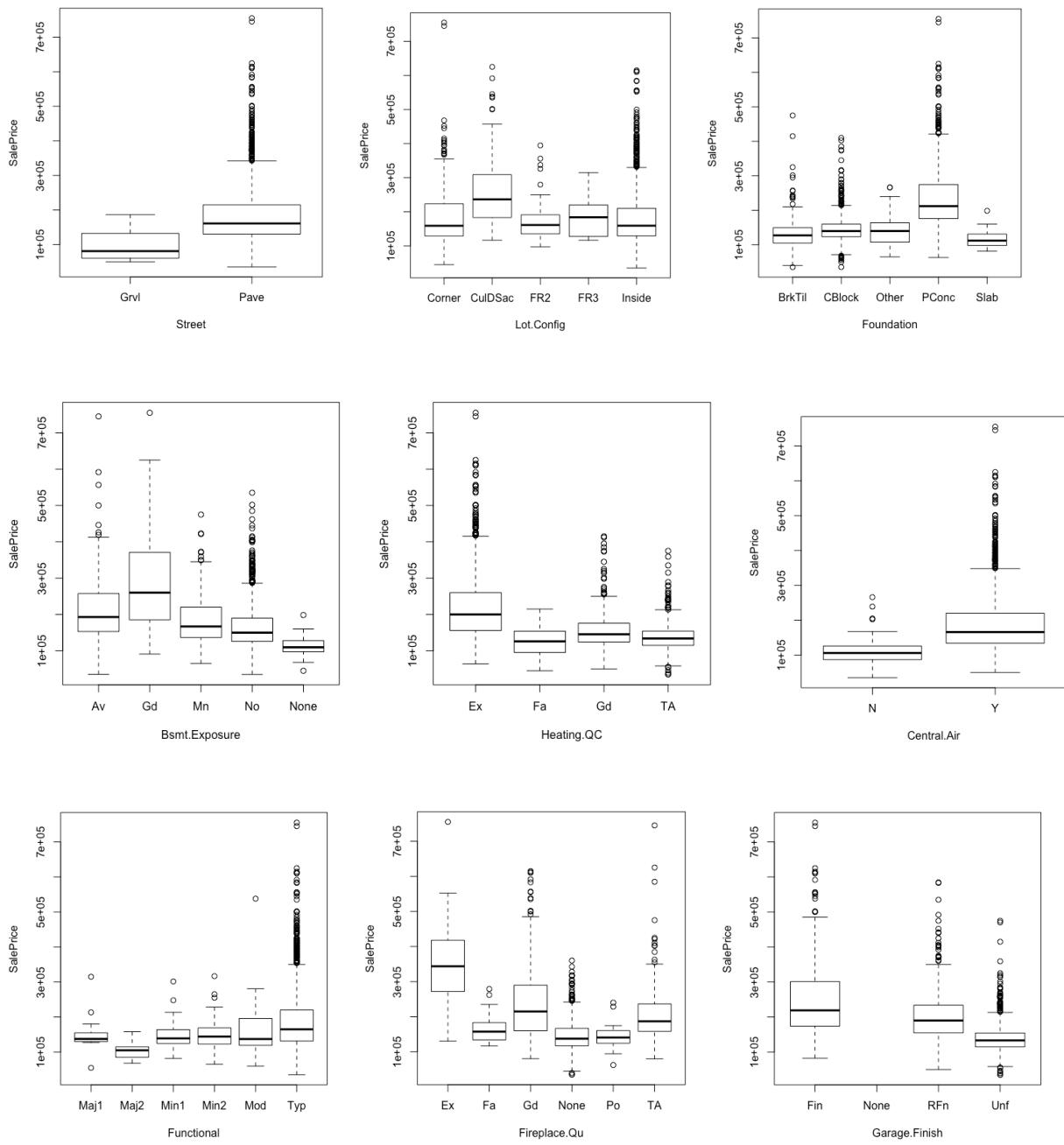


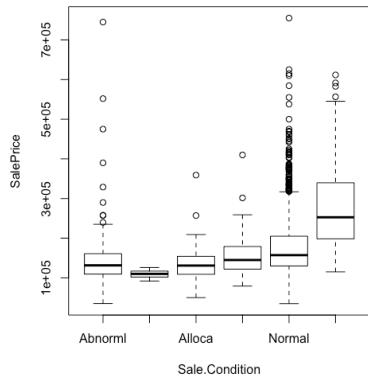
```
Shapiro-Wilk normality test
data: fit$res
W = 0.96388, p-value < 2.2e-16
```

I noticed that there are still outliers in this model. I am ok with that for now though and plan on addressing this at the end once I also have model for my categorical variables.

