

Lamberto Alberici

Professor Yang

May 7th, 2025

Submitted in partial fulfilment of the requirements for BUSA 310: Business Analytics III.

Ames, Iowa Housing Price Prediction Using Machine Learning

Scott Ballentine, Lamberto Alberici, and Connor Agler

Abstract

This project applies multiple linear regression to predict housing prices in Ames, Iowa, using a dataset of 2,919 homes and 80 features. Two models were developed: one based on the most correlated variables and another refined through recursive feature selection. Key predictors included living area and overall quality. Additional analysis examined the impact of the 2008 housing crisis, neighbourhood-level price trends, and how mortgage interest rates influence affordability. The findings highlight how statistical modelling and data analysis can uncover meaningful patterns and improve understanding of housing market dynamics.

Introduction

The housing market plays one of the most important roles in the U.S. economy, shaping not only individuals' financial decisions, but also many broader aspects of growth, investment, and stability. Understanding the dynamics of all of the factors that affect the housing market is critical for buyers, sellers, real estate professionals, and policymakers.

This project aims to analyze and predict housing prices using machine learning and statistical methods, with focus on linear regression modeling. We seek to identify which variables have the most correlation with sales price and how these predictors can be adjusted to improve price forecasting. This analysis begins with careful data collection and using EDA, followed by the development and completion of regression analysis models to search for predictive accuracy. Beyond the core work of the project, it extends into three analyses that are designed to further open our insight into the housing market. These dive into how broader economic variables like mortgage interest rates interact with housing affordability, and how categorical factors like neighborhood can provide perspectives on market segmentation. The

project not only addresses the technical goals of predictive modeling but connects the findings to practical implications for home buyers and sellers seeking informed decision making, real estate firms, who are navigating competitive markets, and policymakers interested in improving affordability trends. Ultimately, the goal of this project is to demonstrate how data-driven approaches can show valuable insights and patterns in the housing market, demonstrating both predictive accuracy and actionable insights. By combining rigorous statistical methods with real-world application, this project highlights the role of business analytics in understanding one of the most dynamic and consequential industries in the United States.

Data and Methodology

This study seeks to answer the research question: How effective are machine learning models at predicting residential prices using multiple linear regression? The analysis pulls from a publicly available dataset from the Kaggle competition platform, including 2,919 observations and 80 variables that are related to residential property sales in Ames, Iowa. The 80 variables capture a wide range of property characteristics and are categorized into 3 main categories, physical features, locational characteristics, and neighborhood/amenities. Some of the main physical feature predictors include dwelling-specific traits like building age and overall quality. Overall quality may seem like a difficult variable to quantify, but it rates the overall quality and finish of the house on a scale from 1-10. Thus, overall quality is a categorical variable with a value of 10 meaning very excellent. The second grouping of predictors, locational characteristics, measures the proximity to services and markets. Neighborhood and amenity predictors reflect community attributes regarding infrastructure and other leisure activities. This

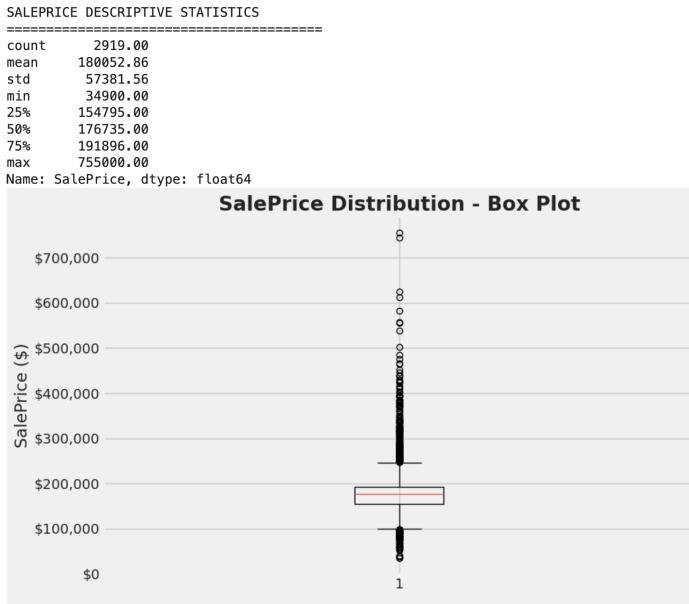
is a representative dataset because many variables are considered, those that directly impact the price of a home and other contextual factors that play a big role in pricing, too.

In general, the machine learning process takes inputs (X) and uses a model (e.g. linear regression) to produce an output (Y). The methodology used establishes a multiple linear regression (MLR) model based on ordinary least squares (OLS) estimates. The formula for MLR models is:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n + \varepsilon$$

Y , sale price, is the dependent and predicted value. $X_1 \dots X_n$ are the explanatory variables. $\beta_0 \dots \beta_n$ are the estimated parameters and ε is the error term. The parameters are estimated using OLS, which minimizes the residual sum of squares (RSS). Essentially, OLS minimizes the predictive error, how far the projections deviate from the actual values. Model performance is evaluated based on R-squared and adjusted R-squared values. R-squared measures the proportion of variance in Y that is explained by the included X 's. This metric is useful as values range from 0 to 1, but it always increases when more predictors are added, even irrelevant ones. Adjusted R-squared accounts for the number predictors relative to the sample size. This metric is much more reliable for model performance because it penalizes the inclusion of unnecessary variables. Adjusted R-squared also allows for model comparison with different numbers of explanatory variables.

The dependent variable in all models is Sale Price, the final transaction amount at which each property was sold. To better understand its distribution and central tendency before modeling, summary statistics and a box plot are provided below.



Results

Before getting into the Models, it is important to get a better understanding of what the data looks like. The initial dataset included comprehensive data for housing in Ames, Iowa, which was cleaned by filling in missing values through mean calculations. The data was then manipulated to run OLS regression; hence, categorical variables were converted into dummy variables. After this brief background of the data, some more data exploration was made to further understand the correlation and distribution of it. Referring to the appendix at the end of this document, Figure 1 portrays a heatmap of the top 10 predictors with the highest correlation to SalesPrice, which serves as a basis for showing which predictors will probably have the greatest impact in predicting SalesPrice. Figure 2 in the appendix seems to prove this by showing a box plot with the Y-axis as SalesPrice and the X-axis as Overall Quality: one of the most important predictors. Figure 2 shows a positive correlation between the two variables. Figure 3 in

the appendix offers a similar example where “GrLivArea” is plotted against SalesPrice to show the positive correlation under a different data visualisation.

Appendix Figures 4 & 5 show the average sales price and number of sales across the years; these histograms perfectly capture how the housing crisis impacted the housing market: the average sales price remained constant, while the number of sales decreased drastically. Appendix Figures 6,7, 8 & 9 are histograms meant to show the distribution of Median SalesPrice & number of sales across Neighbourhood and MSSubClass. Appendix Figure 10 shows how to interpret the MSSubClass categorical data; for instance, 20 = “1-story 1946 & newer all styles.”

