

Linear Regression Analysis on House Price of Ames

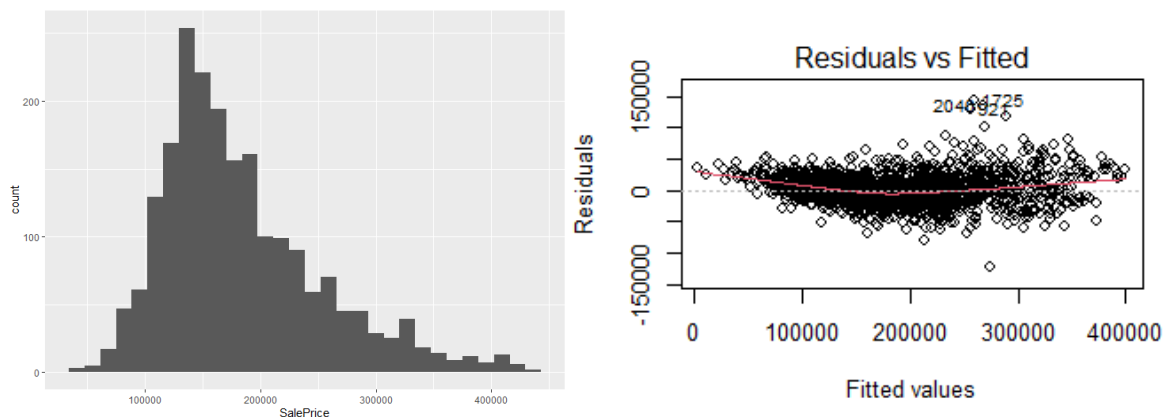
1. Introduction

This report will show the process of linear regression modelling and the interpretation of the final model.

2. Modelling Process

2.1 Data Transformation

The density plot of SalePrice is right-skewed and also if I do the linear regression the Residuals vs Fitted plot shows non-constant variance. Thus, I take log transformation in response SalePrice to eliminate the heteroscedasticity.



Total 28 explanatory variables can be used to predict the house prices. I first do the data transformation for different type of variables.

According to variable descriptions, total 16 qualitative variables can be transformed into factors, which can be further verified by summary function.

There are 4 variables are related to rates or measures but cannot be divided into groups of same values. Their orders matter but values cannot be split into certain levels. Thus, they are not transformed into factors.

For other continuous quantitative variables, scatter plots and density plots or histograms are used to show whether it is necessary to do box-cox transformations. Doing log transformation in Age is more meaningful since we compare the percentage change in house price relative to percentage change in age. House age doesn't matter in units since one house can have hundreds of life length, different from people. Therefore, it is better to take log transformation in the age of house.

Log transformation enlarges the differences between small values and heavily reduce the differences between large values. It makes huge differences between large values ambiguous. Square root transformation also does the same thing but remains the relative differences between large values. Scatterplots of LotFrontage, LotArea and LowQualSinSF shows that most of values are around certain values with quite few huge values of huge difference between most of small values. To eliminate the effects caused by these huge values of small quantities, I use log transformation on these variables. In contrast, scatterplots of GrLivArea, BaseLivArea, and GarageArea also shows there are large differences between small values and large values but with larger amounts in both groups, that is to say, the large differences here would not affect the coefficients because there are also many small values and similar number of large values, different from the before case of thousands of small values and few large values. Therefore, square root transformation is preferred. I use log to eliminate the leverage effects of outliers in LotFrontage, LotArea and LowQualSinSF, and use square root transformation to reduce the units of range and make the coefficients more practical significant in GrLivArea, BaseLivArea, and GarageArea.

Table below shows the transformation done in variables.

| | | |
|------------|--------------------|---|
| Response | Log Transformation | SalePrice |
| Predictors | Factors | Location, Amenities, RoadRail, TwoStory_dum, FlatContour_dum, FlatRoof_dum, Garage_dum, CentralAirNum, KitchenQual_Ex, Zoning_2, Zoning_3, Zoning_4, YrSold_2007, YrSold_2008, YrSold_2009, YrSold_2010 |
| | No change | BedroomAbvGr, Bathrooms, OverallCond, OverallQual, Fireplaces |
| | Transformation | log Age, LotFrontage, LotArea, LowQualFinSF |
| | | sqrt GrLivArea, BaseLivArea, GarageArea |

2.2 Collinearity Detection in Variables

To avoid collinearity in model, I first check the correlation coefficients between variables. Since all correlation coefficients are below 0.7, we need further check in the regression model.

