

Housing in Ames*

Yanzun Jiang, Siyuan Lu, Yi Tang

October 3, 2024

Table of contents

1	Introduction	2
2	Data Description	2
3	Ethics Discussion	4
4	Preliminary Results	5
A	Appendix	9
A.1	Contributions	9
	References	10

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

1 Introduction

The primary research question we aim to answer is: What are the key factors that significantly influence house prices in Ames from 2006 to 2010? Sale price of the house is the response variable. Predictor variables include area, overall quality index, year of construction, house facilities, and value of miscellaneous feature. The hypothesis we state is that there is a statistically significant linear relationship between certain property characteristics (predictors) and the sale price of houses (response). This hypothesis will be tested using linear regression, which would be appropriate in this case. We use residual plots and a Q-Q plot to check the assumption.

Linear regression provides coefficients that quantify the relationship between each predictor and the response variable, making it easier to interpret the impact of each factor on house prices: estimates of how much the sale price is expected to change with a one-unit change in each predictor variable, holding other variables constant. Our primary goal is to understand the impact of each predictor on house prices, so the focus should be on interpretability instead of precision or accuracy.

We found several peer-reviewed articles that focus on similar problems with this paper. “Influencing Factors Analysis of House Prices Based on Multiple Linear Regression” concludes that housing prices are negatively correlated with housing completion costs, land acquisition prices, urban residents’ disposable income, and urban population density (Wang (2013)). This article provides some characteristics other than what we use that can also influence house price in national scope.

In “Flood Risk as a Price-Setting Factor in the Market Value of Real Property”, the analyzed market consider flood risk almost indifferent with house price compare to other factors in the analysis. (Cupal (2015)). This article offers a great example of our research since it uses multiple linear regression with some similar predictors as ours.

The research “House Price Prediction Using Hedonic Pricing Model and Machine Learning Techniques” shows that XGBoosting has higher accuracy in comparison to hedonic pricing model in prediction of property price. (Zaki et al. (2022)). This article provides an alternative method to study the relationship between multiple factors and house price.

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

Table 1: Preview of Data (First Half)

sale_price	lot_area	overall_qual	year_built	roof_style
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 indicate that houses in Ames might differs 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

Our data is collected from Ames City Assessor's Office (De Cock (2011)), then we cleaned the data to only keep some necessary key factors that is highly relavant to house prices. Raw and processed versions of the data from De Cock is published on Journal of Statistics Education in 2011. More detailed information about source of data is described in Section 2. The cleaned data we are using includes some detailed information about housing characteristics, but does not contain personal identifiers.

The analysis can provide deeper insights for stakeholders, including homeowners, potential buyers, real estate agents, and policymakers, to make better decisions about buying, selling and investing in real estate. The Ames housing dataset has gained popularity, especially in the context of academic projects and machine learning competitions. It is often considered a

modern alternative to the Boston Housing dataset. The dataset is well-vetted and trusted by the data science community for its comprehensiveness and relevance.

