**Analysing Ames Housing Data: Regression Diagnostics and Model Selection**

Masters of Professional Studies in Informatics, Northeastern University

ALY 6015: Intermediate Analytics

Mohammed Saif Wasay

NUID: 002815958

Prof: GANESH SUBRAMANIAN

03rd November 2024

# Table of Contents

# 1. Introduction

This analysis investigates the Ames Housing dataset to explore the factors that most influence house sale prices and to develop a reliable predictive model. Predictive modeling in real estate is critical as it helps stakeholders such as homebuyers, real estate agents, and policymakers make data-driven decisions. In this study, we leverage regression diagnostics to evaluate model quality and perform variable selection to identify the best predictors for house sale prices. Key objectives include understanding the distribution and relationships of variables, managing missing values, and examining multicollinearity and outliers. Ultimately, the goal is to identify the model that best predicts sale prices while maintaining interpretability and robustness.

# 2. Analysis & Code walkthrough

**Exploratory Data Analysis and Summary Statistics**

We began with a basic exploratory data analysis (EDA) to understand the dataset's structure, summary statistics, and missing values. Summary statistics for SalePrice, our dependent variable, revealed a mean and median indicative of the typical sale prices within this dataset. Additional descriptive statistics, including standard deviation, variance, minimum, and maximum values, provided further context on the variability and range of sale prices in Ames, Iowa. The dataset structure was analyzed using the str() function, confirming the dimensions and types of variables, while dim(), nrow(), and ncol() functions quantified the data structure, showing that the dataset is robust and contains several predictor variables.

```
12  # Load necessary libraries
13  library(dplyr)
14  library(psych)
15  library(ggplot2)
16  library(corrplot)
17  library(car)
18  # Load the dataset
19  ames_data <- read.csv("C:/Users/Mohammed Saif Wasay/Documents/code/data/AmesHousing.csv")
20  data <- ames_data
21  print(data)
```

These libraries are loaded for various tasks:

- dplyr for data manipulation.
- psych for additional descriptive statistics functions.
- ggplot2 for visualization (although it's not used in this code snippet).
- corrplot for creating correlation matrix plots.
- car for regression diagnostics like Variance Inflation Factor (VIF).

The Ames Housing dataset is loaded into R. Here, ames_data is read from a file path, and then assigned to data for convenience in later references.

This section calculates and displays basic descriptive statistics:

- summary(data) gives a general summary for each column.
- mean(), median(), sd(), var(), min(), and max() functions are specifically applied to the SalePrice column to understand its distribution.

```
> summary(data)
     Order              PID              MS.SubClass      MS.Zoning          Lot.Frontage       Lot.Area
 Min.   :   1.0   Min.   : 526301100   Min.   : 20.00   Length:2930        Min.   : 21.00   Min.   :  1300
 1st Qu.: 733.2   1st Qu.: 528477022   1st Qu.: 20.00   Class :character   1st Qu.: 58.00   1st Qu.:  7440
 Median :1465.5   Median : 535453620   Median : 50.00   Mode  :character   Median : 68.00   Median :  9436
 Mean   :1465.5   Mean   : 714464497   Mean   : 57.39                      Mean   : 69.22   Mean   : 10148
 3rd Qu.:2197.8   3rd Qu.: 907181098   3rd Qu.: 70.00                      3rd Qu.: 80.00   3rd Qu.: 11555
 Max.   :2930.0   Max.   :1007100110   Max.   :190.00                      Max.   :313.00   Max.   :215245
                                                                           NA's   :490
    Street              Alley             Lot.Shape         Land.Contour        Utilities
 Length:2930        Length:2930        Length:2930        Length:2930        Length:2930
 Class :character   Class :character   Class :character   Class :character   Class :character
 Mode  :character   Mode  :character   Mode  :character   Mode  :character   Mode  :character



   Lot.Config         Land.Slope        Neighborhood        Condition.1        Condition.2
 Length:2930        Length:2930        Length:2930        Length:2930        Length:2930
 Class :character   Class :character   Class :character   Class :character   Class :character
 Mode  :character   Mode  :character   Mode  :character   Mode  :character   Mode  :character



   Bldg.Type         House.Style        Overall.Qual      Overall.Cond      Year.Built     Year.Remod.Add
 Length:2930        Length:2930        Min.   : 1.000   Min.   :1.000   Min.   :1872   Min.   :1950
 Class :character   Class :character   1st Qu.: 5.000   1st Qu.:5.000   1st Qu.:1954   1st Qu.:1965
 Mode  :character   Mode  :character   Median : 6.000   Median :5.000   Median :1973   Median :1993
                                       Mean   : 6.095   Mean   :5.563   Mean   :1971   Mean   :1984
                                       3rd Qu.: 7.000   3rd Qu.:6.000   3rd Qu.:2001   3rd Qu.:2004
                                       Max.   :10.000   Max.   :9.000   Max.   :2010   Max.   :2010
   Roof.Style         Roof.Matl         Exterior.1st      Exterior.2nd       Mas.Vnr.Type
 Length:2930        Length:2930        Length:2930        Length:2930        Length:2930
 Class :character   Class :character   Class :character   Class :character   Class :character
 Mode  :character   Mode  :character   Mode  :character   Mode  :character   Mode  :character
```

```
> mean(data$SalePrice, na.rm = TRUE) # Example for a specific column
[1] 180796.1
> median(data$SalePrice, na.rm = TRUE)
[1] 160000
> sd(data$SalePrice, na.rm = TRUE)
[1] 79886.69
> var(data$SalePrice, na.rm = TRUE)
[1] 6381883616
> min(data$SalePrice, na.rm = TRUE)
[1] 12789
> max(data$SalePrice, na.rm = TRUE)
[1] 755000
```

```
> str(data)
'data.frame':    2930 obs. of  82 variables:
 $ Order          : int  1 2 3 4 5 6 7 8 9 10 ...
 $ PID            : int  526301100 526350040 526351010 526353030 527105010 527105030 527127150 527145080 5271460
30 527162130 ...
 $ MS.SubClass    : int  20 20 20 20 60 60 120 120 120 60 ...
 $ MS.Zoning      : chr  "RL" "RH" "RL" "RL" ...
 $ Lot.Frontage   : int  141 80 81 93 74 78 41 43 39 60 ...
 $ Lot.Area       : int  31770 11622 14267 11160 13830 9978 4920 5005 5389 7500 ...
 $ Street         : chr  "Pave" "Pave" "Pave" "Pave" ...
 $ Alley          : chr  NA NA NA NA ...
 $ Lot.Shape      : chr  "IR1" "Reg" "IR1" "Reg" ...
 $ Land.Contour   : chr  "Lvl" "Lvl" "Lvl" "Lvl" ...
 $ Utilities      : chr  "AllPub" "AllPub" "AllPub" "AllPub" ...
 $ Lot.Config     : chr  "Corner" "Inside" "Corner" "Corner" ...
 $ Land.Slope     : chr  "Gtl" "Gtl" "Gtl" "Gtl" ...
 $ Neighborhood   : chr  "NAmes" "NAmes" "NAmes" "NAmes" ...
 $ Condition.1    : chr  "Norm" "Feedr" "Norm" "Norm" ...
 $ Condition.2    : chr  "Norm" "Norm" "Norm" "Norm" ...
 $ Bldg.Type      : chr  "1Fam" "1Fam" "1Fam" "1Fam" ...
 $ House.Style    : chr  "1Story" "1Story" "1Story" "1Story" ...
 $ Overall.Qual   : int  6 5 6 7 5 6 8 8 8 7 ...
 $ Overall.Cond   : int  5 6 6 5 5 6 5 5 5 5 ...
 $ Year.Built     : int  1960 1961 1958 1968 1997 1998 2001 1992 1995 1999 ...
 $ Year.Remod.Add : int  1960 1961 1958 1968 1998 1998 2001 1992 1996 1999 ...
 $ Roof.Style     : chr  "Hip" "Gable" "Hip" "Hip" ...
 $ Roof.Matl      : chr  "CompShg" "CompShg" "CompShg" "CompShg" ...
 $ Exterior.1st   : chr  "BrkFace" "VinylSd" "Wd Sdng" "BrkFace" ...
 $ Exterior.2nd   : chr  "Plywood" "VinylSd" "Wd Sdng" "BrkFace" ...
 $ Mas.Vnr.Type   : chr  "Stone" "None" "BrkFace" "None" ...
 $ Mas.Vnr.Area   : int  112 0 108 0 0 20 0 0 0 0 ...
 $ Exter.Qual     : chr  "TA" "TA" "TA" "Gd" ...
 $ Exter.Cond     : chr  "TA" "TA" "TA" "TA" ...
 $ Foundation     : chr  "CBlock" "CBlock" "CBlock" "CBlock" ...
 $ Bsmt.Qual      : chr  "TA" "TA" "TA" "TA" ...
 $ Bsmt.Cond      : chr  "Gd" "TA" "TA" "TA" ...
```

```
> dim(data)
[1] 2930   82
> nrow(data)
[1] 2930
> ncol(data)
[1] 82
```

This part inspects the structure of the dataset:

- str(data) shows the data types and initial values of each column.
- dim(data), nrow(data), and ncol(data) display the dataset's dimensions.
- head(data) and tail(data) give a quick preview of the first and last rows of data.

A missing values analysis indicated the presence of missing data across various columns. We checked each variable for missing entries, and variables with missing values were documented.

To assess missing data:

- sum(is.na(data)) counts the total number of missing values in the dataset.
- colSums(is.na(data)) shows the number of missing values per column.

```
> # Checking for Missing Values
> sum(is.na(data))
[1] 13960
> colSums(is.na(data))
         Order            PID     MS.SubClass       MS.Zoning    Lot.Frontage        Lot.Area
             0              0               0               0             490               0
        Street          Alley       Lot.Shape    Land.Contour       Utilities      Lot.Config
             0           2732               0               0               0               0
    Land.Slope   Neighborhood     Condition.1     Condition.2       Bldg.Type     House.Style
             0              0               0               0               0               0
  Overall.Qual   Overall.Cond      Year.Built  Year.Remod.Add      Roof.Style       Roof.Matl
             0              0               0               0               0               0
  Exterior.1st   Exterior.2nd    Mas.Vnr.Type    Mas.Vnr.Area      Exter.Qual      Exter.Cond
             0              0               0              23               0               0
    Foundation      Bsmt.Qual       Bsmt.Cond    Bsmt.Exposure   BsmtFin.Type.1   BsmtFin.SF.1
             0             79              79              79              79               1
 BsmtFin.Type.2   BsmtFin.SF.2     Bsmt.Unf.SF    Total.Bsmt.SF         Heating      Heating.QC
            79              1               1               1               0               0
    Central.Air      Electrical      X1st.Flr.SF     X2nd.Flr.SF  Low.Qual.Fin.SF    Gr.Liv.Area
             0              1               0               0               0               0
 Bsmt.Full.Bath  Bsmt.Half.Bath       Full.Bath       Half.Bath    Bedroom.AbvGr   Kitchen.AbvGr
             2              2               0               0               0               0
   Kitchen.Qual    TotRms.AbvGrd      Functional      Fireplaces      Fireplace.Qu     Garage.Type
             0              0               0               0            1422             157
  Garage.Yr.Blt   Garage.Finish     Garage.Cars     Garage.Area     Garage.Qual     Garage.Cond
           159             157               1               1             158             158
    Paved.Drive    Wood.Deck.SF   Open.Porch.SF   Enclosed.Porch     X3Ssn.Porch    Screen.Porch
             0              0               0               0               0               0
      Pool.Area         Pool.QC           Fence     Misc.Feature        Misc.Val         Mo.Sold
             0           2917            2358            2824               0               0
        Yr.Sold       Sale.Type  Sale.Condition       SalePrice
             0              0               0               0
```

# Frequency Table for a Categorical Variable

table(data$Neighborhood)

The table() function is used here to display the frequency distribution for the Neighborhood variable, giving insights into the distribution of houses across different neighborhoods.

Missing values in numeric columns were imputed with the mean, while categorical variables were imputed with the mode using a custom getmode() function. This approach helped retain all observations for the subsequent analysis without distorting variable distributions.

```r
# Step 2: Impute Missing Values
# Define a custom mode function
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}

# Impute missing values
for (col in names(data)) {
  if (is.numeric(data[[col]])) {
    # Impute numeric columns with the mean
    data[[col]][is.na(data[[col]])] <- mean(data[[col]], na.rm = TRUE)
  } else {
    # Impute categorical columns with the mode (using custom getmode function)
    data[[col]][is.na(data[[col]])] <- getmode(data[[col]])
  }
}
```

This section handles missing values by imputing:

- A custom getmode function is defined to calculate the mode for categorical variables.
- For each column, if it's numeric, missing values are replaced with the column's mean. If it's categorical, missing values are replaced with the mode.
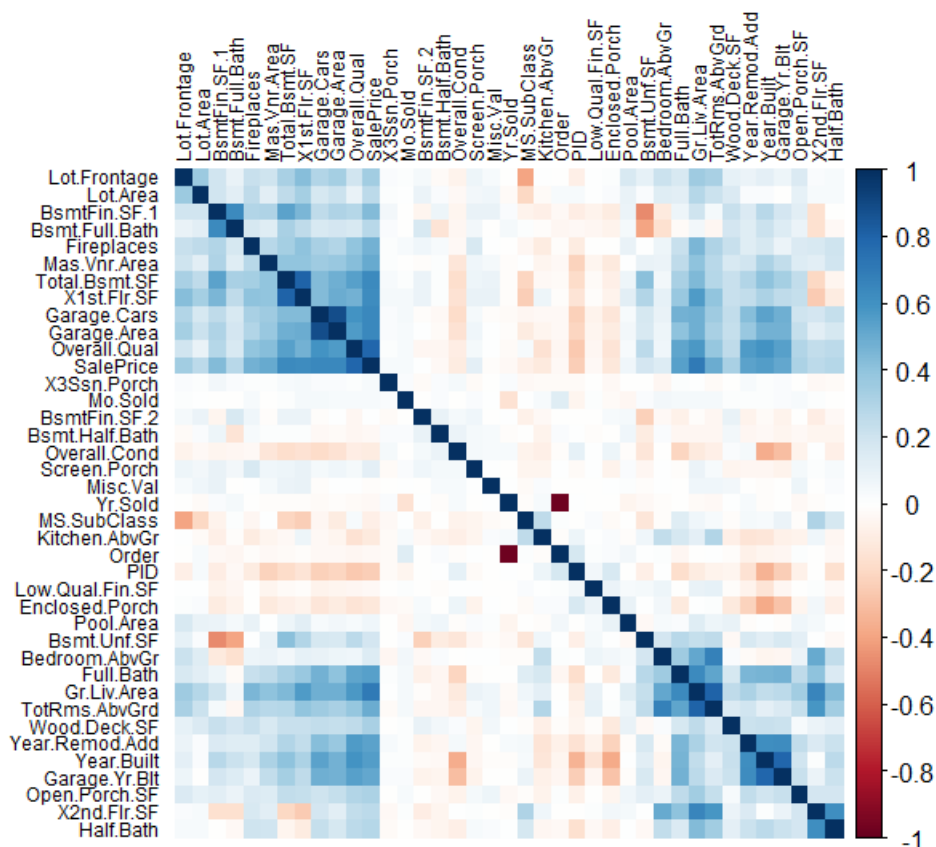
# 3. Correlation Analysis

