

# Housing in Ames\*

Yanzun Jiang, Siyuan Lu, Yi Tang

October 1, 2024

## Table of contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Data Description</b>	<b>2</b>
<b>3</b>	<b>Ethics Discussion</b>	<b>4</b>
<b>4</b>	<b>Preliminary Results</b>	<b>4</b>
<b>A</b>	<b>Appendix</b>	<b>6</b>
A.1	Contributions . . . . .	6
	<b>References</b>	<b>7</b>

---

\*Code and data supporting this proposal is available at: <https://github.com/Stary54264/Housing-in-Ames>

# 1 Introduction

R Core Team (2023)

## 2 Data Description

Table 1: Preview of Data (First Half)

sale_price	lot_area	overall_qual	year_built	roof_style
215000	31770	6	1960	Hip
105000	11622	5	1961	Gable
172000	14267	6	1958	Hip
244000	11160	7	1968	Hip
189900	13830	5	1997	Gable
195500	9978	6	1998	Gable

Table 2: Preview of Data (Second Half)

mas_vnr_area	total_bsmt_sf	central_air	garage_area	misc_val
112	1080	Y	528	0
0	882	Y	730	0
108	1329	Y	312	12500
0	2110	Y	522	0
0	928	Y	482	0
20	926	Y	470	0

The Ames Housing dataset (Table 1, Table 2) was sourced from the **AmesHousing** package (Kuhn (2020)) in R (R Core Team (2023)). It was originally compiled by the Ames City Assessor’s Office through a comprehensive data dump of property tax records from 2006 to 2010, and it aimed to document residential property sales (De Cock (2011)). The dataset was initially designed for property tax assessments and general valuation, focusing on property characteristics such as lot area, the year built, and sale price. In contrast, this research aims to analyze how various property features influence house prices in Ames.

The dataset consists of 2930 observations and 82 variables relevant to understanding housing market dynamics. It was cleaned using **tidyverse** package (Wickham et al. (2019)). After cleaning, we selected 1 response variable, **sale\_price**, and 9 predictor variables: **lot\_area**, **overall\_qual**, **year\_built**, **roof\_style**, **mas\_vnr\_area**, **total\_bsmt\_sf**, **central\_air**, **garage\_area**, and **misc\_val**.

**sale\_price**: Price of the house in dollars  
**lot\_area**: Lot size in square feet  
**overall\_qual**: Rates the overall material and finish of the house  
**year\_built**: Original construction date  
**roof\_style**: Type of roof  
**mas\_vnr\_area**: Masonry veneer area in square feet  
**total\_bsmt\_sf**: Total square feet of basement area  
**central\_air**: Central air conditioning  
**garage\_area**: Size of garage in square feet  
**misc\_val**: Value of miscellaneous feature in dollars

These predictor variables all show the quality of the house, which will affect the price of the house directly. So, we believe there is a linear relationship between these predictor variables and the response variable.

Table 3: Summarize Table of Numerical Data

	Mean	Standard_Deviation	Median
sale_price	180425.31	79811.03	160000
lot_area	10143.13	7898.24	9434
overall_qual	6.09	1.41	6
year_built	1971.13	30.22	1973
mas_vnr_area	101.97	179.15	0
total_bsmt_sf	1050.52	440.66	990
garage_area	472.34	215.23	479
misc_val	51.07	568.76	0

From the summary table (Table 3), we can easily see that **mas\_vnr\_area** and **misc\_val** might be right-skewed since their mean is a lot greater than their median. An interesting point is that the standard deviation of **misc\_val** is quite large, which indicates that houses in Ames might differ significantly in miscellaneous features. By analyzing these variables, we aim to provide insights into how specific property characteristics affect housing prices in Ames, Iowa.