After this thorough understanding of the data, Model 1& Model 2 were created; Model 1 consists of an OLS regression model with the top 10 most correlated variables with Sales Price. Additionally, Model 1 dropped any predictors which were deemed statistically insignificant (p-value ≥ 0.05); this resulted in Model 1 taking the following form:

$$Y = \beta_0 + \beta_1 GrLivArea + \beta_2 OverallQual + \beta_3 TotRmsAbvGrd + \beta_4 GarageCars + \beta_5 TotalBsmtSF + \beta_6 GarageFinish_Unf + \varepsilon$$

Model 2 was created slightly differently: an OLS regression was run with all 388 variables from the dataset, and then, similarly to Model 1, all statistically insignificant predictors were dropped. In Model 2, 342 predictors were dropped, which resulted in a new OLS regression with a few additional insignificant predictors; hence, recursion was needed. The Python code made it so that Model 2 recursively dropped statistically insignificant predictors 3 times before completing the Model. Finally, Model 2 ended with 30 statistically significant predictors.

The Table shown below offers a comparison between Model 1 & Model 2 descriptive statistics. The adjusted R-Squared in both models are pretty similar, but Model 2 has a higher one of 0.508 compared to the 0.447 of Model 1. This means that 50.8% of the variability in SalesPrice is explained by the predictors in Model 2, whereas for Model 1, 44.7% of the

variability in SalesPrice is explained by its predictors. The F-statistics provided are quite different, where Model 1 has a considerably higher one of 394.6 compared to the one in Model 2 of 101.4. This doesn't necessarily matter too much, as both Models are statistically significant, as shown by the P-value of < 0.000; the higher F-stat in Model 1 can probably be explained by the noticeably smaller number of predictors.

	Model 1	Model 2
R-Squared	0.448	0.513
Adjusted R-Squared	0.447	0.508
F-stat	394.6	101.4
P-Value (F-stat)	0.000	0.000

After gaining an understanding of the overall model statistics, it's a good idea to explore some predictors more in depth; therefore, the table below shows the comparison between the 6 predictors in Model 1 with their counterparts in Model 2. Before getting into the details, here are what some of the labels mean:

GrLivArea: Above-grade (ground) living area square feet.

OverallQual: Rates the overall material and finish of the house.

TotRmsAbvGrd: Total rooms above grade (does not include bathrooms).

Garage Cars: Size of garage in car capacity.

TotalBsmtSF: Total square feet of basement area.

GarageFinish_Unf: Interior finish of the garage (Unfinished).

Predictors Model 1	Coefficient	SD	P-Value	Predictors Model 2
Intercept	11.3420	0.027	0.000	
	7.6518	0.447	0.000	Intercept
GrLivArea	0.0001	1.58e-05	0.000	
	0.0001	2.05e-05	0.000	GrLivArea
OverallQual	0.0407	0.004	0.000	
	N/A	N/A	N/A	OverallQual
TotRmsAbvGrd	0.0186	0.004	0.000	
	0.0226	0.005	0.000	TotRmsAbvGrd
GarageCars	0.0489	0.07	0.000	
	N/A	N/A	N/A	GarageCars
TotalBsmtSF	7.676e-05	1.13e-05	0.000	
	N/A	N/A	N/A	TotalBsmtSF
GarageFinish_Unf	(0.0540)	0.010	0.000	
	N/A	N/A	N/A	GarageFinish_Unf

Looking at the table above, it can be seen that some of the highest-correlated predictors to SalesPrice in Model 1 aren't present in Model 2. This is because Model 2 originally ran an OLS regression of the whole dataset, meaning that it might've included variables that were correlated to the ones in Model 1, for example, GarageCars is most likely related to GarageArea (*Size of garage in square feet*), which is present in Model 2. Nevertheless, useful information can still be taken from this table, for instance, a good observation would be that of saying that Model 2's predictors are better because they are regressed while keeping more variables constant. Going off this statement, it seems that, based on the table above, the most impactful predictor is

TotRmsAbvGrd, where SalesPrice increases by roughly 2.26% for every room above grade.

Moving on from the previous table, we can take a closer look at the results from Model 2 in Figure 1, where it's clear to see that the most important predictors are the following:

LandSlope_Sev: Slope of property - Severe Slope

Condition2_PosA: Proximity to various conditions (if more than one is present) - Adjacent to positive off-site feature

RoofMatl_Membran: Roof Material - Membrane

RoofMatl_Metal: Roof Material - Metal

RoofMatl_WdShngl: Roof Material - Wood Shingles

Figure 1

	coef	std err	t	P> t	[0.025	0.975]
const	7.6518	0.447	17.116	0.000	6.775	8.528
Id	4.635e-05	4.51e-06	10.267	0.000	3.75e-05	5.52e-05
LotArea	5.091e-06	6.16e-07	8.263	0.000	3.88e-06	6.3e-06
OverallCond	0.0256	0.004	6.978	0.000	0.018	0.033
TotalBsmtSF	0.0001	1.51e-05	8.454	0.000	9.83e-05	0.000
2ndFlrSF	4.844e-05	1.83e-05	2.653	0.008	1.26e-05	8.42e-05
GrLivArea	0.0001	2.05e-05	6.942	0.000	0.000	0.000
KitchenAbvGr	-0.1261	0.020	-6.406	0.000	-0.165	-0.088
TotRmsAbvGrd	0.0226	0.005	4.986	0.000	0.014	0.031
GarageYrBlt	0.0010	0.000	5.162	0.000	0.001	0.001
GarageArea	0.0001	2.36e-05	6.141	0.000	9.88e-05	0.000
MiscVal	-4.123e-05	8.82e-06	-4.674	0.000	-5.85e-05	-2.39e-05
MSZoning_FV	0.1035	0.021	4.922	0.000	0.062	0.145
MSZoning_RL	0.0869	0.011	7.813	0.000	0.065	0.109
LandSlope_Sev	-0.2546	0.064	-3.958	0.000	-0.381	-0.128
Condition1_Norm	0.0359	0.011	3.206	0.001	0.014	0.058
Condition2_PosA	-0.2992	0.119	-2.512	0.012	-0.533	-0.066
RoofMatl_CompShg	1.9443	0.213	9.132	0.000	1.527	2.362
RoofMatl_Membran	2.4379	0.302	8.086	0.000	1.847	3.029
RoofMatl_Metal	2.3954	0.304	7.886	0.000	1.800	2.991
RoofMatl_Roll	1.6828	0.296	5.685	0.000	1.102	2.263
RoofMatl_Tar&Grv	1.9769	0.217	9.114	0.000	1.552	2.402
RoofMatl_WdShake	1.9933	0.224	8.914	0.000	1.555	2.432
RoofMatl_WdShngl	2.2359	0.225	9.957	0.000	1.796	2.676
Exterior2nd_Other	0.4318	0.205	2.109	0.035	0.030	0.833
MasVnrType_Stone	0.0717	0.015	4.929	0.000	0.043	0.100
Electrical_Mix	-0.5598	0.206	-2.722	0.007	-0.963	-0.157
KitchenQual_Fa	-0.0921	0.025	-3.622	0.000	-0.142	-0.042
MiscFeature_Othr	-0.3967	0.138	-2.869	0.004	-0.668	-0.126
MiscFeature_Shed	-0.3307	0.107	-3.085	0.002	-0.541	-0.121
YearBuilt_1893.0	0.5919	0.237	2.496	0.013	0.127	1.057

To finish off this result section, Figure 2 and Figure 3 show a graph where the predicted SalesPrice and the Actual SalesPrice are plotted together. Without any statistics, it would be hard to determine which one is better, but by knowing in advance that R-squared is better for Model 2, it can be seen that the dots in Figure 3 are more in line with the red-dotted line.