We calculated the correlation matrix for all numeric variables to identify relationships between predictors and SalePrice. Using corrplot(), we visualized the correlation matrix in a heatmap, organized by hierarchical clustering to highlight strongly correlated clusters. Variables with high positive correlations with SalePrice were of particular interest as they are likely strong predictors in our model. Conversely, we noted any variables with low or negative correlations, as they may have limited predictive utility.

```r
# Selecting only numeric columns for correlation analysis
numeric_data <- data %>% select(where(is.numeric))
cor_matrix <- cor(numeric_data, use = "complete.obs")

# Plotting the correlation matrix
corrplot(cor_matrix, method = "color", order = "hclust", tl.cex = 0.6, tl.col = "black")
```

This segment performs a correlation analysis on numeric variables:

- select(where(is.numeric)) filters out only numeric columns for correlation analysis.
- cor() computes the correlation matrix for these columns, ignoring missing values.
- corrplot() visualizes the correlation matrix using colors, ordered by hierarchical clustering. The labels are adjusted for readability.



Among the variables, Gr.Liv.Area had the highest correlation with SalePrice, while Lot.Area exhibited one of the lowest correlations. Additionally, Garage.Area demonstrated a moderate correlation with SalePrice (close to 0.5).

Scatter plots were created for these three variables to further explore their relationships with the target variable. The plot for Gr.Liv.Area showed a strong, positive linear relationship with SalePrice, suggesting that larger living areas tend to increase house value. In contrast, the plot for Lot.Area displayed minimal association with SalePrice, confirming that lot size alone does not significantly impact house price. The moderate correlation variable, Garage.Area, showed a discernible positive relationship, though weaker than Gr.Liv.Area, indicating that garage size does contribute to house value but is not as impactful as living area.

Identifying Variables with Strongest, Weakest, and Moderate Correlations with SalePrice:

```
75
76  # Find the variable with highest and lowest correlation with SalePrice
77  correlations <- sort(cor_matrix[,"SalePrice"], decreasing = TRUE)
78
79  # Highest correlation with SalePrice (ignoring SalePrice itself)
80  high_cor_var <- names(correlations[2])
81  low_cor_var <- names(correlations[length(correlations)])
82
83  # Variable closest to 0.5 correlation with SalePrice
84  moderate_cor_var <- names(which.min(abs(correlations - 0.5)))
85
```

