

Linear Regression of Property Listing Prices in Denver

By: David Heisler

Advisor: Joshua French, Ph.D.

5/18/2017

Abstract:

Linear regression is used to model property prices in the city of Denver as listed on Realtor.com during the month of March. An Ordinary Least Squares linear regression model will be used, with a brief comparison to a Robust Linear model. The model will use the number of bedrooms, bathrooms, car spaces in a garage, property size in square feet, and the zip code as regressors. Before fitting the model, it is confirmed that these regressors are not excessively collinear with one another by computing their Variance Inflation Factors. It is demonstrated that these regressors are both significant and logically appropriate for explaining the response. This paper also introduces the assumptions pertaining to the linear regression models, and subsequently conducts diagnostics on the model to show these assumptions are satisfied.

Linear Regression of Property Listing Prices in Denver

Introduction

For this independent study project, I have chosen to create and evaluate a linear model of property listing prices in the city of Denver. The data for this project were gathered from the website Realtor.com during the month of March. A brief description of each variable in the data set is as follows:

Data Description

- price: The listing price of a property as given on Realtor.com. This variable is numeric and the natural log of this will be the response for the model.
- zip code: This is the zip code each property is listed in. This variable is different depending on where the observation is located in the city of Denver. As such, it is a factor (categorical) variable.
- beds: The number of bedrooms the property has as listed on Realtor.com. This is a numeric variable.
- baths: The number of bathrooms the property has as listed on Realtor.com. This is a numeric variable.
- garage: The number of cars that can fit in the garage of the property or the number of parking spaces reserved for the property. Specifically, this distinction is made when there is no physical garage, but the property still has reserved parking spaces available to it. This value is 0 if the property does not have a garage or the property does not have a parking space. For example, some homes in Denver only offer street parking. This is a numeric variable.
- property_size: This is the size of the property in square feet. This is a numeric variable
- lot_size: This is the size of the lot the property is on. This variable is NA when there is no lot associated with a property. For example, some apartments and condos don't offer lots. This is a numeric variable.

Motivation

I chose this as my independent study project, because in the future I want to start a real estate investment company. One of the main pillars of my investment philosophy will be the statistical analysis of potential investment opportunities. I one day hope to develop this model even further and turn it into a production quality tool to aid me in determining what properties represent the best investment opportunities.

Methodology

The model I chose for this project is a linear regression model fitted using Ordinary Least Squares. Not all variables in my data set will be included in my final model. When selecting regressors, I attempted to balance logical explanation with respect to the response and

observation completeness. One example of a variable that would be logical in explaining the response, but reduced my data set too much due to incompleteness, is the `lot_size` variable. Over 900 observations in my data set have NA values for the `lot_size` variable, and thus my data set would have been reduced dramatically if I included this in my model. The regressors included in the model will be `baths`, `log(property_size)`, `garage`, `beds`, and `zipcode`. Before determining these regressors were acceptable, I verified they weren't too collinear with one another. I did this by finding their Variance Inflation Factors (see Appendix A). Before going any further, some background of the OLS model will be helpful.

Ordinary Least Squares Regression

The Ordinary Least Squares (OLS) model is one of the most basic linear regression models available to statisticians. The model may be represented in matrix form as:

$$y = X\beta + \epsilon$$

The components of this model are:

- y : an $n \times 1$ column vector of response values. This vector is random, and observable. In my model, the dimensions are 1876×1 , since some of the observations from the data set were removed due to incomplete regressor values.
- X : this is an $n \times p$ matrix of the observed values for the regressors. These values are fixed, and the dimensions in this model are 1876×38 . The 38 comes from the expansion of the zip code variable, which is a dummy variable.
- β : this is the $p \times 1$ vector of regression coefficients. In this model, the dimensions are 38×1 .
- ϵ : This is the $n \times 1$ column vector of errors, which are random and unobservable. For this model, the dimension is 1876×1 .

A linear regression model is developed by estimating the β parameters that minimize the residual sums of squares - this is where the "least squares" comes in. By the Gauss-Markov theorem, so long as the errors have constant variance, have mean zero, and are uncorrelated, $\hat{\beta}$ is the best linear unbiased estimator of β .