Figure 2

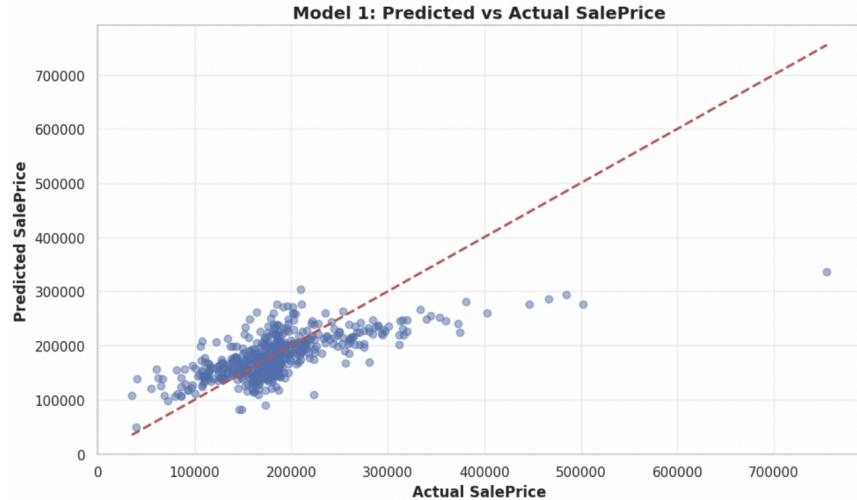
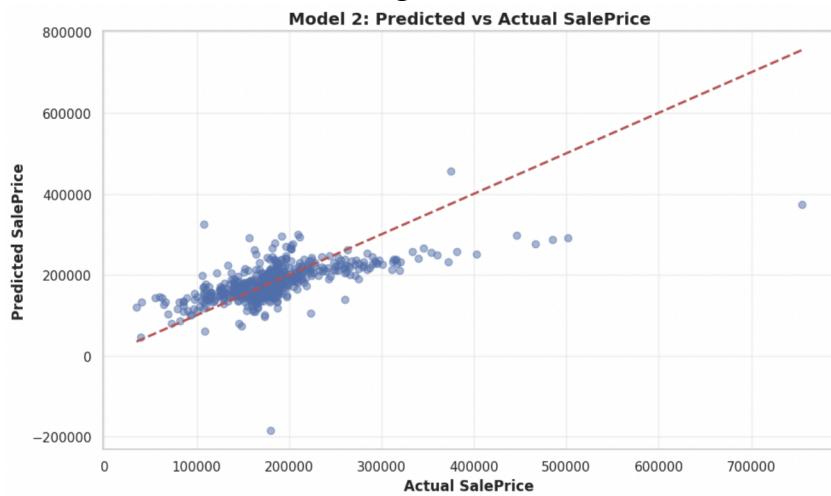


Figure 3



Additional Analysis #1

With 5 years of data from 2006-2010, segmenting sales price into 3 different time periods across the entire dataset provided insights into macroeconomic cycles during the period. The 3 periods are pre-recession during the housing boom, during the recession from Dec 2007-Jun 2009, and post-recession representing economic recovery. The following table shows the summary statistics of sales price broken down into the 3 periods.

MarketPeriod	count	mean	std	min	\
Pre-Recession (Boom)	1255.0	182076.558566	59112.484991	35311.0	
Recession (Dec 2007–Jun 2009)	998.0	178867.118236	53482.260365	37900.0	
Post-Recession (Recovery)	666.0	178016.258258	59635.789368	34900.0	
	25%	50%	75%	max	
MarketPeriod					
Pre-Recession (Boom)	156761.0	178403.0	192589.0	755000.0	
Recession (Dec 2007–Jun 2009)	155661.5	176175.0	191397.5	582933.0	
Post-Recession (Recovery)	150354.5	174432.0	190190.5	611657.0	

As a real estate professional, the above numbers are consistent with the broader reflections of the entire housing market of the US. The volume in sales dropped from 1,255 pre-recession to roughly 1,000 during the recession and much more aggressively during the recovery phase. The decrease in volume of sales reflects tighter credit conditions during the recession and home buyer confidence decreasing, as well. It is safe to say that Ames, Iowa is a relatively stable market compared to the broader market during the late 2000s. This is reflected in the median selling price of homes. The selling price decreased across all 3 phases, consistent with the national market. The decline was mild, indicating that Ames and similar, smaller midwestern markets were not as hit as hard as some during the recession. The high-end selling homes saw a notable drawback. The maximum selling price during the boom was \$755,000 and then this dipped to roughly \$583,000 during the recession. This suggests that there was a drastic decrease in demand

for luxurious properties. It can be expected that this trend in high-end properties was magnified in larger markets. During the post-recession period, sales volume remained very low in Ames and prices had not yet regained the pre-recession highs. These are indicators of a sluggish recovery from the recession, but there could be regional differences depending on which city or state is examined.

Additional Analysis #2

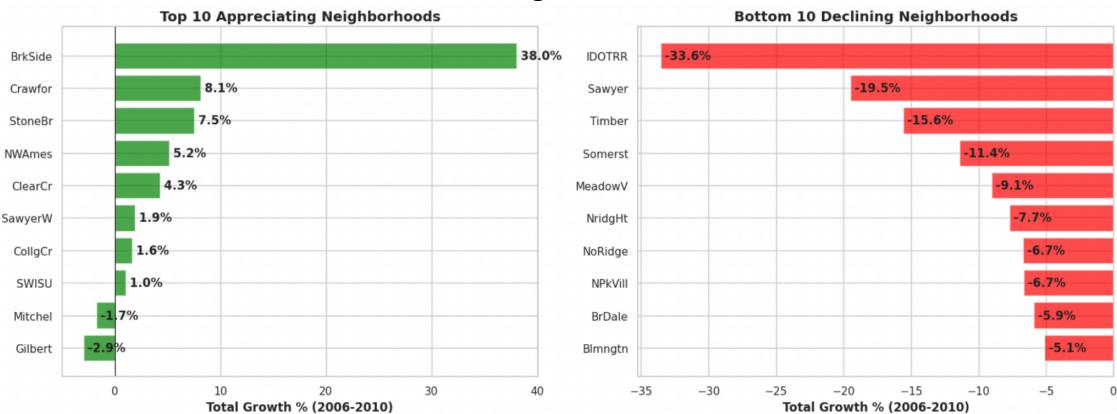
After running a general analysis which explores how all predictors affected SalePrice during the years 2006-20010, it's worth diving deeper into some predictors; this section of the paper will focus on analysing how different neighbourhoods were impacted by the housing crisis. The key detail that's important to remember is the fact that real estate markets aren't homogeneous, hence neighbourhoods will perform differently. Furthermore, understanding neighbourhood trajectories helps identify investment opportunities.

For this type of analysis, the data were grouped by neighbourhood and year, where the median sale price was calculated. Median sale price was chosen over the mean, so that outliers wouldn't sway the true representation of each neighbourhood. Appreciation rates were then calculated with the following formula: $\frac{\text{Median Price 2010} - \text{Median Price 2006}}{\text{Median Price 2006}}$, which returned the percentage change for each neighbourhood from before to after the housing crisis. Lastly, two neighbourhoods were excluded (Blueste and Veenker) due to some missing data, leaving a final count of 23 neighbourhoods.

