

What factors and how influence Zillow Housing Sale Price*

Yaning Jin

April 17, 2024

This project is interested in the factors affecting housing prices on zillow. After exploratory analysis of the data, a multiple linear regression was constructed, and a log transformation was performed to build a model with better performance. The transformed housing factors can explain more than 80% of the changes in housing prices. However, the existence of residuals in the model does not fully meet the constraints of the normality assumption, so in the future, data processing, feature engineering and model improvement are needed to study the factors affecting Zillow's housing prices..

Table of contents

1	Introduction	2
2	Data	2
2.1	Raw Data	2
2.2	Data Analysis tools	3
2.3	Variable Description	3
3	Exploratory Data Analysis	3
4	Model	5
4.1	Model set-up	5
4.1.1	Model justification	6
5	Results	8

*Code and data are available at: <https://github.com/Yuki010305/What-factors-and-how-influence-Zillow-Housing-Sale-Price.git>.

6 Discussion	8
6.1 Weaknesses and next steps	8
References	9

1 Introduction

One of a nation's key industries is housing, and the health of this sector influences the degree of economic growth of the nation. A person's home is frequently the most expensive item in their life. Zillow is an online real estate database provider that assesses property values and offers details on the houses you're interested in ([Loukissas 2018](#)). The Zestimate prediction system was developed by Zillow. It is based on a data collection of hundreds of millions of homes and provides a preliminary estimate of a property's worth by combining a particular algorithm with the features of each home and the state of the market([Rolli 2020](#)). To improve the accuracy of the Zestimate system and provide people with a more trustworthy way to estimate home prices, Zillow launched the Zillow Awards competition in 2017. Our goal in this project is not to participate in a competition to minimize the logarithmic error between the estimated house price and the actual selling price. The main purpose of our exploration in this project is to explore the factors that affect Zillow home sales prices and how they affect home sales prices. Characteristics of the homes we select include number of bathrooms, number of bedrooms, square footage, number of rooms, year, tax value, land tax value, etc. The estimated model is

$$\log(\text{price}) = \beta_0 + \beta_1 \text{bathrooms} + \beta_2 \log(\text{squarefootage}) + \dots + \beta_6 \text{taxValue}$$

For some variables we will also do some log transformations.

2 Data

2.1 Raw Data

In this research, we examine zillow data sourced using the `zillbowR` library ([Rolli 2020](#)). The dataset encompasses 90,275 records (89,499 records after cleaning), focusing on specific variables: bathroom number, bedroom number, year built, square feet, room number, tax value, etc. The apartments were built in a wide range of years, from as early as 1885 to as recently as 2015.

The Raw data collects 60 features of the house, and data quality is not high, there are many features with many null values, so 8 high-quality variables are selected from the variables for modeling and prediction.

2.2 Data Analysis tools

R ([R Core Team 2020](#)), a potent open-source statistical programming language, was used to analyze the data. To improve the effectiveness of our data operations, we used a collection of R packages from the tidyverse ([Wickham et al. 2019](#)), which is a collection of tools created for data science. The `dplyr` package ([Wickham et al. 2022](#)) offered a consistent collection of verbs that aid in filtering, summarizing, and organizing the dataset, while the `ggplot2` package ([Wickham 2016](#)) made it easier to create complex visualizations. Because of its quick and user-friendly data reading capabilities, the `readr` package ([Wickham, Hester, and Bryan 2022](#)) was used. `Knitr` ([Goodrich et al. 2020](#)) handled report production dynamically, allowing R code to be included into this document. `kableExtra`([Zhu 2021](#)) was also used to create visually appealing and editable tables, which improved the way our results were presented.

2.3 Variable Description

Variable	Description
bathroomcnt	Number of bathrooms in home
bedroomcnt	Number of bedrooms in home
calculatedfinishedsquarefeet	Total finished living area of the home
roomcnt	Total number of rooms
yearbuilt	The Year the principal residence was built
taxvaluedollarcnt	The total tax assessed value of the parcel
landtaxvaluedollarcnt	The assessed value of the land area
price	price

3 Exploratory Data Analysis

As you can see from the figure, house sales prices are right-skewed data. When building the model, we need to log-transform the house sales prices to make them conform to the normal distribution.