### 3 Ethics Discussion

### 4 Preliminary Results

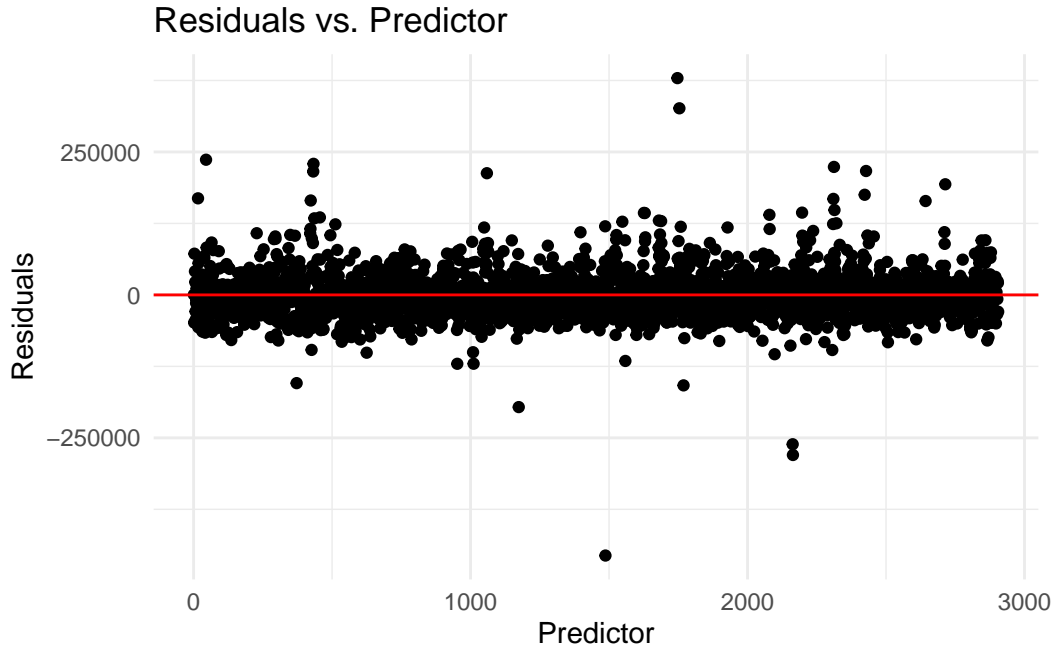


Figure 1: Uncorrelated errors testing

I plotted residuals against observation order to test uncorrelated errors. If errors are uncorrelated, residuals should scatter randomly around the horizontal axis without any clear pattern. From Figure 1, we can see that residuals appear to be randomly distributed across observations, which satisfy the assumption of independence in the linear regression model. But there still have some larger residuals, we can regard these points as potential outliers.

The constant variance assumption means we assume that the variance of the marginal distribution of responses does not change with the value of the predictor. The red loess curve added to the plot highlights this trend, showing a wider spread of residuals as the fitted values rise. From Figure 2, all data points lead to a fanning effect which against the assumption of constant variance(also known as homoskedasticity). The red loess curve added to the plot highlights this trend, showing a wider spread of residuals as the fitted values rise.

For linearity, the residuals are scattered relatively randomly around the horizontal line at zero, without showing any distinct or systematic patterns, such as curves or trends. Based on the Figure 2 plot, most of the residuals are centered around the horizontal line and no obvious non-linear patterns like U-shaped or inverted U-shaped trends in the plot, which satisfies the