Figure 4 shows the overall housing crisis effect on most of the neighbourhoods, ranging from BrkSide, at a 38% growth, to IDOTRR, at a -33.6% decline. Figure 4 clearly shows that even though BrkSide has a higher growth rate, compared to the lowest loss rate shown in

IDOTRR, the overall distribution of the neighbourhoods shows a prominent negative aftermath. Only 8 out of the 23 neighbourhoods had positive appreciation, while the rest all experienced decline, showing how the housing crisis really impacted most neighbourhoods.

Figure 4



Although it is important to understand how all the neighbourhoods did, the part that is the most interesting to real estate professionals is what differentiates neighbourhoods which rose in value, compared to neighbourhoods which declined in value. Figure 5 attempts to provide an answer to this conundrum; 3 out of the top 5 most-grown neighbourhoods after the housing crisis have a higher initial median sale price than 3 out of the top 5 most depreciated neighbourhoods after the housing crisis. This fact leads to the conclusion that higher-end neighbourhoods will tend to do better compared to lower-end neighbourhoods, at least that's what the market in Ames, Iowa, showed after the housing crisis.

Figure 5

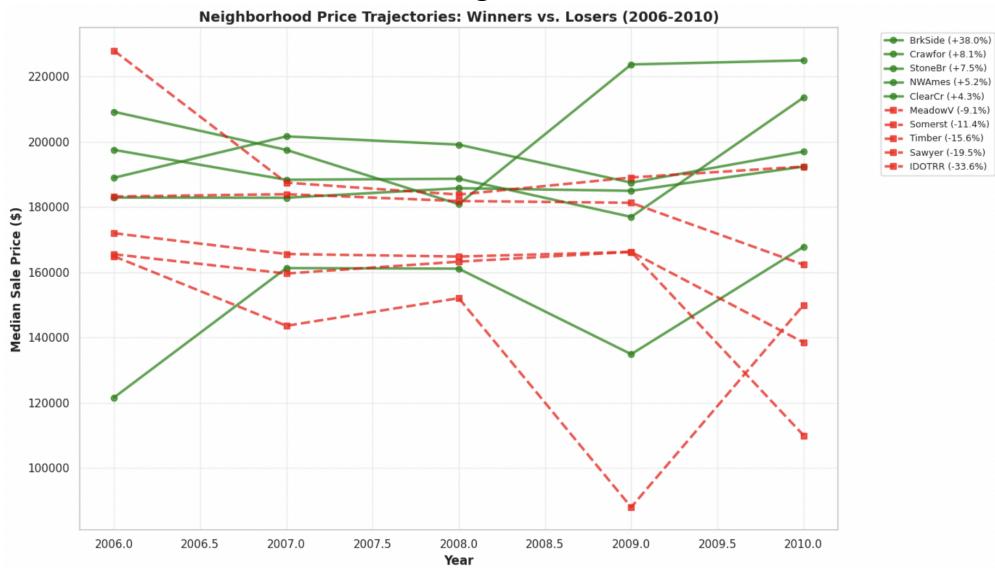


Figure 6 gave us a hint as to what drove a neighbourhood's growth or loss during the housing crisis. Figure 6 builds off of that by plotting the Initial Median Price in 2006 vs the Total Growth from 2006-2010; this figure will help immensely in understanding the relationship between the two. At first glance, the scatter plot shows almost no relationship between the two variables, which isn't promising; on the other hand, it raises another interesting observation. From Figure 6, it seems as if bigger neighbourhoods suffer from less change, whether that's positive or negative. Before moving on to testing that theory, an OLS regression with Total Growth as its Y parameter and Initial Median Price as its X parameter proved to be not statistically significant (shown in Figure 7), hence it cannot be assumed that higher-end neighbourhoods have any advantage in surviving a housing crisis.

Figure 6

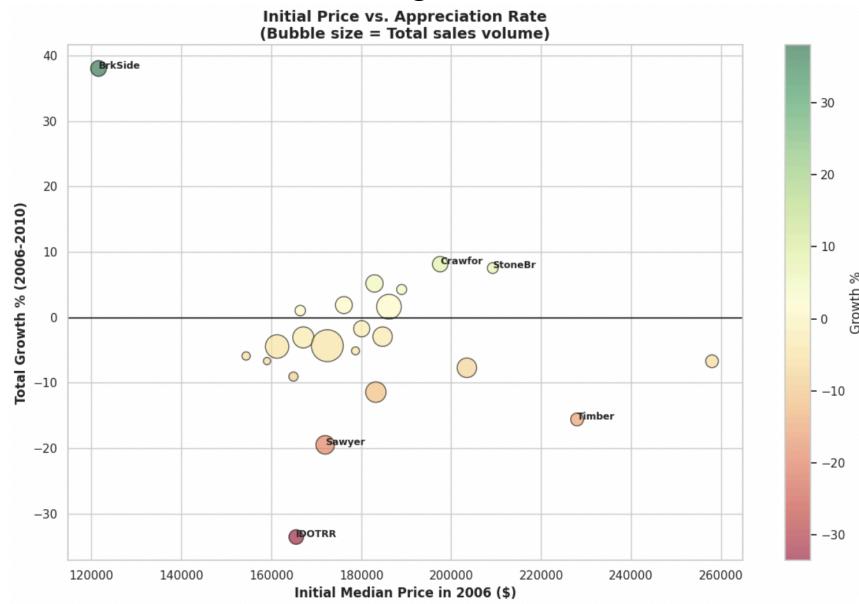


Figure 7

	coef	std err	t	P> t	[0.025	0.975]
const	19.7959	18.339	1.079	0.293	-18.343	57.935
InitialPrice_2006	-0.0001	0.000	-1.259	0.222	-0.000	8.23e-05

After determining that the initial median price wasn't a suitable predictor for total growth, the hypothesis "Bigger neighbourhoods suffer lower absolute changes" can be tested. Figure 8 shows that this hypothesis is indeed correct, showing that the correlation between size and absolute growth is -0.1348. This means that larger neighbourhoods have a lower volatility during a housing crisis.

Figure 8

Correlation between size and ABSOLUTE growth: -0.1348			
Negative correlation: Larger neighborhoods DO have less extreme outcomes			
Large neighborhoods (>200 sales):			
Neighborhood Total_Sales Total_Growth_Pct			
NAMES 443 -4.328314			
CollCr 267 1.624842			
OldTown 239 -4.451167			
Small neighborhoods (<50 sales):			
Neighborhood Total_Sales Total_Growth_Pct			
NPkVill 23 -6.671193			
Blmgtn 28 -5.118049			
BrDale 30 -5.903268			
MeadowW 37 -9.050067			
ClearCr 44 4.256968			
SWISU 48 1.033003			

To conclude this additional analysis, the “animation link” below leads to an animation done with Tableau, which shows how all the neighbourhoods changed with the passing of the years between 2006-2010.

