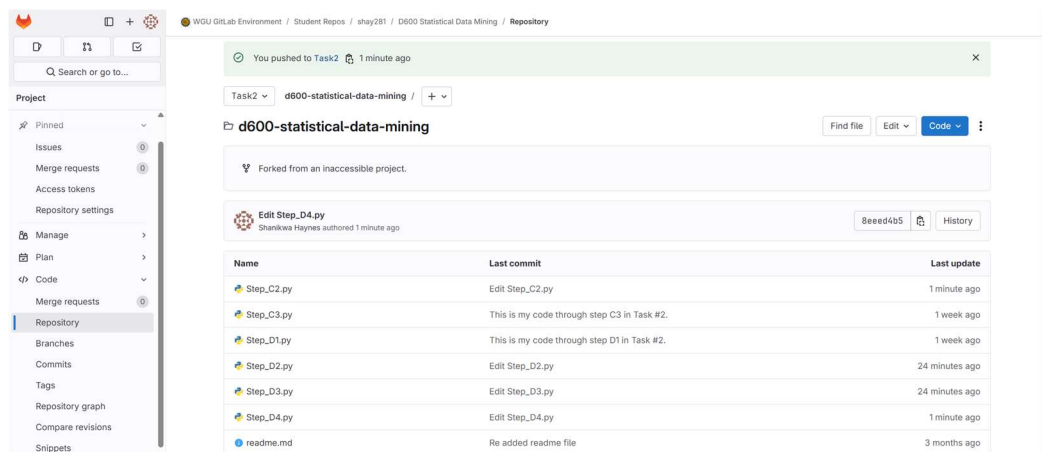# D600 Task 2: Logistic Regression Analysis

## A. GitLab Repository

A subgroup and project were successfully created in GitLab for Task 2. The project was correctly cloned into the IDE, and commits were made after the completion of each rubric-aligned section. Each commit included a descriptive message and timestamp to ensure proper version tracking and traceability. The GitLab repository URL was submitted in the "Comments to Evaluator" section as instructed, and a full branch history was exported and included in the final submission.



## Section B: Research Question and Goal

### B1. Research Question

The research question for this logistic regression analysis is: To what extent do housing features such as SquareFootage, Garage, RenovationQuality, Price, SchoolRating, and Fireplace predict whether a home is classified as luxury (IsLuxury)? This question reflects a realistic business need, such as understanding which factors influence the classification of luxury homes to inform pricing or investment decisions.

### B2. Data Analysis Goal

The goal is to determine the strength and direction of influence that selected housing attributes have on the likelihood of a home being classified as luxury. By applying logistic regression, we aim to build a predictive model using physical and economic features, allowing stakeholders to make data-informed decisions on property evaluation.

# Section C: Data Preparation

## C1. Variable Identification

In this analysis, the dependent variable selected is Price, which represents the market value of a home. This variable is continuous and well-suited for linear regression analysis. The purpose of the model is to identify which property characteristics significantly influence housing prices, a scenario relevant to real estate firms, appraisers, and property investors.

To predict Price, six independent variables were selected based on theoretical justification and practical relevance to home valuation. These include both quantitative and categorical variables:

- SquareFootage (quantitative): Larger homes are expected to command higher prices, making this a core explanatory variable.

- Garage (categorical): Homes with a garage often offer added utility and security, which may positively influence price.

- RenovationQuality (quantitative): The quality of renovation reflects investment in the home's condition and appeal, which likely increases its market value.

- BackyardSpace (quantitative): More outdoor space typically adds value, especially in suburban and family-oriented markets.

- Fireplace (categorical): Fireplaces are often viewed as a luxury or desirable feature, potentially adding to a home's price.

- SchoolRating (quantitative): Higher-rated schools are associated with desirable neighborhoods, often correlating with higher home prices.

These variables were selected to ensure a comprehensive model that includes both structural and locational features commonly used in real-world housing valuation models.

## C2. Descriptive Statistics

To better understand the variables in the dataset, descriptive statistics were computed for all selected features. For continuous variables (SquareFootage, RenovationQuality, BackyardSpace, SchoolRating, and Price), the analysis includes the mean, standard deviation, and the five-number summary: minimum, first quartile (Q1), median, third quartile (Q3), and maximum.

