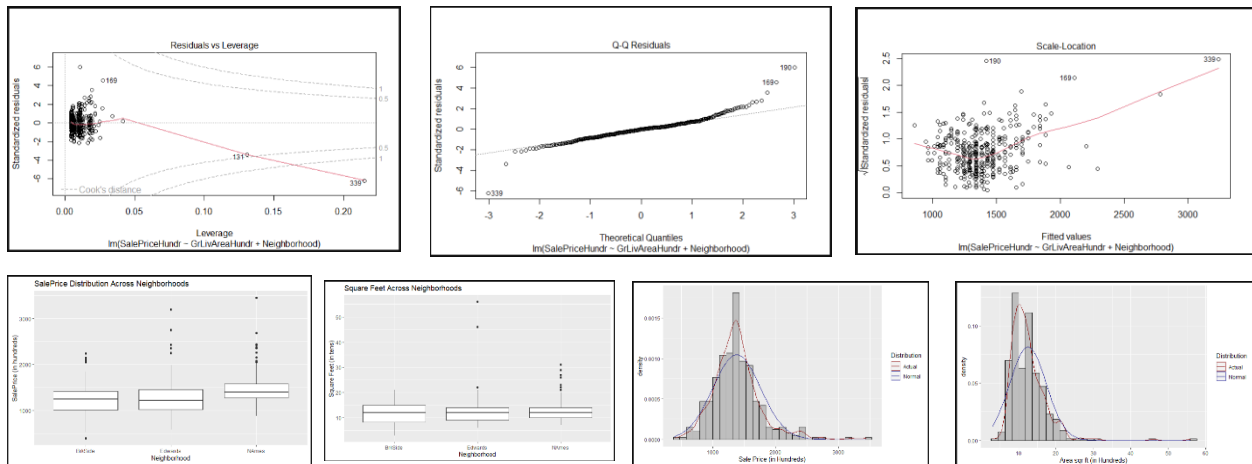


Analysis 1 Summary: Century 21 Ames is seeking to analyze the relationship of the square foot living area and the sale prices of houses in the North Ames, Edwards, and Brookside neighborhoods in Ames, Iowa to improve their customer service and competitive edge.

Problem: Using multiple linear regression analysis estimate if there is a difference in the price compared to the living space square footage in the North Ames, Edwards, and Brookside neighborhoods and how those neighborhoods statistically compare with each other.

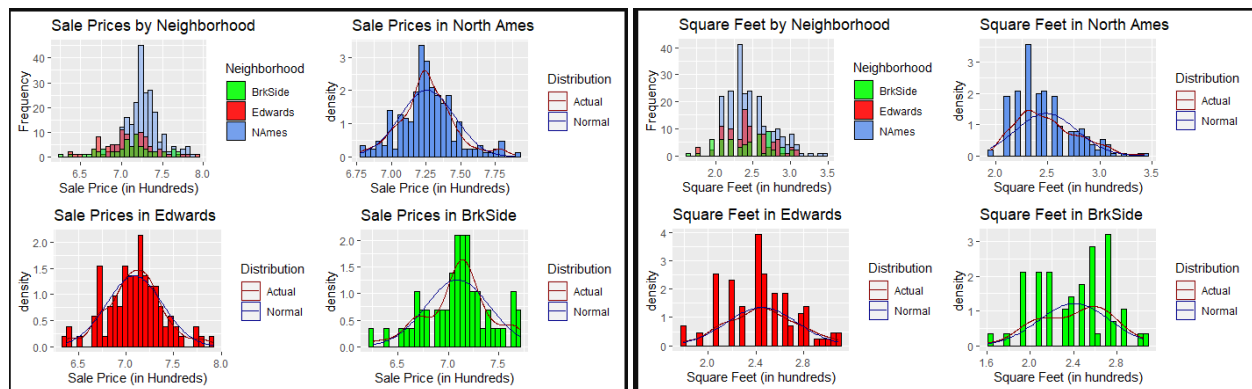
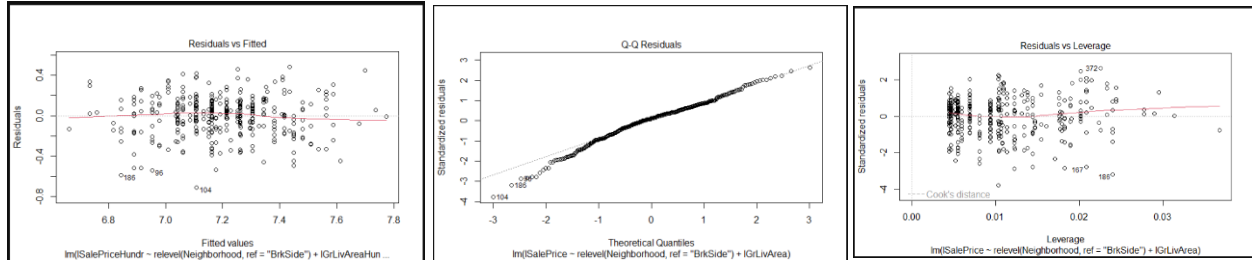
Assumptions Check: An initial visual check of **normality** using histograms and QQ-plots show some right skewedness, probably resulting from outliers with strong leverage (See graphics below). Because the sample size is greater than 30 the Central Limit Theorem (CLT) is expected to resolve inconsistencies in normality. Nonetheless, we will log transform some of the data to check for improved fit. A visual review of the residuals also shows **some variance across the line** despite a good variance inflation factor (VIF) for sale price and living space (1.007636, 1.003811). The variables appear to be **independent** and have good **linearity**. The primary issue with this test is the **outliers**. Because the outliers have so much leverage, we opted to review these extremes to check for human error and other anomalies (see leverage and boxplots below).



Influential Outliers: A review of the outliers showed that IDs 524 and 1299 were new construction and in an incomplete status. We use Zillow to review the properties in Edwards associated with these IDs and concluded that the square footage exceeded any properties in that neighborhood and the low prices were probably related to lot price rather than a newly constructed home. Additionally, ID 534 in Brookside was assessed to be missing a zero on both sales (39300) and square footage (324 sqft). Lastly, ID 725 in Edwards had an extraordinary sale price (320000) for the living space (1698) and condition (5). All these outliers were determined to be the result of human error or reporting requirements for planned construction and were removed from the model.

Comparing Competing models: All models in analysis 1 used one generalized formula ($\text{salePrice} \sim \text{GrLivArea} + \text{Neighborhood}$) adjusted for the logged variable. We tested a multiple linear regression ($r^2 = 0.5081$), log-linear ($r^2 = 0.4785$), linear-log ($r^2 = 0.4897$), log-log ($r^2 = 0.4903$), and a logged interaction ($r^2 = 0.5094$). The internal CV press for the log-log model returned a fairly good RMSE of 0.189923, suggesting strong performance, an r^2 of 0.49082, suggesting decent predictive performance, and a MAE of 0.145467 indicating only minor difference between the predicted and actual values.

Log-Log transformations: We opted to go with the Log-Log model with no interaction, primarily because it best met the assumptions. As you can see below, the log-log transformation and removed outliers improved *normality*, *linearity*, the *residual distribution*, and *variance* across most neighborhoods, despite some right skewedness for living space in North Ames.