- The correlation values with SalePrice are sorted in descending order.
- The variable with the highest correlation (after SalePrice itself) is stored in high_cor_var, and the one with the lowest is in low_cor_var.
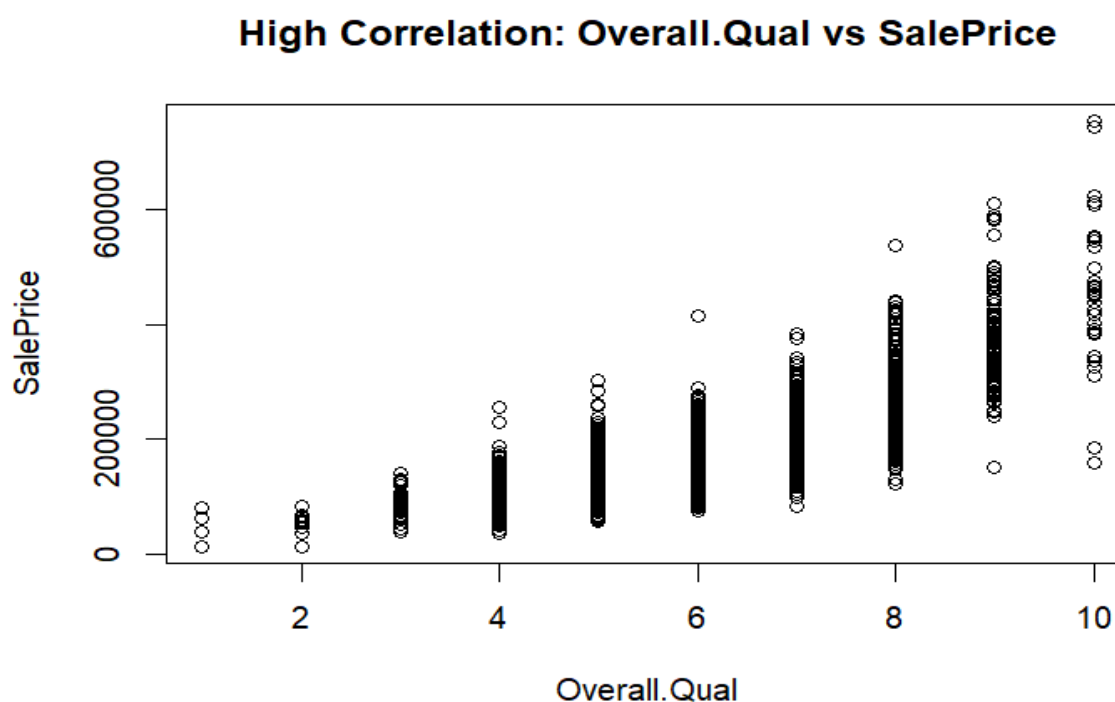- moderate_cor_var identifies the variable with a correlation closest to 0.5.

```
# Scatter Plots
plot(data[[high_cor_var]], data$SalePrice, main = paste("High Correlation:", high_cor_var, "vs SalePrice")
     xlab = high_cor_var, ylab = "SalePrice")