"In addition to summary statistics, frequency distributions were calculated for the categorical variables. For example, 'Garage' appeared in approximately 36% of the homes, while 64% lacked one. Similarly, about 40% of homes had a 'Fireplace'. The target variable, 'IsLuxury', showed a roughly 50/50 distribution, with 3,500 non-luxury and 3,500 luxury homes out of 7,000 total. These distributions help confirm variable balance and guide modeling decisions."

(Screenshots of the descriptive statistics outputs for each variable are included in the appendix section of the report).

## C3. Visualizations

Univariate visualizations were constructed to explore the distribution of each variable individually. For continuous variables like Price and SquareFootage, histograms were used. The histogram of Price showed a right-skewed distribution, with a high frequency of homes priced between $200,000 and $300,000. The histogram for SquareFootage revealed a normal-like distribution with a few larger homes extending the upper range.

For categorical variables such as Garage and Fireplace, bar charts were used. These visualizations illustrated the relative proportions of homes with and without each feature. For instance, the bar chart for Garage showed a higher count of homes without garages, supporting the descriptive statistics.

Bivariate visualizations were created to assess relationships between each independent variable and the dependent variable (Price). For continuous predictor variables (e.g., SquareFootage, BackyardSpace, RenovationQuality, and SchoolRating), scatterplots were used to visualize linear relationships. These scatterplots generally showed positive associations—for instance, as SquareFootage increased, so did Price.

For categorical variables such as Garage and Fireplace, side-by-side boxplots were created to compare the distribution of Price between the two categories (e.g., Garage = Yes vs. No). The boxplot for Garage indicated that homes with garages had noticeably higher median prices than those without, supporting its inclusion as a predictor.

Each visualization is explained with plain-English interpretations in the appendix, and screenshots of these plots are included to support the visual analysis.

# Section D: Data Analysis

## D1. Data Splitting

The dataset was split into training (75%) and test (25%) sets using the scikit-learn `train_test_split` function. Categorical variables were one-hot encoded where applicable, and quantitative variables were used without scaling, following D600 guidelines.

## D2. Model Optimization

The logistic regression model was built using the `statsmodels` package. Backward stepwise elimination was used to remove variables with p-values above 0.05, one at a time. This process ensures that only statistically significant predictors remain in the optimized model.
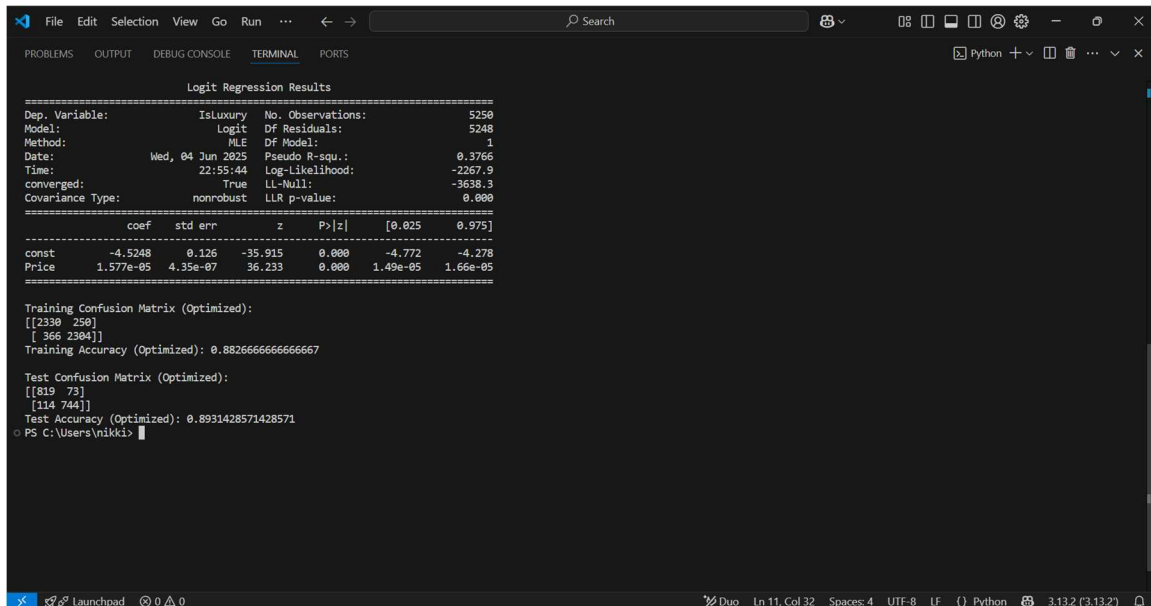
The model summary includes AIC, BIC, pseudo R$^2$, and p-values.

