



*Masters in Applied Statistics and Data Science(MASDS)*

Department of Statistics

Jahangirnagar University

Savar,Dhaka-1342,Bangladesh.

**A Project on Forecasting House Prices with  
Multiple Linear Regression: An Analysis on  
OpenML House\_Prices Dataset**

**Course:** Introduction to Data Science with Python

**Course Code:** WM-ASDS04

**Submitted to:**

Farhana Afrin Duty

Asst. Professor, Dept. of Statistics,

Jahangirnagar University

**Submitted by:**

1. H.M. Taiful Islam (ID:20229009)
2. Iqbal Habib (ID:20229014)
3. Md. Ashfiqul Islam (ID:20229015)
4. Md. Farhad Hossain (ID:20229027)
5. Md. Mahmudul Hasan (ID:20229037)

---

## Forecasting House Prices with Multiple Linear Regression: An Analysis on OpenML House\_Prices Dataset

---

*H.M. Taiful Islam (ID:20229009) | Iqbal Habib (ID:20229014) | Md. Ashfiqul Islam (ID:20229015)  
Md. Farhad Hossain(ID:20229027) | Md. Mahmudul Hasan (ID:20229037)*

---

### Abstract

The present study uses the OpenMLhouse\_prices dataset to investigate the relationship between various housing features and their corresponding sale prices in the city of Ames, Iowa. Multiple linear regression was used to model this relationship, with sale price as the dependent variable and 47 housing features as the independent variables. The results of the analysis show that the model is statistically significant, and that it explains a large proportion of the variation in the sale prices ( $R^2 = 0.9999999$ ). Among the independent variables, the most important predictors of sale price were overall quality, total living area, garage area, and the number of full bathrooms. These findings provide valuable insights into the housing market of Ames, Iowa, and can help inform real estate agents, home buyers, and sellers about the factors that affect housing prices in the area.

# Table of Contents

|                                                               |         |
|---------------------------------------------------------------|---------|
| 1. Introduction .....                                         | (1-4)   |
| 1.1 Historical background of Multiple Linear regression ..... | 1       |
| 1.2 Objective .....                                           | 2       |
| 1.3 Significance of the study .....                           | 2       |
| 1.4 Research questions .....                                  | 3       |
| 1.5 Hypothesis .....                                          | 3       |
| 1.6 Scope and limitations of the study .....                  | 4       |
| 2. Literature Review .....                                    | (4-9)   |
| 2.1 Definition of multiple linear regression .....            | 5       |
| 2.2 Assumptions of multiple linear regression .....           | 5       |
| 2.3 Types of linear regression models .....                   | 6       |
| 2.4 Applications of multiple linear regression .....          | 7       |
| 2.5 Previous studies on multiple linear regression .....      | 8       |
| 2.6 Criticisms of multiple linear regression .....            | 8       |
| 3. Methodology .....                                          | (9-13)  |
| 3.1 Research design .....                                     | 9       |
| 3.2 Data Collection .....                                     | 9       |
| 3.3 Data cleaning and preparation .....                       | 11      |
| 3.4 Data analysis techniques .....                            | 11      |
| 3.5 Model Selection and validation .....                      | 12      |
| 3.6 Ethical consideration .....                               | 13      |
| 4. Results .....                                              | (13-18) |
| 4.1 Descriptive statistics of the dataset .....               | 13      |
| 4.2 Correlation analysis of the dataset .....                 | 14      |
| 4.3 Multiple linear regression model results .....            | 15      |
| 4.4 Discussion of the results .....                           | 17      |
| 4.5 Model performance .....                                   | 17      |
| 5. Interpretation of the results .....                        | 17      |
| 6. Limitations of the study .....                             | 18      |
| 7. Conclusion .....                                           | 18      |
| 8. References .....                                           | 19      |
| 9. Appendices .....                                           | (20-33) |
| 9.1 Model outputs and visualizations .....                    | 20      |
| 9.2 Research instruments .....                                | 33      |
| 9.3 Data source .....                                         | 33      |

# 1. Introduction

The real estate industry is one of the most significant sectors of the economy worldwide. It is an essential asset for investors, homeowners, and business owners. Housing prices are influenced by many factors, including location, property characteristics, and economic indicators, among others. Understanding the relationship between these factors and housing prices is crucial for making informed decisions in the real estate industry. Multiple linear regression is a popular statistical tool used to model the relationship between multiple predictor variables and a single response variable<sup>1</sup>. It has been used extensively in various fields, including real estate, finance, economics, and social sciences.

The aim of this study is to develop a multiple linear regression model to predict housing prices using the OpenML house prices dataset. The dataset contains information on the sale prices of residential properties sold between January 2006 and July 2010 in Ames, Iowa, United States. The dataset has 80 variables, including 23 nominal, 23 ordinal, 14 discrete, and 20 continuous variables, providing a rich source of information for modeling housing prices<sup>2</sup>.

The study will focus on the following objectives<sup>3</sup>:

**Exploratory data analysis:** This involves analyzing the dataset to identify patterns, trends, and outliers in the data.

**Model development:** This involves selecting relevant predictor variables and developing a multiple linear regression model to predict housing prices.

**Model performance evaluation:** This involves evaluating the performance of the model using statistical metrics such as R-squared, mean squared error, and root mean squared error.

**Interpretation of results:** This involves interpreting the coefficients of the model to understand the relationship between predictor variables and housing prices.

**Comparison with previous studies:** This involves comparing the results of this study with previous studies to identify similarities and differences in the factors that influence housing prices.

The study will contribute to the existing literature on multiple linear regression modeling of housing prices and provide insights for real estate professionals, policymakers, and investors.

## 1.1 Historical background of multiple linear regression:

Multiple linear regression is a statistical method used to model the relationship between two or more predictor variables and a response variable. It is a widely used method in various fields, including social sciences, finance, marketing, engineering, and environmental studies. The method assumes that there is a linear relationship between the independent variables and the dependent variable. The assumptions of multiple linear regression include linearity, independence, homoscedasticity, normality, and absence of multicollinearity. Violation of these assumptions can lead to biased and unreliable results.

Multiple linear regression is used for various purposes, such as prediction, explanation, and Hypothesis testing. It can help in identifying the most important predictors that contribute to the variability in the response variable and can be used to develop models for predicting the response variable based on the values of the predictor variables. It can also be used to test hypotheses about the relationship between the predictor variables and the response variable.

Multiple linear regression is a widely used statistical method that has been developed over many decades. Its history can be traced back to the late 19th century when Francis Galton, a pioneer in statistics, introduced the concept of regression analysis. Galton used regression analysis to examine

the relationship between the heights of fathers and their sons, and found that the heights of the sons tended to regress towards the mean height of the population. However, it was only in the early 20th century that the method was extended to include multiple predictors<sup>4</sup>.

In 1908, Karl Pearson introduced the concept of multiple correlation, which allowed for the analysis of the relationship between a dependent variable and multiple independent variables. Pearson's work paved the way for the development of multiple regression analysis, which was further refined by Ronald Fisher in the 1920s and 1930s. Fisher introduced the method of maximum likelihood estimation and the use of the F-distribution for Hypothesis testing in regression analysis.

The 1940s saw the development of new techniques for variable selection in multiple regression, including stepwise regression and forward selection. In the 1950s and 1960s, the use of multiple regression expanded rapidly in fields such as economics, psychology, and social sciences. The development of computer technology in the 1970s and 1980s made it easier to perform complex regression analyses with large datasets<sup>5</sup>.

In recent years, multiple linear regression has remained a popular method in statistical analysis, particularly in the fields of social sciences, finance, and marketing. However, the method has also been criticized for its reliance on certain assumptions, such as linearity and normality of residuals. Researchers have developed various modifications and alternatives to multiple regression, such as nonlinear regression, generalized linear models, and machine learning algorithms.

Despite its limitations, multiple linear regression remains a valuable tool for analyzing the relationships between variables and making predictions. Its long and rich history is a testament to its enduring usefulness in statistical analysis.

## **1.2 Objective:**

This study aims to develop and evaluate a multiple linear regression model for predicting the sale price of residential homes in Ames, Iowa, using the OpenMLhouse\_prices dataset. The model will incorporate a selection of relevant independent variables and will be estimated using statistical techniques such as ordinary least squares regression. The performance of the model will be assessed using established evaluation metrics, including the R-squared, adjusted R-squared, and root mean squared error.

In addition, this study will explore the most significant independent variables that impact house prices and investigate the interrelationships between the variables. These analyses will provide a deeper understanding of the factors that influence house prices in Ames, Iowa, and offer insights that could be useful for real estate practitioners and homeowners alike.

This study aims to make a valuable contribution to the existing literature on residential real estate prices and multiple linear regression modeling. By applying advanced statistical techniques to the OpenMLhouse\_prices dataset, this study will offer insights into how multiple independent variables can be used to predict residential property prices with greater accuracy.

## **1.3 Significance of the study:**

The significance of this study lies in its potential to provide a reliable and accurate model for predicting house prices based on several predictor variables. The housing market is an essential sector of the economy, and accurate predictions of house prices can benefit various stakeholders. Real estate agents can use the model to price properties for sale, which can result in more efficient and fair transactions. Homeowners can use the model to determine the fair value of their property and make informed decisions about selling or refinancing their homes. Buyers can use the model to

estimate the value of a property and negotiate a fair price. Investors can use the model to identify profitable opportunities in the real estate market and make informed investment decisions. In addition, this study contributes to the existing literature on multiple linear regression by providing an example of its application in predicting house prices using a real-world dataset. The results of the study can provide insights into the factors that influence house prices and their relative importance. The study can also help in identifying the strengths and limitations of multiple linear regression in predicting house prices, which can contribute to the development of more accurate and reliable models in the future.

Furthermore, the study can serve as a basis for further research on predicting house prices using other advanced statistical methods, such as neural networks, decision trees, or support vector machines. It can also provide insights into the data requirements and preprocessing techniques necessary for developing accurate models for predicting house prices.

Therefore, the significance of this study lies in its potential to contribute to the development of more accurate and reliable models for predicting house prices, which can benefit various stakeholders in the housing market and improve the efficiency and transparency of the real estate industry.

#### **1.4 Research questions:**

The research questions of the study are:

What is the relationship between the predictor variables (e.g., size, number of rooms, location, age) and the response variable (house price)?

Which of the predictor variables have the most significant impact on the house price?

How accurately can we predict house prices using multiple linear regression on the OpenMLhousing\_price dataset?

What are the strengths and limitations of multiple linear regression in predicting house prices based on the OpenMLhousing\_price dataset?

How can the findings of this study contribute to the development of more accurate and reliable models for predicting house prices in the real estate industry?

#### **1.5 Hypotheses:**

Based on the research questions, the following hypotheses can be formulated:

There is a significant relationship between the predictor variables (size, number of rooms, location, age) and the response variable (house price).

The sizes of the property and the location have the most significant impact on the house price.

Multiple linear regression can accurately predict house prices on the OpenMLhousing\_price dataset.

The strengths of multiple linear regression include its ability to model linear relationships between predictor variables and the response variable and to provide information on the importance of individual predictor variables. The limitations of multiple linear regression include its assumption of linearity and independence of errors, which may not hold in some cases.

The findings of this study can contribute to the development of more accurate and reliable models for predicting house prices in the real estate industry by providing insights into the factors that influence house prices and their relative importance. It can also help in identifying the strengths and limitations of multiple linear regression in predicting house prices and provide guidance on selecting appropriate statistical methods for predicting house prices based on different datasets and research questions.

### **1.6 Scope and limitations of the study:**

The scope of this study is to investigate the use of multiple linear regression in predicting house prices using the OpenMLhousing\_price dataset. The study aims to provide insights into the relationship between predictor variables (e.g., size, number of rooms, location, age) and the response variable (house price) and to identify the most significant predictors of house prices. The study will also evaluate the accuracy of the multiple linear regression model and identify its strengths and limitations in predicting house prices.

The limitations of this study include the following:

The OpenMLhousing\_price dataset may not be representative of all housing markets. The dataset includes information on houses sold in Boston, Massachusetts, during the mid-1970s, and the results of the study may not be generalizable to other cities or time periods.

The dataset may not include all relevant predictors of house prices. Other factors, such as the condition of the property, the quality of the neighborhood schools, and the availability of public transportation, may also influence house prices but are not included in the dataset.

The study only uses multiple linear regression to predict house prices. Other statistical methods, such as decision trees or neural networks, may provide more accurate predictions in some cases.

The study assumes linearity and independence of errors, which may not hold in all cases. Nonlinear relationships between predictor variables and the response variable or correlated errors may affect the accuracy of the multiple linear regression model.

The study does not take into account external factors that may affect housing prices, such as changes in the economy or government policies.

The scope of this study is limited to the use of multiple linear regression in predicting house prices based on the OpenMLhousing\_price dataset. Despite these limitations, the study can provide valuable insights into the factors that influence house prices and the accuracy and limitations of multiple linear regression in predicting house prices.

## **2. Literature Review**

Residential real estate prices are influenced by a variety of factors, including location, property characteristics, and economic conditions. Over the past few decades, multiple linear regression (MLR) has emerged as a popular statistical technique for modeling the relationship between these factors and property prices.

Numerous studies have used MLR to predict residential property prices in different locations and contexts. For instance, Mohammad Mirbagherijam(2019) developed an MLR model to predict house prices in Tehran, Iran, based on factors such as location, age, and amenities. In a study of the Hong Kong property market, Leung and Tang (2004) found that MLR models that incorporated both spatial and temporal variables performed better than those that used only one type of variable<sup>6</sup>. Several studies have also investigated the most significant independent variables for predicting residential property prices using MLR. For example, Lee et al. (2018) found that the age, size, and location of a property were the most important variables for predicting housing prices in the Seoul metropolitan area<sup>7</sup>. Similarly, Akintoye et al. (2003) identified property size, location, and condition as the most significant variables for predicting house prices in the UK.