4 Preliminary Results

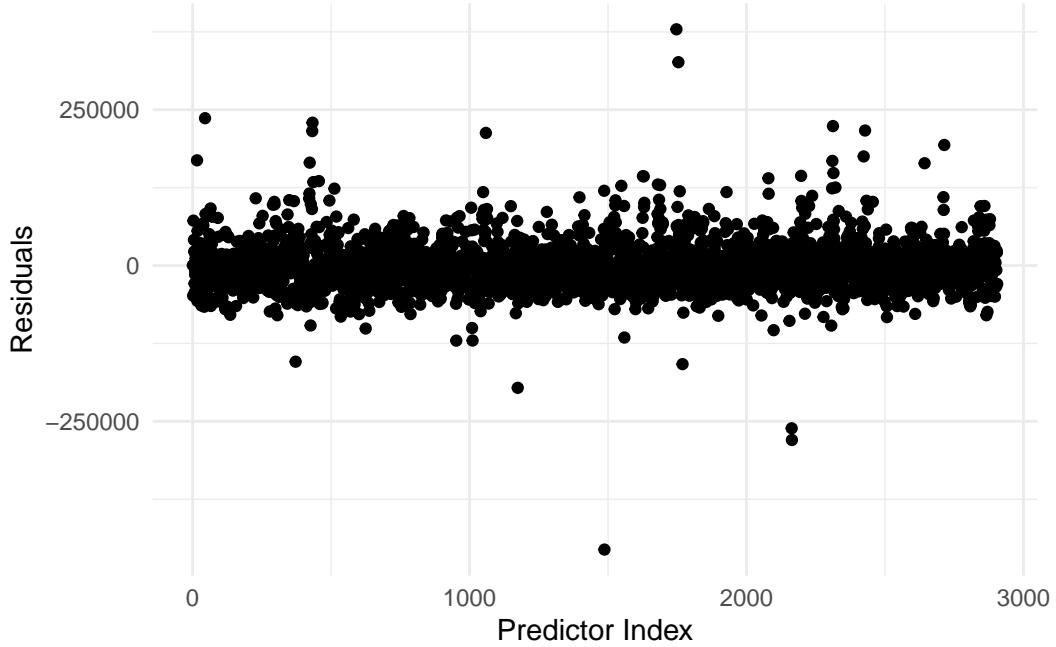


Figure 1: Residuals vs. Observation Index

From Figure 1, we can see that residuals appear to be randomly scattered along the x-axis without any special patterns, which means the dataset satisfies the uncorrelated errors assumption.

From Figure 2, all points are randomly scattered along a straight line, showing a linear relationship. Thus, it satisfies the linearity assumption. However, the graph shows a wider spread in errors with the increase of the fitted values. This means the dataset has violated the constant variance assumption. To solve this, we can use variance stabilizing transformation or box-cox transformation in the later analysis.

From Figure 3, most of the points lies on the diagonal, which gives evidence that the distribution of the errors are normally distributed. Therefore, it satisfies the normality assumption.

From Figure 4, we can see that all points are randomly scattered along the diagonal, which implies that the mean responses are a single function of a linear combination.