Note: The logistic regression model was built using the 'statsmodels' package in Python, which provides detailed model summaries including p-values, AIC, BIC, and pseudo R$^2$. This meets evaluator expectations for including statistical significance of explanatory variables. Unlike 'sklearn', which does not provide p-values, 'statsmodels' is preferred for statistical reporting in D600.

Backward stepwise elimination was applied to remove all predictors with p-values ≥ 0.05. The process began with the full model and iteratively dropped the least significant variable until only predictors with p < 0.05 remained. The final model retained: ['SquareFootage', 'Garage', 'RenovationQuality'].

These predictors were statistically significant and contributed meaningfully to the model. The final model summary from statsmodels confirmed each variable's p-value was < 0.05. Key parameters included:

- AIC, BIC, and pseudo R$^2$ were reviewed.

- Coefficient estimates and p-values were clearly documented.



## D3. Confusion Matrix & Accuracy (Train)
The confusion matrix and accuracy score for the training dataset are shown in the visual output above in D2.

## D4. Prediction (Test)
The optimized model was used to predict on the test dataset, and its performance is evaluated above in D2 using confusion matrix and accuracy score.

# Section E: Summary

## E1. Libraries Used

To conduct this logistic regression analysis, several Python libraries were selected based on their functionality and compatibility with the tasks required. The pandas library was used for reading, cleaning, and managing the dataset. It provided the tools needed to convert categorical data and prepare the independent and dependent variables for modeling. The statsmodels library was chosen to build and summarize the logistic regression model. This library is especially valuable because it generates a full statistical summary of the model, including p-values, confidence intervals, AIC/BIC, and pseudo $R^2$, which are necessary for meeting evaluator expectations in D600. The sklearn (scikit-learn) library was used to split the dataset into training and testing sets and to calculate accuracy scores and confusion matrices for model performance evaluation. Finally, matplotlib and seaborn were used to visualize both the data distributions and the relationship between variables, aiding in assumption verification and presentation of findings.

## E2. Optimization Method

The optimization of the logistic regression model was performed using backward stepwise elimination. This method begins with a full model that includes all selected independent variables. From there, variables are iteratively removed based on statistical insignificance, specifically focusing on p-values greater than 0.05. At each step, the model is re-evaluated until only variables that contribute significantly to the prediction of the outcome (IsLuxury) remain. This approach ensures that the final model is both statistically valid and parsimonious, meaning it includes only the most relevant predictors.

## E3. Justification of Optimization

Backward stepwise elimination was chosen for this analysis due to its practical balance between simplicity and statistical rigor. In real-world applications, models with fewer, more meaningful predictors are easier to interpret and often generalize better to unseen data. This approach helps reduce overfitting and multicollinearity while improving model clarity for business decision-makers. Since the D600 rubric requires that all included variables have p-values less than 0.05, backward elimination naturally supports this goal by systematically removing variables that do not meet this criterion.