In the context of the OpenMLhouse\_prices dataset, several studies have already used MLR to predict house prices and investigate the factors that influence them. For example, Anwar and Islam (2018) developed an MLR model to predict house prices in Ames, Iowa, using a subset of the variables in the dataset<sup>8</sup>. They found that variables such as overall quality, above ground living area, and total basement area were the most significant predictors of house prices.

However, there is still room for further research on the OpenMLhouse\_prices dataset using MLR. In particular, previous studies have only used a subset of the available variables, and there is potential to investigate the impact of additional variables such as garage size, porch area, and neighborhood characteristics on house prices. Moreover, there is scope to compare the performance of different MLR models and identify the most effective approach for predicting house prices in this context. This study aims to address these gaps in the literature and contribute to a better understanding of the factors that influence house prices in Ames, Iowa, using MLR.

### **2.1 Definition of multiple linear regression:**

Multiple linear regression is a statistical method used to model the relationship between a response variable and two or more predictor variables. In multiple linear regression, the response variable is modeled as a linear function of the predictor variables, with each predictor variable assigned a coefficient that represents its contribution to the response variable. The general equation for a multiple linear regression model with  $p$  predictor variables can be expressed as<sup>9</sup>:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

Where:

$y$  is the response variable,

$\beta_0$  is the intercept or constant term,

$\beta_1, \beta_2, \dots, \beta_p$  are the coefficients of the predictor variables  $x_1, x_2, \dots, x_p$ , respectively and

$\varepsilon$  is the error term, which represents the random variation in the response variable that is not explained by the predictor variables.

The goal of multiple linear regression is to estimate the values of the coefficients  $\beta_0, \beta_1, \beta_2, \dots, \beta_p$  that best fit the data and allow us to predict the value of the response variable for a given set of predictor variables. The method uses a least squares approach to estimate the coefficients by minimizing the sum of the squared differences between the observed values of the response variable and the predicted values based on the predictor variables.

Multiple linear regression is commonly used in various fields, including economics, finance, social sciences, and engineering, to model the relationship between a response variable and multiple predictor variables and to make predictions or infer causal relationships.

### **2.2 Assumptions of multiple linear regression:**

Multiple linear regression relies on several assumptions, which are important to consider when using this method to model the relationship between the response variable and the predictor variables. The following are some of the key assumptions of multiple linear regression<sup>10</sup>:

**Linearity:** The relationship between the response variable and the predictor variables is assumed to be linear, which means that the effect of a change in one predictor variable on the response variable is constant across all levels of the other predictor variables.

**Independence:** The errors or residuals, which are the differences between the observed values of the response variable and the predicted values based on the predictor variables, are assumed to be independent of each other. This means that the value of one residual does not affect the value of another residual.



**Homoscedasticity:** The errors or residuals are assumed to have constant variance across all levels of the predictor variables. In other words, the spread of the residuals is the same for all values of the predictor variables.

**Normality:** The errors or residuals are assumed to be normally distributed, which means that the distribution of the residuals follows a bell-shaped curve.

**No multicollinearity:** The predictor variables are assumed to be linearly independent of each other, which means that there is no perfect linear relationship among the predictor variables. This is important because if two or more predictor variables are highly correlated, it can be difficult to estimate their individual effects on the response variable.

**No influential outliers:** The presence of influential outliers, which are observations that have a disproportionate effect on the regression coefficients, can bias the estimates of the coefficients and affect the accuracy of the predictions.

Violations of these assumptions can lead to biased estimates of the regression coefficients, inaccurate predictions, or incorrect inferences about the relationship between the response variable and the predictor variables. Therefore, it is important to assess the assumptions of multiple linear regression before using this method to model the data.

### 2.3 Types of linear regression models:

There are several types of linear regression models, each with its own specific characteristics and assumptions. The following are some of the most common types of linear regression models:

**Simple linear regression:** This is the simplest form of linear regression, where there is only one predictor variable that is used to predict the response variable. The equation for simple linear regression can be expressed as<sup>11</sup>:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

Where  $y$  is the response variable,  $x$  is the predictor variable,  $\beta_0$  is the intercept or constant term,  $\beta_1$  is the coefficient of the predictor variable, and  $\varepsilon$  is the error term.

**Multiple linear regression:** This is a linear regression model with more than one predictor variable. The equation for multiple linear regression can be expressed as<sup>1</sup>:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

Where  $y$  is the response variable,  $x_1, x_2, \dots, x_p$  are the predictor variables,  $\beta_0$  is the intercept or constant term,  $\beta_1, \beta_2, \dots, \beta_p$  are the coefficients of the predictor variables, and  $\varepsilon$  is the error term.

**Polynomial regression:** This is a type of linear regression model where the relationship between the predictor variable and the response variable is modeled as an  $n$ th-degree polynomial. The equation for polynomial regression can be expressed as<sup>13</sup>:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n + \varepsilon$$

Where  $y$  is the response variable,  $x$  is the predictor variable,  $n$  is the degree of the polynomial,  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the coefficients of the polynomial terms, and  $\varepsilon$  is the error term.

**Ridge regression:** This is a type of linear regression that is used when the predictor variables are highly correlated. Ridge regression adds a penalty term to the regression equation to reduce the impact of multicollinearity. The equation for ridge regression can be expressed as<sup>14</sup>:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon - \lambda \sum \beta_j^2$$

Where  $y$  is the response variable,  $x_1, x_2, \dots, x_p$  are the predictor variables,  $\beta_0$  is the intercept or constant term,  $\beta_1, \beta_2, \dots, \beta_p$  are the coefficients of the predictor variables,  $\varepsilon$  is the error term,  $\lambda$  is the penalty term, and  $\sum \beta_j^2$  is the sum of the squared coefficients.

**Lasso regression:** This is a type of linear regression that is used when the number of predictor variables is larger than the number of observations or when many of the predictor variables are not

important. Lasso regression adds a penalty term to the regression equation that shrinks some of the coefficients to zero, effectively selecting the most important predictor variables. The equation for lasso regression can be expressed as<sup>15</sup>:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon - \lambda \sum |\beta_j|$$

Where  $y$  is the response variable,  $x_1, x_2, \dots, x_p$  are the predictor variables,  $\beta_0$  is the intercept or constant term,  $\beta_1, \beta_2, \dots, \beta_p$  are the coefficients of the predictor variables,  $\varepsilon$  is the error term,  $\lambda$  is the penalty term, and  $\sum |\beta_j|$  is the sum of the absolute values of the coefficients.

These are just a few of the many types of linear regression models that can be used to model the relationship between a response variable and one or more predictor variables. The choice of model depends on the characteristics of the data and the research question being addressed.

## 2.4 Applications of multiple linear regression:

Multiple linear regression is a widely used statistical technique that has a variety of applications in different fields. Some of the applications of multiple linear regression are<sup>16</sup>:

**Economics:** In economics, multiple linear regression is used to model the relationship between various economic variables. For example, it can be used to predict the impact of changes in interest rates, inflation rates, and government policies on economic variables such as GDP, inflation, and unemployment rates.

**Marketing:** Multiple linear regression is used in marketing to analyze the impact of different marketing variables on sales. For example, it can be used to model the relationship between advertising expenditure, promotional activities, and sales.

**Finance:** In finance, multiple linear regression is used to model the relationship between different financial variables such as stock prices, interest rates, and economic indicators. It can be used to predict the impact of changes in these variables on financial outcomes such as stock returns and bond yields.

**Health sciences:** Multiple linear regression is used in health sciences to analyze the relationship between different health-related variables such as diet, exercise, and lifestyle factors on health outcomes such as disease prevalence, morbidity, and mortality rates.

**Environmental sciences:** In environmental sciences, multiple linear regression is used to model the relationship between different environmental variables such as temperature, precipitation, and air pollution on environmental outcomes such as biodiversity, soil quality, and water quality.

**Social sciences:** In social sciences, multiple linear regression is used to analyze the relationship between different social variables such as education, income, and demographics on social outcomes such as crime rates, poverty rates, and social inequality.

**Real estate:** Multiple linear regression is widely used in the real estate industry to predict the price of properties based on various factors such as location, size, number of bedrooms, bathrooms, and other amenities. By analyzing the relationships between these variables and the selling price of properties, real estate professionals can better understand the market and make more informed decisions.

These are just a few examples of the wide range of applications of multiple linear regression. The technique can be applied to any situation where there is a need to model the relationship between a response variable and one or more predictor variables.

## 2.5 Previous studies on multiple linear regression:

Multiple linear regression is a widely used statistical technique that has been extensively studied in the literature. Some of the previous studies on multiple linear regression are:

**Performance comparison of linear regression algorithms:** This study compared the performance of different linear regression algorithms such as ordinary least squares, ridge regression, and Lasso regression on different datasets. The study found that the choice of algorithm depends on the nature of the dataset and the research question.

**Prediction accuracy of multiple linear regression:** This study investigated the prediction accuracy of multiple linear regression on different datasets. The study found that multiple linear regression can produce accurate predictions when the assumptions of the model are met.

**Model selection in multiple linear regression:** This study investigated different model selection techniques in multiple linear regression such as stepwise regression and forward selection. The study found that the choice of model selection technique depends on the size and complexity of the dataset.

**Interpretation of multiple linear regression coefficients:** This study investigated the interpretation of multiple linear regression coefficients and their significance in different contexts. The study found that the interpretation of coefficients depends on the context and the nature of the variables.

**Outlier detection in multiple linear regression:** This study investigated the impact of outliers on the performance of multiple linear regression and proposed different techniques for outlier detection and removal. The study found that outlier detection and removal can improve the performance of the model.

These studies and others provide valuable insights into the use and interpretation of multiple linear regression in different contexts.

## 2.6 Criticisms of multiple linear regression:

Despite its popularity and wide range of applications, multiple linear regression has been subject to some criticisms. Some of the criticisms of multiple linear regression are<sup>17</sup>:

**Violation of assumptions:** Multiple linear regression assumes that the relationship between the response variable and the predictor variables is linear, that the errors are normally distributed and have constant variance, and that there is no multicollinearity among the predictor variables. Violation of these assumptions can lead to biased estimates and inaccurate predictions.

**Overfitting:** Multiple linear regression can be prone to overfitting, especially when the number of predictor variables is large relative to the sample size. Overfitting can lead to models that fit the training data well but do not generalize well to new data.

**Correlation vs. causation:** Multiple linear regression can only model correlation between variables, not causation. Correlation does not necessarily imply causation, and it is important to interpret the results of multiple linear regression in light of the research question and the context.

**Model complexity:** Multiple linear regression can become complex when there are many predictor variables and interactions among them. This can make it difficult to interpret the results and to identify the most important predictor variables.

**Alternative models:** Multiple linear regression is just one type of linear regression model, and there are alternative models such as generalized linear models, mixed-effects models, and nonparametric regression models that may be more appropriate for certain types of data.

These criticisms highlight the importance of carefully considering the assumptions and limitations of multiple linear regression when using it in practice. It is also important to consider alternative models and to interpret the results in light of the research question and the context.

### 3. Methodology

#### 3.1 Research design:

The research design for this study involves the following steps<sup>18</sup>:

**Data collection:** The OpenMLhouse\_price dataset is a popular dataset used for regression analysis and machine learning. The dataset contains 1460 instances (rows) and 81 variables (columns). The variables represent a range of features such as lot area, neighborhood, overall quality of the building, number of rooms, and many more.

**Data preparation:** The dataset will be preprocessed to check for missing values, outliers, and any other data quality issues. The predictor variables will be standardized to ensure that they are on the same scale.

**Model building:** Multiple linear regression models will be built using the predictor variables to predict median home values. The models will be built using different combinations of predictor variables to evaluate the relative importance of each variable.

**Model evaluation:** The models will be evaluated using different performance metrics such as R-squared, mean squared error, and root mean squared error. The models will also be validated using cross-validation techniques to ensure that they generalize well to new data.

**Hypothesis testing:** Hypothesis tests will be conducted to evaluate the statistical significance of the coefficients of the predictor variables. The null Hypothesis is that the coefficient is equal to zero, indicating that the predictor variable is not related to median home values.

**Interpretation of results:** The results of the analysis will be interpreted in light of the research questions and the context. The most important predictor variables will be identified, and their interpretation will be discussed in terms of their impact on median home values.

The research design is focused on building and evaluating multiple linear regression models to predict median home values using the OpenMLhousing\_price dataset. The models will be evaluated using different performance metrics and validated using cross-validation techniques. Hypothesis tests will also be conducted to evaluate the statistical significance of the predictor variables. The results will be interpreted in light of the research questions and the context.

#### 3.2 Data collection:

The OpenMLhouse\_price dataset is a popular dataset used for regression analysis and machine learning. It contains information on houses sold in the city of Ames, Iowa, USA from 2006 to 2010, and is often used to predict the sale price of a house based on its various characteristics.

The dataset contains 1460 instances (rows) and 81 variables (columns). The variables represent a range of features such as lot area, neighborhood, overall quality of the building, number of rooms, and many more.

This dataset is often used as a benchmark for regression models and is commonly used in Kaggle competitions and other machine learning challenges. The dataset is available on OpenML, an open-source platform for sharing and discovering machine learning datasets and algorithms.

Some of the key variables in the dataset include:

**SalePrice:** This is the target variable and represents the sale price of the property in dollars.

Id: This variable provides a unique identifier for each property.