Parameters and Interpretation of the model:

An analysis of the North Ames, Edwards, and Brookside neighborhoods shows relatively no statistically significant difference in the relationship of sales price and living area between Brookside (the reference neighborhood) and Edwards neighborhoods (p -value = 0.382). Although a purchase in the Edwards neighborhood is associated with a 0.02755 (or \$97 non-logged) decrease in log-transformed price, this difference is not significant. We are 95 percent confident that the true logged sales price difference is found between (-0.0894956 and 0.0344004).

```
-----Neighborhood no interaction-----
Call:
lm(formula = lSalePrice ~ relevel(Neighborhood, ref = "BrkSide") +
    lGrLivArea, data = AmesHousing_Data3)

Residuals:
    Min       1Q   Median       3Q      Max
-0.72335 -0.10484  0.02247  0.11660  0.48696

Coefficients:
(Intercept)              7.66029    0.24041  31.863  < 2e-16 ***
relevel(Neighborhood, ref = "BrkSide")Edwards -0.02755    0.03150  -0.874   0.382
relevel(Neighborhood, ref = "BrkSide")NAmes    0.12412    0.02809   4.419  1.3e-05 ***
lGrLivArea              0.57229    0.03388  16.891  < 2e-16 ***

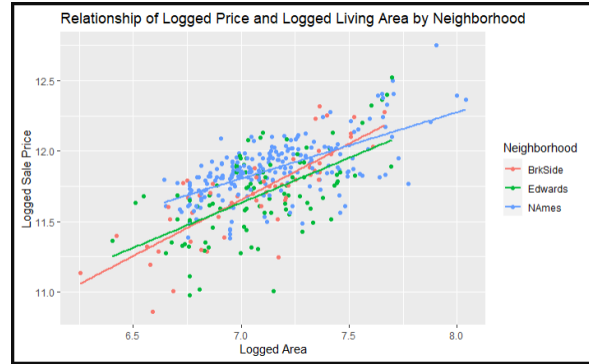
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1886 on 375 degrees of freedom
Multiple R-squared:  0.4903,    Adjusted R-squared:  0.4863
F-statistic: 120.3 on 3 and 375 DF,  p-value: < 2.2e-16

              2.5 %    97.5 %
(Intercept)  7.18756764  8.1330200
relevel(Neighborhood, ref = "BrkSide")Edwards -0.08949560  0.0344004
relevel(Neighborhood, ref = "BrkSide")NAmes    0.06889019  0.1793571
lGrLivArea    0.50566510  0.6389101
```

However, there is strong evidence (p -value < .0001) that the relationship between sales price and living area is different between Brookside and North Ames. For North Ames, there is evidence, holding all other variables constant, that purchasing in that neighborhood is associated with an estimated 0.12412 (\$113 non-logged) increase in the log-transformed sale price for each unit of living space compared to Brookside. We are 95 percent confident that the true logged difference in the purchasing price between North Ames and Brookside is between (0.06889 and 0.179357).

Regarding the relationship between sales price and living area across all values, there is overwhelming evidence (p -value < 0.0001) of a positive linear relationship. For each log-transformed unit of living space, the log-transformed sales price is increased on average by 0.57229 units. In real world terms, that is for every 100 square feet of living space there is an estimated increase of \$171 in the base purchase price.

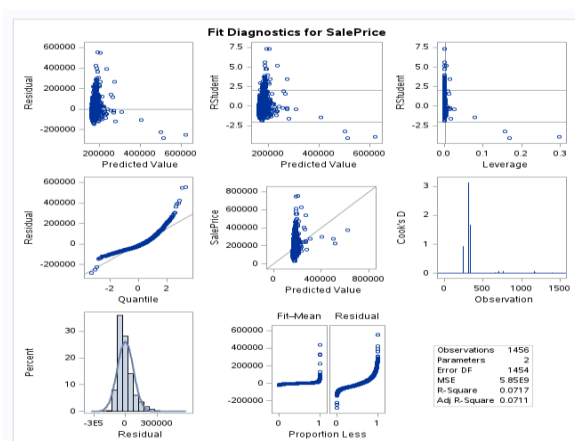
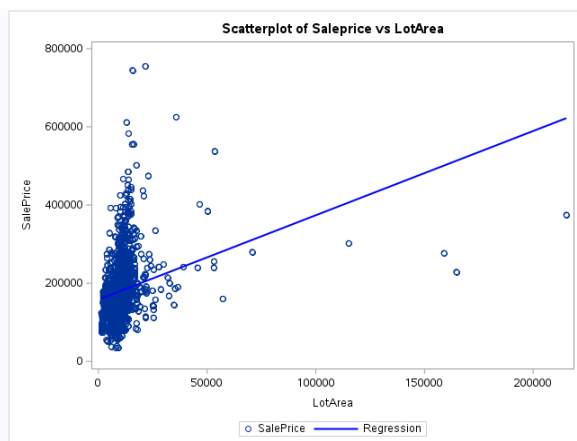


Conclusion: There is overwhelming evidence (p -value < 0.0001) of a positive linear relationship between the sales price and the square footage of living space. However, this relationship does not mean there are significant differences between North Ames, Edwards, and Brookside neighborhoods in purchasing price and living space. A breakdown of the price to square footage between Edwards and Brookside shows no significant difference. Our best estimate was that square footage cost \$97 dollars less in Edwards than Brookside, but this was not statistically significant. In North Ames, however, there was strong evidence that square footage cost \$113 more per hundred square feet than the same living space in Brookside. This model explained only 49.03 percent of the variation and other variables should be considered to strengthen the model. (See Appendix 1)

Analysis 2 Summary: Expanding on the request from the customer, we will seek to build the most predictive model for the sale prices of homes in all of the neighborhoods in Ames Iowa.

The Problem: Build three models; a simple linear regression model, a multiple linear regression model comparing sales price, above ground living Area (GrLivArea), and number of full baths, and a custom multiple linear regression model that will show the relationship of the most predictive variables. These models will be compared to each other using metrics such as adjusted R squared, CV press and Kaggle score.

Model 1: Simple Regression – Examining the relationship between Salesprice and lot area.

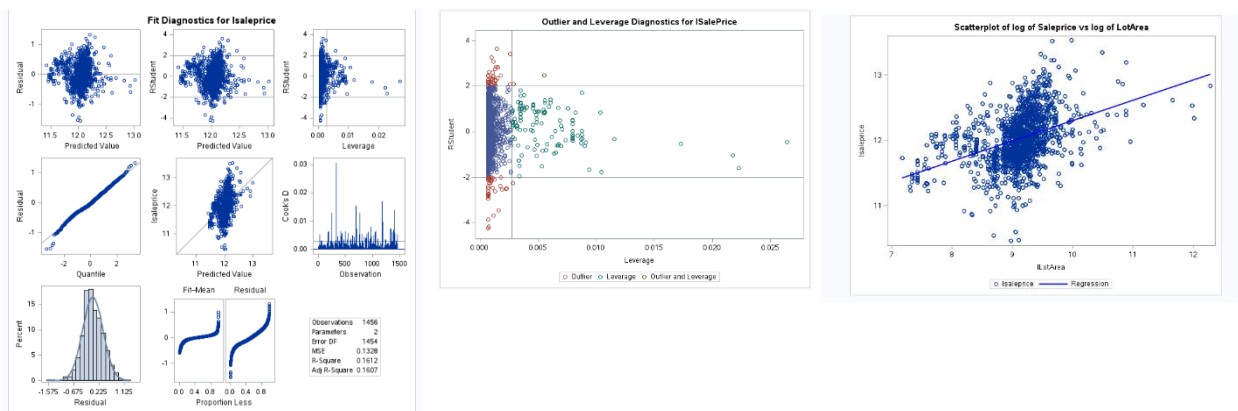


Assumptions Check

An initial visual review of the plots between saleprice and lot area shows some right skewedness in the histogram and q-q plot, suggesting the assumption for normal distribution is not met. Additionally, there appears to be clustering of the residuals and several outliers with heavy leverage, suggesting that there is not equal spread. In regard to the sample, we assume that the observations are independent and that the data represents the entire single family housing population in Ames.