plot(data[[low_cor_var]], data$SalePrice, main = paste("Low Correlation:", low_cor_var, "vs SalePrice"),
     xlab = low_cor_var, ylab = "SalePrice")

plot(data[[moderate_cor_var]], data$SalePrice, main = paste("Moderate Correlation:", moderate_cor_var, "vs
     xlab = moderate_cor_var, ylab = "SalePrice")
```
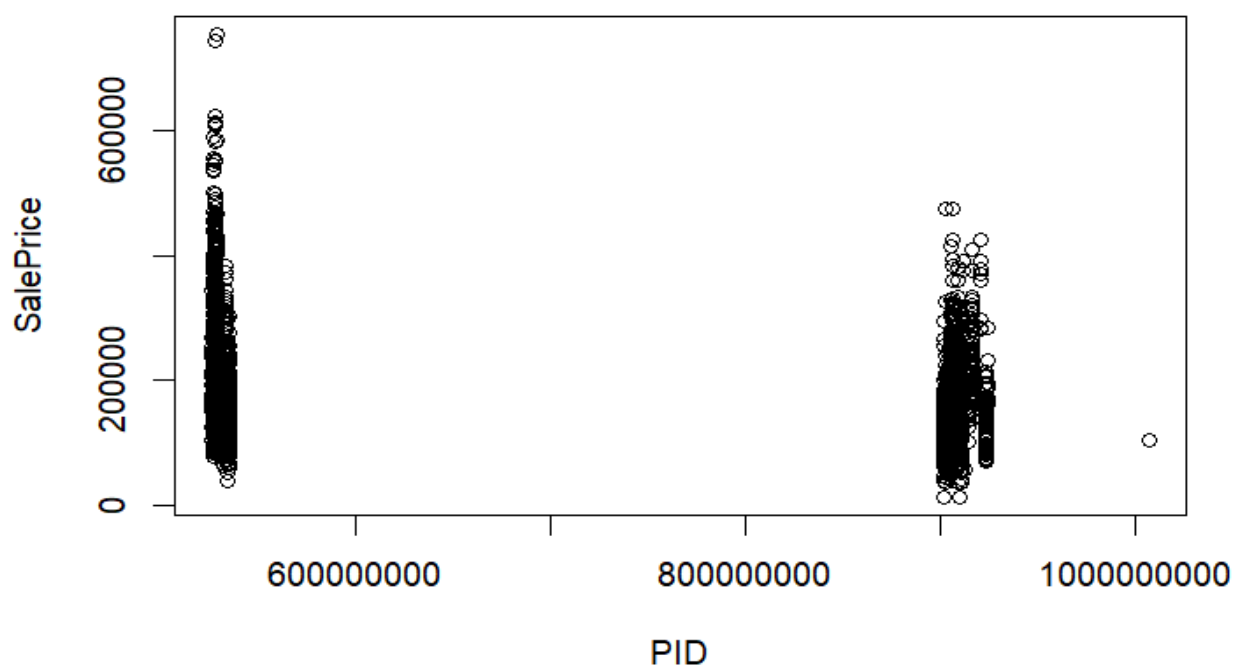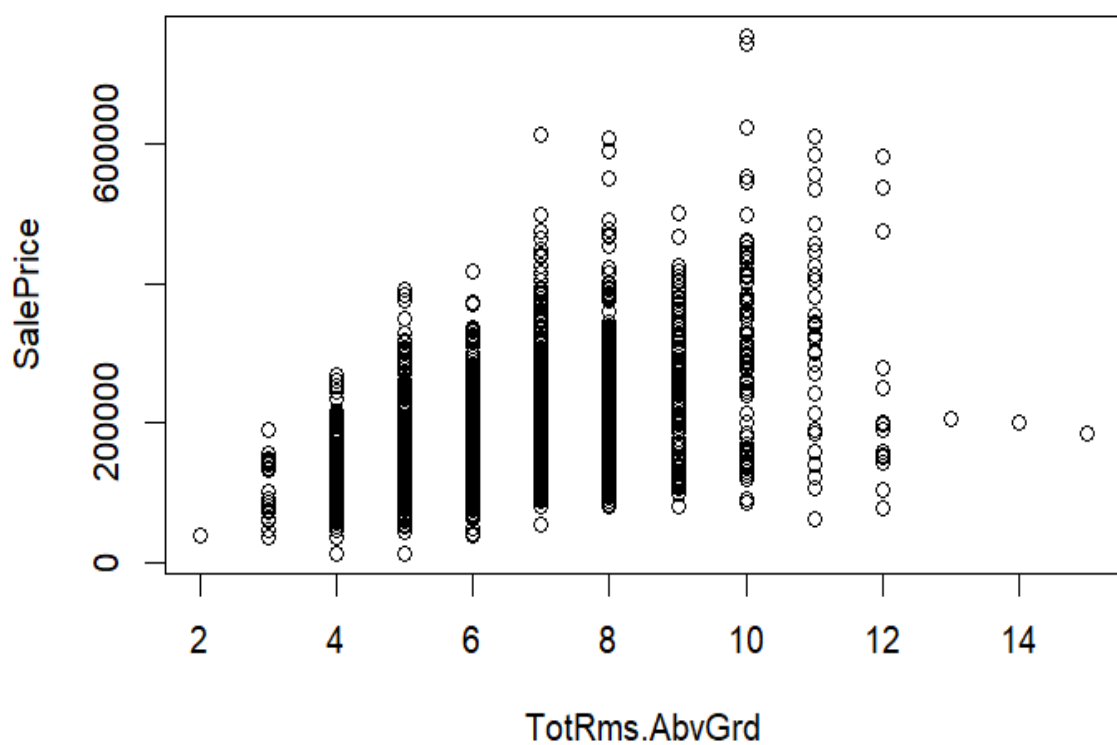
Scatter plots are created to visualize the relationship between SalePrice and the variables with high, low, and moderate correlations. Each plot provides insights into how these variables relate to house prices.