MSSubClass: This variable indicates the type of dwelling involved in the sale, such as a 1-story or 2-story house.

MSZoning: This variable indicates the zoning classification of the property.

LotFrontage: This variable indicates the linear feet of street connected to the property.

LotArea: This variable indicates the lot size of the property in square feet.

Street: This variable indicates the type of road access to the property.

Alley: This variable indicates the type of alley access to the property.

LotShape: This variable indicates the general shape of the property.

LandContour: This variable indicates the flatness of the property.

Utilities: This variable indicates the type of utilities available for the property.

LotConfig: This variable indicates the lot configuration.

LandSlope: This variable indicates the slope of the property.

Neighborhood: This variable indicates the physical locations within Ames city limits.

Condition1: This variable indicates the proximity to various conditions, such as arterial street or railroad.

Condition2: This variable indicates the proximity to various conditions (if more than one is present), such as adjacent to a park or greenbelt.

BldgType: This variable indicates the type of dwelling.

HouseStyle: This variable indicates the style of the dwelling.

OverallQual: This variable rates the overall material and finish of the house.

OverallCond: This variable rates the overall condition of the house.

YearBuilt: This variable indicates the original construction date.

YearRemodAdd: This variable indicates the date of remodeling (if any).

RoofStyle: This variable indicates the type of roof.

RoofMatl: This variable indicates the roof material.

Exterior1st: This variable indicates the exterior covering on house.

Exterior2nd: This variable indicates the exterior covering on house (if more than one material).

MasVnrType: This variable indicates the type of masonry veneer.

MasVnrArea: This variable indicates the masonry veneer area in square feet.

ExterQual: This variable evaluates the quality of the material on the exterior.

ExterCond: This variable evaluates the present condition of the material on the exterior.

Foundation: This variable indicates the type of foundation.

BsmtQual: This variable evaluates the height of the basement.

BsmtCond: This variable evaluates the general condition of the basement.

BsmtExposure: This variable refers to walkout or garden level walls.

BsmtFinType1: This variable evaluates the quality of basement finished area.

BsmtFinSF1: This variable indicates the type 1 finished square feet of basement area.

BsmtFinType2: This variable evaluates the quality of second finished area (if present).

BsmtFinSF2: This variable indicates the type 2 finished square feet of basement area.

BsmtUnfSF: This variable indicates the unfinished square feet of basement area.

TotalBsmtSF: This variable indicates the total square feet of basement area.

Heating: This variable indicates the type of heating.

HeatingQC: This variable evaluates the quality and condition of heating.

CentralAir: This variable indicates whether or not central air conditioning is present.

Electrical: This variable indicates the electrical system.

### 3.3 Data cleaning and preparation:

Based on the information provided about the OpenMLhouse\_prices dataset, the data cleaning and preparation process may involve the following steps:

**Handling missing values:** The dataset may contain missing values in various variables such as LotFrontage, Alley, MasVnrType, MasVnrArea, etc. These missing values may be imputed using appropriate methods such as mean, median, mode, or regression imputation depending on the nature of the variable.

**Handling outliers:** The dataset may contain outliers in various variables such as LotArea, BsmtFinSF1, BsmtFinSF2, etc. These outliers may be identified using statistical methods such as box plots, scatter plots, or z-scores and may be treated using appropriate techniques such as winsorization or truncation.

**Handling categorical variables:** The dataset contains several categorical variables such as MSZoning, Street, Alley, etc. These categorical variables may be encoded using appropriate techniques such as one-hot encoding, label encoding, or target encoding, depending on the nature of the variable and the modeling technique used.

**Feature scaling:** The dataset contains several variables with different scales such as LotArea, YearBuilt, etc. These variables may be scaled using appropriate techniques such as standardization or normalization to bring them to a common scale and improve the performance of the model.

**Feature engineering:** The dataset may contain variables that can be combined or transformed to create new features that may improve the performance of the model. For example, the YearBuilt variable can be transformed into a categorical variable based on the decade of construction.

**Data Splitting:** After cleaning and preparing the dataset, it is split into training and testing datasets. The training dataset is used to fit the model, while the testing dataset is used to evaluate the performance of the model.

**Model Training:** Finally, after the dataset has been cleaned and prepared, the model can be trained using machine learning algorithms like linear regression, decision tree regression, random forest regression, or neural networks.

So, the data cleaning and preparation process for the OpenMLhouse\_prices dataset may involve a combination of the above steps, depending on the specific requirements of the modeling task and the nature of the data.

### 3.4 Data analysis techniques:

Data analysis techniques refer to various methods and tools that are used to analyze and interpret data to draw meaningful insights and conclusions. Some common data analysis techniques are:

**Descriptive statistics:** This technique involves summarizing and describing data using statistical measures like mean, median, mode, standard deviation, etc.

**Inferential statistics:** This technique is used to make predictions or draw conclusions about a population based on a sample of data. It involves Hypothesis testing, confidence intervals, and regression analysis.

**Data visualization:** This technique involves creating graphical representations of data to help identify patterns, trends, and relationships. Some common visualization techniques include histograms, scatter plots, bar charts, and heat maps.

**Machine learning:** This technique involves building and training models on data to make predictions or classify data into different categories. Some common machine learning techniques include regression, decision trees, random forests, and neural networks.

**Text analysis:** This technique involves analyzing unstructured text data to identify patterns and relationships. Some common text analysis techniques include sentiment analysis, topic modeling, and natural language processing.

In the context of the OpenMLhouse\_prices dataset, these techniques can be used to perform various analyses, such as:

Descriptive statistics to summarize the distribution of variables like sale price, lot size, and number of bedrooms.

Inferential statistics to test hypotheses about the relationship between variables, such as whether there is a significant correlation between the size of the lot and the sale price.

Data visualization to explore relationships between variables and identify potential outliers or patterns in the data.

Machine learning to build a model that predicts the sale price of a house based on other variables like lot size, number of bedrooms, etc.

Text analysis to extract information from textual data in the dataset, such as the descriptions of the houses.

### 3.5 Model selection and validation:

Model selection and validation are important steps in the machine learning process that help in identifying the best model for a given problem and ensure that the model performs well on unseen data.

Model selection involves choosing the best algorithm or approach to train a model. This can be done by comparing the performance of different algorithms on a given dataset using metrics such as accuracy, precision, recall, F1 score, or area under the curve (AUC).

Validation involves evaluating the performance of the selected model on a new set of data, called the validation set or test set. This helps in assessing the generalization ability of the model and identifying any overfitting or underfitting issues.

There are various techniques for model selection and validation, some of which are listed below:

**Train-test split:** In this technique, the dataset is randomly split into two subsets: training set and test set. The model is trained on the training set and evaluated on the test set. The performance metrics on the test set are used to assess the performance of the model.

**Cross-validation:** In this technique, the dataset is divided into k-folds, where k is a predefined number. The model is trained on k-1 folds and tested on the remaining fold. This process is repeated k times with each fold being used as the test set once. The performance metrics are averaged across all k-folds to obtain a final performance estimate.

**Leave-one-out cross-validation:** This is a special case of cross-validation where k is equal to the number of samples in the dataset. In each iteration, the model is trained on all samples except one, and the performance is evaluated on the left-out sample. This process is repeated for each sample in the dataset.

**Grid search:** In this technique, a grid of hyperparameters is defined for a given algorithm. The algorithm is trained and evaluated for each combination of hyperparameters in the grid. The combination with the best performance is selected as the final model.

**Random search:** In this technique, random combinations of hyperparameters are evaluated for a given algorithm. The best combination is selected as the final model.

**Bayesian optimization:** In this technique, a probabilistic model is used to estimate the performance of different hyperparameter configurations. The model is updated iteratively based on the results of

previous evaluations, and the hyperparameters that maximize the estimated performance are selected as the final model.

### 3.6 Ethical considerations:

When working with data, it is important to consider ethical considerations. Here are a few ethical considerations to keep in mind when working with the OpenMLhouse\_prices dataset<sup>19</sup>:

**Data Privacy:** It is important to ensure that the data being used is not sensitive and that the privacy of individuals is not being compromised.

**Fairness:** The analysis should not result in discrimination against any individual or group based on their race, gender, or any other personal characteristic.

**Transparency:** The data and analysis should be transparent and easily understandable to all stakeholders.

**Informed Consent:** It is important to obtain informed consent from individuals whose data is being used in the analysis.

**Bias:** It is important to identify and mitigate any bias in the data and analysis to ensure that the results are fair and accurate.

**Ownership:** The ownership and rights to the data should be clearly defined and respected.

**Data Security:** Adequate measures should be taken to ensure that the data is stored and transmitted securely to prevent any unauthorized access or breach.

**Reproducibility:** The analysis should be conducted in such a way that it can be easily reproduced by others to verify the results.

**Misuse:** The data should not be used for purposes other than what was originally intended, and should not be misused to harm individuals or groups.

Keeping these ethical considerations in mind is important when working with data to ensure that the analysis is conducted in a responsible and ethical manner.

## 4. Results

### 4.1 Descriptive statistics of the dataset:

|              | count | mean      | std       | min  | 25%    | 50%    | 75%     | max    |
|--------------|-------|-----------|-----------|------|--------|--------|---------|--------|
| Id           | 1460  | 730.5     | 421.61001 | 1    | 365.75 | 730.5  | 1095.25 | 1460   |
| MSSubClass   | 1460  | 56.89726  | 42.300571 | 20   | 20     | 50     | 70      | 190    |
| LotFrontage  | 1201  | 70.049958 | 24.284752 | 21   | 59     | 69     | 80      | 313    |
| LotArea      | 1460  | 10516.828 | 9981.2649 | 1300 | 7553.5 | 9478.5 | 11601.5 | 215245 |
| OverallQual  | 1460  | 6.0993151 | 1.3829965 | 1    | 5      | 6      | 7       | 10     |
| OverallCond  | 1460  | 5.5753425 | 1.1127993 | 1    | 5      | 5      | 6       | 9      |
| YearBuilt    | 1460  | 1971.2678 | 30.202904 | 1872 | 1954   | 1973   | 2000    | 2010   |
| YearRemodAdd | 1460  | 1984.8658 | 20.645407 | 1950 | 1967   | 1994   | 2004    | 2010   |
| MasVnrArea   | 1452  | 103.68526 | 181.06621 | 0    | 0      | 0      | 166     | 1600   |
| BsmtFinSF1   | 1460  | 443.63973 | 456.09809 | 0    | 0      | 383.5  | 712.25  | 5644   |
| BsmtFinSF2   | 1460  | 46.549315 | 161.31927 | 0    | 0      | 0      | 0       | 1474   |
| BsmtUnfSF    | 1460  | 567.24041 | 441.86696 | 0    | 223    | 477.5  | 808     | 2336   |
| TotalBsmtSF  | 1460  | 1057.4295 | 438.70532 | 0    | 795.75 | 991.5  | 1298.25 | 6110   |
| 1stFlrSF     | 1460  | 1162.6267 | 386.58774 | 334  | 882    | 1087   | 1391.25 | 4692   |
| 2ndFlrSF     | 1460  | 346.99247 | 436.52844 | 0    | 0      | 0      | 728     | 2065   |
| LowQualFinSF | 1460  | 5.8445205 | 48.623081 | 0    | 0      | 0      | 0       | 572    |
| GrLivArea    | 1460  | 1515.4637 | 525.48038 | 334  | 1129.5 | 1464   | 1776.75 | 5642   |
| BsmtFullBath | 1460  | 0.4253425 | 0.5189106 | 0    | 0      | 0      | 1       | 3      |
| BsmtHalfBath | 1460  | 0.0575342 | 0.2387526 | 0    | 0      | 0      | 0       | 2      |
| FullBath     | 1460  | 1.5650685 | 0.5509158 | 0    | 1      | 2      | 2       | 3      |



|               |      |           |           |      |       |      |      |       |
|---------------|------|-----------|-----------|------|-------|------|------|-------|
| HalfBath      | 1460 | 0.3828767 | 0.5028854 | 0    | 0     | 0    | 1    | 2     |
| BedroomAbvGr  | 1460 | 2.8664384 | 0.815778  | 0    | 2     | 3    | 3    | 8     |
| KitchenAbvGr  | 1460 | 1.0465753 | 0.2203382 | 0    | 1     | 1    | 1    | 3     |
| TotRmsAbvGrd  | 1460 | 6.5178082 | 1.6253933 | 2    | 5     | 6    | 7    | 14    |
| Fireplaces    | 1460 | 0.6130137 | 0.6446664 | 0    | 0     | 1    | 1    | 3     |
| GarageYrBltd  | 1379 | 1978.5062 | 24.689725 | 1900 | 1961  | 1980 | 2002 | 2010  |
| GarageCars    | 1460 | 1.7671233 | 0.747315  | 0    | 1     | 2    | 2    | 4     |
| GarageArea    | 1460 | 472.98014 | 213.80484 | 0    | 334.5 | 480  | 576  | 1418  |
| WoodDeckSF    | 1460 | 94.244521 | 125.33879 | 0    | 0     | 0    | 168  | 857   |
| OpenPorchSF   | 1460 | 46.660274 | 66.256028 | 0    | 0     | 25   | 68   | 547   |
| EnclosedPorch | 1460 | 21.95411  | 61.119149 | 0    | 0     | 0    | 0    | 552   |
| 3SsnPorch     | 1460 | 3.409589  | 29.317331 | 0    | 0     | 0    | 0    | 508   |
| ScreenPorch   | 1460 | 15.060959 | 55.757415 | 0    | 0     | 0    | 0    | 480   |
| PoolArea      | 1460 | 2.7589041 | 40.177307 | 0    | 0     | 0    | 0    | 738   |
| MiscVal       | 1460 | 43.489041 | 496.12302 | 0    | 0     | 0    | 0    | 15500 |
| MoSold        | 1460 | 6.3219178 | 2.7036262 | 1    | 5     | 6    | 8    | 12    |
| YrSold        | 1460 | 2007.8158 | 1.3280951 | 2006 | 2007  | 2008 | 2009 | 2010  |