Preliminary Analysis for Categorical Variables:

I plotted all of categorical variables to see which ones had a relationship with SalePrice. The following plots are the categorical variables I thought had a relationship with SalePrice. There are some variables that was a bit difficult to be able to tell if there was or wasn't a relationship. I went ahead and threw them into my model for now. I will remove them later if needed.





Model Comparison:

I will now compare add one categorical variable to the model I created with only the numeric variables. Then use the anova function to compare them both.

```
Analysis of Variance Table

Model 1: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF
Model 2: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF + Street
Res.Df   RSS Df Sum of Sq    F Pr(>F)
1    1790 33.818
2    1789 33.732  1  0.085418 4.5302 0.03344 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Analysis of Variance Table

Model 1: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF
Model 2: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config
Res.Df   RSS Df Sum of Sq    F Pr(>F)
1    1790 33.818
2    1785 33.455  5  0.36261 3.8694 0.00172 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Analysis of Variance Table

Model 1: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF
Model 2: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config + Foundation
Res.Df   RSS Df Sum of Sq    F Pr(>F)
1    1790 33.818
2    1781 33.336  9  0.48108 2.8558 0.002403 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

Analysis of Variance Table

Model 1: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin_SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF
Model 2: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin_SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config + Foundation +
  Bsmt.Exposure
Res.Df   RSS Df Sum of Sq    F    Pr(>F)
1  1790 33.818
2  1777 32.959 13   0.85834 3.5598 1.522e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

After Running the `nova` function a few more times. I ended up with this final model.

```

> anova(model1,model2)
Analysis of Variance Table

Model 1: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin_SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF
Model 2: num.ames$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
  BsmtFin_SF.1 + Total.Bsmt.SF + X1st.Flr.SF + Gr.Liv.Area +
  Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config + Foundation +
  Bsmt.Exposure + Heating.QC + Central.Air + Functional
Res.Df   RSS Df Sum of Sq    F    Pr(>F)
1  1790 33.818
2  1768 29.938 22   3.8792 10.413 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