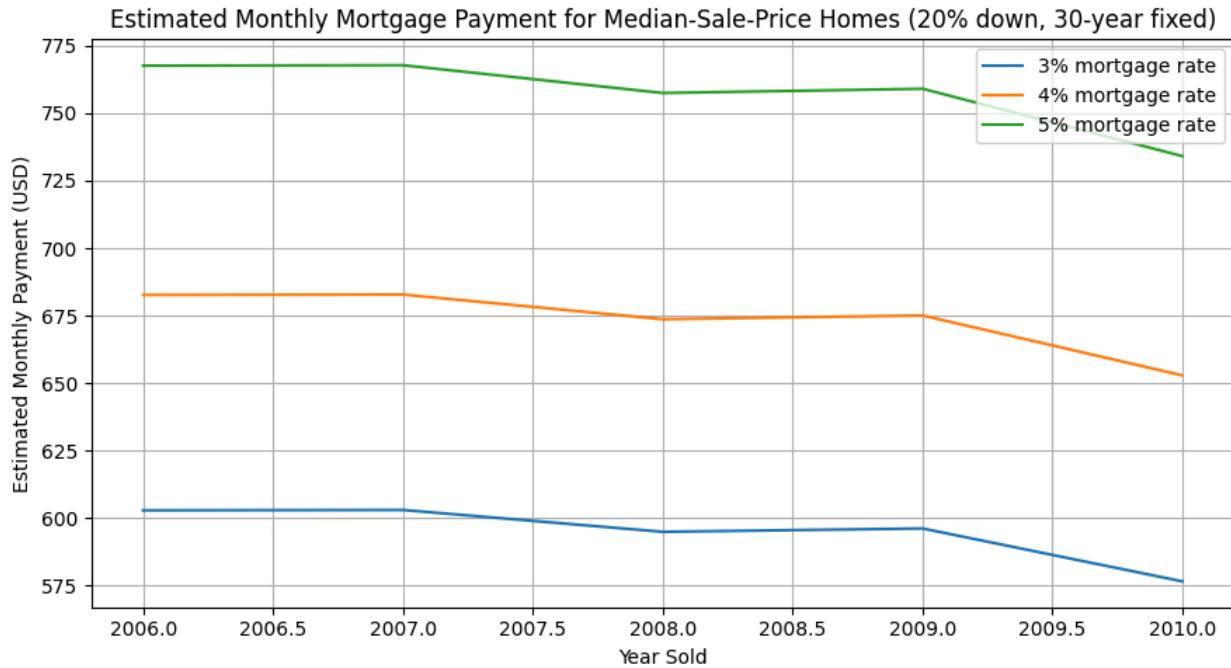
[*Animation Link*](#)



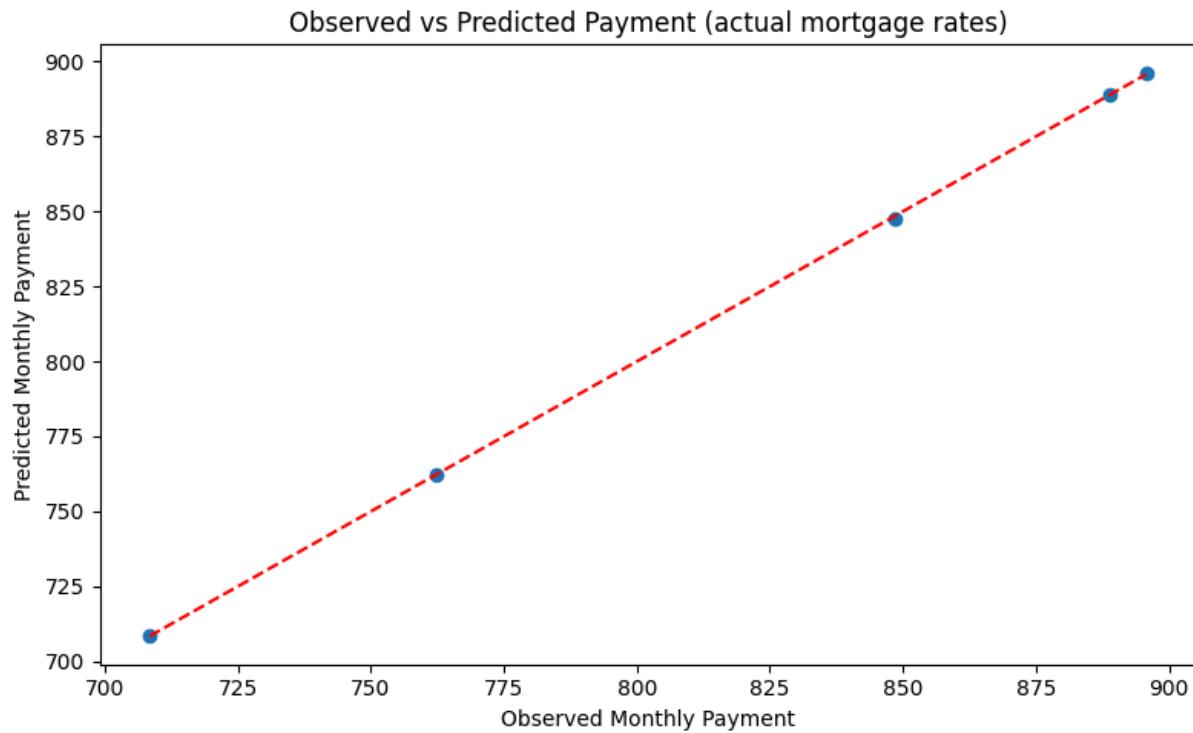
Additional Analysis #3

The goal of this analysis is to set up how the U.S. housing affordability is shaped by the interaction between median home sale prices and mortgage interest rates. Affordability is captured by the monthly mortgage payments required to purchase a median-priced home with a standard 30-year fixed mortgage. The methodology behind the analysis is based on the median annual home sale price calculated from the housing dataset; the 30-year fixed mortgage rates were retrieved from the Federal Reserve Economic Data(FRED). A log-linear regression was executed to quantify the impact of both home prices and interest rates on affordability:

$$\ln(\text{Monthly Payment}) = \beta_0 + \beta_1 \ln(\text{Median Price}) + \beta_2 (\text{Mortgage Rate}) + \epsilon$$