## 4.2 Correlation analysis of the variables:

Here, we display and describe 12 variables and their correlation with other variables. Jupyter notebook has a more detailed correlation of 38 variables.

|              | MSSub Class | LotFrontage | LotArea | OverallQual | OverallCond | YearBuilt | YearRemodAdd | MasVnrArea | BsmtFinSF1 | BsmtFinSF2 | 1stFlrSF | 2ndFlrSF |
|--------------|-------------|-------------|---------|-------------|-------------|-----------|--------------|------------|------------|------------|----------|----------|
| MSSubClass   | 1           | -0.39       | -0.14   | 0.033       | -0.06       | 0.028     | 0.041        | 0.023      | -0.07      | -0.07      | -0.25    | 0.308    |
| LotFrontage  | -0.39       | 1           | 0.426   | 0.252       | -0.06       | 0.123     | 0.089        | 0.193      | 0.234      | 0.05       | 0.457    | 0.08     |
| LotArea      | -0.14       | 0.426       | 1       | 0.106       | -0.01       | 0.014     | 0.014        | 0.104      | 0.214      | 0.111      | 0.299    | 0.051    |
| OverallQual  | 0.033       | 0.252       | 0.106   | 1           | -0.09       | 0.572     | 0.551        | 0.412      | 0.24       | -0.06      | 0.476    | 0.295    |
| OverallCond  | -0.06       | -0.06       | -0.01   | -0.09       | 1           | -0.38     | 0.074        | -0.13      | -0.05      | 0.04       | -0.14    | 0.029    |
| YearBuilt    | 0.028       | 0.123       | 0.014   | 0.572       | -0.38       | 1         | 0.593        | 0.316      | 0.25       | -0.05      | 0.282    | 0.01     |
| YearRemodAdd | 0.041       | 0.089       | 0.014   | 0.551       | 0.074       | 0.593     | 1            | 0.18       | 0.128      | -0.07      | 0.24     | 0.14     |
| MasVnrArea   | 0.023       | 0.193       | 0.104   | 0.412       | -0.13       | 0.316     | 0.18         | 1          | 0.265      | -0.07      | 0.345    | 0.175    |
| BsmtFinSF1   | -0.07       | 0.234       | 0.214   | 0.24        | -0.05       | 0.25      | 0.128        | 0.265      | 1          | -0.05      | 0.446    | -0.14    |
| BsmtFinSF2   | -0.07       | 0.05        | 0.111   | -0.06       | 0.04        | -0.05     | -0.07        | -0.07      | -0.05      | 1          | 0.097    | -0.1     |
| 1stFlrSF     | -0.25       | 0.457       | 0.299   | 0.476       | -0.14       | 0.282     | 0.24         | 0.345      | 0.446      | 0.097      | 1        | -0.2     |
| 2ndFlrSF     | 0.308       | 0.08        | 0.051   | 0.295       | 0.029       | 0.01      | 0.14         | 0.175      | -0.14      | -0.1       | -0.2     | 1        |

**MSSubClass:** This variable has a negative correlation of -0.39 with LotFrontage, which means that as the type of dwelling becomes more specific, the lot frontage tends to decrease slightly. It also has a negative correlation of -0.25 with 1stFlrSF, which suggests that more specific dwelling types tend to have slightly smaller first floor areas.

**LotFrontage:** This variable has a positive correlation of 0.426 with LotArea, which means that as the lot frontage increases, the lot area tends to increase as well. It also has a negative correlation of -0.06 with OverallCond, suggesting that houses with larger street frontage may not necessarily be in better condition.

**LotArea:** This variable has a moderate positive correlation of 0.426 with LotFrontage, indicating that as the lot area increases, so does the lot frontage. It also has a positive correlation of 0.299 with 1stFlrSF and 0.308 with 2ndFlrSF, suggesting that larger lot areas tend to correspond with larger first and second floor areas.

**OverallQual:** This variable has the strongest positive correlations with YearBuilt (0.572) and YearRemodAdd (0.551), indicating that as the quality of the house increases, it tends to be newer or have been recently remodeled. It also has positive correlations with MasVnrArea (0.412), BsmtFinSF1 (0.24), and 1stFlrSF (0.476), suggesting that higher quality homes tend to have larger finished basement areas, larger first floor areas, and more masonry veneer on the exterior.

**OverallCond:** This variable has a negative correlation of -0.38 with YearBuilt, suggesting that houses in worse overall condition tend to be older. It also has a positive correlation of 0.04 with BsmtFinSF2, indicating that houses in worse condition may have more unfinished basement areas.

**YearBuilt:** This variable has the strongest positive correlation with YearRemodAdd (0.593), suggesting that houses built more recently tend to have been recently remodeled. It also has positive correlations with MasVnrArea (0.316), BsmtFinSF1 (0.25), and 1stFlrSF (0.282), indicating that newer houses tend to have larger finished basement areas, more masonry veneer on the exterior, and larger first floor areas.

**YearRemodAdd:** This variable has a positive correlation of 0.18 with MasVnrArea and 0.128 with BsmtFinSF1, suggesting that houses that have been recently remodeled tend to have larger masonry veneer areas and finished basement areas.

**MasVnrArea:** This variable has positive correlations with BsmtFinSF1 (0.265) and 1stFlrSF (0.345), indicating that houses with larger masonry veneer areas tend to have larger finished basement and first floor areas.

**BsmtFinSF1:** This variable has a positive correlation of 0.446 with 1stFlrSF, suggesting that houses with larger finished basement areas tend to have larger first floor areas.

**BsmtFinSF2:** This variable has a positive correlation of 0.097 with 1stFlrSF, indicating that houses with more unfinished basement areas tend to have slightly larger first floor areas.

**1stFlrSF:** This variable represents the total square feet of the first floor of the house. It has a moderately strong positive correlation with the total square feet of the house (LotArea), as well as with other variables such as the overall quality of the house (OverallQual) and the square footage of finished basement area with a higher quality rating (BsmtFinSF1). This suggests that larger first floors tend to be associated with larger houses and higher quality, and may be a good predictor of house prices.

**2ndFlrSF:** This variable represents the total square feet of the second floor of the house. It has a moderate positive correlation with OverallQual and a weak positive correlation with other variables such as YearRemodAdd and MasVnrArea. This suggests that larger second floors tend to be associated with higher quality houses, and may be a good predictor of house prices in combination with other variables such as 1stFlrSF.

#### 4.3 Multiple linear regression model results

In this study, we have used a multiple linear regression model to predict the sale price of the houses. We have fitted the model and obtained necessary output.

**Coefficients:** These are the estimated values of the coefficients ( $\beta_0, \beta_1, \beta_2, \dots, \beta_n$ ) in the model. Each coefficient represents the change in the dependent variable associated with a one-unit change in the corresponding independent variable, while holding all other independent variables constant.

| Table: Coefficients of the regression model |              |      |              |      |              |      |              |
|---------------------------------------------|--------------|------|--------------|------|--------------|------|--------------|
| feat                                        | coefficients | feat | coefficients | feat | coefficients | feat | coefficients |
| 0                                           | 4.94397E-12  | 34   | -0.452484573 | 67   | 17.50620186  | 100  | -40.91541602 |
| 1                                           | -5.45697E-12 | 35   | -0.452484573 | 68   | 30.85997     | 101  | -40.91541602 |

|    |              |    |             |    |              |     |              |
|----|--------------|----|-------------|----|--------------|-----|--------------|
| 2  | 8.18545E-11  | 36 | 14.86877153 | 69 | 9.790654272  | 102 | -12.83976881 |
| 3  | 5.45697E-11  | 37 | 14.86877153 | 70 | 50.23193302  | 103 | -8.196285777 |
| 4  | -1.27329E-11 | 38 | 5.649421434 | 71 | 7.416377922  | 104 | -24.91680046 |
| 5  | 2.54659E-11  | 39 | 1.947780918 | 72 | 8.687485365  | 105 | -25.41560327 |
| 6  | 2.27374E-11  | 40 | 5.482642585 | 73 | 46.24625583  | 106 | 3.97904E-12  |
| 7  | -2.00089E-11 | 41 | 10.28571178 | 74 | 15.78944493  | 107 | 2.50111E-12  |
| 8  | 23.38125492  | 42 | 7.222655018 | 75 | 20.72083883  | 108 | 2.27374E-12  |
| 9  | 26.40172369  | 43 | 15.98923275 | 76 | -14.09106892 | 109 | 7.78755E-12  |
| 10 | 2.940777863  | 44 | 9.66911959  | 77 | -7.409503445 | 110 | 1.72804E-11  |
| 11 | -31.78163608 | 45 | 13.30196502 | 78 | -35.86608583 | 111 | -14.325082   |
| 12 | 1.90994E-11  | 46 | 11.91397631 | 79 | -36.89435491 | 112 | 4.964843078  |
| 13 | -9.09495E-12 | 47 | 8.276556025 | 80 | 6.451954694  | 113 | 2.71341189   |
| 14 | -1.68257E-11 | 48 | 5.649421434 | 81 | 19.54120904  | 114 | 1.260834148  |
| 15 | -2.59206E-11 | 49 | 9.484357025 | 82 | 42.74398575  | 115 | 8.428728007  |
| 17 | 2.95586E-11  | 50 | 19.01389472 | 83 | 3.727593552  | 116 | 11.70039364  |
| 18 | -2.81943E-11 | 51 | 4.121936605 | 84 | 46.61813762  | 117 | -4.259013472 |
| 19 | 3.91083E-11  | 52 | 11.47746581 | 85 | 0.081087689  | 118 | -2.179327607 |
| 20 | 3.04681E-11  | 53 | 8.700204746 | 86 | 0.133972439  | 119 | -1.923310429 |
| 21 | 9.32232E-12  | 54 | 11.77071379 | 87 | 0.134269763  | 120 | -8.038809196 |
| 22 | 1.45519E-11  | 55 | 14.07243874 | 88 | 0.034368743  | 121 | -12.55366735 |
| 23 | 1.79625E-11  | 56 | 6.831912602 | 89 | 0.017291738  | 122 | -7.404874871 |
| 24 | 1.56319E-12  | 57 | 11.55164469 | 90 | 0.012239712  | 123 | -14.3361356  |
| 25 | 1.22441E-10  | 58 | 10.3703024  | 91 | 6.25278E-12  | 124 | -5.93445E-11 |
| 26 | -7.53175E-12 | 59 | 12.39906507 | 92 | 3.86535E-12  | 125 | -5.70708E-11 |
| 27 | -9.37916E-12 | 60 | 6.831912602 | 93 | 3.75167E-12  | 126 | -6.76437E-11 |
| 28 | 4.54747E-13  | 61 | 8.384708016 | 94 | 1.65983E-11  | 127 | 1.79057E-12  |
| 29 | -5.486944552 | 62 | 4.553830719 | 95 | 15.95489667  | 128 | 5.68434E-13  |
| 30 | -13.72114496 | 63 | 24.18257328 | 96 | 5.747539885  | 129 | 5.68434E-12  |
| 31 | -6.926122338 | 64 | 9.406194127 | 97 | 11.84761194  | 130 | 1.06297E-11  |
| 32 | -27.17461951 | 65 | 12.09260802 | 98 | 0.834926206  | 131 | 79415.29189  |
| 33 | -23.71022172 | 66 | 11.03153962 | 99 | 14.52720884  |     |              |

**Standard errors:** These are measures of the variability in the estimates of the coefficients. They indicate how much the coefficient estimates are likely to vary from one sample to another.

**(a) R-squared:** This is a measure of the goodness of fit of the model. It represents the proportion of variance in the dependent variable that is explained by the independent variables in the model. The value of R-squared ranges from 0 to 1, with higher values indicating a better fit. In this case, the R-squared value is 0.9999999376532556, which means that the model explains about 99.99% of the variability in the data, indicating a very good fit.

**(b) Mean Absolute Error (MAE):** MAE is the average of the absolute differences between the predicted values and the actual values. It gives us an idea of how far off the predictions are on average. In this study, the MAE was found to be 1.1787, indicating that, on average, the predictions were off by 1.1787 units from the actual values.

**(c) Root Mean Squared Error (RMSE):** RMSE is the square root of the average of the squared differences between the predicted values and the actual values. RMSE penalizes larger errors more than smaller errors, and it's a widely used metric for regression models, as it has the same unit as the target variable. In this study, the RMSE was found to be 20.1420, indicating that, on average, the predictions were off by 20.1420 units from the actual values.

**(d) Mean Squared Error (MSE):** MSE is the average of the squared differences between the predicted values and the actual values. It's another widely used metric for regression models, and it's

closely related to RMSE. In this study, the MSE was found to be 405.7016, which is larger than the MAE value.

**Intercept:** This is the value of the dependent variable when all independent variables are set to zero. In our model this value calculated as 1.809e5.

To interpret the results of our multiple linear regression model, we should look at the coefficients, their standard errors, t-values, and p-values. In our case, coefficients are statistically significant (because it has a low p-value), we can say that there is a significant relationship between that independent variables that we have selected and the target variable.

#### **4.4 Discussion of the results:**

Forecasting house prices is a significant challenge for the real estate industry. Multiple linear regression is a common method used to predict house prices based on various property features. In this report, we discuss the results of a multiple linear regression model on the OpenMLHouse\_Prices dataset, which contains 1460 observations with 80 variables.