Log-Log Transformation

To mitigate for the lot size outliers (cases with 2 to 5 acres) we will perform a log-log transformation on the sales price and lot area and examine the relationship between the two logged variables. (Note: Limited information on the cost associated with acreage precluded an analysis that would have allowed us to drop these most egregious outliers. Instead of ranged exclusion, we opted for log-log transformations.) By first plotting the log - log transformed data, there is a better fitting linear relationship between the logged saleprice (response variable) and the explanatory variable (Lot area). Furthermore, a first look at the graphics seem to meet the assumptions. There does not appear to be any potential influential point. Judging from the scatter plot, q-q plot and histogram of the residuals, there is no evidence that the residuals do not follow a normal distribution with constant variance. We continue to assume here that the observations are independent.



$$\text{Pred log(Saleprice)} = \beta_0 + \beta_1 \text{log(LotArea)}$$

$$\text{Pred log(Saleprice)} = 9.193 + 0.311\text{log(LotArea)}$$

We are 95% confident that for each doubling of the lot area, the median Saleprice rate will increase between approximately $(2^{0.27})$ 20.6% and $(2^{0.35})$ 27.5%. Our best estimate is an increase of $(2^{0.31})$ 24%.

The REG Procedure
Model: MODEL1
Dependent Variable: lsaleprice

Number of Observations Read	1456
Number of Observations Used	1456

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	37.13003	37.13003	279.52	<.0001
Error	1454	193.14466	0.13284		
Corrected Total	1455	230.27468			

Root MSE	0.36447	R-Square	0.1612
Dependent Mean	12.02455	Adj R-Sq	0.1607
Coeff Var	3.03103		

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits
Intercept	1	9.19260	0.16996	54.18	<.0001	8.85980 9.52540
lLotArea	1	0.31091	0.01880	16.72	<.0001	0.27443 0.34739

Model 2: Multiple Regression with SalePrice ~ GrLivArea + Full Bath

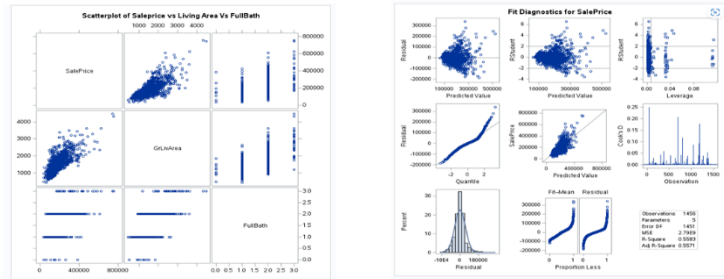
By plotting the data first, we can tell there is some evidence of a positive relationship between the explanatory variables and the response variable. Please see appendix 2 for additional visual plots.

Assumptions check:

There does appear to be an influential point in the Rstudent plot. That observation will require further investigation. Judging from the scatter plot, q-q plot and histogram of the residuals, there is no evidence that the residuals do not follow a normal distribution but there is not enough evidence to indicate a constant variance. We assume here that the observations are independent.

Log – Log Transformation

We will do a log transformation of Saleprice and living area and plot those variables against each other to see if there is a linear relationship between them. Please see appendix 2 for plot.



Looking at the scatterplot matrix there is a positive linear relationship between the explanatory variables (Living area and full bath) and the response variable (Sales price). The log transformed variables show a more positive linear relationship between lsaleprice and lGrLivArea than the original data. We will build and fit a model with the log transformed data.

Assumptions Check

There do appear to be three influential points. That observation will require further investigation. Judging from the scatter plot, q-q plot and histogram of the residuals, there is no evidence that the residuals do not follow a normal distribution and a constant variance. We assume here that the observations are independent.

We fit our model with ran on our explanatory (lsaleprice) and dependent (lgrlivArea and Full bath). This model has an Adjusted R square of 0.5637 which means this model estimates about 56.37% of the variation in salesprice is explained by the explanatory variables, a SBC of -3860.2586, and a CV press of 101.20 . Please see appendix for more details on the model.

The next step is to build a model with the log of sale price as the explanatory variable and fit the model.

$$\begin{aligned} \text{pred}\{\text{Saleprice}\} &= \beta_0 + \beta_1 \text{ lgrlivArea} + \beta_2 \text{ FB0} + \beta_3 \text{ FB1} + \beta_4 \text{ FB2} \text{ Reference} = \text{FB3} \\ \text{pred}\{\text{Saleprice}\} &= 6.9797 + 0.723 \text{ lgrlivArea} - 0.2424 \text{ FB0} - 0.2992 \text{ FB1} - 0.142 \text{ FB2} \\ \text{pred}\{\text{Saleprice} \mid \text{ lGrLivArea, FullBath}=0\} &= 6.9797 + 0.723 \text{ lgrlivArea} - 0.2424 \text{ FB0} = 6.9209 + 0.723 \text{ lgrlivArea} \\ \text{pred}\{\text{Saleprice} \mid \text{ lGrLivArea, FullBath}=1\} &= 6.9797 + 0.723 \text{ lgrlivArea} - 0.2992 \text{ FB1} = 6.6805 + 0.723 \text{ lgrlivArea} \\ \text{pred}\{\text{Saleprice} \mid \text{ lGrLivArea, FullBath}=2\} &= 6.9797 + 0.723 \text{ lgrlivArea} - 0.142 \text{ FB2} = 6.8377 + 0.723 \text{ lgrlivArea} \\ \text{pred}\{\text{Saleprice} \mid \text{ lGrLivArea, FullBath}=3\} &= 6.9797 + 0.723 \text{ lgrlivArea} \end{aligned}$$

We are 95% confident that for each doubling of the living area, the median Saleprice rate will increase between approximately 59% ($2^{0.67} = 1.59$) and 71% ($2^{0.78} = 1.71$). Our best estimate is an increase of 65% ($2^{0.72} = 1.647$) after holding all other variables (Full Bath) constant. To get the regression

equation for a specific full bath number, we will need to adjust that full bath's coefficient with the intercept in our regression model.

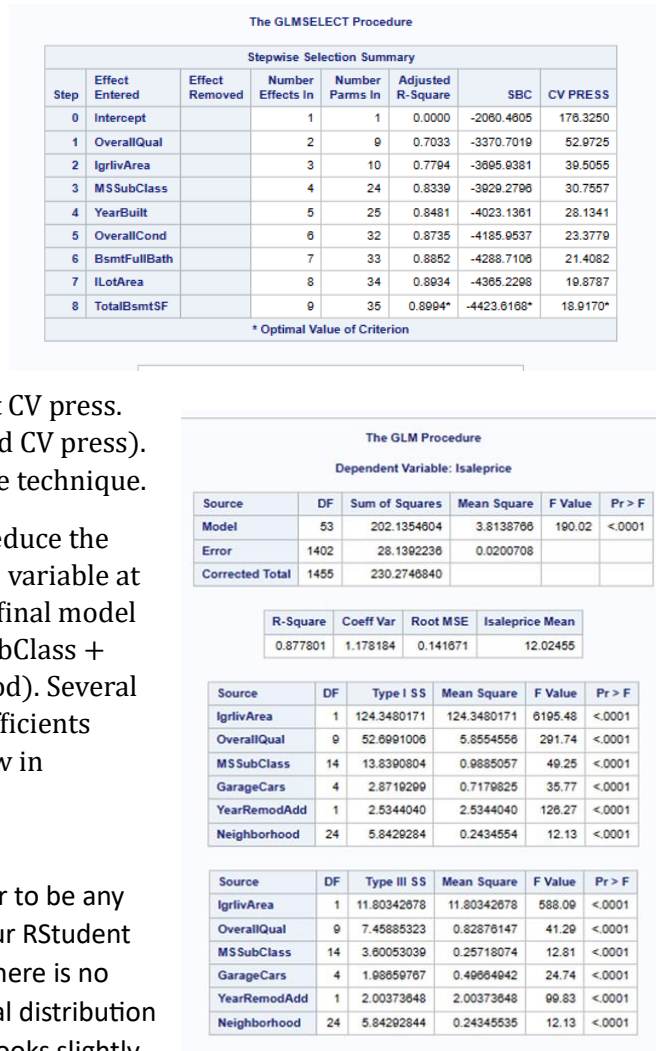
Custom Model

For the custom model, we decided to go with an automatic selection technique and then finetune the model with manual selection techniques. We did this by running stepwise, forward and backward selection models on the explanatory variables in the housing dataset and our response variable (Saleprice). The stepwise method made the most sense with the highest R score and lowest CV press. (See the graphic for Adjusted R square, SBC, and CV press). So, for our custom model we chose the stepwise technique.