The residuals of the model are:

$$\hat{\epsilon} = y - \hat{y}.$$

or the difference between the observed response values and our model's fitted values. We want to minimize the residual sum of squares, or:

$$\hat{\epsilon}^T \hat{\epsilon} = (y - X\hat{\beta})^T (y - X\hat{\beta}).$$

We want to find $\hat{\beta}$ by differentiating with respect to $\hat{\beta}$ and setting equal to 0. This leaves us with,

$$X^T X \hat{\beta} = X^T y.$$

Using the normal equations, and so long as $X^T X$ is invertible,

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

multiplying both sides by X,

$$X \hat{\beta} = X(X^T X)^{-1} X^T y.$$

Once we find $\hat{\beta}$, our fitted model takes the form:

$$\hat{y} = X \hat{\beta}$$

In the fitted model, the \hat{y} is the 1876 x 1 column vector of fitted (or predicted) values for an observation given the regressors in X. The $\hat{\beta}$ is the 38 x 1 column vector of estimates for the regression coefficients. These are the values that minimize the residual sums of squares.

Exploratory Analysis

Before I begin discussing the assumptions related to this model, I want to perform some exploratory analysis to learn more about the data I have.

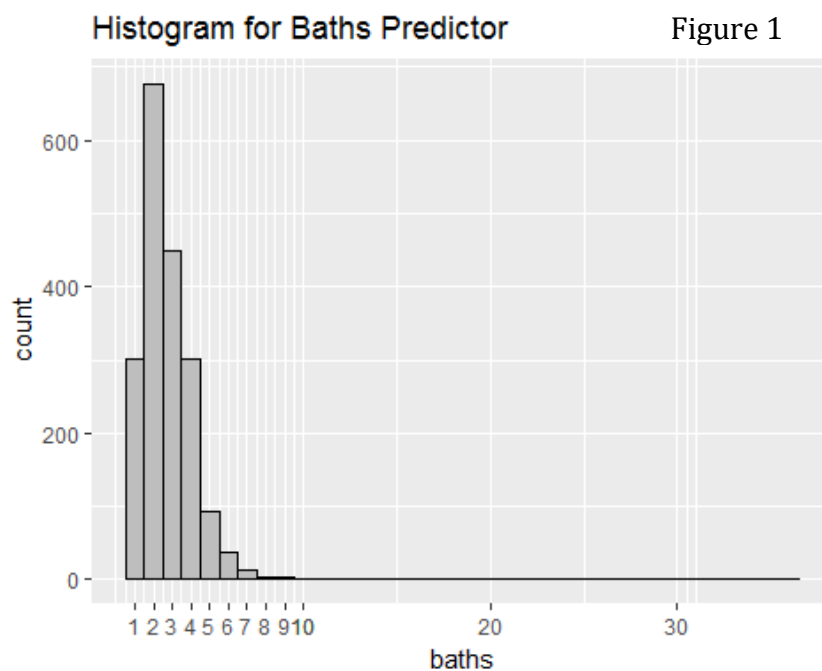


Figure 1

Figure 1 is a histogram for the `baths` regressor. From this plot, it is evident that almost all of the properties in my data set have 1-5 bathrooms. However, there does seem to be some positive skewing in this data.