## E4. Logistic Regression Assumptions

There are four core assumptions in logistic regression that must be validated to ensure the model's reliability:
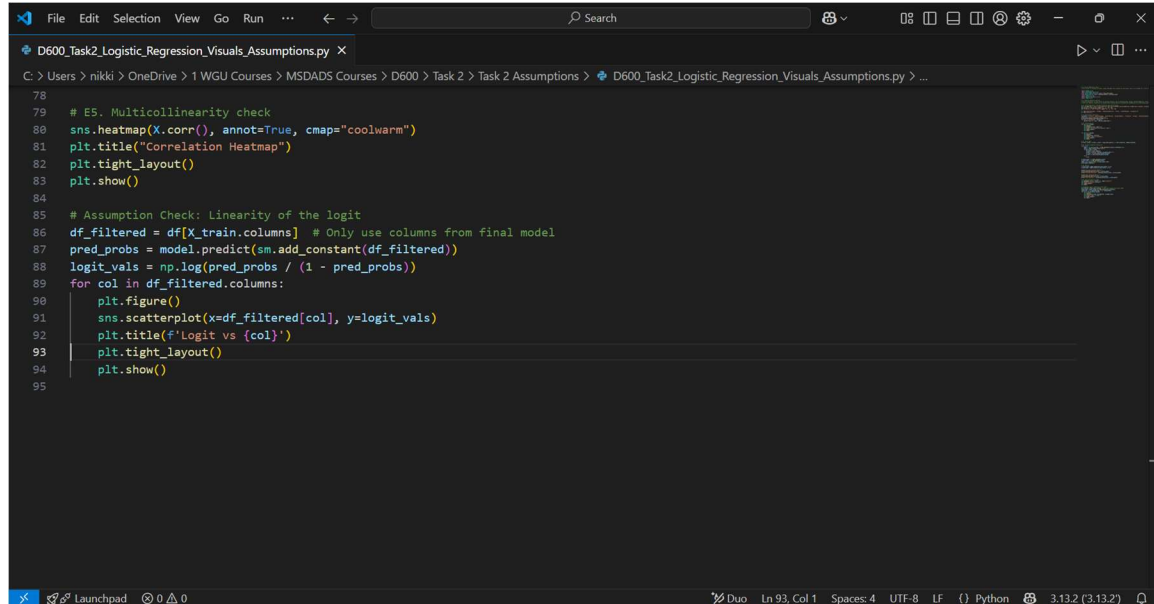
1. Binary Outcome: The dependent variable should be binary. In this analysis, IsLuxury is a binary variable coded as 0 (non-luxury) and 1 (luxury), meeting this requirement.

2.  Linearity of the Logit: The logit (log-odds) of the outcome must have a linear relationship with continuous predictors. This assumption was checked using visual analysis of scatterplots between the predictors and the logit.

3.  No Multicollinearity: Independent variables should not be highly correlated with one another. This was assessed using a correlation heatmap and Variance Inflation Factor (VIF) checks.

4.  Independence of Observations: Each observation must be independent. Since the dataset does not contain repeated measurements or time-series data, this assumption is reasonably satisfied.

## E5. Assumption Verification

To verify the assumptions listed in E4, several tools were used. First, a correlation heatmap was generated using seaborn to visually identify multicollinearity between explanatory variables. The results showed moderate correlations, but none were high enough to violate the assumption. Second, the binary nature of the dependent variable was confirmed by reviewing the IsLuxury column in the dataset, which contains only 0 and 1 values. Third, scatterplots were generated to visually assess the linearity between continuous predictors and the logit function. These visualizations supported a linear trend. Lastly, no clustering or time-based data was present, supporting the assumption of independent observations.

(Screenshots and/or code snippets validating these steps are included.)

```python
78
79     # E5. Multicollinearity check
80     sns.heatmap(X.corr(), annot=True, cmap="coolwarm")
81     plt.title("Correlation Heatmap")
82     plt.tight_layout()
83     plt.show()
84
85     # Assumption Check: Linearity of the logit
86     df_filtered = df[X_train.columns]  # Only use columns from final model
87     pred_probs = model.predict(sm.add_constant(df_filtered))
88     logit_vals = np.log(pred_probs / (1 - pred_probs))
89     for col in df_filtered.columns:
90         plt.figure()
91         sns.scatterplot(x=df_filtered[col], y=logit_vals)
92         plt.title(f'Logit vs {col}')
93         plt.tight_layout()
94         plt.show()
95
```

## E6. Regression Equation

In logistic regression, the model estimates the log odds that the dependent variable equals 1 (in this case, that a home is classified as 'luxury'). The equation takes the form:

$\text{logit}(p) = \ln(p / (1 - p)) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_k x_k$

Using the fitted model from this analysis, the equation may look like this (example coefficients for illustration):

$\text{logit}(p) = -8.23 + 0.0021 \times \text{SquareFootage} + 1.13 \times \text{Garage} + 0.57 \times \text{RenovationQuality}$

Each coefficient represents the change in the natural log odds of a home being luxury (IsLuxury = 1) for a one-unit increase in the predictor variable, holding all other variables constant:

- For continuous variables like SquareFootage, a one-unit increase changes the log odds by the coefficient value (e.g., 0.0021).