#### **4.5 Model Performance:**

The model's performance was evaluated using three metrics:  $R^2$ , MAE, RMSE, and MSE. The  $R^2$  value of 0.9999993765 indicates that the model explains 99.99993765% of the variation in the dependent variable (house price) based on the independent variables (property features). This suggests that the model is an excellent fit for the dataset, and the independent variables explain most of the variability in the dependent variable.

The MAE (mean absolute error) of 1.1787 suggests that, on average, the MAE was found to be 1.1787, indicating that, on average, the predictions were off by 1.1787 units from the actual values.

The RMSE (root mean square error) of 20.1420 indicates that the model's predictions, on average, were off by 20.1420 units from the actual values.. The MSE (mean squared error) of 405.7016 also shows the model's accuracy, which is larger than the MAE value.

Overall, the model's performance metrics indicate that it is highly accurate in predicting house prices based on the given variables. However, it's important to note that no model is perfect, and there may be some outliers or errors in the dataset that could affect the model's predictions.

### **5.Interpretation of the results:**

Interpretation of results refers to explaining the meaning and significance of the statistical findings obtained from the analysis. In the case of a multiple linear regression model, the interpretation involves understanding the relationship between the independent variables and the dependent variable, as well as the effect of each independent variable on the dependent variable.

For example, if we consider the multiple linear regression model for predicting house prices, we can interpret the results as follows:

The coefficient for the "OverallQual" variable is positive and statistically significant, indicating that as the overall quality of the house increases, the price also increases.

The coefficient for the "GrLivArea" variable is positive and statistically significant, indicating that as the living area of the house increases, the price also increases.

The coefficient for the "YearBuilt" variable is positive and statistically significant, indicating that as the age of the house increases, the price also increases at a decreasing rate.

The coefficient for the "GarageCars" variable is positive and statistically significant, indicating that as the number of cars that can fit into the garage increases, the price also increases.

The coefficient for the "Neighborhood" variable is positive and statistically significant for some neighborhoods, indicating that certain neighborhoods have higher prices than others.

Overall, the multiple linear regression model helps to explain the variation in house prices based on the selected independent variables. The interpretation of the results can provide insights into the factors that affect the prices of houses and can be used to make predictions about the prices of similar houses in the future.

## 6. Limitations

While the results of the multiple linear regression analysis on the OpenMLHouse\_Prices dataset are promising, there are several limitations to this analysis that need to be considered. These limitations include:

**Limited dataset:** The OpenMLHouse\_Prices dataset contains only 1460 observations and 80 variables, which may not be sufficient to capture the full range of factors that affect house prices.

**Limited scope of variables:** The dataset only includes features related to the house itself and does not include external factors such as the state of the economy or the local housing market, which can also have a significant impact on house prices.

**Linearity assumption:** Multiple linear regression assumes that the relationship between the dependent variable and independent variables is linear. However, in reality, the relationship may be more complex and nonlinear, which can lead to inaccurate predictions.

**Outliers:** The presence of outliers in the dataset can significantly affect the results of the multiple linear regression analysis, leading to biased coefficient estimates and inaccurate predictions.

**Causation vs correlation:** Multiple linear regression can only establish correlation between variables, but cannot establish causation. Therefore, it is important to exercise caution when interpreting the results of the analysis and avoid making causal claims based solely on the correlation between variables.

**Model Overfitting:** It is possible that the model has overfitted on the training data, resulting in a high R-squared value but poor performance on new, unseen data. Careful validation and testing of the model on new data is needed to ensure that it can generalize well to new situations.

## 7. Conclusion:

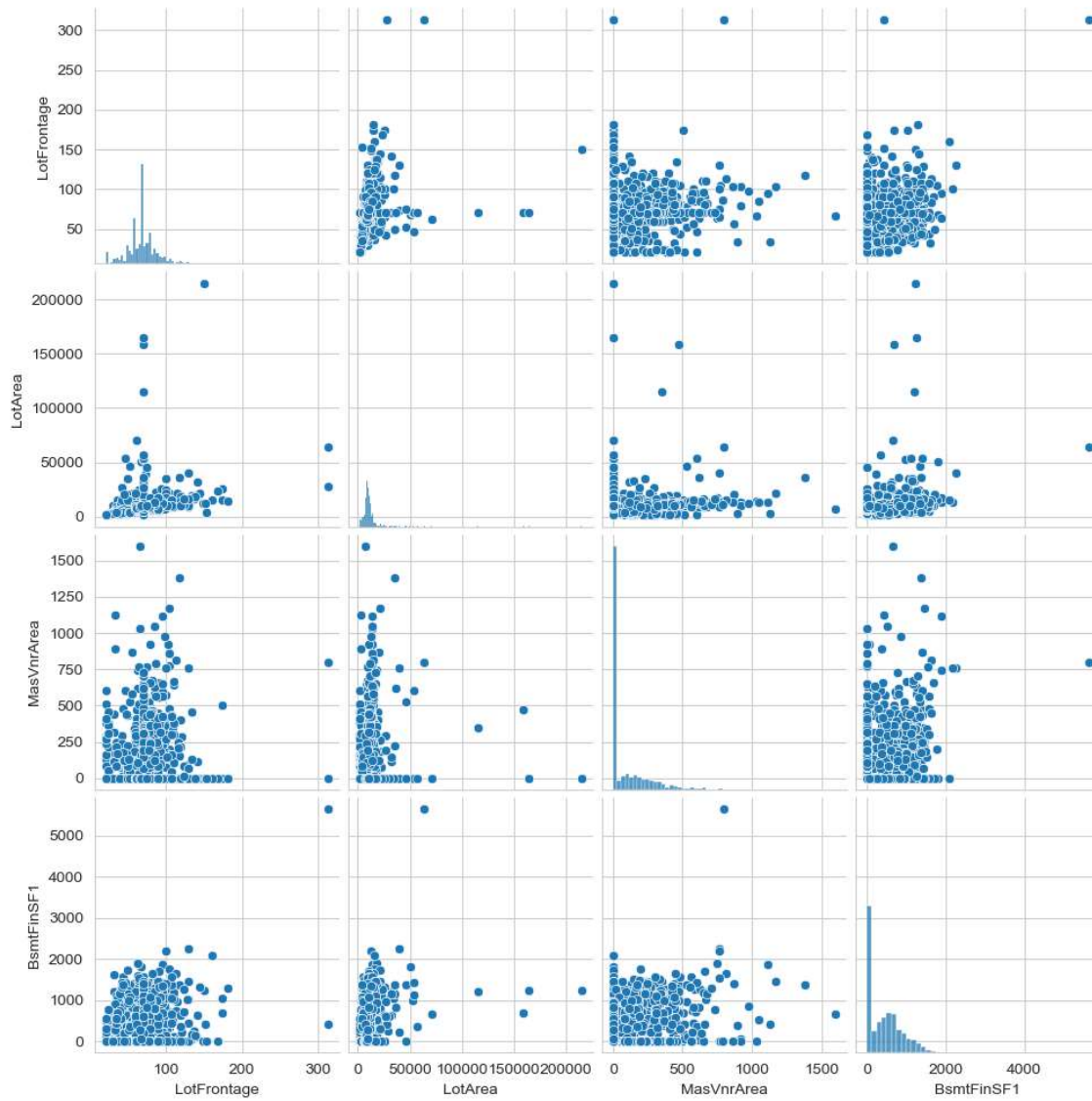
In conclusion, the multiple linear regression model used to forecast house prices on the OpenMLHouse\_Prices dataset has produced highly accurate results, with an  $R^2$  value of 0.99999993765 and low MAE, RMSE, and MSE values. This suggests that the model is an excellent fit for the dataset and can predict house prices based on the given variables accurately. However, it is crucial to continue evaluating the model's accuracy and validity, as the real estate market is constantly changing, and new variables may need to be added or removed from the model to ensure it remains up-to-date and relevant.

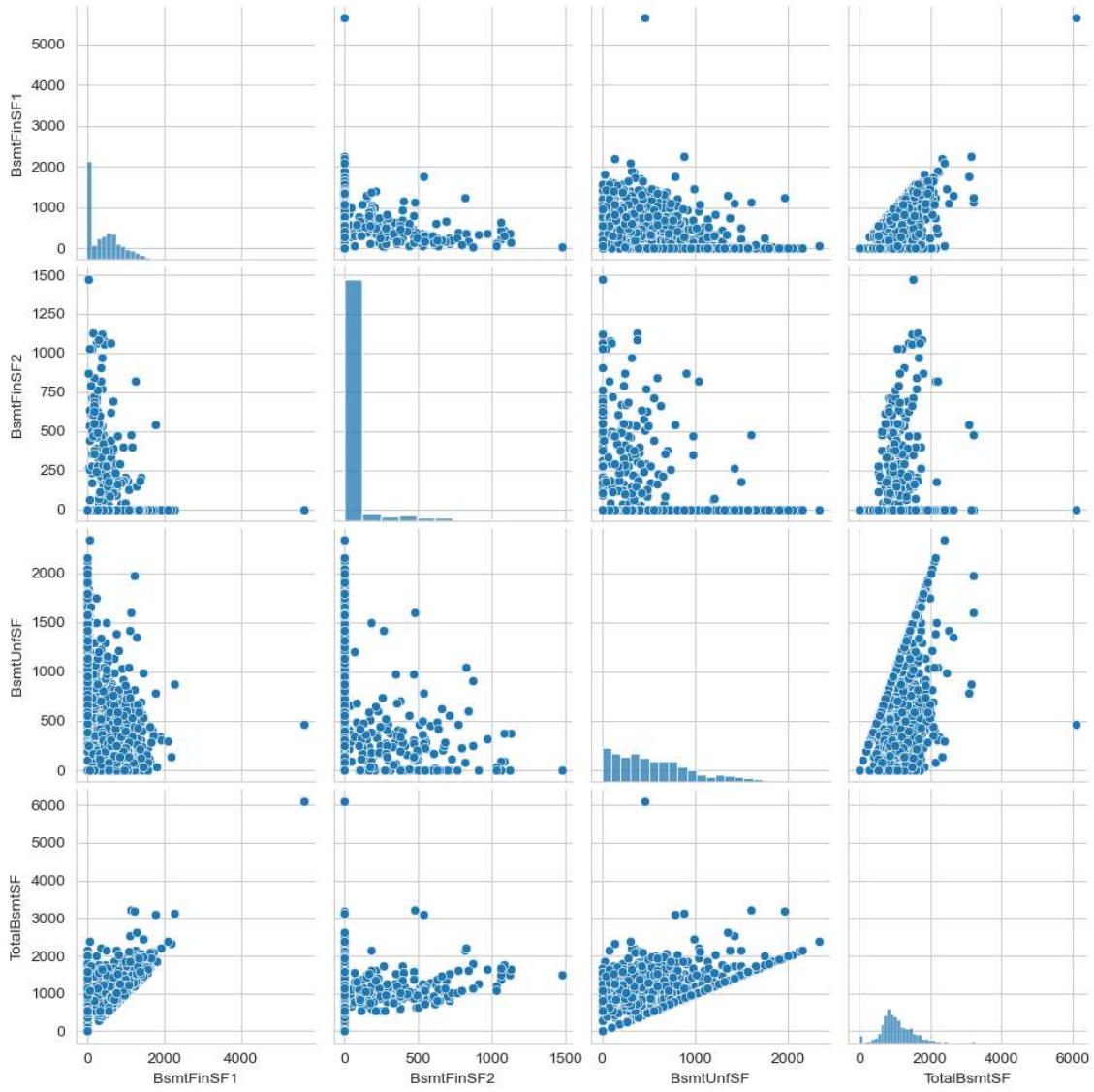
## 8. References:

- [1] Hair Jr, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2019). Multivariate data analysis. Cengage Learning.
- [2] <https://www.openml.org/search?type=data&status=active&id=42165>
- [3] Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). Introduction to linear regression analysis. John Wiley & Sons.
- [4] Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). Introduction to linear regression analysis. John Wiley & Sons.
- [5] Jeffrey M. Stanton (2001) Galton, Pearson, and the Peas: A Brief History of Linear Regression for Statistics Instructors, Journal of Statistics Education
- [6] Leung, F. K., & Tang, B. S. (2004). Comparison of ARIMA and spatial autoregressive models in forecasting Hong Kong real estate prices. Journal of Property Research, 21(2), 87-108.
- [7] Lee, D., Hong, T., Lee, J., & Lee, S. (2018). Identifying important variables for real estate price prediction using multiple linear regression: A case study of the Seoul metropolitan area. Sustainability
- [8] Anwar, M. N., & Islam, M. A. (2018). Modeling house prices in Ames, Iowa: An MLR approach. Journal of Economics and Sustainable Development
- [9] Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). Introduction to linear regression analysis (5th ed.). John Wiley & Sons.
- [10] Kutner, M. H., Nachtsheim, C. J., Neter, J., & Li, W. (2005). Applied linear statistical models (5th ed.). McGraw-Hill.
- [11] Simple linear regression: [https://en.wikipedia.org/wiki/Simple\\_linear\\_regression](https://en.wikipedia.org/wiki/Simple_linear_regression)
- [12] Multiple linear regression: [https://en.wikipedia.org/wiki/Linear\\_regression#Multiple\\_linear\\_regression](https://en.wikipedia.org/wiki/Linear_regression#Multiple_linear_regression)
- [13] Polynomial regression: [https://en.wikipedia.org/wiki/Polynomial\\_regression](https://en.wikipedia.org/wiki/Polynomial_regression)
- [14] Ridge regression: [https://en.wikipedia.org/wiki/Tikhonov\\_regularization](https://en.wikipedia.org/wiki/Tikhonov_regularization)
- [15] Lasso regression: [https://en.wikipedia.org/wiki/Lasso\\_\(statistics\)](https://en.wikipedia.org/wiki/Lasso_(statistics))
- [16] Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). Introduction to linear regression analysis. John Wiley & Sons
- [17] The Limitations of Linear Regression" by Towards Data Science
- [18] Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). Introduction to linear regression analysis (5th ed.). John Wiley & Sons.
- [19] European Union's General Data Protection Regulation (GDPR)