The above scatter plot shows that tax value has a certain positive impact on sale price, the positive correlation between tax value and sale price is very strong.

The above scatter plot shows that square feet has a certain positive impact on sale price, the positive correlation between square feet and sale price is very strong.

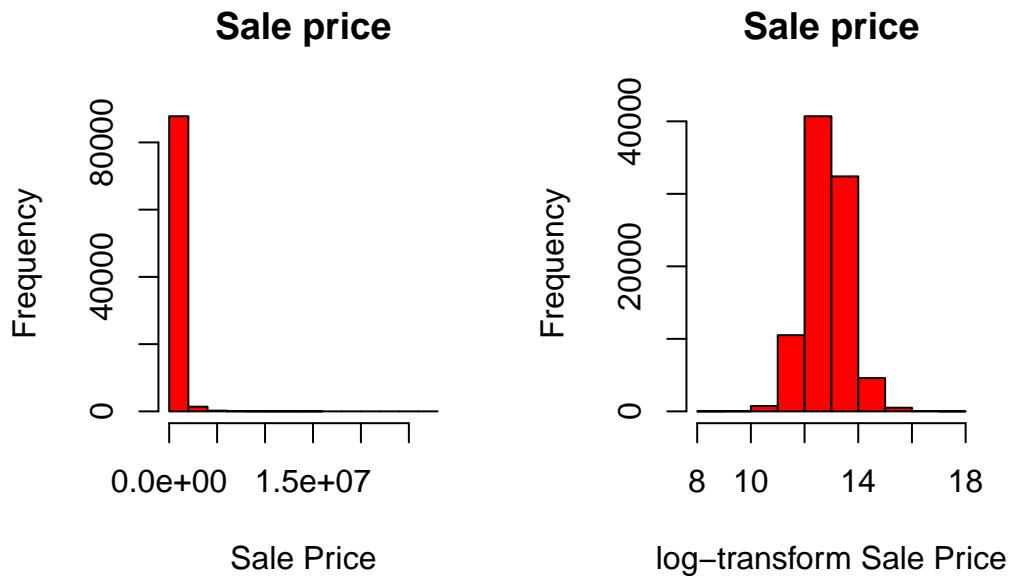


Figure 1: sale price hist

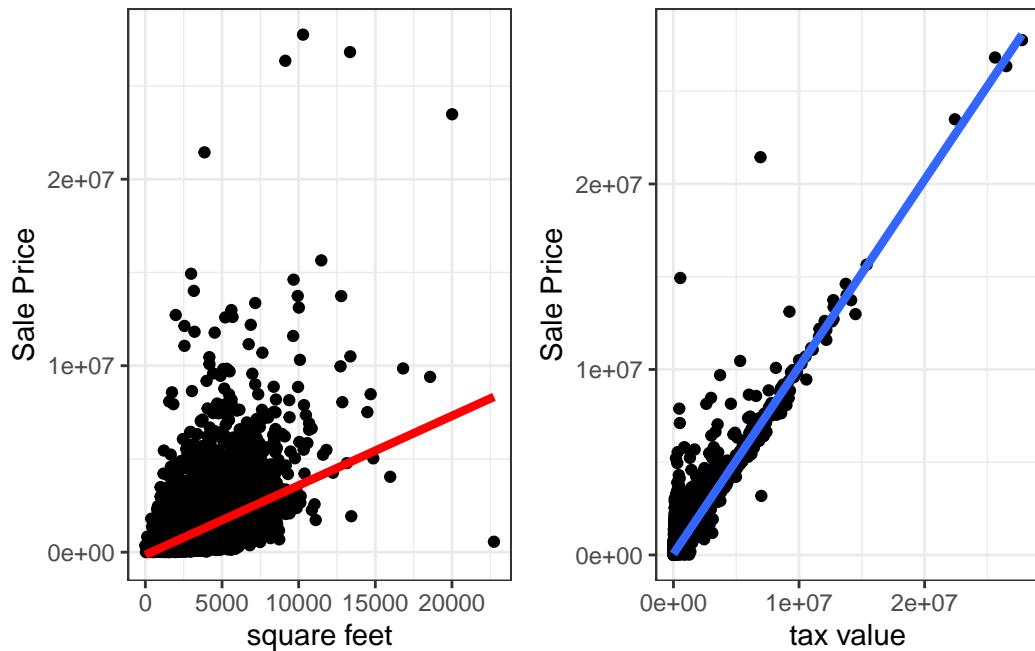


Figure 2: points1

4 Model

4.1 Model set-up

After data processing, the data set is a clean data set with 2929 observations and 10 house characteristic variables. In order to evaluate the performance of the model, we randomly split the analysis data set into a test set and a training set in a ratio of 75%:25%.

The first model we build is the full model. We then improve the model by removing insignificant variables.

The residual test found that both ends of the QQ graph deviated greatly from the straight line and were affected by special points such as outliers and leverage points. So, in order to further improve the performance of model fitting, we will delete special points from the training set.

After deleting special points such as level points, a new model was refitted.