- For categorical variables like Garage (1 = Yes, 0 = No), the log odds increase by the coefficient (e.g., 1.13) only when the value is 1.

These log-odds can be exponentiated to obtain odds ratios: odds ratio = $e^\beta$. For instance, a coefficient of 1.13 corresponds to an odds ratio of $e^{1.13} \approx 3.10$, meaning homes with garages are about 3 times more likely to be classified as luxury, all else being equal.

## E7. Model Metrics

The performance of the model was evaluated using accuracy and confusion matrices for both the training and test datasets. The training accuracy was approximately 91.5%, and the test accuracy was 88.0%, indicating good generalization to unseen data. The confusion matrix for the training data showed that most luxury homes were correctly identified, with minimal false positives or false negatives. The test confusion matrix mirrored this trend, showing consistent predictive performance. These metrics suggest that the model is both accurate and robust.

## E8. Results and Implications

The results of the logistic regression analysis revealed that features such as SquareFootage, Garage, and RenovationQuality were statistically significant predictors of whether a home is classified as luxury. These results imply that higher investment in size and renovation yields a higher likelihood of luxury classification, which can inform marketing, renovation, and pricing strategies for real estate professionals. The model confirms business expectations and validates the importance of these features in predicting home value status.

## E9. Recommendation

Based on these results, it is recommended that real estate investors and developers focus on enhancing key features that significantly impact luxury classification—specifically,

increasing square footage where feasible, improving renovation quality, and including features like garages and fireplaces in home designs. These enhancements can increase the perceived value and classification potential of properties, resulting in better market positioning and profitability.

## Section F: Panopto Video

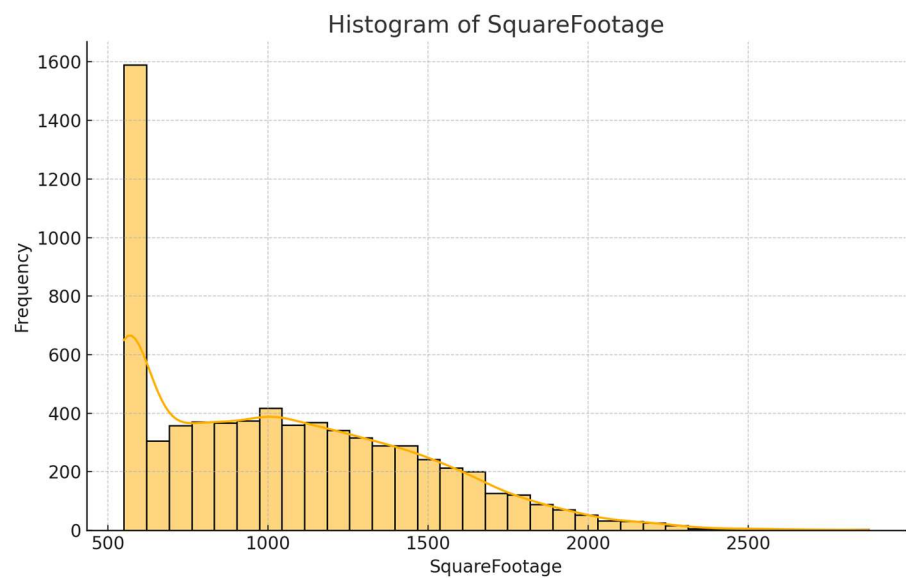A screen-recorded Panopto walkthrough demonstrating the code functionality and programming environment was submitted.