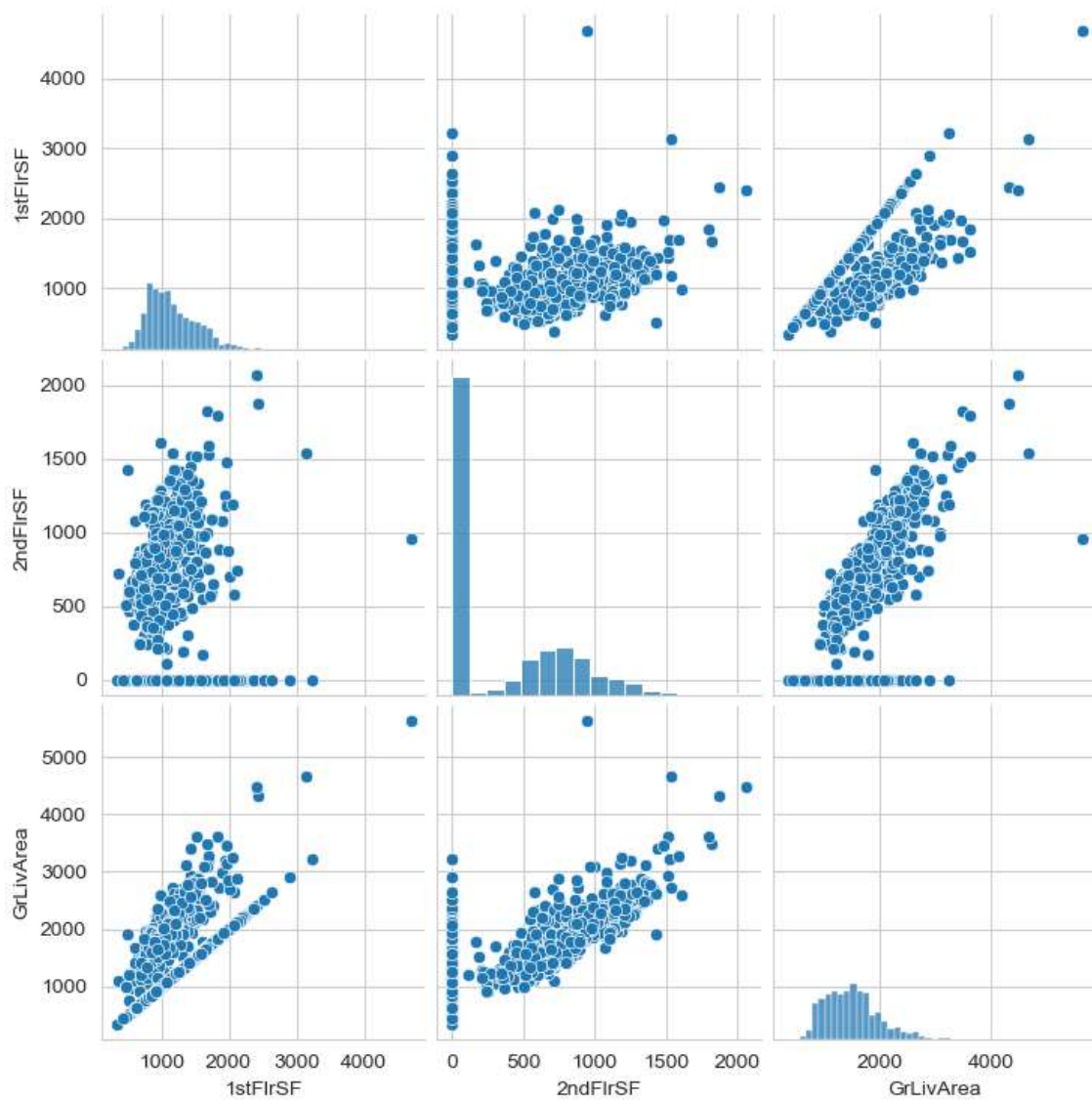
## 9. Appendices

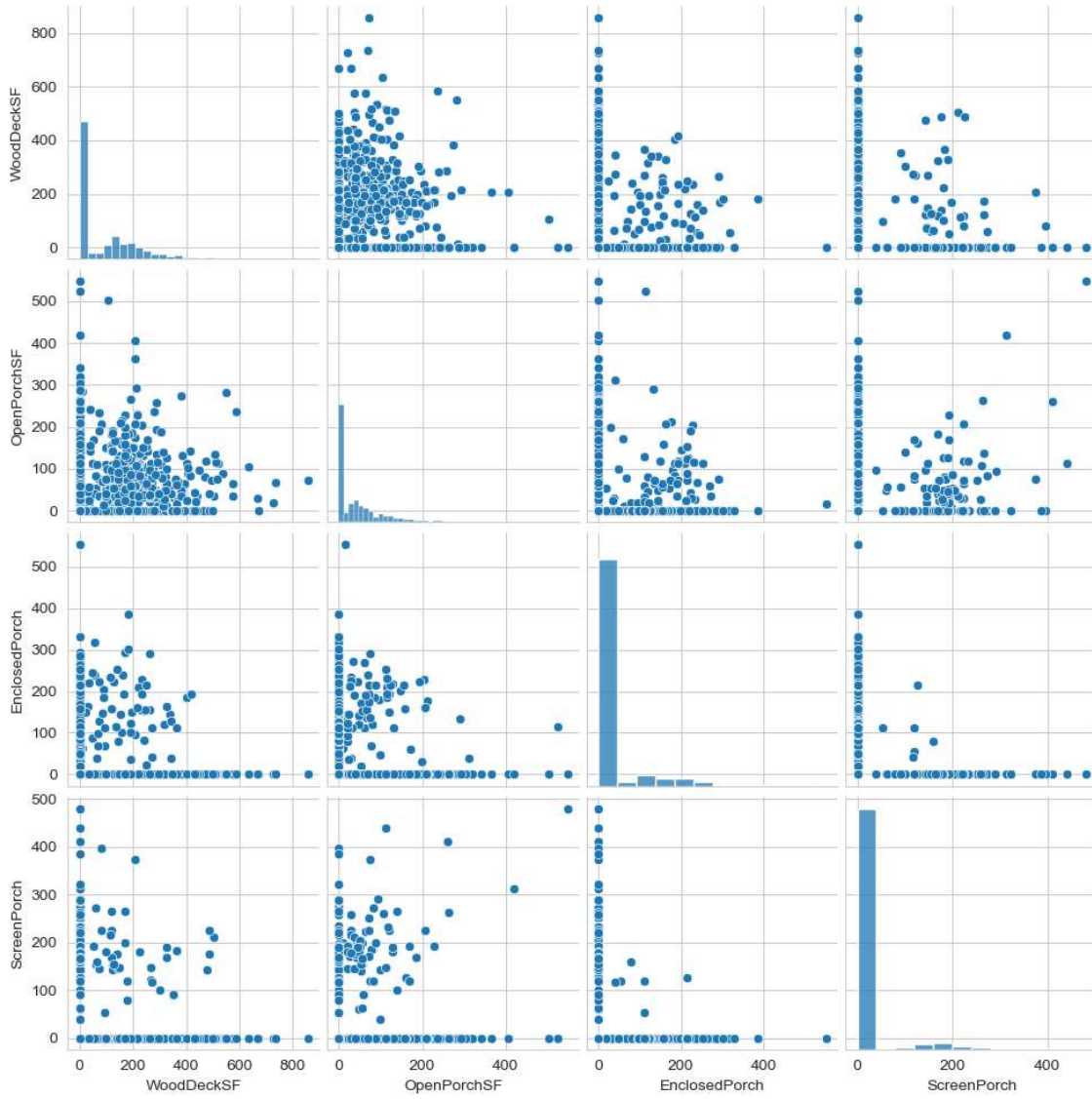
### 9.1 Model outputs and visualizations

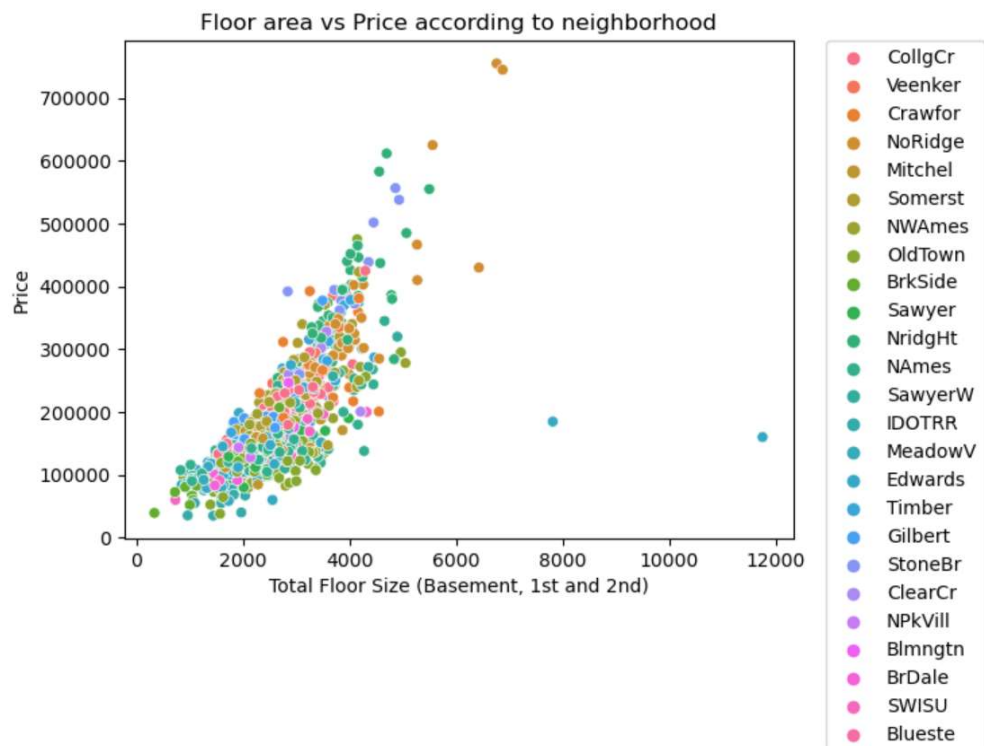
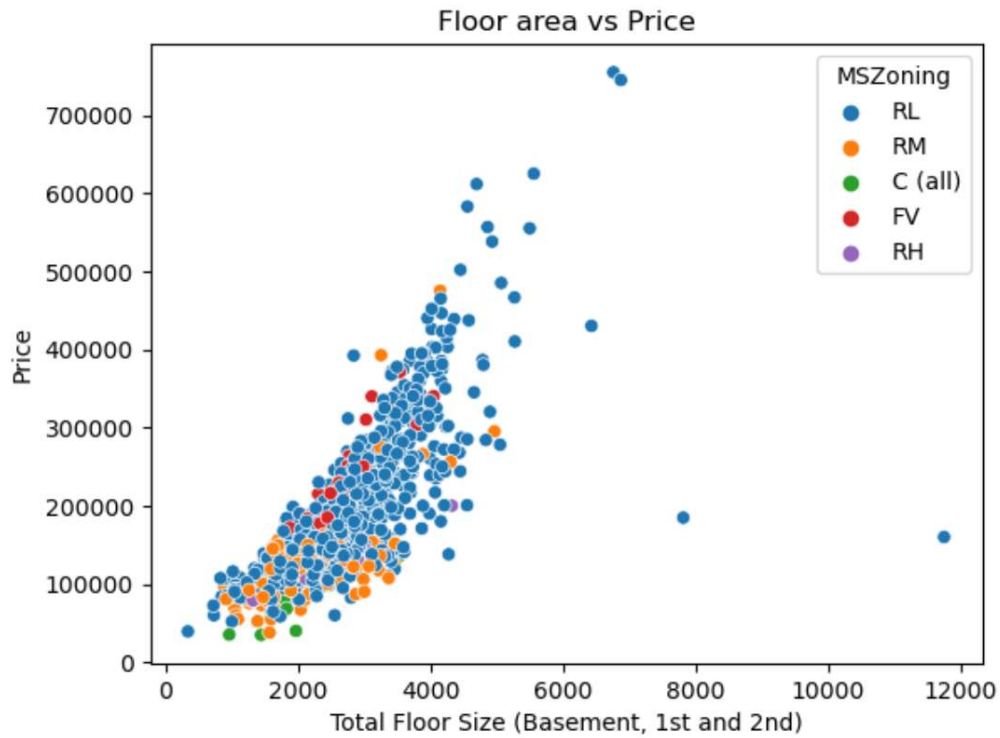