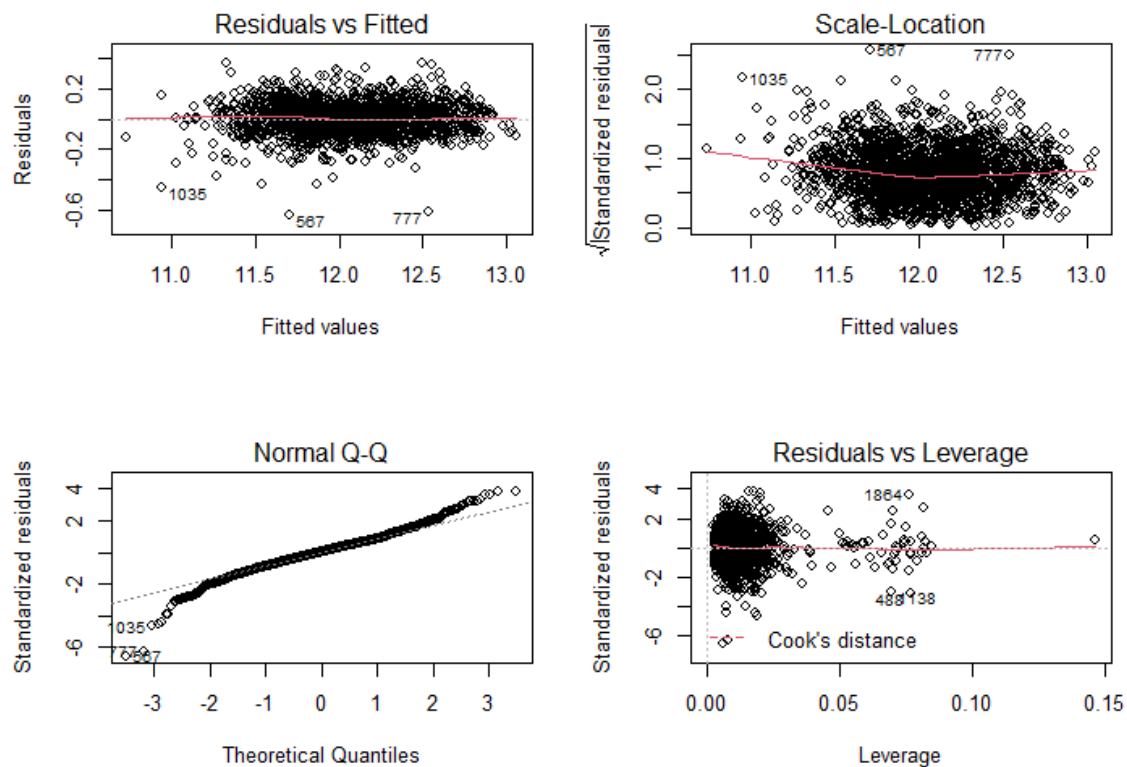
2.3 Model Determination

I use the backward selection to determine the model by first regressing all explanatory variables and then deleting variables of highest p-values until all variables are significant.

The final linear regression model is as below.

$$\begin{aligned} \log(\text{SalePrice}) = & \log(\text{Age}) + \sqrt{\text{GrLivArea}} + \sqrt{\text{BaseLivArea}} + \text{Location} + \text{RoadRail} \\ & + \text{BedroomAbvGr} + \text{OverallCond} + \text{OverallQual} + \log(\text{LotArea}) \\ & + \text{TwoStory_dum} + \text{FlatContour_dum} + \text{FlatRoof_dum} + \sqrt{\text{GarageArea}} \\ & + \text{CentralAirNum} + \log(\text{LowQualFinSF}) + \text{Fireplaces} + \text{KitchenQual_Ex} \\ & + \text{Zoning_2} + \text{Zoning_3} + \text{Zoning_4} + \text{YrSold_2008} + \text{YrSold_2010} \end{aligned}$$

F statistic is significant and also for t statistics of each variable. Thus, predictors are related to response and there is no collinearity. Adjusted R-squared of 0.9292 shows that the model captures most of the characteristics of sale price and 92.92% of variation in log of sale price can be explained by the model.



The almost horizontal line in Residuals vs Fitted plot shows that there is not obvious trend in residuals thus the variance is constant. The Normal Q-Q plot shows that the residuals have a heavier tailed distribution in small values but it does not cause problems since we have a large sample size. Also, there are no outliers with high leverage in Residuals vs Leverage plot.

3. Interpretation

$$\begin{aligned}\log(\text{SalePrice}) = & 9.5193137 - 0.0901882 \log(\text{Age}) + 0.0253959 \sqrt{\text{GrLivArea}} \\ & + 0.0042050 \sqrt{\text{BaseLivArea}} + 0.0655204 \text{Location.L} + 0.0078998 \text{Location.Q} \\ & + 0.0031158 \text{Location.C} - 0.0649119 \text{RoadRail} - 0.0090965 \text{BedroomAbvGr} \\ & + 0.0476878 \text{OverallCond} + 0.0658852 \text{OverallQual} + 0.0872371 \log(\text{LotArea}) \\ & - 0.0494014 \text{TwoStory_dum} - 0.0264875 \text{FlatContour_dum} \\ & + 0.0801957 \text{FlatRoof_dum} + 0.0042992 \sqrt{\text{GarageArea}} \\ & + 0.0854610 \text{CentralAirNum} - 0.0168232 \log(\text{LowQualFinSF}) \\ & + 0.0345650 \text{Fireplaces} + 0.0581327 \text{KitchenQual_Ex} - 0.0522310 \text{Zoning_2} \\ & + 0.0409902 \text{Zoning_3} + 0.0455631 \text{Zoning_4} + 0.0161815 \text{YrSold_2008} \\ & + 0.0316844 \text{YrSold_2010}\end{aligned}$$

Since I do log transformations in most of variables and others are of small values, the coefficients are all relatively small. Log can be interpreted as all else being same, one percent increase in age causes 0.0901882 percentage decrease in house price. Coefficients of factors represent the relative value added to the intercept. All else being the same, sale prices with excellent kitchen quality are 5.81327 percentage higher than those which do not have excellent quality.

4. Prediction

Predicting log transformed response uses following formula,

$$\hat{y} = \exp(\widehat{\log y}) * \left(\frac{\sum_{i=1}^n \exp \hat{u}_i}{n} \right),$$

where \hat{u}_i is the residual.