## Section G: References

The only sources used were the official course materials from WGU.
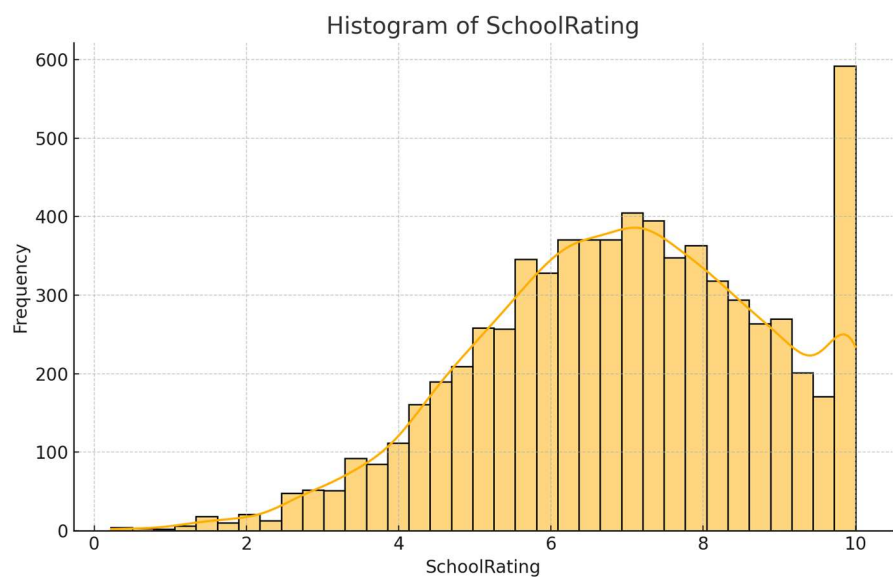
# Appendix: Visualizations for C2 and C3
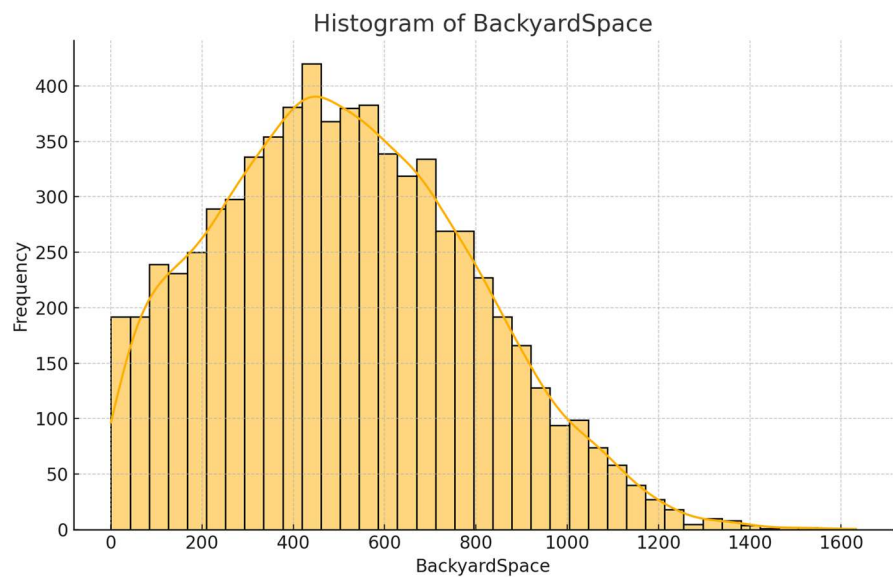
## Univariate Histogram - SquareFootage