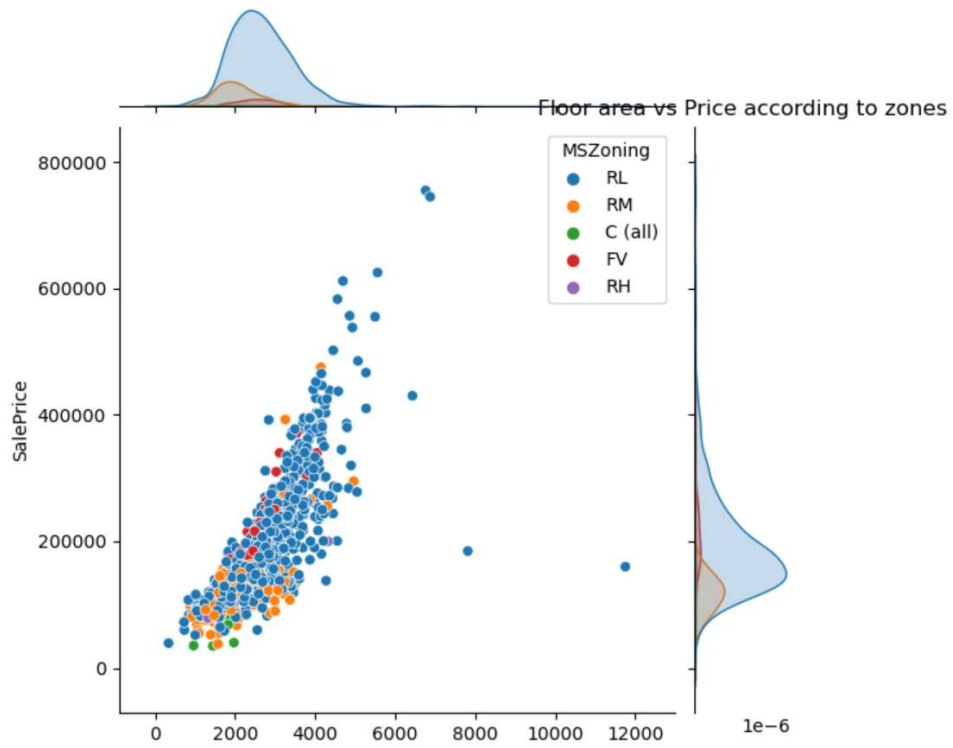


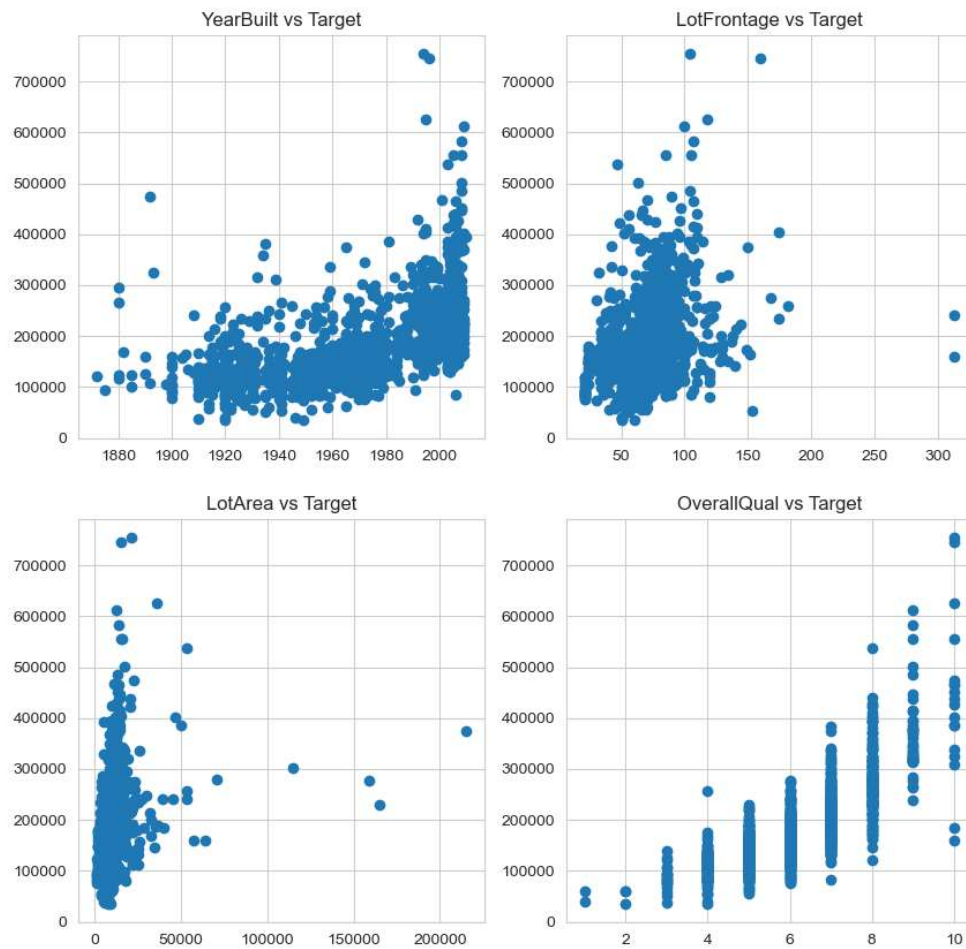


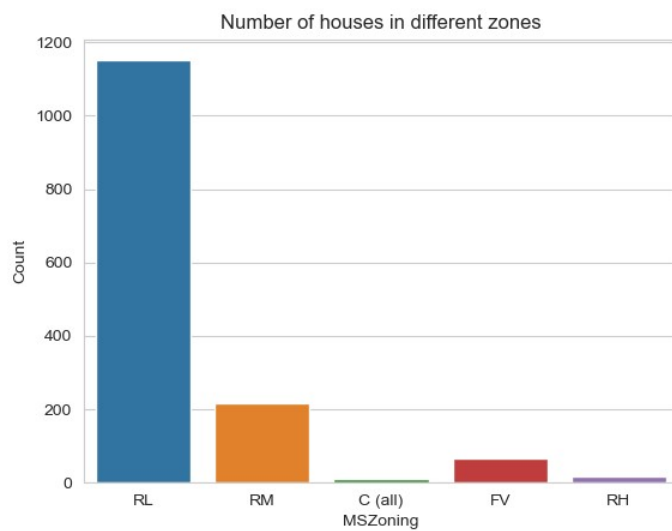
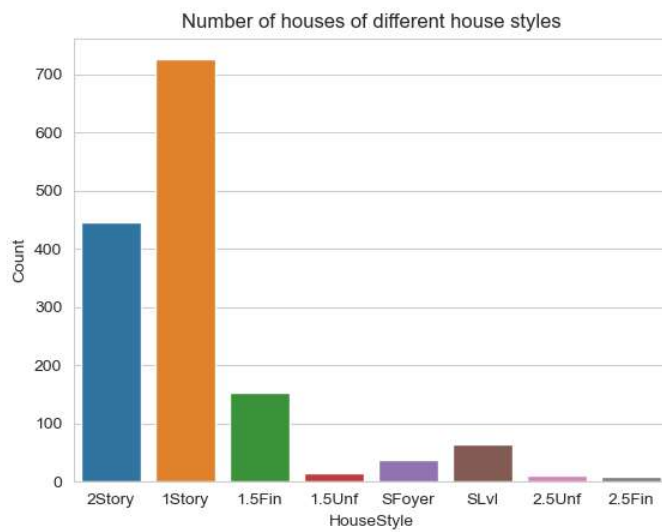
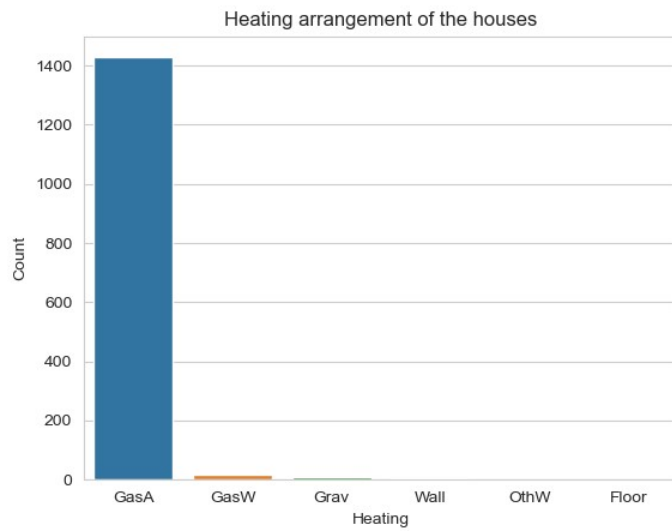


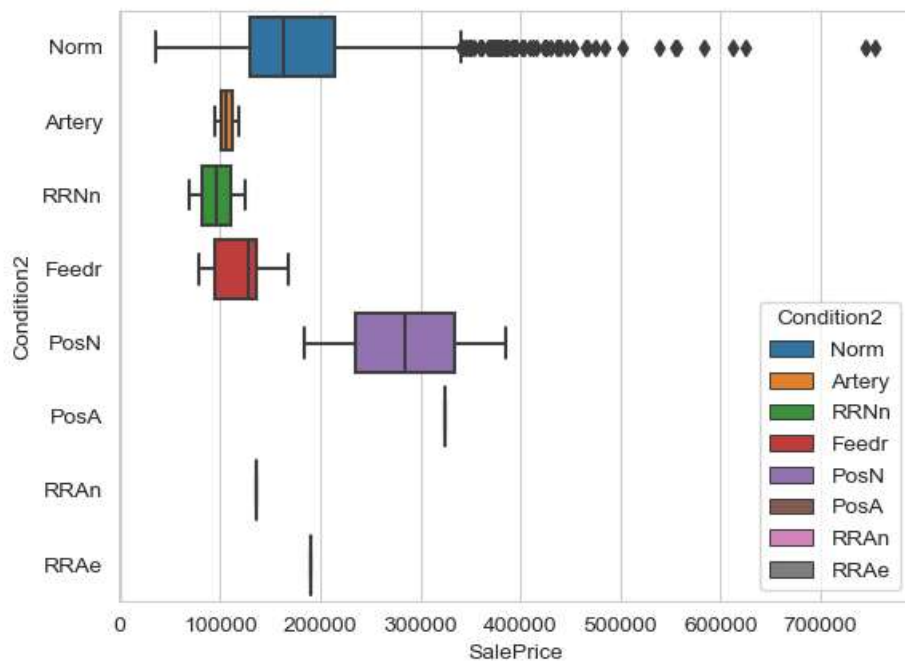
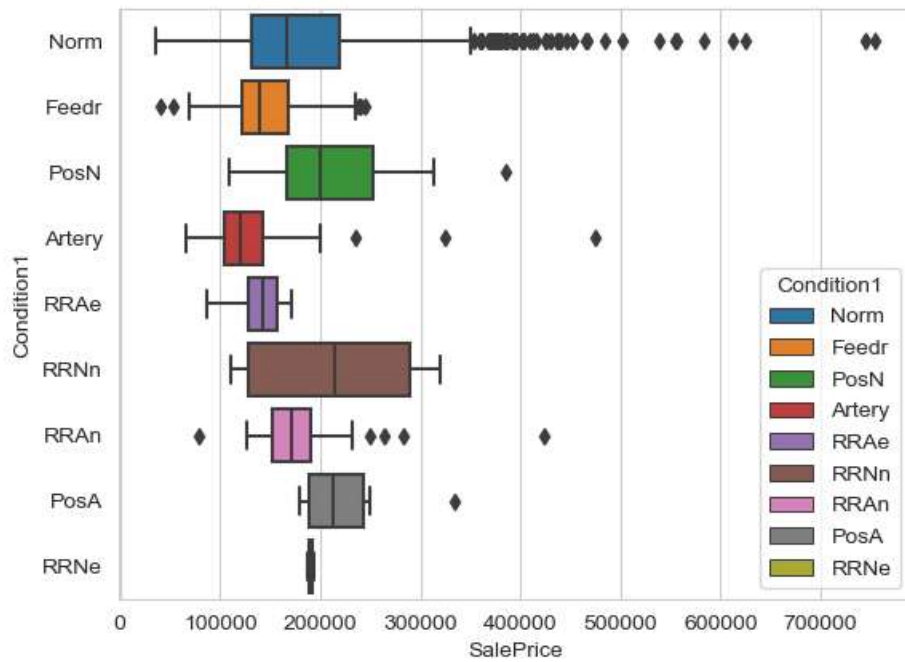