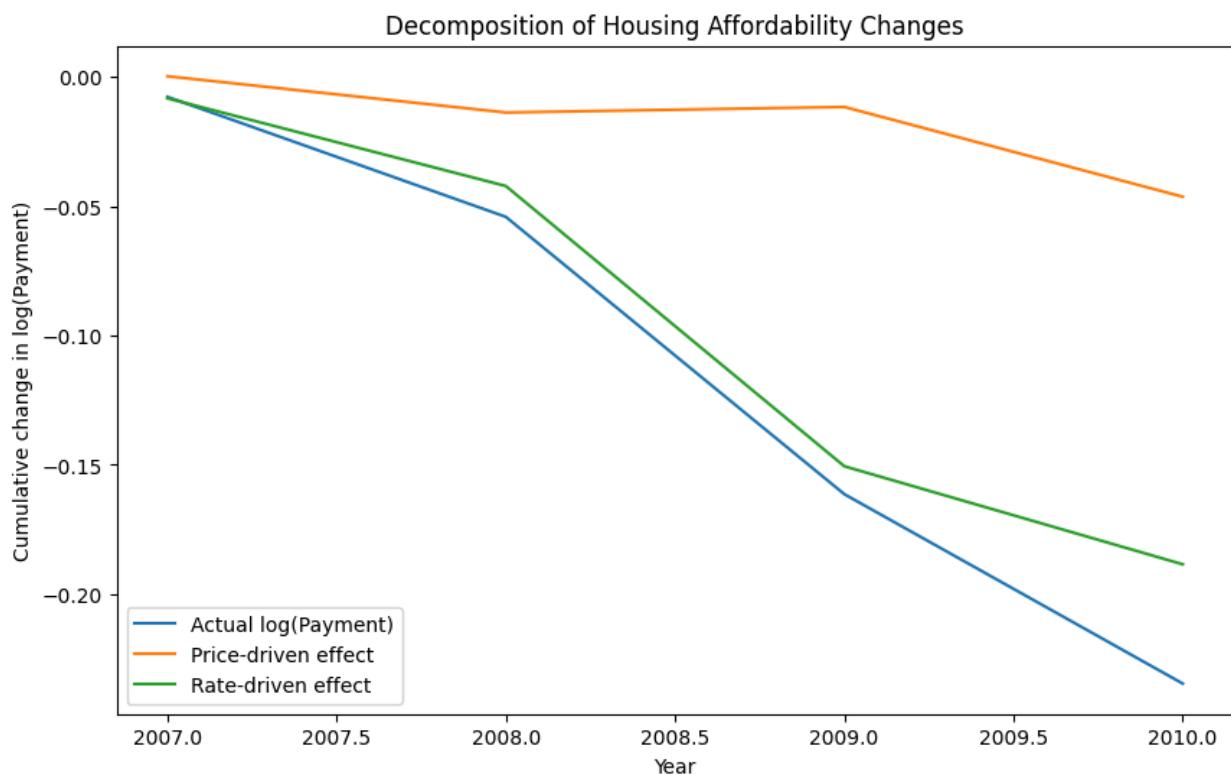


This line chart shows the trend from 2006 to 2010 in estimated payments, capturing how affordability has shifted over the period of time from before the recession, through the recession, and after the recession as both prices and rates changed.



This chart validates the regression method by showing how closely predicted affordability aligns with actual calculated payments.

The results of the Regression analysis shows that the coefficient of the log simulating price is about 1.0. This means that a 1% rise in home prices increases monthly payments by ~1%. The coefficient on mortgage rates is ~0.15-0.25. This means that a 1% increase in mortgage rates drives monthly payments up by 15-25%. R^2 is about 0.99, which shows that all variation in affordability is explained by price and interest rate movements. To better understand the affordability dynamics, monthly payments were decomposed into two components: Price-driven effects, which holds the rates constant, and Rate-driven effects, which holds the price constant.



This chart highlights when affordability shifts were primarily due to prices (steady growth) versus mortgage rates.

The interpretation of this analysis is that home prices gradually increase monthly payments, nearly one-to-one with rises in median home sales price. Mortgage rates act as an affordability accelerator. Even a small jump in rates creates sharp rises in monthly payments, largening housing cost burdens. Timing is super important through these two topics. After 2008, falling rates softened the blow of rising prices. In opposition to this, in 2021 and 2022, a rush in rates combined with rising prices produced the worst affordability decline in decades. Housing affordability is formed by a dual pressure gauge: long-term price growth and volatile interest rates. The regression confirms that while prices provide the baseline, interest rates determine the severity and pace of affordability shocks. Both must be analyzed together to understand the U.S. housing market dynamics.

Conclusion

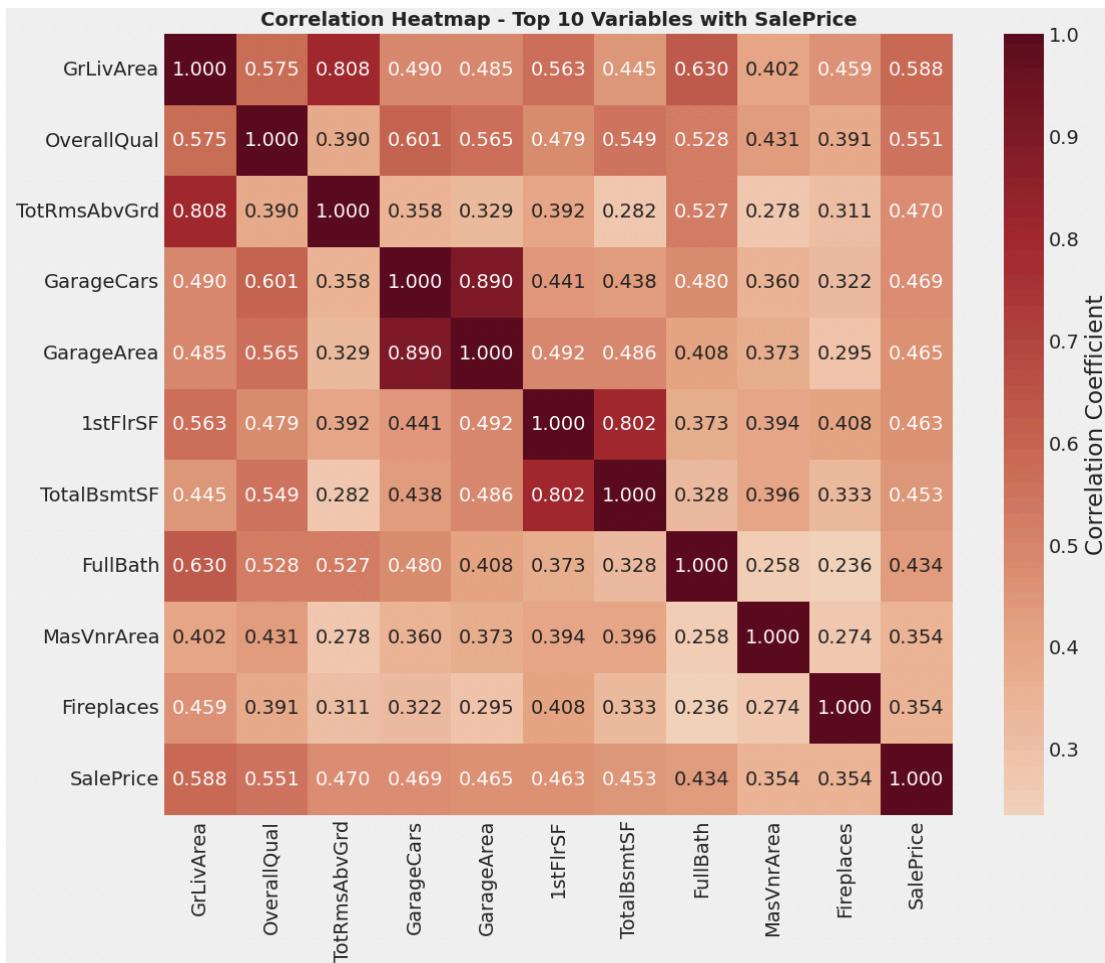
This project demonstrates how predictive analytics can be applied to one of the most important sectors of the U.S. economy: The housing market. By conducting thorough data cleaning, exploratory data analysis, and regression modeling, we identified key variables such as living area and neighborhood. These factors strongly influence housing prices. The comparison of different regression models highlights the importance of balancing simplicity with explanatory power. Beyond technical modeling, our extension analyses illustrate how broader economic factors, such as interest rates, shape affordability and influence sales trends over time. Incorporating these perspectives helps connect data driven insights to the lived experience of homebuyers, sellers, and industry professionals who navigate these dynamics. Visuals such as distributions of sales price and correlation of predictors provide a clear picture of the market shifts and patterns that enhance the interpretation of the results. Our findings reinforce the value

of both prescriptive and predictive models in real estate decision making. For sellers and homebuyers, the strongest price drivers offer a more informed insight on property valuation. These insights support competitive strategies for professionals, even in challenging market conditions. For researchers and policy makers, the results underscore the potential of quantitative methods to shed light on questions of equity, affordability, and long-term market health.