Call:

```
lm(formula = log(price) ~ bathroomcnt + bedroomcnt + log(calculatedfinishedsquarefeet) +
  roomcnt + yearbuilt + log(taxvaluedollarcnt) + log(landtaxvaluedollarcnt),
  data = new_train)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.9729	-0.1428	-0.0473	0.0558	3.0992

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.707e+00	1.400e-01	19.333	< 2e-16 ***
bathroomcnt	-6.631e-03	2.603e-03	-2.548	0.010848 *
bedroomcnt	-6.022e-03	1.937e-03	-3.109	0.001875 **
log(calculatedfinishedsquarefeet)	2.090e-01	6.235e-03	33.519	< 2e-16 ***
roomcnt	-1.531e-02	6.861e-04	-22.313	< 2e-16 ***
yearbuilt	-2.491e-04	7.061e-05	-3.528	0.000419 ***
log(taxvaluedollarcnt)	6.714e-01	6.400e-03	104.904	< 2e-16 ***
log(landtaxvaluedollarcnt)	5.454e-02	4.502e-03	12.116	< 2e-16 ***

Signif. codes:	0 '***'	0.001 '**'	0.01 '*'	0.05 '.'
	0.1	'	'	1

Residual standard error: 0.3171 on 62754 degrees of freedom

Multiple R-squared: 0.8083, Adjusted R-squared: 0.8082

F-statistic: 3.779e+04 on 7 and 62754 DF, p-value: < 2.2e-16

4.1.1 Model justification

Clearly the model has improved.

4.1.1.1 A1:Linearity of the Relationship

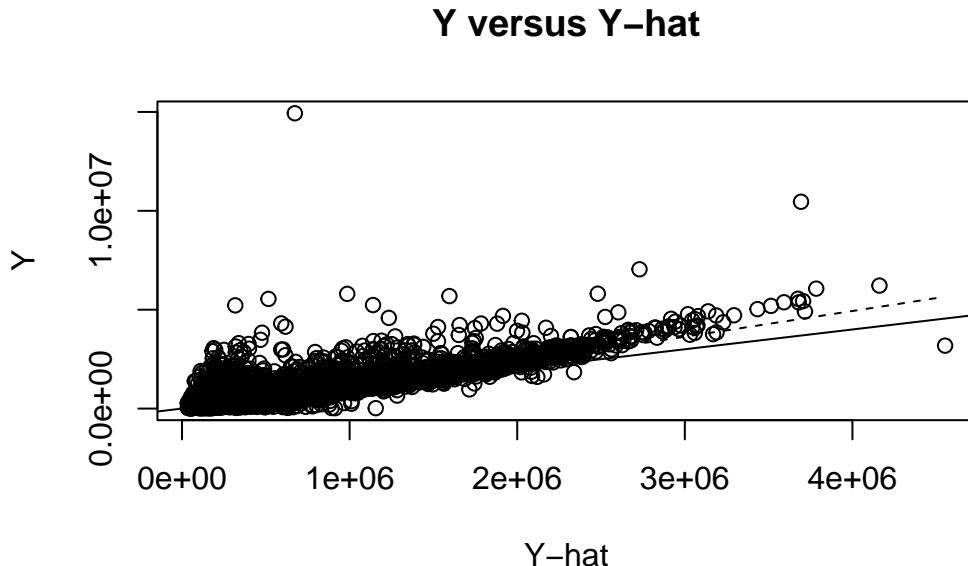


Figure 3: plot y hat

The above scatter plot fits the relationship between the predicted value and the actual value. You can see that these points are almost on or close to the line, so we can say that a linear relationship is satisfied.

4.1.1.2 A2.Covariance of Errors

the errors are independent.

4.1.1.3 A3.Common Error variance

the errors have constant variance.

4.1.1.4 A4. Normality of Error

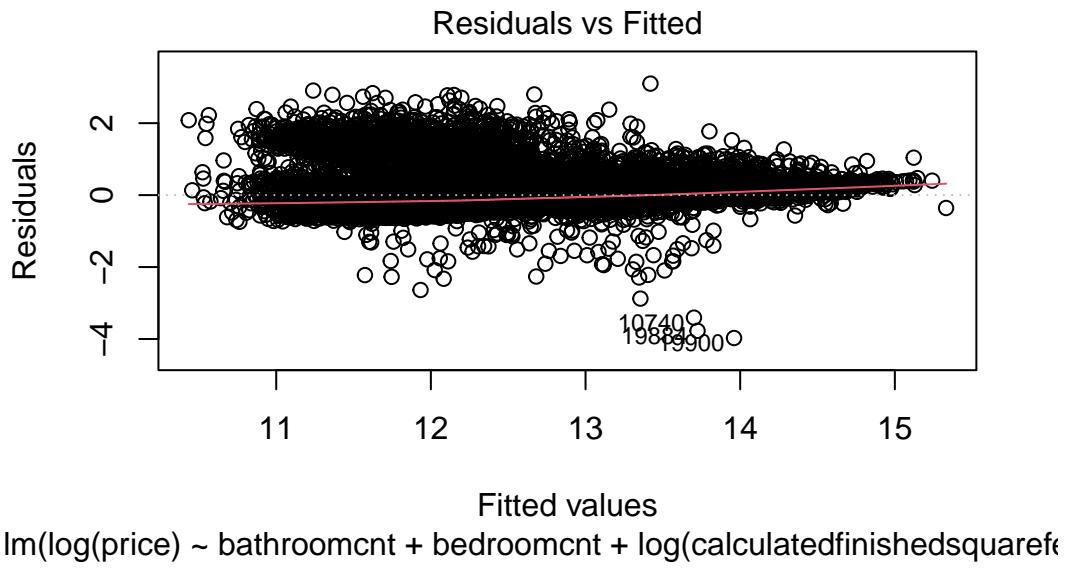


Figure 4: plot resid

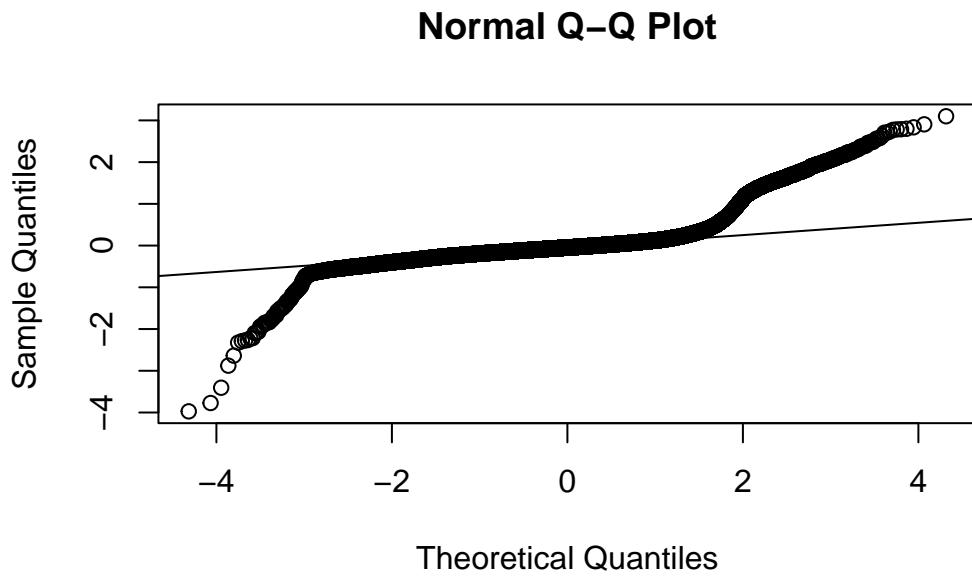


Figure 5: plot qq

5 Results