Because our model was so large, we opted to reduce the number of variables manually by removing one variable at a time and retesting the model. We will fit our final model ($ISalePrice \sim IGrLivArea + OverallQual + MSSubClass + GarageCars + YearRemodelAdd + Neighborhood$). Several of our variables were categorical and their coefficients require individual assessment, which we review in appendix 2.

Assumptions Check

Looking at our fit residuals, there does not appear to be any influential point on the residual plots based on our RStudent Plot. Judging from the scatter plot, q-q plot and there is no evidence that the residuals do not follow a normal distribution and a constant variance although our histogram looks slightly skewed to the left due to the presence of some outliers in the data. We will assume that all observations are independent of each other. (Please see appendix for fit residuals plot)



Comparing Competing Models

Predictive Models	Adjusted R2	CV PRESS	Kaggle Score
Simple Linear Regression	.16	103.5	.983
Multiple Linear Regression	.56	101.2	.864
Custom MLR Model	.899	31.53	.552
Other Models

Conclusion

Comparing the three models, that is our simple Linear regression, our multiple linear regression (MLR), and the Custom MLR model, it is quite obvious that the custom MLR model with $\text{IGrLivArea} + \text{OverallQual} + \text{MSSubClass} + \text{GarageCars} + \text{YearRemodelAdd} + \text{Neighborhood}$ as predictors for salesprice is the most useful model to predict the sale price of the homes in Iowa. It has the highest adjusted R2 of 0.899, and the lowest CV press of 31.53, our chosen metric for comparison. Approximately 89.9% of the variation in the salesprice is explained by explanatory variables in the model. Interpreting the slope for our regression equation will depend on what slope we want to interpret. For instance, interpreting the slope for living area will mean that we are 95% confident that for each doubling of the living area, the median Saleprice will increase between approximately 38% ($2^{0.46} = 1.38$) and 45% ($2^{0.54}$). Our best estimate is an increase of 41% ($2^{0.5} = 1.41$) after holding all other variables constant. There is strong evidence ($p\text{-value} < .0001$) of a relationship with sales price.

Additionally, several categorical variables showed a statistically significant relationship with sales price (See appendix 2), including several neighborhood, Model type, and the number of cars to the garage, but varies based on the category. To get the regression equation for a specific categorical variables, we will need to adjust that variable's coefficient with the intercept in our regression model.

Project RShiny App: [Statistical Insights \(weiprecht.github.io\)](https://weiprecht.github.io)

Project RShiny Page: [SFDS Final Project \(weiprecht.github.io\)](https://weiprecht.github.io)

Interactive Graphic: [Ames Housing App \(shinyapps.io\)](https://shinyapps.io)

Project RShiny Github: [Ames Housing App \(shinyapps.io\)](https://shinyapps.io)

Appendix 1: Analysis 1 Code –

Initial analysis of original data

```

# Initial dataset for Analysis question 1: check assumptions
r housing_regression

# Fit a multiple linear regression model
lm_model <- lm(SalePrice ~ GrLivArea + Neighborhood, data = AmesHousing_Data2)

# Summarize the model
summary(lm_model)

# Calculate vif for the variables in the model
vif_result <- vif(lm_model)
vif_result

# Calculate mean and standard deviation
mean_sale_price <- mean(AmesHousing_Data2$SalePriceHundr)
sd_sale_price <- sd(AmesHousing_Data2$SalePriceHundr)
mean_sqr_ft <- mean(AmesHousing_Data2$GrLivAreaHundr)
sd_sqr_ft <- sd(AmesHousing_Data2$GrLivAreaHundr)

# Check Assumptions Visually
plot(lm_model)
avPlots(lm(SalePrice ~ GrLivArea + Neighborhood, data = AmesHousing_Data2))

# For histograms
ggplot(AmesHousing_Data2, aes(x = SalePriceHundr)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "grey") +
  geom_density(aes(color = "Actual")) +
  stat_function(aes(color = "Normal"),
               fun = dnorm,
               args = list(mean = mean_sale_price, sd = sd_sale_price)) +
  xlab("Sale Price (in Hundreds)") +
  scale_colour_manual("Distribution", values = c("darkred", "darkblue"))

ggplot(AmesHousing_Data2, aes(x = GrLivAreaHundr)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "grey") +
  geom_density(aes(color = "Actual")) +
  stat_function(aes(color = "Normal"),
               fun = dnorm,
               args = list(mean = mean_sqr_ft, sd = sd_sqr_ft)) +
  xlab("Area sqr ft (in Hundreds)") +
  scale_colour_manual("Distribution", values = c("darkred", "darkblue"))

# Boxplot for SalePrice
ggplot(AmesHousing_Data2, aes(x = Neighborhood, y = SalePriceHundr)) +
  geom_boxplot() +
  labs(title = "SalePrice Distribution Across Neighborhoods",
       x = "Neighborhood",
       y = "SalePrice (in hundreds)")

# Boxplot for square feet
ggplot(AmesHousing_Data2, aes(x = Neighborhood, y = GrLivAreaHundr)) +
  geom_boxplot() +
  labs(title = "Square Feet Across Neighborhoods",
       x = "Neighborhood",
       y = "Square Feet (in hundreds)")