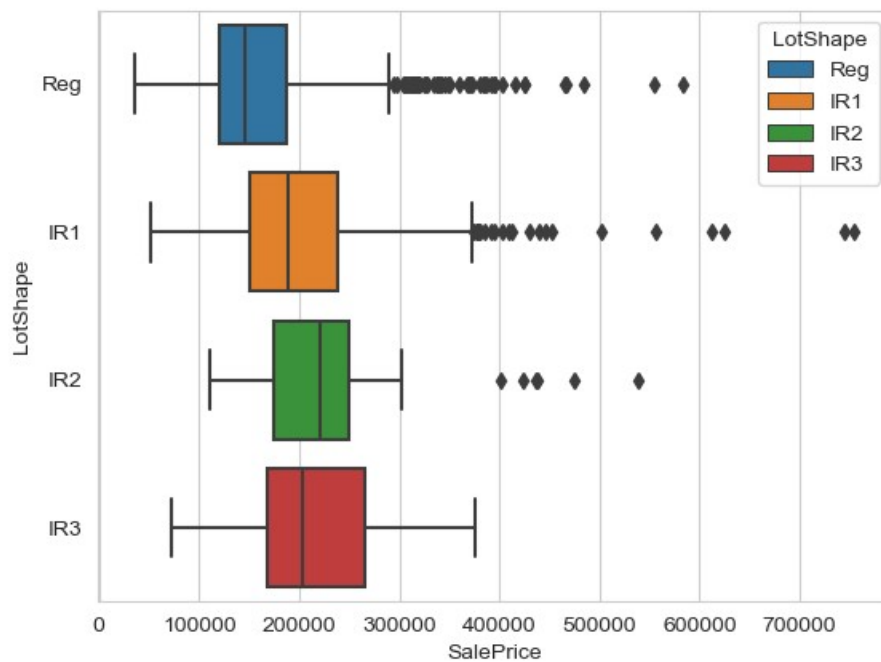
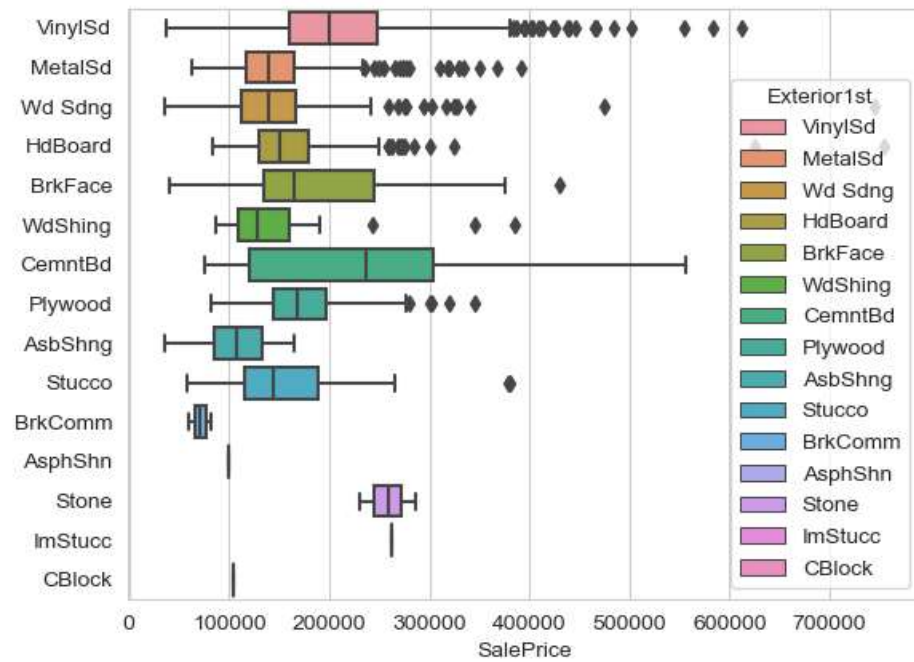




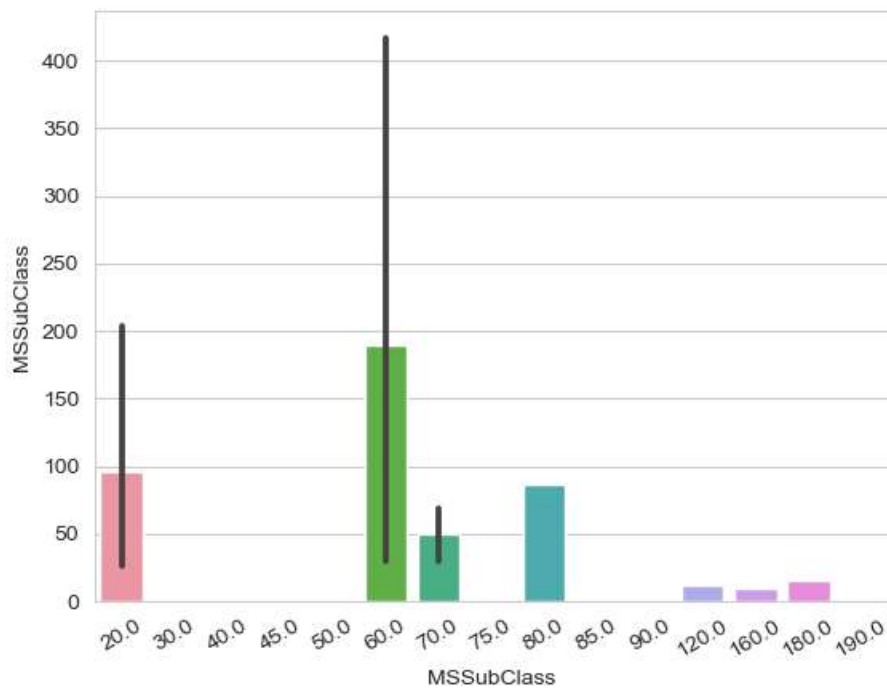
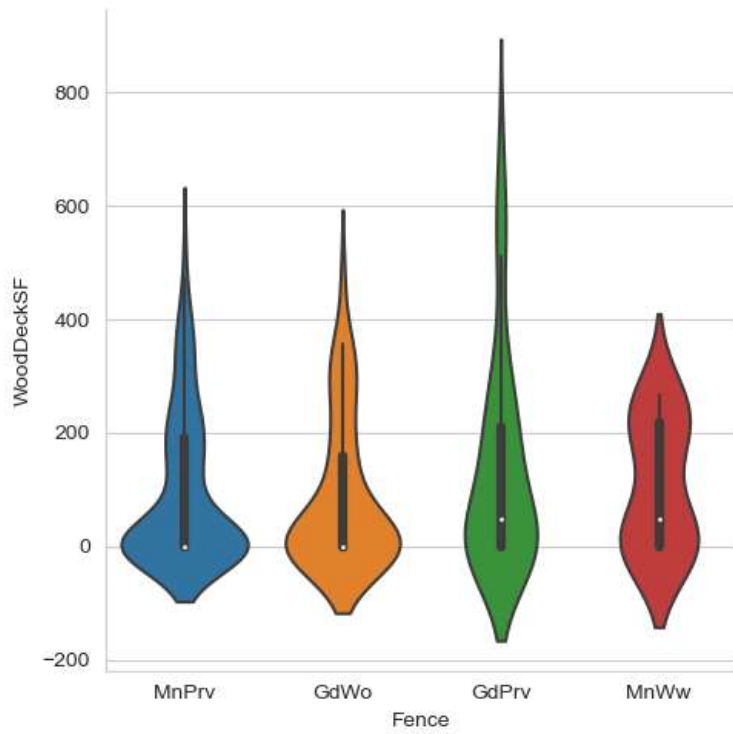


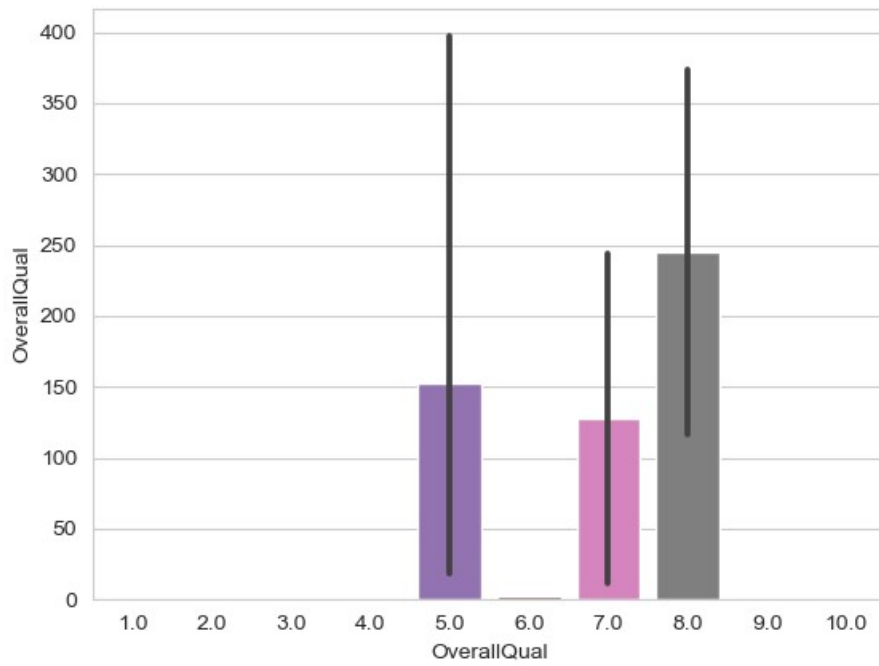


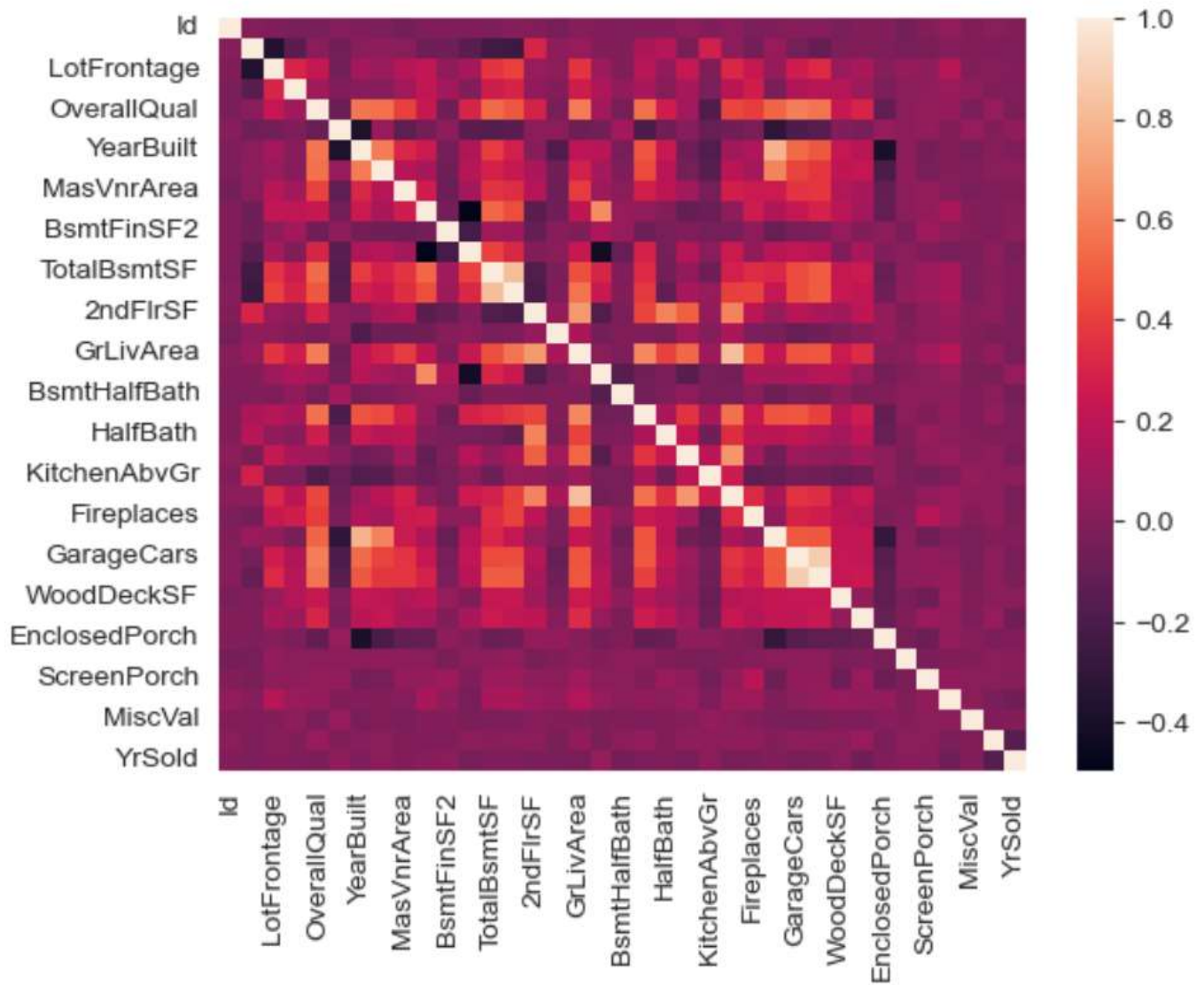


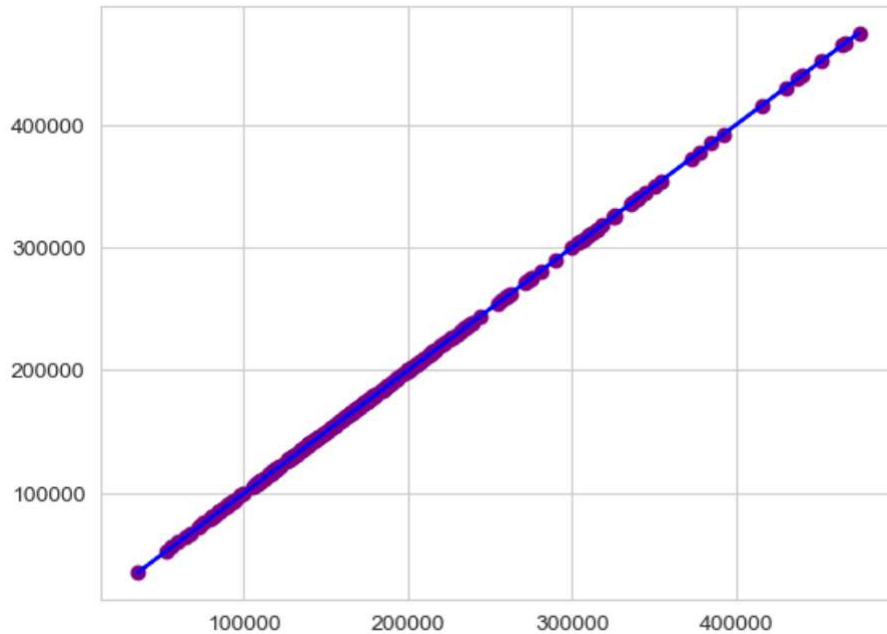












## 9.2 Research instruments

Jupyter notebook, pandas, numpy, matplotlib, seaborn, sklearn and various other resources.

## 9.3 Data sources

Data was obtained from Open ML online database. Here is the link of the dataset:  
<https://www.openml.org/search?type=data&status=active&id=42165>