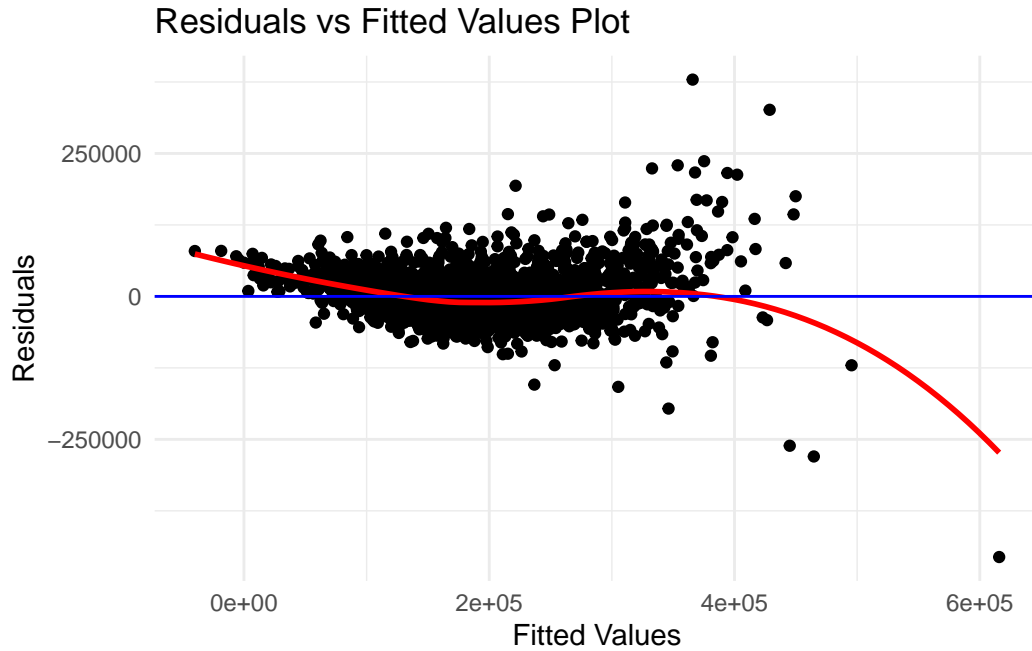


Figure 2: Constant variance testing

assumption of linearity, although there are a few outliers. This suggests that the relationship between the predictors (like `lot_area`, `overall_qual`, and `year_built`) and the response variable (sale price) is mostly linear.

In a well-behaved model that follows a normal distribution, the points in the Q-Q plot should align closely with the red reference line. Based on the Figure 3, most of the residuals are near the red line, which suggests that the central part of the residual distribution generally follows a normal distribution. Although there are some deviations at the tails, these are unlikely to have a major impact on the model's overall performance. Therefore, the Q-Q plot indicates that the data mostly satisfies the normality assumption.

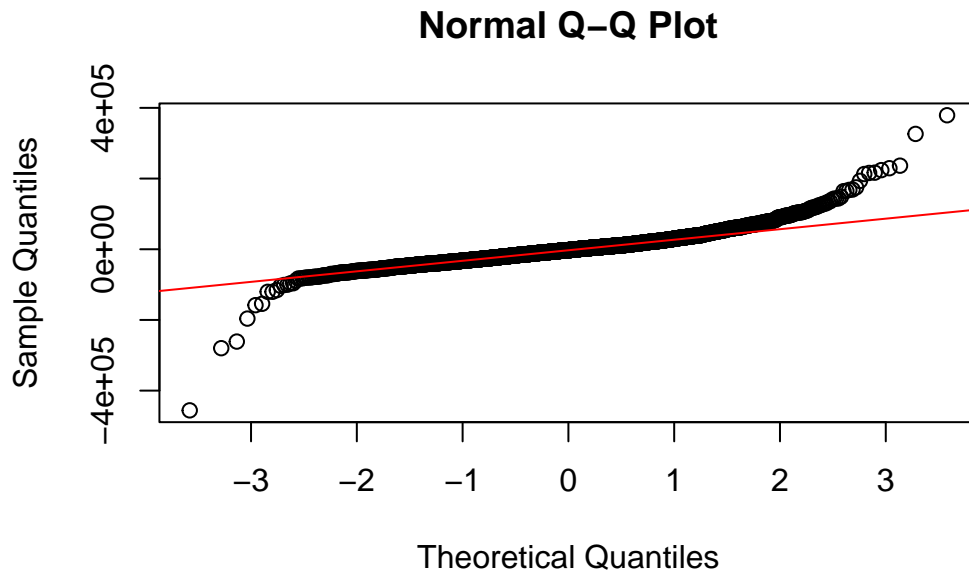


Figure 3: Q-Q Plot

## A Appendix

### A.1 Contributions

Group contribution is available at <https://github.com/Stary54264/Housing-in-Ames/graphs/contributors>. Below is a more specific version of group contribution.

Yanzun Jiang: Organized discussions and meetings; assigned tasks to group members; set up Github workspace for collaborating; downloaded data for setting up the linear regression model; cleaned data to make further analysis easier; completed Section 2 in the proposal; made the reference list; revised and combined group member's work together.

Siyuan Lu:

Yi Tang: Built linear regression model to predict house sale prices by using five key predictors in cleaned data. It assisted to understand the relationship between variables and ensure data meets key assumptions for statistical validity.

## References

- De Cock, Dean. 2011. “Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project.” *Journal of Statistics Education* 19 (3).
- Kuhn, Max. 2020. *AmesHousing: The Ames Iowa Housing Data*. <https://CRAN.R-project.org/package=AmesHousing>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Golemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.