```

```

call:
lm(formula = SalePrice ~ GrLivArea + Neighborhood, data = AmesHousing_Data2)

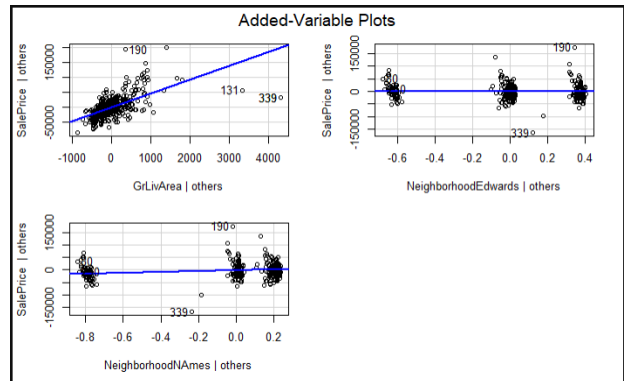
Residuals:
    Min       1Q   Median       3Q      Max
-165078  -16215    281    13578  175400

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  69781.538   5442.400  12.822 < 2e-16 ***
GrLivArea     45.760     3.149   14.533 < 2e-16 ***
NeighborhoodEdwards -2882.155   4930.632   -0.585 0.559204
NeighborhoodAmes  16105.621   4395.352   3.664 0.000283 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

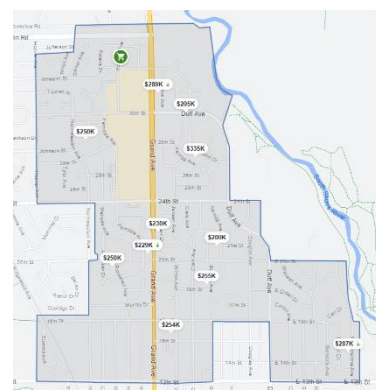
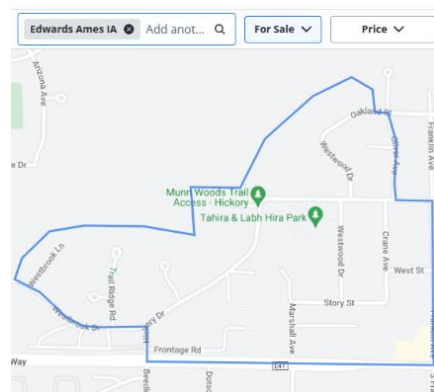
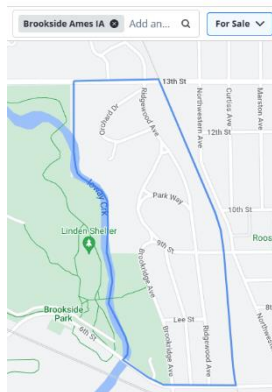
Residual standard error: 29760 on 379 degrees of freedom
Multiple R-squared:  0.3965, Adjusted R-squared:  0.3917
F-statistic: 83 on 3 and 379 DF, p-value: < 2.2e-16

              GVIF Df GVIFA(1/(2*Df))
GrLivArea    1.008149  1      1.004066
Neighborhood 1.008149  2      1.002031

```



Zillow maps of Brookside, Edwards, and North Ames neighborhoods.



Internal Cross Validation:

```

## {r internal_cross_validation}
# define the model formula
formula <- SalePrice ~ GrLivArea + Neighborhood

# Create a train control object for cross-validation
ctrl <- trainControl(method = "cv", # cross-validation method
                    number = 5, # Number of folds (you can change this)
                    verboseIter = TRUE, # Display iteration progress
                    summaryFunction = defaultSummary) # use default summary function

# Train the model with cross-validation
lm_model <- train(formula,
                 data = AmesHousing_data2,
                 method = "lm", # Linear regression method
                 trcontrol = ctrl)

# View results
lm_model

```

```

383 samples
2 predictor

No pre-processing
Resampling: Cross-validated (5 fold)
Summary of sample sizes: 307, 305, 307, 306, 307
Resampling results:

RMSE      Rsquared  MAE
29756.1   0.4114803  20897.76

Tuning parameter 'intercept' was held constant at a value of TRUE

```

Interaction code:

```

## {r housing_regression_logged}
# ref is always the element that does not show up in the model
cat("In-----Neighborhood no interaction-----\n")
pricefit <- lm(lSalePrice ~ relevel(neighborhood, ref = "BrkSide") + lGrLivArea, data = AmesHousing_data3)
summary(pricefit)
confint(pricefit)

# Get coefficients and statistical summary in a tidy format
tidy_summary2 <- tidy(pricefit)
print(tidy_summary2)
plot(pricefit)

cat("\n-----Neighborhood interaction-----\n")

# log-pricing (log all) Regression
pricefit2 <- lm(lSalePrice ~ relevel(neighborhood, ref = "BrkSide") + lGrLivArea, data = AmesHousing_data3)
summary(pricefit2)
confint(pricefit2)

# Get coefficients and statistical summary in a tidy format
tidy_summary2 <- tidy(pricefit2)
print(tidy_summary2)
plot(pricefit2)

```

```

-----Neighborhood interaction-----
[call]:
lm(formula = lSalePrice ~ relevel(neighborhood, ref = "BrkSide") +
    lGrLivArea, data = AmesHousing_data3)
Residuals:
    Min       1Q   Median       3Q      Max
-0.72588 -0.10046  0.02184  0.10631  0.52051
Coefficients:
(Intercept)                6.06485    0.54988   11.029 < 2e-16 ***
relevel(neighborhood, ref = "BrkSide")Edwards  2.42788    0.63272    3.837 0.000146 ***
relevel(neighborhood, ref = "BrkSide")Names    0.78836    0.07784   10.257 < 2e-16 ***
lGrLivArea                -0.15917    0.10353   -1.537 0.125044
relevel(neighborhood, ref = "BrkSide")Edwards:lGrLivArea -0.32534    0.08933   -3.642 0.000309 ***
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1855 on 373 degrees of freedom
Multiple R-squared:  0.5094, Adjusted R-squared:  0.5029
F-statistic: 77.47 on 5 and 373 DF, p-value: < 2.2e-16

(Intercept)                2.5 %      97.5 %
4.9835007  7.14609601
relevel(neighborhood, ref = "BrkSide")Edwards
-0.5487246  2.53461692
relevel(neighborhood, ref = "BrkSide")Names
1.1837337  3.67202489
lGrLivArea
0.6453021  0.55142031
relevel(neighborhood, ref = "BrkSide")Edwards:lGrLivArea
-0.3627466  0.04440939
relevel(neighborhood, ref = "BrkSide")Names:lGrLivArea
-0.5009937 -0.14968137

```

Scatterplot code:

```

## {r}
library(ggplot2)
ggplot(AmesHousing_data3, aes(x = lGrLivArea, y = lSalePrice, color = Neighborhood)) +
  geom_point() +
  geom_smooth(method = "lm", formula = y ~ x, se = FALSE) +
  labs(title = "Relationship of Logged Price and Logged Living Area by Neighborhood",
       x = "Logged Area", y = "Logged Sale Price", color = "Neighborhood")

```

Conversion code (transform coefficients back to original.)

```

## {r conversion}
# coefficients from the regression model
intercept <- 7.66029
coeff_Edwards <- -0.02755
coeff_NAMES <- 0.12412
coeff_GrLivArea <- 0.57229

# Transforming coefficients back to original scale
original_intercept <- exp(intercept)
original_coeff_Edwards <- exp(coeff_Edwards)
original_coeff_NAMES <- exp(coeff_NAMES)
original_coeff_GrLivArea <- exp(coeff_GrLivArea)

# Interpretation in dollars and square feet
# Assuming GrLivArea is measured in square feet
# Assuming the SalePrice was measured in dollars
# Interpretation of coefficients in original scale
interpretation_intercept <- paste("Intercept (in hundreds of dollars):", round(original_intercept, 2))
interpretation_Edwards <- paste("Effect of Edwards (in hundreds of dollars):", round(original_coeff_Edwards, 2))
interpretation_NAMES <- paste("Effect of NAMES (in hundreds of dollars):", round(original_coeff_NAMES, 2))
interpretation_GrLivArea <- paste("Effect of GrLivArea (in hundreds of square feet):", round(original_coeff_GrLivArea, 2))

# Printing interpretations
print(interpretation_intercept)
print(interpretation_Edwards)
print(interpretation_NAMES)
print(interpretation_GrLivArea)