I then conducted a diagnostics checks on my new model which includes both the numeric and non-numeric values: Everything Looks good

```

Call:
lm(formula = num.ames$logSalePrice ~ Overall.Qual + Year.Built +
  Year.Remod.Add + BsmtFin.SF.1 + Total.Bsmt.SF + Xlst.Flr.SF +
  Gr.Liv.Area + Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config +
  Foundation + Bsmt.Exposure + Heating.QC + Central.Air + Functional,
  data = Ames2, subset = train)

Residuals:
    Min      1Q  Median      3Q     Max 
-0.80369 -0.06770  0.00426  0.07541  0.46330 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 4.268e+00  4.909e-01   8.694 < 2e-16 ***
Overall.Qual 8.439e-02  3.826e-03  22.055 < 2e-16 ***
Year.Built   1.185e-03  1.982e-04   5.982 2.66e-09 ***
Year.Remod.Add 1.923e-03  2.210e-04   8.704 < 2e-16 ***
BsmtFin.SF.1 1.102e-04  8.361e-06  13.179 < 2e-16 ***
Total.Bsmt.SF 8.064e-05  1.851e-05   4.357 1.40e-05 ***
Xlst.Flr.SF   8.403e-05  1.920e-05   4.375 1.28e-05 ***
Gr.Liv.Area   2.855e-04  9.005e-06  31.699 < 2e-16 ***
Wood.Deck.SF  1.645e-04  2.719e-05   3.845 0.000125 ***
Open.Porch.SF 1.889e-04  5.193e-05   3.638 0.000282 ***
StreetPave    1.179e-01  5.438e-02   2.169 0.030231 *  
Lot.ConfigCulDSac 4.720e-02  1.734e-02   2.723 0.006536 ** 
Lot.ConfigFR2  -4.339e-02  2.187e-02  -1.984 0.047390 *  
Lot.ConfigFR3  -8.240e-02  4.693e-02  -1.756 0.079277 .  
Lot.ConfigInside -1.201e-03  8.381e-03  -0.143 0.886093  
FoundationCBlock 5.292e-03  1.302e-02   0.406 0.684478  
FoundationOther -1.895e-02  4.308e-02  -0.419 0.675247  
FoundationPConc  1.780e-02  1.570e-02   1.133 0.257276  
FoundationSlab   2.326e-02  4.749e-02   0.490 0.624098  
Bsmt_ExposureGd 2.937e-02  1.336e-02   2.198 0.028091 *  
Bsmt_ExposureMn -2.027e-02  1.381e-02  -1.468 0.142197  
Bsmt_ExposureNo -2.609e-02  9.752e-03  -2.676 0.007530 ** 
Bsmt_ExposureNone -4.252e-02  4.399e-02  -0.966 0.333939  
Heating_QCFa  -4.888e-02  2.112e-02  -2.315 0.026725 *  
Heating_QCGd  -1.891e-02  9.712e-03  -1.948 0.051616 .  
Heating_QCTA  -6.149e-02  9.218e-03  -6.670 3.41e-11 *** 
Central.AirY   1.146e-01  1.590e-02   7.207 8.43e-13 *** 
FunctionalMaj2  4.080e-02  7.637e-02   0.535 0.592543  
FunctionalMin1  1.798e-01  4.323e-02   4.158 3.36e-05 *** 
FunctionalMin2  1.847e-01  4.287e-02   4.309 1.73e-05 *** 
FunctionalMod   1.419e-01  4.888e-02   2.952 0.003204 ** 
FunctionalTyp   2.270e-01  3.796e-02   5.981 2.67e-09 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1301 on 1768 degrees of freedom
Multiple R-squared:  0.898,    Adjusted R-squared:  0.8962 
F-statistic:  502 on 31 and 1768 DF,  p-value: < 2.2e-16

```

After Looking at the model, It looks like the only variable that needs to be removed is Foundation categorical variable. There are other categorical variables that may have a P-Value over 0.05 but they're still going to be left in our model since they're attached to other variables in their subsidiary that "DO" have P-Values less than 0.05.