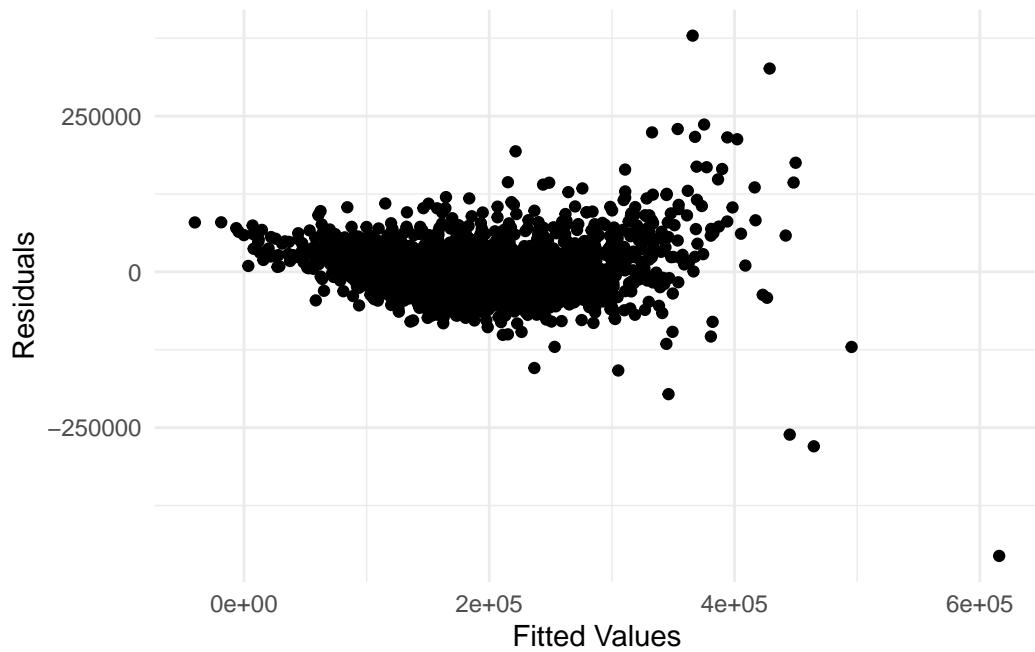


Figure 2: Residuals vs. Fitted Values

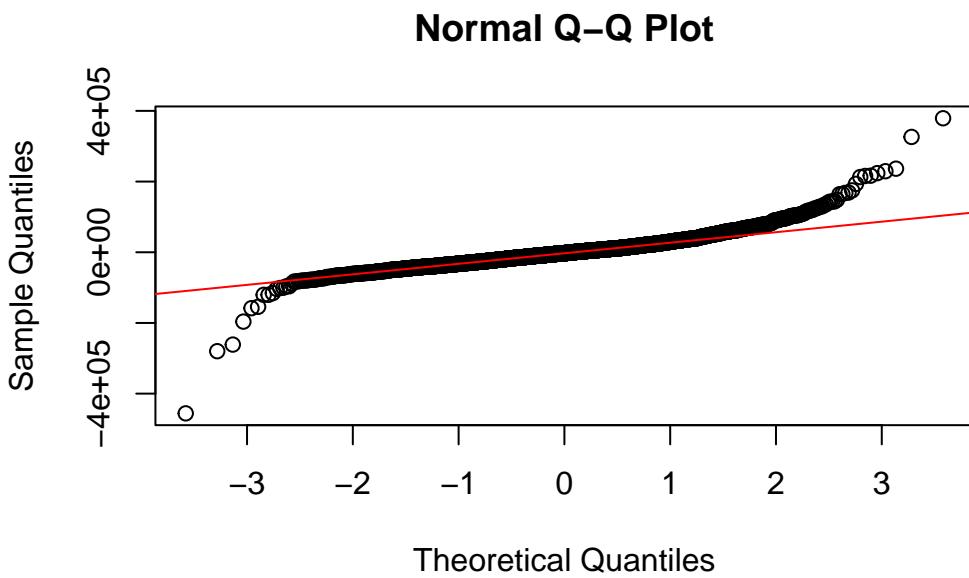


Figure 3: Q-Q Plot

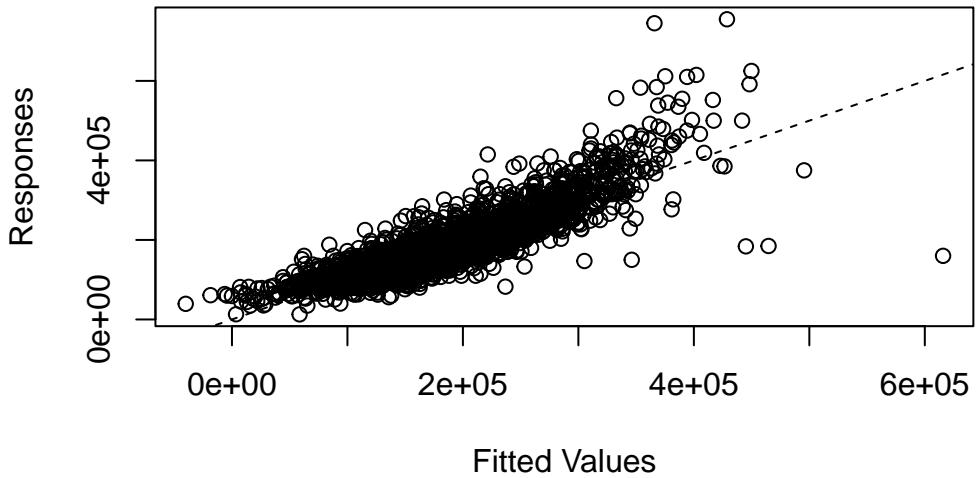


Figure 4: Responses vs. Fitted Values

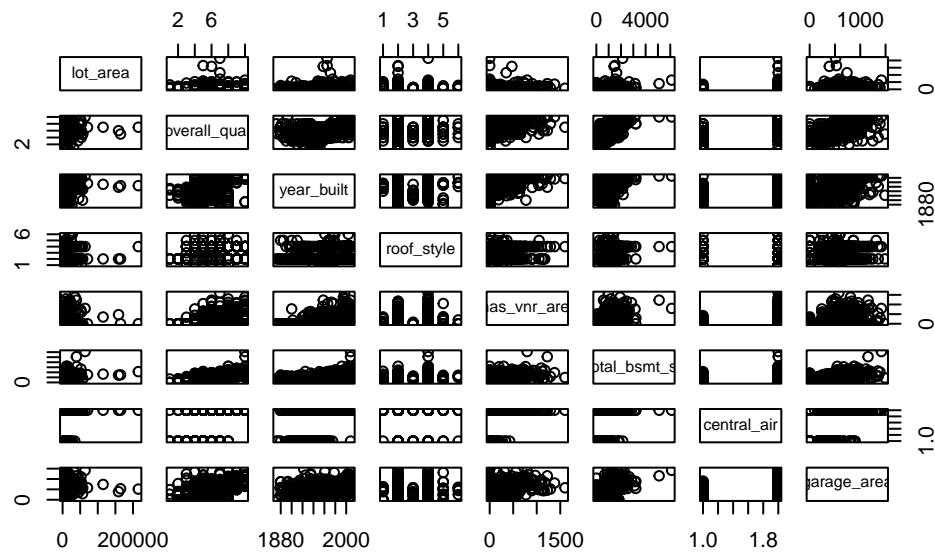


Figure 5: Pairwise Scatterplots of Predictors

Figure 5 shows no significant curves, which indicates the predictors are linearly or weakly associated. Thus, every predictor is related to each other in no more complicated way than linearly.

Although the constant variance assumption is violated, we could still gain insights from the analysis. From the analysis above, we can find out that when every predictor is 0, `roof_style` is “Flat”, and `central_air` is “N”, the expected value of `sale_price` would be -426801.9. A one-unit increase in `lot_area`, `overall_qual`, `year_built`, `mas_vnr_area`, `total_bsmt_sf`, `garage_area`, and `misc_val` would result in the expected value of `sale_price` to increase by 1.24, 28596.94, 189.49, 53.04, 28.98, 62.75, and -6.88 respectively. With `central_air`, the expected value of `sale_price` would decrease by 3355.89. Different `roof_type` would give different increase in the expected value of `sale_price`.

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; introduced the dataset; made the summary table; created file for R code; made the reference list; revised group member's work; combined group member's work together.
- Siyuan Lu: Set research question and hypothesis; searched and read peer-reviewed articles; introduced the project; checked data ethics.
- Yi Tang: Built linear regression model; checked conditions for performing linear regression; checked extra conditions for performing multiple linear regression; showed the results of the linear regression model.

References

- Cupal, Martin. 2015. “Flood Risk as a Price-Setting Factor in the Market Value of Real Property.” *Procedia Economics and Finance* 23: 658–64.
- 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/>.
- Wang, Cheng. 2013. “Influencing Factors Analysis of House Prices Based on Multiple Linear Regression.” *Applied Mechanics and Materials* 415: 722–25.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Zaki, John, Anand Nayyar, Surjeet Dalal, and Zainab H Ali. 2022. “House Price Prediction Using Hedonic Pricing Model and Machine Learning Techniques.” *Concurrency and Computation: Practice and Experience* 34 (27): e7342.