```

Appendix 2: Initial dataset for Analysis Question 2

```

AmesHousing_data$Neighborhood <- as.factor(AmesHousing_data$Neighborhood)

# create a new column that holds the 'salePrice', 'GrLivArea', and 'LotArea' in hundreds
AmesHousing_data$SalePriceHundr <- floor(AmesHousing_data$SalePrice /100)
AmesHousing_data$GrLivAreaHundr <- floor(AmesHousing_data$GrLivArea /100)
AmesHousing_data$LotAreaHundr <- floor(AmesHousing_data$LotArea /100)

# Drop previously identified outliers rows 131, 136, and 339
rows_to_drop <- c(131, 136, 190, 339)
AmesHousing_data4 <- AmesHousing_data[-rows_to_drop, ]

AmesHousing_data4
summary(AmesHousing_data4)

# Fit a multiple linear regression model
lm_model <- lm(SalePrice ~ LotArea, data = AmesHousing_data4)

# Summarize the model
summary(lm_model)

# Calculate mean and standard deviation
mean_sale_price <- mean(AmesHousing_data4$SalePriceHundr)
sd_sale_price <- sd(AmesHousing_data4$SalePriceHundr)
lotmean_sqr_ft <- mean(AmesHousing_data4$LotAreaHundr)
lotsd_sqr_ft <- sd(AmesHousing_data4$LotAreaHundr)

# Check Assumptions Visually
plot(lm_model)

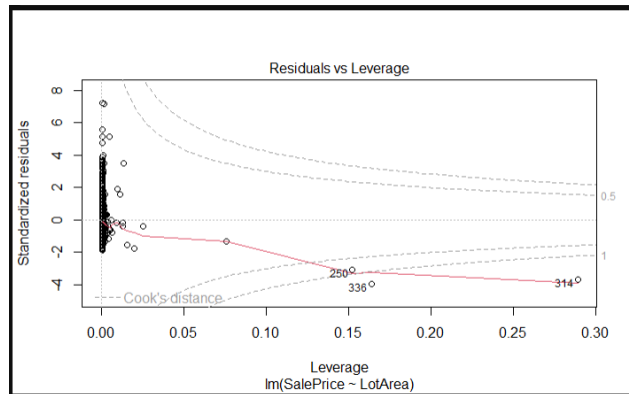
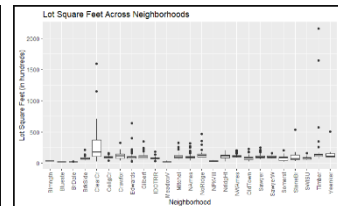
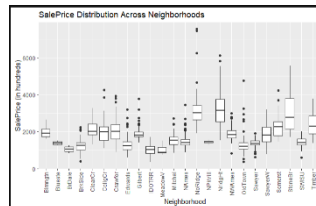
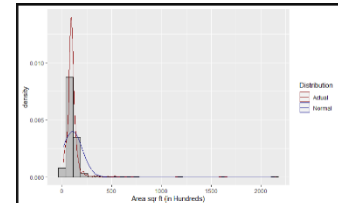
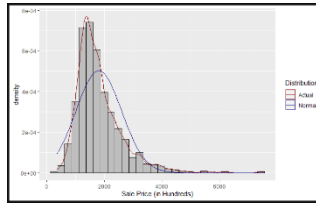
# For histograms
ggplot(AmesHousing_data4, aes(x = SalePriceHundr)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "grey") +
  geom_density(aes(color = "Actual")) +
  stat_function(aes(color = "Normal"),
               fun = dnorm,
               args = list(mean = mean_sale_price, sd = sd_sale_price)) +
  xlab("Sale Price (in Hundreds)") +
  scale_colour_manual("Distribution", values = c("darkred", "darkblue"))

ggplot(AmesHousing_data4, aes(x = LotAreaHundr)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "grey") +
  geom_density(aes(color = "Actual")) +
  stat_function(aes(color = "Normal"),
               fun = dnorm,
               args = list(mean = lotmean_sqr_ft, sd = lotsd_sqr_ft)) +
  xlab("Area sq ft (in Hundreds)") +
  scale_colour_manual("Distribution", values = c("darkred", "darkblue"))

# Boxplot for saleprice
ggplot(AmesHousing_data4, aes(x = Neighborhood, y = SalePriceHundr)) +
  geom_boxplot() +
  labs(title = "SalePrice Distribution Across Neighborhoods",
       x = "Neighborhood",
       y = "SalePrice (in hundreds)") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))

# Boxplot for square feet
ggplot(AmesHousing_data4, aes(x = Neighborhood, y = LotAreaHundr)) +
  geom_boxplot() +
  labs(title = "Lot Square Feet Across Neighborhoods",
       x = "Neighborhood",
       y = "Lot Square Feet (in hundreds)") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))

```



Additional Data for Model 1: This is the code for the sales price by the lot areas. The code includes a cv press details and then the simple linear model.

```

proc glmselect data = housing3;
/* where LotArea not in (63887,40094); */
model lSalePrice = llotArea / selection= Stepwise(stop = cv) cvmethod = random(10) CVDETAILS stats = adjrsq;
run;

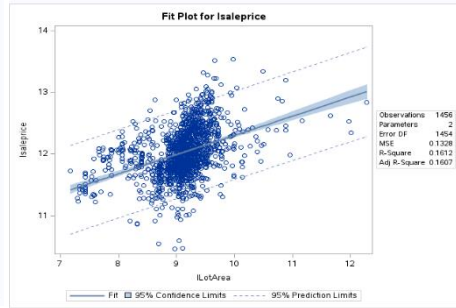
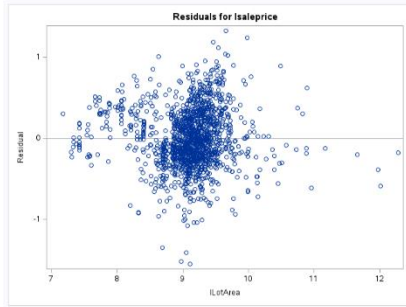
```

```

/* Regression Code to get plots for SLR model */
proc reg data = housing3;

```

model lsaleprice = llotarea;
run;



The GLMSELECT Procedure

Stepwise Selection Summary						
Step	Effect Entered	Effect Removed	Number of Effects in Model	Adjusted R-Square	SBC	CV PRESS
0	Intercept		1	0.0000	-2877.8359	230.3871
1	lLotArea		2	0.1607*	-2828.9558*	193.4870*

*Optimal Value of Criterion

Selection stopped because all effects are in the first model.

The GLMSELECT Procedure
Selected Model

The selected model is the model at the last step (Step 1).

Effects: Intercept lLotArea

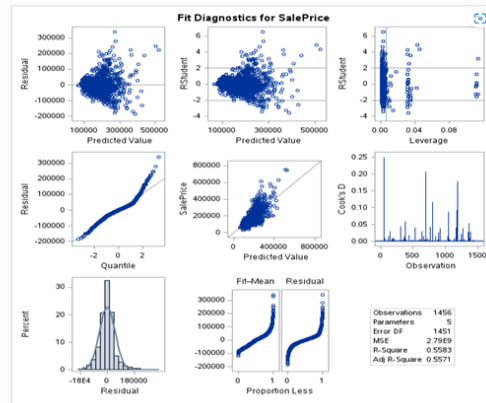
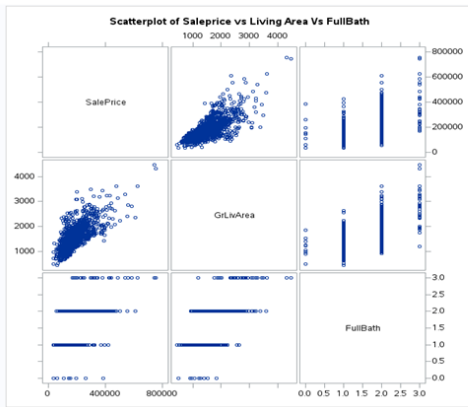
Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	1	37.10263	37.10263	279.82
Error	1454	183.14455	0.12624	
Corrected Total	1455	220.27468		