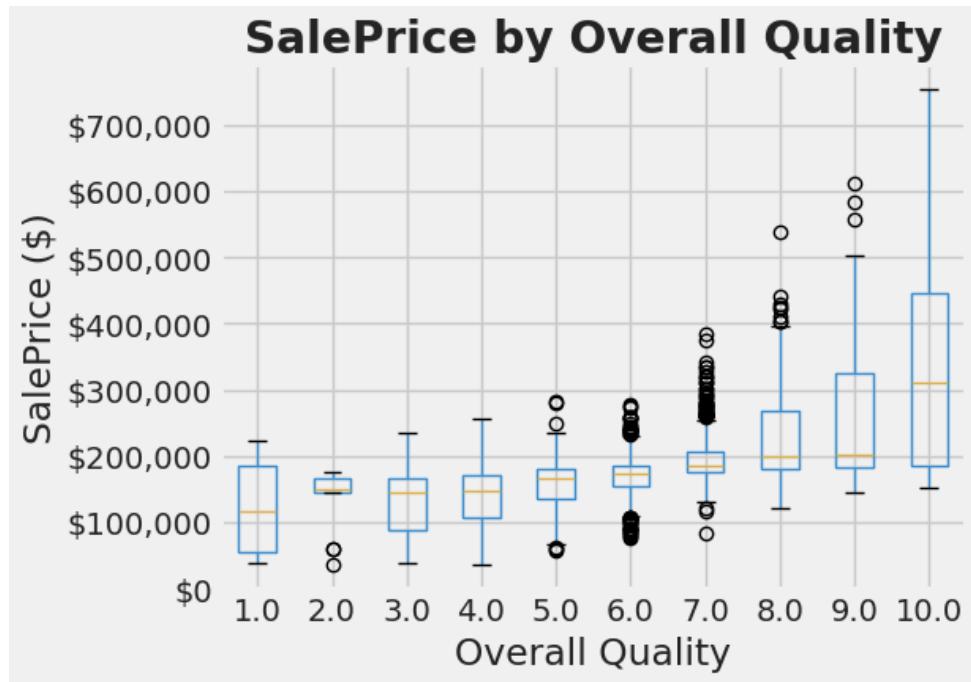
Ultimately, this project highlights how data analytics not only explains past trends but also equips shareholders with tools to anticipate and respond to future developments. By bridging technical rigor with practical application, we demonstrate how predictive modeling can play a vital role in guiding decisions within the housing market. The housing market is one of the most dynamic and consequential areas of the U.S. economy.