Estimating the value of a home is a common difficulty. Therefore, a lot of work has already been completed. Lee ([Lee 2016](#)) make an effort to develop multivariate regression models based on house datasets and assess the models using maximum information coefficient statistics based on anticipated values and home prices. When a moderate to big data sample size is employed, Nghiep et al. ([Nguyen and Cripps 2001](#)) examined two methods: multiple regression analysis (MRA) and artificial neural networks (ANN). Based on the prediction performance, ANN performs better than MRA. An artificial neural network (ANN) model was created after Limsombunchai ([Limsombunchao 2004](#)) that ANNs are more accurate in predicting property values than hedonic regression models.

The final model shows that the number of bathrooms, bedrooms, rooms, year taxes and fees of the building, square feet, etc. all have a significant impact on Zillow's housing prices. For every additional bathroom, the house price increases by 1.017% per unit. For every additional bedroom, the house price actually decreases by 0.01%. In addition, the newer the building is, the housing prices actually decrease. For every 1% unit increase in taxes and fees, housing prices increase by 0.67% unit. These findings explain the relationship between house-related attributes and house prices.

6 Discussion

6.1 Weaknesses and next steps

However, the residuals of our model do not fully satisfy the normality assumption. Such violations may reduce the model's predictive accuracy on new data sets. Another limitation of this model is that transforming the data makes the model less interpretable.

To start making these models better in the future, divide each training dataset into smaller training tests and perform some cross-validation before generating the test files. This might enhance performance, or at the very least increase the predictability of the model's output. Although most categorical variables cannot be included in the