Final Checks Look Good.

```

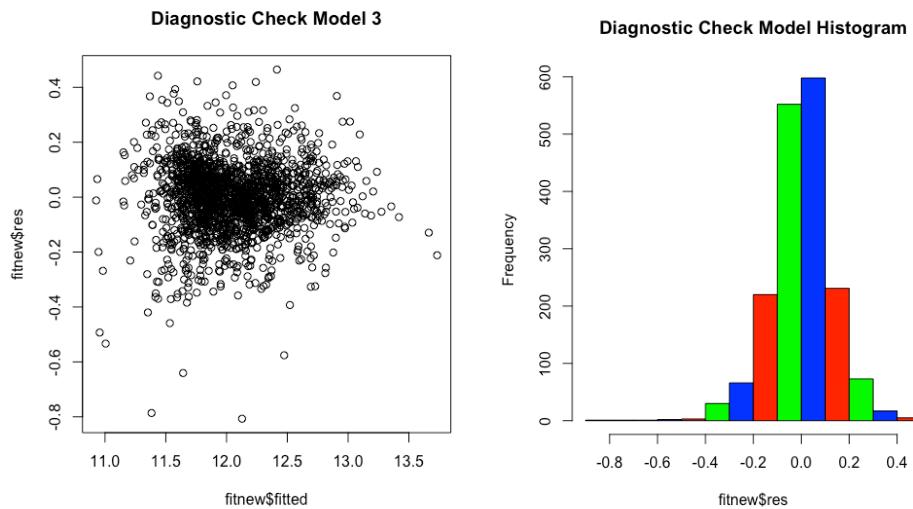
lm(formula = num.ames$logSalePrice ~ Overall.Qual + Year.Built +
  Year.Remod.Add + BsmtFin.SF.1 + Total.Bsmt.SF + X1st.Flr.SF +
  Gr.Liv.Area + Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config +
  Bsmt.Exposure + Heating.QC + Central.Air + Functional, data = Ames2,
  subset = train)

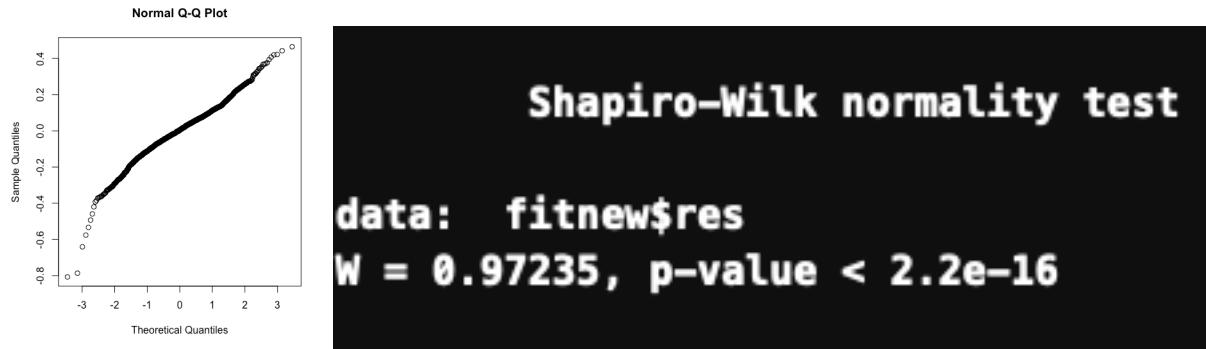
Residuals:
    Min      1Q  Median      3Q     Max 
-0.80698 -0.06840  0.00451  0.07490  0.46453 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 3.919e+00 4.227e-01  9.271 < 2e-16 ***
Overall.Qual 8.447e-02 3.748e-03 22.539 < 2e-16 ***
Year.Built   1.325e-03 1.553e-04  8.533 < 2e-16 ***
Year.Remod.Add 1.964e-03 2.149e-04  9.138 < 2e-16 ***
BsmtFin.SF.1 1.699e-04 8.294e-06 13.251 < 2e-16 ***
Total.Bsmt.SF 8.286e-05 1.839e-05  4.464 8.57e-06 ***
X1st.Flr.SF   8.236e-05 1.902e-05  4.330 1.57e-05 ***
Gr.Liv.Area   2.869e-04 8.897e-06 32.252 < 2e-16 ***
Wood.Deck.SF 1.019e-04 2.711e-05  3.759 0.000176 *** 
Open.Porch.SF 1.904e-04 5.178e-05  3.678 0.000242 *** 
StreetPave    1.216e-01 5.392e-02  2.256 0.024218 *  
Lot.ConfigCulDSac 4.618e-02 1.731e-02  2.668 0.007705 ** 
Lot.ConfigFR2 -4.368e-02 2.184e-02 -2.000 0.045681 *  
Lot.ConfigFR3 -8.309e-02 4.687e-02 -1.773 0.076451 . 
Lot.ConfigInside -1.964e-03 8.357e-03 -0.235 0.814175  
Bsmt.ExposureGd 2.924e-02 1.335e-02  2.190 0.028634 * 
Bsmt.ExposureMn -2.007e-02 1.380e-02 -1.455 0.145921  
Bsmt.ExposureNo -2.584e-02 9.744e-03 -2.652 0.008069 ** 
Bsmt.ExposureNone -3.021e-02 3.091e-02 -0.977 0.328563  
Heating.QCFa   -5.043e-02 2.099e-02 -2.403 0.016362 * 
Heating.QCGd   -2.633e-02 9.586e-03 -2.121 0.034947 * 
Heating.QCTA   -6.488e-02 8.780e-03 -7.298 4.39e-13 *** 
Central.AirY   1.129e-01 1.540e-02  7.332 3.44e-13 *** 
FunctionalMaj2  3.156e-02 7.563e-02  0.417 0.676484  
FunctionalMin1  1.790e-01 4.318e-02  4.144 3.57e-05 *** 
FunctionalMin2  1.835e-01 4.281e-02  4.287 1.91e-05 *** 
FunctionalMod  1.423e-01 4.796e-02  2.967 0.003049 ** 
FunctionalTyp  2.269e-01 3.794e-02  5.982 2.66e-09 *** 
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.1301 on 1772 degrees of freedom
Multiple R-squared:  0.8978,    Adjusted R-squared:  0.8963 
F-statistic: 576.8 on 27 and 1772 DF,  p-value: < 2.2e-16