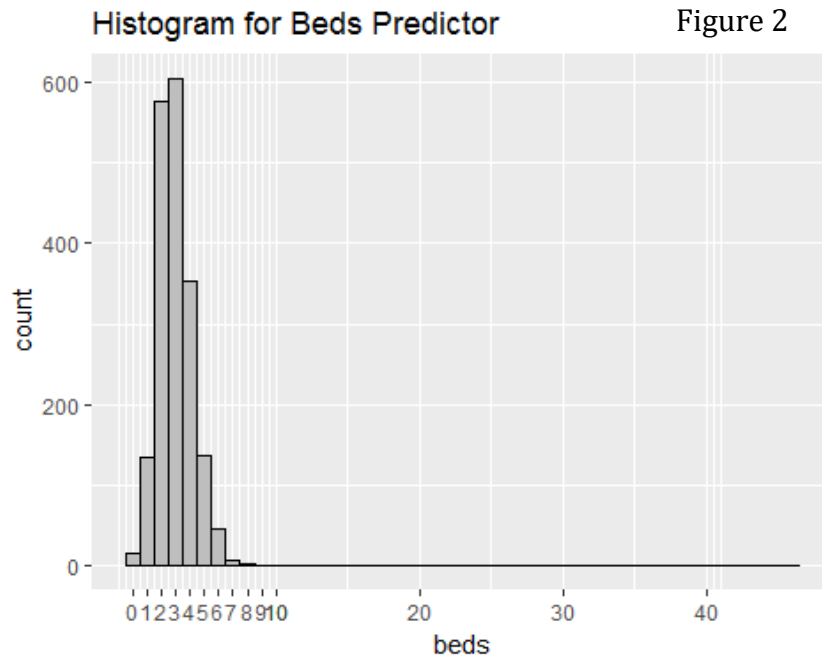


Figure 2 is a histogram for the beds regressor. This plot shares similarities with the baths histogram. Most of the data is between 1-5, but there is some positive skewing. The major difference between this histogram and the baths histogram is that there are some properties that have a value of 0 for the beds regressor. I looked up the properties with 0 beds listed in my data set on Realtor.com, and they are a mix of properties that are

classified as just land, and properties which appear to have the correctly listed number of beds. Other than the properties which are land, I believe the properties in my data set with 0 beds were caused by the algorithm I used to collect the data.

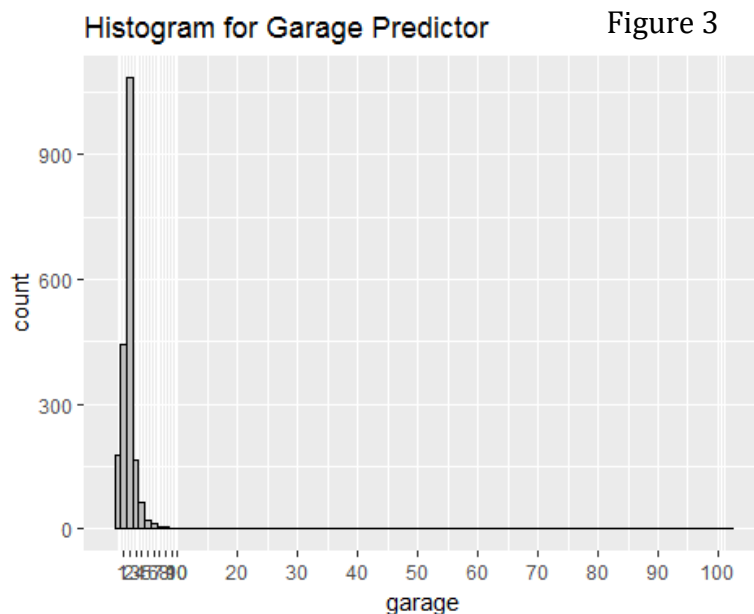


Figure 3 is a histogram for the garage regressor. The pattern in this data follows closely to what was seen in the baths and beds histograms. However, there is a greater positive skewing affect, because one property is listed as having a value of 102 for the garage regressor. After looking this property up on Realtor.com, it appears that this property is in an apartment complex, and the listing shows all of the parking available at the complex. This observation actually became a potential outlier, so I will address it later.



Figure 4 is a scatter plot of the price versus the property_size. This scatter plot appears to show nonconstant variance. Specifically, an increase in variance of the response, price, as the regressor, property_size, increases. As such, I will apply a natural log transformation to the response, price, and examine the relationship between these two variables.



After applying the transformation, Figure 5 appears to show nonlinear relationship. In order to try to change this, I will apply a natural log transformation to the regressor, property_size, and reexamine the relationship again.

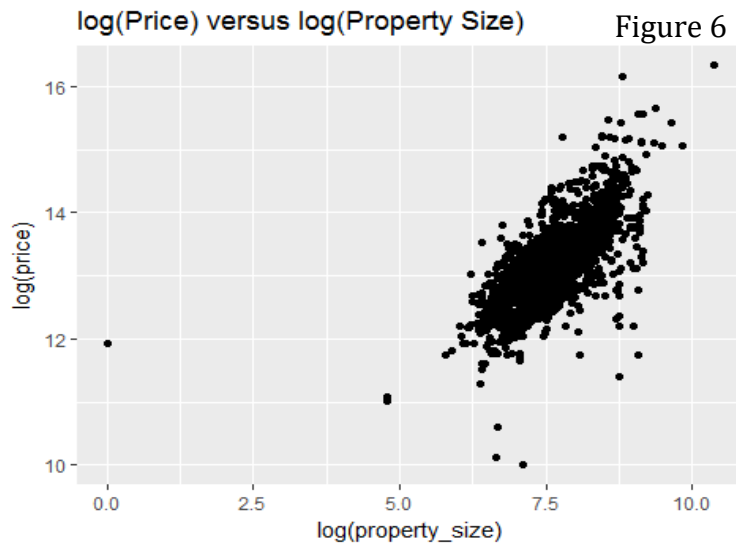


Figure 6

After applying this final transformation, it is evident that there is now a linear relationship between the $\log(\text{price})$ and $\log(\text{property_size})$ variables. As such, I will move forward with these two transformations used in my model.