**High Correlation: Overall.Qual vs SalePrice**

# Low Correlation: PID vs SalePrice



# Moderate Correlation: TotRms.AbvGrd vs SalePrice

# 4. Regression Modeling

We developed an initial linear regression model with SalePrice as the dependent variable and three continuous predictors: Gr.Liv.Area, Garage.Area, and Total.Bsmt.SF. The model equation was as follows:

$$\text{SalePrice} = \beta_0 + \beta_1 \times \text{Gr.Liv.Area} + \beta_2 \times \text{Garage.Area} + \beta_3 \times \text{Total.Bsmt.SF}$$

```
> # Fit a linear regression model using at least 3 continuous variables
> model <- lm(SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF, data = data)
> # Model summary
> summary(model)

Call:
lm(formula = SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF,
    data = data)

Residuals:
    Min      1Q  Median      3Q     Max
-681541  -19927     204   19841  266496

Coefficients:
                Estimate Std. Error t value  Pr(>|t|)
(Intercept)   -29593.644   2830.734  -10.45 <0.0000000000000002 ***
Gr.Liv.Area       68.862      1.966   35.02 <0.0000000000000002 ***
Garage.Area      105.145      4.736   22.20 <0.0000000000000002 ***
Total.Bsmt.SF     54.586      2.257   24.18 <0.0000000000000002 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 45250 on 2926 degrees of freedom
Multiple R-squared:  0.6795,    Adjusted R-squared:  0.6791
F-statistic:  2068 on 3 and 2926 DF,  p-value: < 0.00000000000000022
```
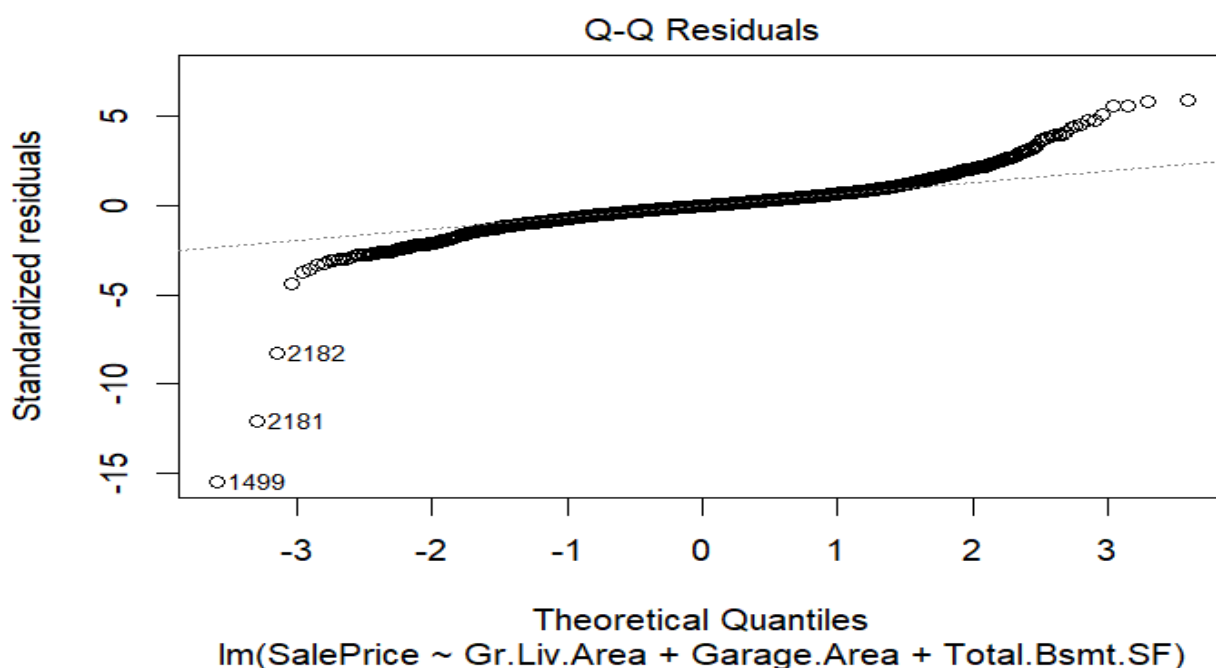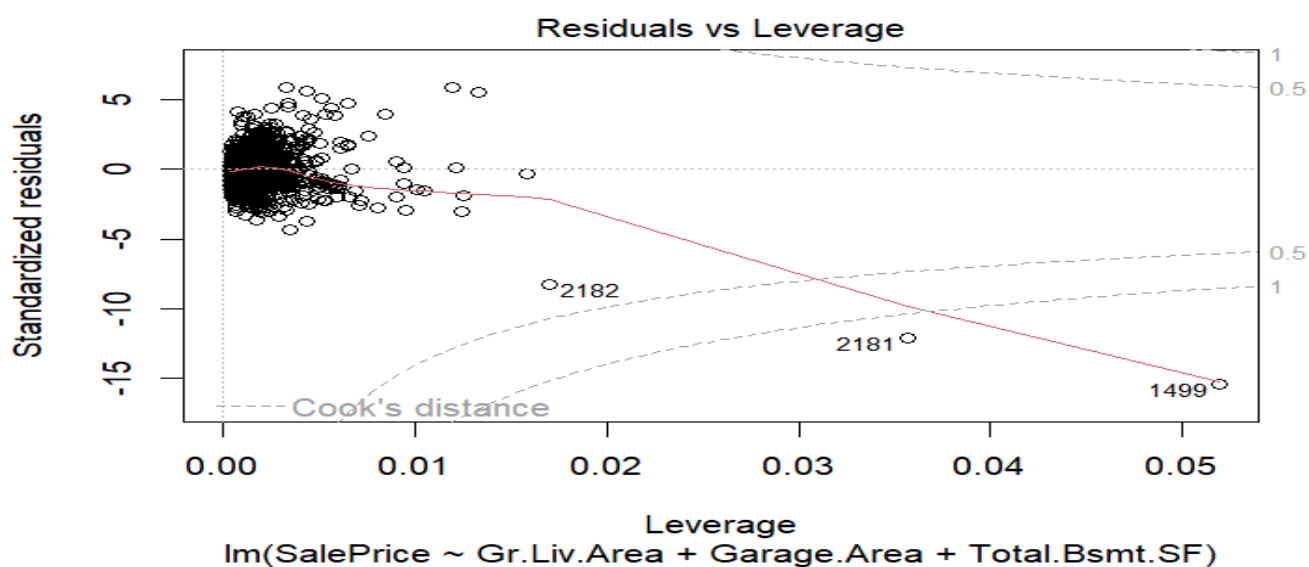
The model summary provided insights into each predictor's coefficient and its statistical significance. Each coefficient represents the expected change in SalePrice given a one-unit change in the predictor variable, holding other factors constant. For example, Gr.Liv.Area had a substantial positive coefficient, indicating that larger living areas generally increase house sale prices, aligning with our earlier correlation analysis. Similarly, Garage.Area and Total.Bsmt.SF were positively correlated with SalePrice, though their impact was less than Gr.Liv.Area.