```





Removing Outliers:

I first had to put the missing data back into my dataframe. By using the na.action=na.excluded function at the end of my linear model.

Afterwards I received the following output once I ran my code:

Misc.Val	Yr.Sold	Sale.Type	Sale.Condition	SalePrice	sres
307	0	2010	WD	Alloca	50138 -3.440812
372	0	2009	WD	Family	82500 -6.349877
561	0	2009	WD	Normal	62383 -4.903205
695	0	2009	WD	Abnorml	66500 -3.998236
790	0	2009	COD	Abnorml	60000 -5.317468
1471	0	2008	WD	Abnorml	185000 -3.072142
1553	0	2008	WD	Normal	40000 -6.248267
1780	0	2007	New	Partial	147000 -4.549078
1947	0	2007	WD	Normal	64500 -3.600840
1996	0	2007	WD	Normal	133900 -3.032386
2598	0	2006	COD	Abnorml	80000 -3.070253
2839	0	2006	Con	Normal	35000 -3.929471
2876	0	2006	WD	Abnorml	35311 -4.251353

Misc.Val	Yr.Sold	Sale.Type	Sale.Condition	SalePrice	sres
919	0	2009	WD	Normal	205000 3.245110
929	0	2009	WD	Normal	158000 3.230606
1535	0	2008	WD	Alloca	257076 3.162753
1638	0	2007	New	Partial	392000 3.659438
2005	0	2007	WD	Normal	144000 3.559679
2208	0	2007	WD	Normal	153500 3.005213
2214	0	2007	WD	Normal	315000 3.373858
2854	0	2006	WD	Normal	195000 3.202883

The following 21 Variables I ended up removing. My Model was already in good shape and I figured removing these 21 outliers wouldn't hurt my model since it is such a big dataset.