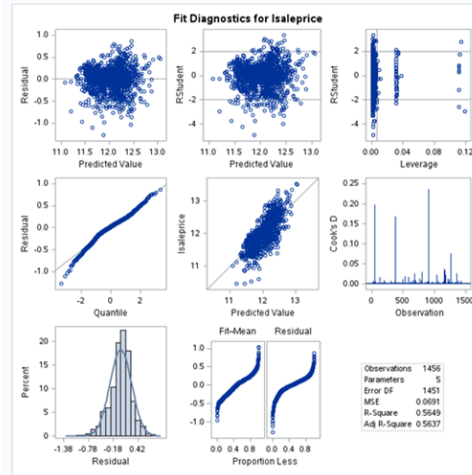
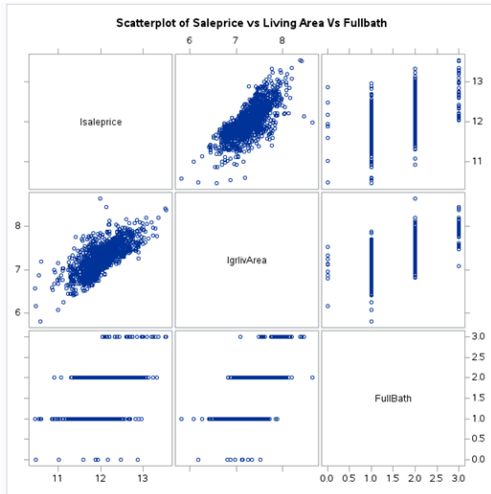
Root MSE	0.36487
Dependent Mean	12.02455
R-Square	0.1612
Adj R-Sq	0.1607
AIC	-4479.02261
AEC	-1478.11625
SBC	-2828.95582
CV PRESS	193.48697

Code for Model 2: Includes CV Press and output.

```
/* Backward Selection for 2nd Model */ proc glmselect data = housing3; Class Fullbath; model lSalePrice = lGrLivArea Fullbath / selection= backward(stop = cv) cvmethod = random(5) CVDETAILS stats = adjrsq; /* selection = stepwise(stop = SL SLE = 0.05 SLS =0.05) STATS=adjrsq; */ run;
```

```
/*Proc glm model to fit 2nd Model*/  
proc glm data = housing3 plots= all;  
Class Fullbath;  
model lSalePrice = lGrLivArea Fullbath / solution clparm;  
run;
```





The GLMSELECT Procedure

Backward Selection Summary

Step	Effect	Number of Effects	Number of Parameters	Adjusted R-Square	SBC	CV PRESS
0		1	1	0.5637*	-3890.2586*	101.2437*
1	lgrivArea	2	2	0.5637*	-3890.2586*	101.2437*
2	FullBath	3	3	0.5637*	-3890.2586*	101.2437*

Selection stopped at a local minimum of the cross validation PRESS.

Step Details

Candidate For Removal	Effect	Candidate CV PRESS	Compare CV PRESS
FullBath	FullBath	100.2716	101.2437

The GLMSELECT Procedure Selected Model

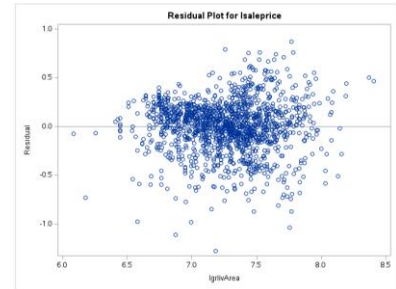
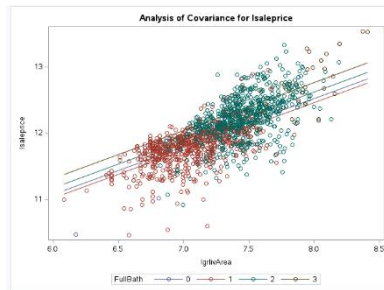
The selected model is the model at the last step (Step 0).

Effects: Intercept lgrivArea FullBath

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value
Model	4	130.07651	32.51913	470.91
Error	1451	100.19017	0.06906	
Corrected Total	1455	230.27468		

Root MSE 0.26276
Dependent Mean 12.0493
R-Square 0.5649
Adj R-Sq 0.5637
AIC -2428.97662
AICC -2428.91763
SBC -3890.2586
CV PRESS 101.24399



Parameter estimates for model 2:

Parameter	Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits	
Intercept	6.979739075	0.22497561	31.02	<.0001	6.538426854	7.421051297
lgrivArea	0.723070081	0.02794304	25.88	<.0001	0.668257016	0.777883147
FullBath 0	-0.242434141	0.10202764	-2.38	0.0176	-0.442571579	-0.042296703
FullBath 1	-0.299228013	0.05304306	-5.64	<.0001	-0.403277260	-0.195178737
FullBath 2	-0.142160707	0.04898498	-2.90	0.0038	-0.238249861	-0.048071754
FullBath 3	0.000000000

Code for Model 3:

Stepwise Selection For custom model

proc glmselect data = housing3;

Class MsSubclass Utilities LotConfig LandSlope Neighborhood Condition1 Condition2 BldgType HouseStyle OverallQual OverallCond RoofStyle RoofMatl Exterior1st Exterior2nd MasVnrType ExterQual ExterCond Foundation BsmtQual BsmtCond BsmtExposure BsmtFinType1 BsmtFinType2 Heating HeatingQC CentralAir Electrical KitchenQual Functional FireplaceQu GarageType GarageFinish GarageQual GarageCond PavedDrive PoolQC Fence MiscFeature SaleType SaleCondition;

model lSalePrice = MSSubClass LotFrontage lLotArea Utilities Lotconfig LandSlope Neighborhood Condition1 Condition2 BldgType HouseStyle OverallQual OverallCond YearBuilt YearRemodAdd RoofStyle RoofMatl Exterior1st Exterior2nd MasVnrType MasVnrArea ExterQual ExterCond Foundation BsmtQual BsmtCond BsmtExposure BsmtFinType1 BsmtFinType2 BsmtUnfSF TotalBsmtSF Heating HeatingQC CentralAir Electrical fstFlrSF scndFlrSF

```

LowQualFinSF IGrLivArea BsmtFullBath BsmtHalfBath FullBath HalfBath BedroomAbvGr KitchenAbvGr KitchenQual
TotRmsAbvGrd Functional Fireplaces FireplaceQu GarageType GarageYrBlt GarageFinish GarageCars GarageArea
GarageQual GarageCond PavedDrive WoodDeckSF OpenPorchSF EnclosedPorch SsnPorch3 ScreenPorch PoolArea
PoolQC Fence MiscFeature MiscVal MoSold YrSold SaleType SaleCondition / selection=stepwise (stop = cv)
cvmethod = random(10) CVDETAILS stats = adjrsq;

```

```

/* selection=backward (stop = cv) cvmethod = random(10) CVDETAILS stats = adjrsq; */ /* selection =
backward(stop = SL SLS = .01) stats = adjrsq; */

```

```
run;
```

Building and fitting custom model with proc glm

```

proc glm data=housing3 alpha=0.05 plots=all; class Neighborhood OverallQual MsSubclass GarageCars ;
model lSalePrice = IGrLivArea OverallQual MsSubclass GarageCars YearRemodAdd Neighborhood / solution clparm;
run;

```

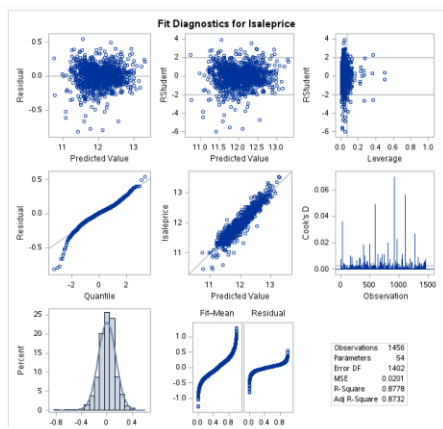
Proc Glmselect on custom model to get the CV press score

```

proc glmselect data=housing3;
class Neighborhood OverallQual MsSubclass GarageCars ;
model lSalePrice = IGrLivArea OverallQual MsSubclass GarageCars YearRemodAdd Neighborhood /
selection=stepwise (stop = cv) cvmethod = random(10) CVDETAILS stats = adjrsq;
run;

```

Plots and Parameter Estimates for the custom model 3:



The GLMSELECT Procedure
Selected Model

The selected model is the model at the last step (Step 6).