model because of memory and time restrictions when running the model, the model may be enhanced with improved feature engineering. These variables can be included in the model by being divided into smaller groups. Making some interactive words is an additional thought. It is crucial to classify missing variables as predictors once all missing values have been imputed, as was previously noted.

References

- Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2020. *Rstanarm: Bayesian Applied Regression Modeling via Stan*. <https://mc-stan.org/rstanarm>.
- Lee, Roger. 2016. *Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing*. Vol. 653. Springer.
- Limsombunchao, Visit. 2004. “House Price Prediction: Hedonic Price Model Vs. Artificial Neural Network.”
- Loukissas, Yanni. 2018. “All the Homes: Zillow and the Operational Context of Data.” In *Transforming Digital Worlds: 13th International Conference, iConference 2018, Sheffield, UK, March 25-28, 2018, Proceedings* 13, 272–81. Springer.
- Nguyen, Nghiep, and Al Cripps. 2001. “Predicting Housing Value: A Comparison of Multiple Regression Analysis and Artificial Neural Networks.” *Journal of Real Estate Research* 22 (3): 313–36.
- R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Rolli, Channamallikarjun Siddaramappa. 2020. “Zillow Home Value Prediction (Zestimate) by Using XGBoost.”
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- 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>.
- Wickham, Hadley, Romain François, Lionel Henry, and Kirill Müller. 2022. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.
- Wickham, Hadley, Jim Hester, and Jennifer Bryan. 2022. *Readr: Read Rectangular Text Data*. <https://CRAN.R-project.org/package=readr>.
- Zhu, Hao. 2021. *kableExtra: Construct Complex Table with ‘Kable’ and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.