The following is the new summary output with the outliers now removed:

```

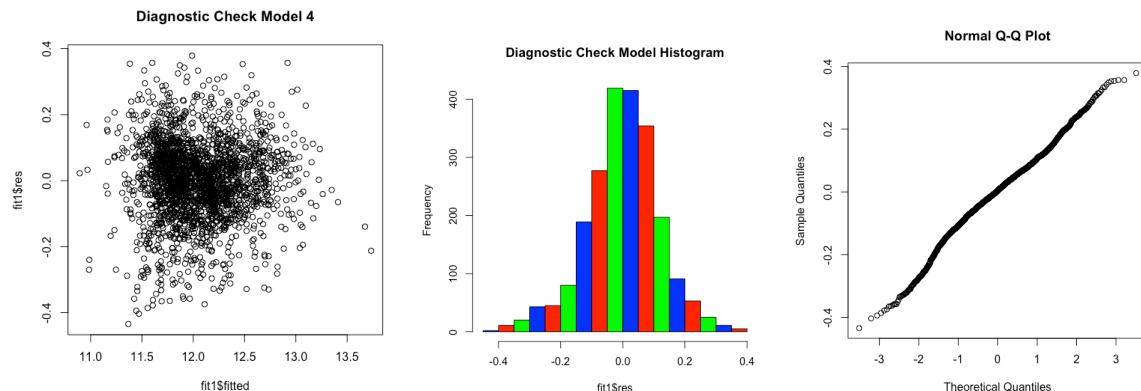
Call:
lm(formula = x$logSalePrice ~ Overall.Qual + Year.Built + Year.Remod.Add +
    BsmtFin.SF.1 + Total.Bsmt.SF + Xist.Flr.SF + Gr.Liv.Area +
    Wood.Deck.SF + Open.Porch.SF + Street + Lot.Config + Bsmt.Exposure +
    Heating.QC + Central.Air + Functional, data = Ames, na.action = na.exclude)

Residuals:
    Min      1Q  Median      3Q     Max 
-0.43427 -0.06700  0.00412  0.07138  0.37859 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 3.615e+00  3.411e-01 10.600 < 2e-16 ***
Overall.Qual 7.955e-02  3.084e-03 25.797 < 2e-16 ***
Year.Built   1.434e-03  1.279e-04 11.291 < 2e-16 ***
Year.Remod.Add 1.985e-03  1.728e-04 11.482 < 2e-16 ***
BsmtFin.SF.1 1.015e-04  6.809e-06 14.913 < 2e-16 ***
Total.Bsmt.SF 1.065e-04  1.520e-05 7.007 3.22e-12 ***
Xist.Flr.SF   6.752e-05  1.573e-05 4.291 1.85e-05 ***
Gr.Liv.Area   2.939e-04  7.301e-06 40.253 < 2e-16 ***
Wood.Deck.SF  8.423e-05  2.216e-05 3.800 0.000148 ***
Open.Porch.SF 1.454e-04  4.231e-05 3.438 0.000597 *** 
StreetPave    1.310e-01  4.883e-02 2.683 0.007351 ** 
Lot.ConfigCulDSac 3.709e-02  1.424e-02 2.665 0.009242 ** 
Lot.ConfigFR2  -4.076e-02  1.656e-02 -2.461 0.013930 *  
Lot.ConfigFR3  -7.558e-02  4.006e-02 -1.887 0.059327 .  
Lot.ConfigInside -2.044e-03  6.799e-03 -0.301 0.763691  
Bsmt.ExposureGd 3.599e-02  1.081e-02 3.329 0.000885 *** 
Bsmt.ExposureMn -1.888e-02  1.117e-02 -1.690 0.091085 . 
Bsmt.ExposureNo -2.419e-02  7.781e-03 -3.188 0.001905 ** 
Bsmt.ExposureNone -2.604e-02  2.489e-02 -1.046 0.295702  
Heating.QCFa   -7.449e-02  1.704e-02 -4.372 1.29e-05 *** 
Heating.QCGd   -2.148e-02  7.737e-03 -2.777 0.005534 ** 
Heating.QCTA   -5.807e-02  7.134e-03 -8.140 6.51e-16 *** 
Central.Airy    8.985e-02  1.262e-02 7.117 1.48e-12 *** 
FunctionalMaj2  4.488e-02  6.430e-02 0.698 0.485235  
FunctionalMin1  2.619e-01  3.967e-02 6.602 5.05e-11 *** 
FunctionalMin2  2.559e-01  3.978e-02 6.433 1.53e-10 *** 
FunctionalMod   2.561e-01  4.316e-02 5.935 3.41e-09 *** 
FunctionalTyp  3.065e-01  3.594e-02 8.527 < 2e-16 *** 
--- 
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.1181 on 2209 degrees of freedom 
Multiple R-squared:  0.9105, Adjusted R-squared:  0.9094 
F-statistic: 832.1 on 27 and 2209 DF, p-value: < 2.2e-16

```



```

Shapiro-Wilk normality test

data: fit1$res
W = 0.98937, p-value = 8.332e-12

```

After removing the outliers from our model, I ran the diagnostic checks one more time on my final model and everything looks great! Fitted vs Residuals has a consistent variance. Our histogram is a normal distribution. Normal Q-Q Plot and Shapiro-Wilk Look good.

AmesHousing_Predict Data:

Now that I have developed my model, I will now predict SalePrice using the AmesHousing_Predict data:

```
> pred = predict(fit1,Predict.Data,interval="prediction")
> pred
      fit      lwr      upr
1 11.40149 11.15050 11.65248
2 11.77428 11.53968 12.00888
```

Thus, I predict that 95% of SalePrice for:

Housing 1, has a Price between 11.2(logofSalePrice) and 11.66(logofSalePrice)

Housing 2, has a Price between 11.5(logofSalePrice) and 12.0(logofSalePrice)