Effects: Intercept IgrLivArea OverallQual MsSubClass GarageCars YearRemodAdd Neighborhood

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	53	202.13548	3.81368	190.02
Error	1402	28.13922	0.02007	
Corrected Total	1455	230.27468		

Root MSE	0.14197
Dependent Mean	12.02455
R-Square	0.8778
Adj R-Sq	0.8732
AIC	-4179.78917
AICC	-4175.38917
SBC	-5352.46298
CV PRESS	31.53895

Cross Validation Details				
Index	Observations			CV PRESS
	Fitted	Left Out		
1	1322	134	2.0810	
2	1296	158	3.8217	
3	1309	147	5.1021	
4	1307	149	3.4346	
5	1313	143	2.5285	
6	1299	157	3.0733	
7	1312	144	2.9239	
8	1323	133	2.7120	
9	1310	146	3.5130	
10	1311	145	2.6481	
Total			31.5390	

Final Custom Model Coefficients:

Parameter	Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits
Intercept	4.280743248	B 0.52245785	8.19	<.0001	3.255859605 5.30626591
IgrflArea	0.486418591	B 0.02047039	24.25	<.0001	0.456282705 0.536574478
OverallQual 1	-0.916865062	B 0.15065716	-6.09	<.0001	-1.212402816 -0.621327308
OverallQual 2	-0.965093114	B 0.09622808	-10.03	<.0001	-1.153855727 -0.776330501
OverallQual 3	-0.809112472	B 0.05374896	-15.05	<.0001	-0.914549616 -0.703675428
OverallQual 4	-0.640796589	B 0.04512535	-14.20	<.0001	-0.729317083 -0.552276114
OverallQual 5	-0.562658965	B 0.04264835	-13.19	<.0001	-0.646314521 -0.478996469
OverallQual 6	-0.500773143	B 0.04183898	-11.97	<.0001	-0.582846883 -0.418699403
OverallQual 7	-0.425310777	B 0.04032303	-10.55	<.0001	-0.504410742 -0.348210812
OverallQual 8	-0.328527673	B 0.03884786	-8.46	<.0001	-0.404733866 -0.252321480
OverallQual 9	-0.145029325	B 0.04238874	-3.42	0.0008	-0.228181507 -0.061877144
OverallQual 10	0.000000000	B
MSubClass 20	0.124482872	B 0.03019460	4.12	<.0001	0.065251413 0.183714331
MSubClass 30	0.017673844	B 0.03303850	0.53	0.5927	-0.047132449 0.082480136
MSubClass 40	0.043169378	B 0.07655780	0.56	0.5729	-0.107010808 0.193349664
MSubClass 45	0.053570117	B 0.05050201	1.06	0.2890	-0.045497523 0.152837758
MSubClass 50	-0.001424324	B 0.02928683	-0.05	0.9612	-0.058835818 0.055987170
MSubClass 60	0.058997913	B 0.03138839	1.91	0.0564	-0.001636129 0.121431955
MSubClass 70	-0.045014524	B 0.03296705	-1.36	0.1726	-0.109723815 0.019694768
MSubClass 75	-0.040022595	B 0.04510565	-0.89	0.3751	-0.128504432 0.048459242
MSubClass 80	0.114089289	B 0.03520347	3.24	0.0012	0.045032148 0.183146430
MSubClass 85	0.180177468	B 0.04432706	4.06	<.0001	0.093222966 0.267131970
MSubClass 90	-0.035811811	B 0.03438749	-1.04	0.3007	-0.102998416 0.031836055
MSubClass 120	0.074112962	B 0.03658759	2.03	0.0430	0.002340642 0.145865281
MSubClass 180	-0.123874888	B 0.03988541	-3.11	0.0019	-0.202116410 -0.045633367
MSubClass 180	0.112673690	B 0.06346229	1.78	0.0780	-0.011817579 0.237164960
MSubClass 190	0.000000000	B
GarageCars 0	-0.336307896	B 0.06868764	-5.03	<.0001	-0.467496504 -0.205119289
GarageCars 1	-0.229013060	B 0.06586194	-3.48	0.0005	-0.358211628 -0.099814492
GarageCars 2	-0.182450281	B 0.06522846	-2.80	0.0062	-0.310406177 -0.054494385
GarageCars 3	-0.112052432	B 0.06677217	-1.68	0.0935	-0.243036554 0.018931690
GarageCars 4	0.000000000	B
YearRemodAdd	0.002476537	B 0.00024786	9.99	<.0001	0.001990320 0.002962753
Neighborhood Blmngtn	-0.193593512	B 0.05764461	-3.36	0.0008	-0.306672486 -0.080514539
Neighborhood Blueste	-0.140708997	B 0.11214701	-1.25	0.2098	-0.360703022 0.079285028
Neighborhood BrDale	-0.240913052	B 0.08222150	-3.87	0.0001	-0.362970332 -0.118855772
Neighborhood BrkSide	-0.172058189	B 0.04985307	-3.45	0.0006	-0.269852837 -0.074263541
Neighborhood ClearCr	-0.021640270	B 0.05120345	-0.42	0.6726	-0.122083892 0.078803352
Neighborhood CollgCr	-0.145262395	B 0.04497194	-3.23	0.0013	-0.233481944 -0.057042846
Neighborhood Crawfor	-0.013758384	B 0.04917083	-0.28	0.7797	-0.110214315 0.082697547
Neighborhood Edwards	-0.234898557	B 0.04671790	-5.03	<.0001	-0.326543080 -0.143254035
Neighborhood Gilbert	-0.180440449	B 0.04678802	-3.86	0.0001	-0.272222525 -0.086658372
Neighborhood IDOTRR	-0.345221880	B 0.05194359	-6.65	<.0001	-0.447117420 -0.243326341
Neighborhood Meadow	-0.260048959	B 0.06314308	-4.12	<.0001	-0.383914045 -0.136183873
Neighborhood Mitchal	-0.167281571	B 0.04810205	-3.48	0.0005	-0.261621321 -0.072901821
Neighborhood NAmes	-0.186076386	B 0.04482673	-4.15	<.0001	-0.274011049 -0.098141683
Neighborhood NPKVlll	-0.122759836	B 0.06599543	-1.86	0.0631	-0.252220273 0.006700601
Neighborhood NWAmes	-0.181263362	B 0.04656076	-3.89	0.0001	-0.272599624 -0.089927100
Neighborhood NoRidge	-0.031855827	B 0.04978914	-0.64	0.5224	-0.129625164 0.065813310
Neighborhood NrtidgHt	-0.055224300	B 0.04718702	-1.17	0.2421	-0.147789077 0.037340477
Neighborhood OldTown	-0.309081793	B 0.04758637	-6.50	<.0001	-0.402429661 -0.215733625
Neighborhood SW SU	-0.209232442	B 0.05487654	-3.83	0.0001	-0.316489095 -0.101975788
Neighborhood Sawyer	-0.205185796	B 0.04702151	-4.36	<.0001	-0.297425901 -0.112945690