Appendix

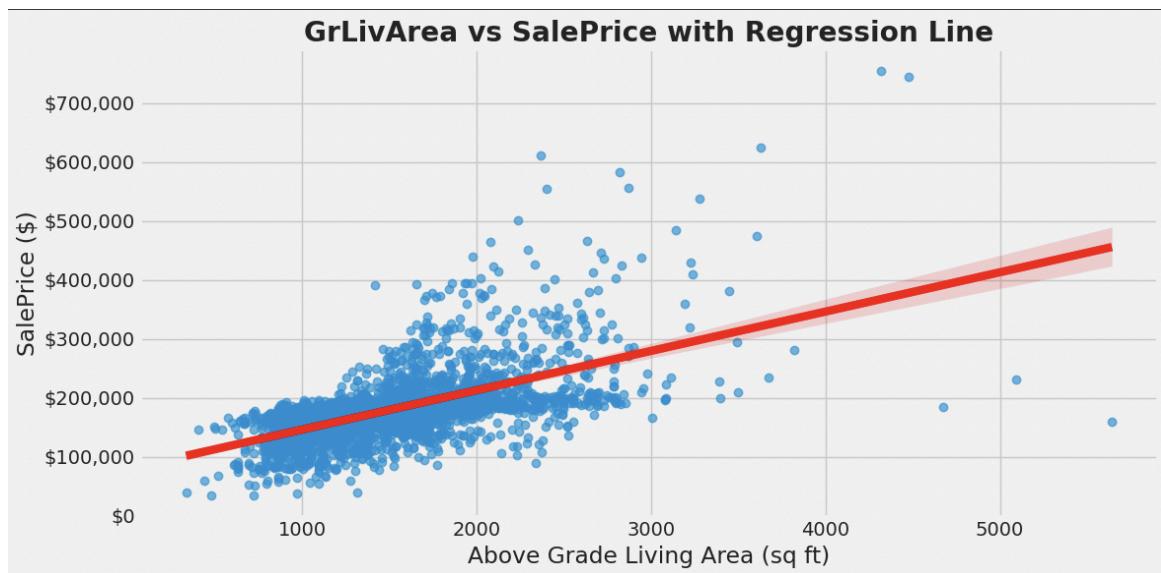
Appendix Figure 1



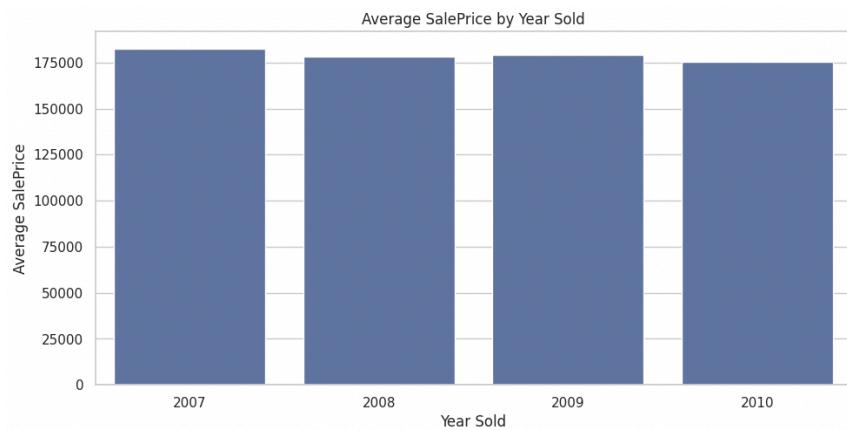
Appendix Figure 2



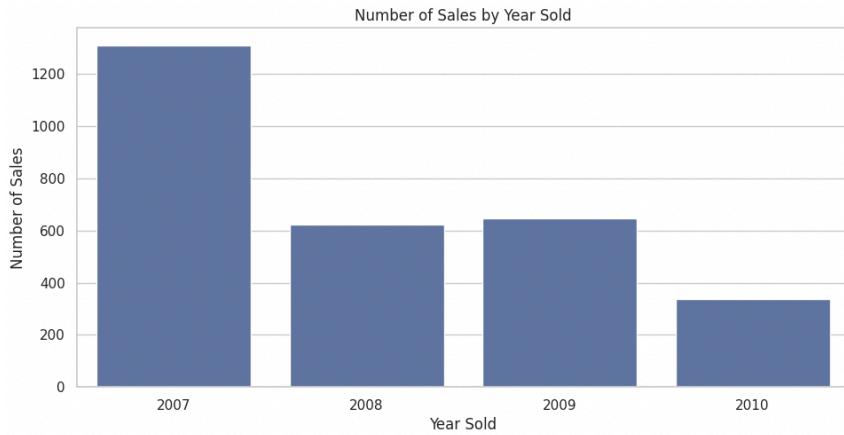
Appendix Figure 3



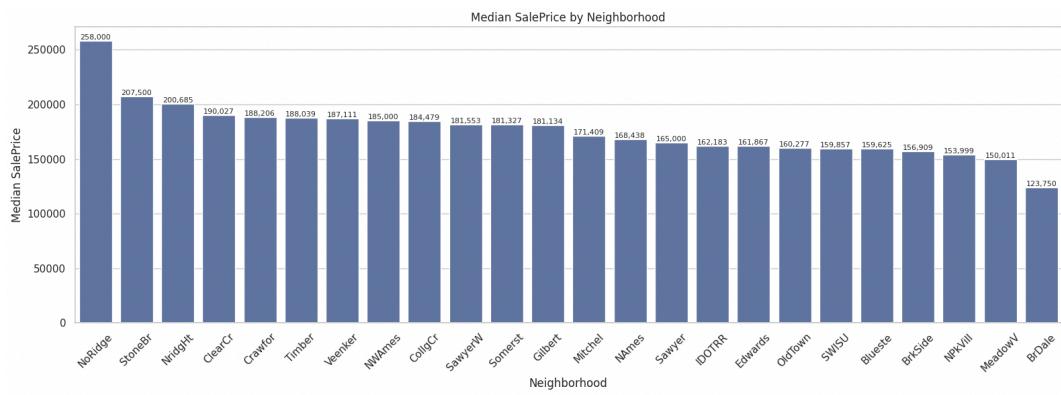
Appendix Figure 4



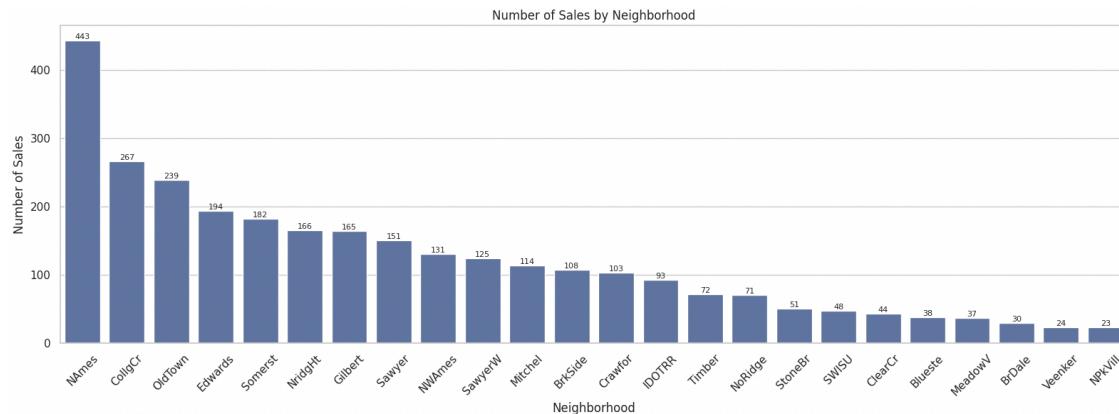
Appendix Figure 5



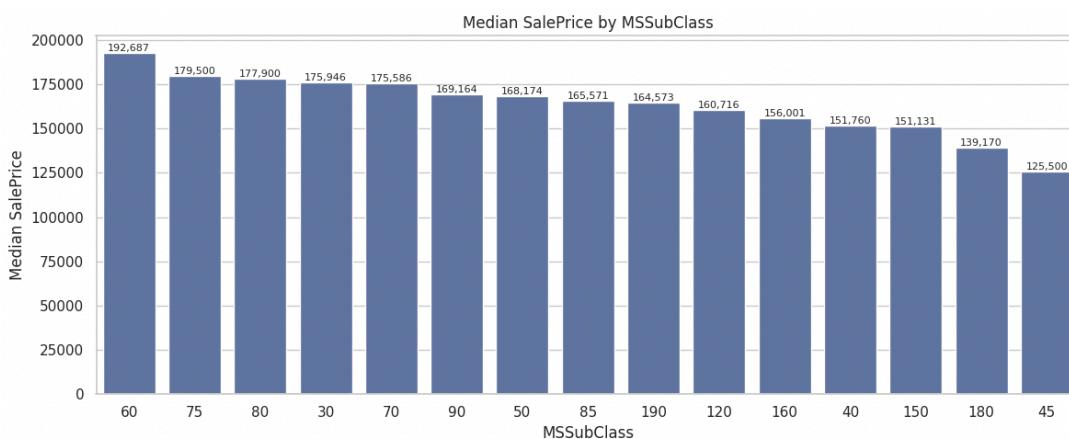
Appendix Figure 6



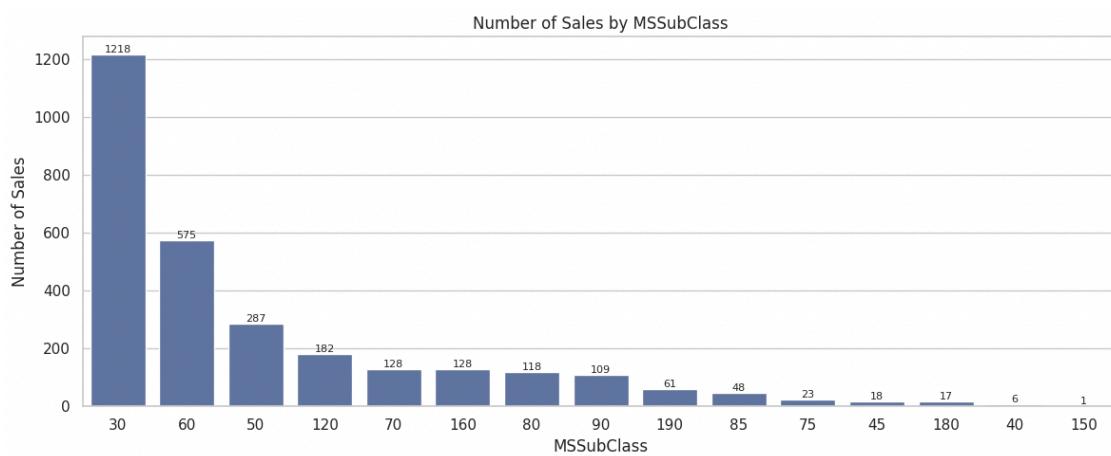
Appendix Figure 7



Appendix Figure 8



Appendix Figure 9



Appendix Figure 10

MSSubClass: Identifies the type of dwelling involved in the sale.

20	1-STORY 1946 & NEWER ALL STYLES
30	1-STORY 1945 & OLDER
40	1-STORY W/FINISHED ATTIC ALL AGES
45	1-1/2 STORY - UNFINISHED ALL AGES
50	1-1/2 STORY FINISHED ALL AGES
60	2-STORY 1946 & NEWER
70	2-STORY 1945 & OLDER
75	2-1/2 STORY ALL AGES
80	SPLIT OR MULTI-LEVEL
85	SPLIT FOYER
90	DUPLEX - ALL STYLES AND AGES
120	1-STORY PUD (Planned Unit Development) - 1946 & NEWER
150	1-1/2 STORY PUD - ALL AGES
160	2-STORY PUD - 1946 & NEWER
180	PUD - MULTILEVEL - INCL SPLIT LEV/FOYER
190	2 FAMILY CONVERSION - ALL STYLES AND AGES