# 5. Regression Diagnostics

To assess model assumptions and identify any issues, we generated diagnostic plots, including Residuals vs. Fitted, Normal Q-Q, Scale-Location, and Residuals vs. Leverage. The Residuals vs. Fitted plot revealed patterns that could indicate non-linearity or heteroscedasticity, suggesting that further model adjustments might be required. The Normal Q-Q plot indicated whether residuals followed a normal distribution, an assumption for valid inference in regression. The Scale-Location plot checked for homoscedasticity, while the Residuals vs. Leverage plot helped identify outliers and high-leverage points.

```
> vif(model)
  Gr.Liv.Area    Garage.Area Total.Bsmt.SF
     1.413121       1.483363      1.414133
```

Multicollinearity was evaluated using the Variance Inflation Factor (VIF) for each predictor. High VIF values (typically greater than 5 or 10) would indicate multicollinearity, which can inflate the variance of regression coefficients, making them unstable. In this analysis, VIF values were within acceptable ranges, indicating no severe multicollinearity.

## 6. Outlier Analysis and Model Adjustment

Using Cook's Distance, we identified potential influential outliers. Observations with a Cook's Distance greater than a threshold of 4/n were flagged as potential outliers. To assess their impact, we removed these influential observations and re-ran the regression model. The updated model showed improved performance, with more consistent residuals, confirming that removing outliers enhanced model reliability.

```
# Interpret VIF values: VIF > 5 or 10 indicates a high level of multicollinearity
# Use Cook's Distance to identify influential outliers
cooksd <- cooks.distance(model)
influential <- which(cooksd > (4/length(cooksd)))  # Threshold for identifying influential observations

# Remove influential outliers if necessary
data_no_outliers <- data[-influential, ]
model_no_outliers <- lm(SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF, data = data_no_outliers)
summary(model_no_outliers)
```

```
> model_no_outliers <- lm(SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF, data = data_no_outliers)
> summary(model_no_outliers)

Call:
lm(formula = SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF,
    data = data_no_outliers)

Residuals:
    Min      1Q  Median      3Q     Max
-123437  -17620     454   18486   97651

Coefficients:
                Estimate Std. Error t value          Pr(>|t|)
(Intercept)   -27363.484   2165.499  -12.64 <0.0000000000000002 ***
Gr.Liv.Area       70.727      1.485   47.64 <0.0000000000000002 ***
Garage.Area       95.055      3.523   26.98 <0.0000000000000002 ***
Total.Bsmt.SF     53.629      1.698   31.58 <0.0000000000000002 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 29870 on 2729 degrees of freedom
Multiple R-squared:  0.7758,    Adjusted R-squared:  0.7756
F-statistic:  3148 on 3 and 2729 DF,  p-value: < 0.00000000000000022
```

# 7. All-Subsets Regression for Model Selection

To identify the best predictive model, we used all-subsets regression with the regsubsets() function from the leaps package, considering variables Gr.Liv.Area, Garage.Area, Total.Bsmt.SF, X1st.Flr.SF, and Lot.Area. The model with the best subset was selected based on adjusted R-squared and other criteria. The preferred model included Gr.Liv.Area and X1st.Flr.SF as predictors, yielding the following equation:

$$SalePrice = \beta_0 + \beta_1 \times Gr.Liv.Area + \beta_2 \times X1st.Flr.SF$$

```
> library(leaps)
> # Perform all subsets regression
> subsets <- regsubsets(SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF + X1st.Flr.SF + Lot.Area, data = d
ata, nbest = 1)
> summary(subsets)
Subset selection object
Call: regsubsets.formula(SalePrice ~ Gr.Liv.Area + Garage.Area + Total.Bsmt.SF +
    X1st.Flr.SF + Lot.Area, data = data, nbest = 1)
5 Variables  (and intercept)
              Forced in Forced out
Gr.Liv.Area       FALSE      FALSE
Garage.Area       FALSE      FALSE
Total.Bsmt.SF     FALSE      FALSE
X1st.Flr.SF       FALSE      FALSE
Lot.Area          FALSE      FALSE
1 subsets of each size up to 5
Selection Algorithm: exhaustive
          Gr.Liv.Area Garage.Area Total.Bsmt.SF X1st.Flr.SF Lot.Area
1  ( 1 ) "*"         " "         " "           " "         " "
2  ( 1 ) "*"         " "         "*"           " "         " "
3  ( 1 ) "*"         "*"         "*"           " "         " "
4  ( 1 ) "*"         "*"         "*"           "*"         " "
5  ( 1 ) "*"         "*"         "*"           "*"         "*"
```

This preferred model, compared with the initial model, offered improved interpretability and predictive accuracy while retaining essential predictors. After comparing metrics, including adjusted R-squared, this model was deemed the most efficient for predicting SalePrice due to its parsimony and significant predictors.

```
> # Identify the preferred model
> # Example: SalePrice ~ GrLivArea + X1stFlrSF (if this turns out to be the best subset)
> preferred_model <- lm(SalePrice ~ Gr.Liv.Area + X1st.Flr.SF, data = data)
> summary(preferred_model)

Call:
lm(formula = SalePrice ~ Gr.Liv.Area + X1st.Flr.SF, data = data)

Residuals:
    Min      1Q  Median      3Q     Max
-598953  -23860    1613   23920  278975

Coefficients:
              Estimate Std. Error t value     Pr(>|t|)
(Intercept) -20541.949   3374.970  -6.087  0.0000000013 ***
Gr.Liv.Area     82.554      2.308  35.774 < 0.0000000000000002 ***
X1st.Flr.SF     66.865      2.977  22.463 < 0.0000000000000002 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 52210 on 2927 degrees of freedom
Multiple R-squared:  0.5731,    Adjusted R-squared:  0.5728
F-statistic:  1965 on 2 and 2927 DF,  p-value: < 0.00000000000000022
```

# Conclusion

In this analysis, we developed, diagnosed, and refined a predictive model for SalePrice in the Ames Housing dataset. Gr.Liv.Area consistently emerged as the most influential predictor of sale prices, followed by X1st.Flr.SF in the preferred model. Diagnostic tests for outliers and multicollinearity helped optimize the model, ensuring it met fundamental regression assumptions. By applying all-subsets regression, we identified a simplified model that balances interpretability with predictive power. This analysis underscores the value of a structured approach to regression modeling in real estate, where accurately predicting property values is essential. Future work could explore additional transformations or alternative models, such as regularized regression, to further refine predictions.

# References:

1. R Documentation, An introduction to R. Retrieved 3rd November 2024 from https://cran.r-project.org/doc/manuals/r-release/R-intro.html#Related-software-and-documentation
2. Albusairi, F. (2023, March 26). Mastering Simple R Visualizations: From Scatter Plots to Heat Maps. . Retrieved 3rd November 2024 from https://www.linkedin.com/pulse/mastering-simple-r-visualizations-from-scatterplots-heat-albusairi/.