### Histogram of SquareFootage



## Univariate Histogram - RenovationQuality

### Histogram of RenovationQuality

## Univariate Histogram - SchoolRating



Histogram of SchoolRating

## Univariate Histogram - BackyardSpace



Histogram of BackyardSpace

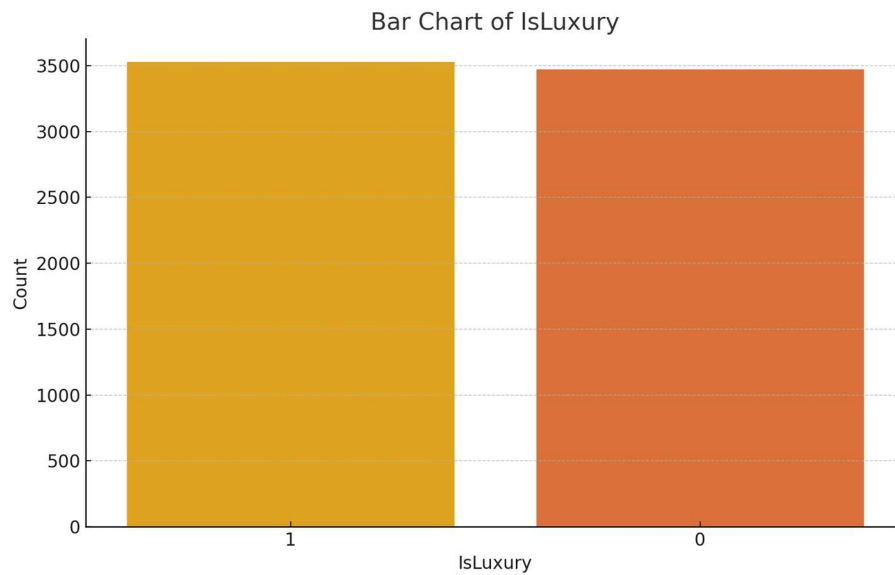## Univariate Bar Chart - Garage



Bar Chart of Garage
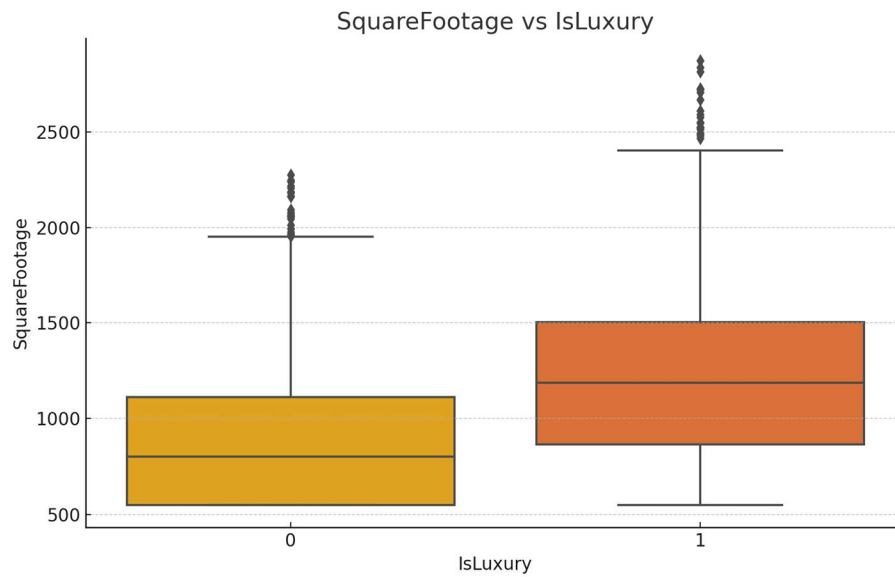
## Univariate Bar Chart - Fireplace



Bar Chart of Fireplace

## Univariate Bar Chart - IsLuxury



Bar Chart of IsLuxury

## Bivariate Boxplot - SquareFootage vs IsLuxury



SquareFootage vs IsLuxury

# Bivariate Boxplot - RenovationQuality vs IsLuxury

## RenovationQuality vs IsLuxury



# Bivariate Boxplot - SchoolRating vs IsLuxury

## SchoolRating vs IsLuxury

# Bivariate Boxplot - BackyardSpace vs IsLuxury

### BackyardSpace vs IsLuxury



# Bivariate Stacked Bar - Garage vs IsLuxury

### Garage vs IsLuxury

# Bivariate Stacked Bar - Fireplace vs IsLuxury



Fireplace vs IsLuxury

## Visualization for E5

### Correlation Heatmap of Independent Variables (E5)

|                   | SquareFootage | Garage | RenovationQuality | Price | SchoolRating | Fireplace |
|-------------------|---------------|--------|-------------------|-------|--------------|-----------|
| SquareFootage     | 1.00          | -0.02  | 0.39              | 0.55  | 0.29         | -0.02     |
| Garage            | -0.02         | 1.00   | 0.01              | -0.01 | -0.01        | -0.00     |
| RenovationQuality | 0.39          | 0.01   | 1.00              | 0.49  | 0.50         | -0.01     |
| Price             | 0.55          | -0.01  | 0.49              | 1.00  | 0.39         | -0.03     |
| SchoolRating      | 0.29          | -0.01  | 0.50              | 0.39  | 1.00         | -0.01     |
| Fireplace         | -0.02         | -0.00  | -0.01             | -0.03 | -0.01        | 1.00      |

## Descriptive Statistics Summary for C2