Model Validation

In order to properly implement a linear regression model, there are a few assumptions that must be satisfied. Without these assumptions, any inferences, conclusions, or predictions will be misleading. I will now discuss some of these assumptions generally, and later will address these within the context of my model in the Diagnostics section of this paper. There are three main categories of assumptions that must be satisfied. In no particular order, we have:

1. Assumptions of the errors in the model,
2. Assumptions of the structure of the model, and
3. Unusual observations

Assumptions of the errors

In order for most inferential methods to be applied after a linear regression model has been fit, it is important that a few assumptions regarding the errors are satisfied. Specifically, the errors, ϵ , must have constant variance, be described by a normal distribution, and must be uncorrelated. Since my data is spatial in nature due to the `zipcode` regressor, I will also check for spatial correlation among the errors. A brief description of how the satisfaction of these assumptions is determined is described below.

Constant Variance

In order to check the errors have constant variance, a plot of the residuals, $\hat{\epsilon}$, versus the fitted values, \hat{y} , can be used. Using this plot, a model that satisfies this assumption will have points scattered with even thickness in the y-direction. This plot is also useful to show that the errors are centered at 0, a key assumption of the Gauss-Markov Theorem.

Normality

In order to determine if errors are consistent with a normal distribution, a Q-Q plot of the errors can be used. It is then possible to determine how well the errors fit the Q-Q plot. In this plot, there is a line which represents data from a normal distribution. The ideal plot to show the satisfaction of this assumption would have points scattered along this line.

Geospatial Correlation

Since one of the regressors in my model is spatial in nature (`zipcode`), I will want to check for any spatial correlation in the errors. To check this assumption, I will create boxplots for each of the residuals based on `zipcode`. If there is spatial correlation, this visualization will show many of the boxplots dramatically skewed from zero. If there is no spatial correlation, then this visualization will show all of the boxplots centered near 0. A little deviation from 0 is acceptable, a problem is present when there is systemic deviation from 0.

Structural Assumptions

Now that the assumptions regarding the errors, ϵ , have been discussed, the assumption concerning the model's structure can be introduced. Of all assumptions, this is the most important. Without this assumption being satisfied, the effectiveness of the implemented OLS model is nullified.

This assumption assesses whether:

$$E(y) = X\beta.$$

This means that the regressors in X are appropriate for describing the typical structure or average pattern of the response. In order to show this assumption has been satisfied, there are a couple of plots that can be used. I will be using the Component Plus Residual plots and Added Variable plots.

The Component Plus Residual plot is used to visualize the relationship between a specified regressor, given the other regressors are present in the model. Specifically, the plot takes the residuals from the model plus the i th regression coefficient, $\hat{\beta}_i$, multiplied by X_i and compares it to X_i . Mathematically, we have:

$$\hat{\epsilon} + \hat{\beta}_i X_i \sim X_i$$

The plot looks like a simple scatter plot, where the slope of the points is the estimated coefficient, $\hat{\beta}_i$, for a specified regressor, X_i . This plot is useful for determining if there is a linear relationship between the response and the specified regressor. In this plot, one should ensure there is no underlying nonlinear structure in the points.

The Added Variable plot is similar to the Component Plus Residual plot; however, this plot is most beneficial for detecting points that may be influential. The plot is most easily described using mathematical notation in a few steps, where (i) denotes every regressor but the i th:

$$\hat{\epsilon}_{(i)} = Y - X_{(i)}\hat{\beta}_{(i)} \quad (1)$$

$$\tilde{\epsilon} = X_i - \hat{X}_{(i)} \quad (2)$$

$$\hat{\epsilon}_{(i)} \sim \tilde{\epsilon}$$

First, we regress the response on all of the regressors except the regressor we are creating the Added Variable plot for. Next, we regress this regressor on all other regressors in our model. Finally, we regress the residuals from the equation (1) against the residuals from equation (2).

To determine if an observation is influential, I will look to see if the line in the plot, which represents the estimated coefficient, $\hat{\beta}_i$, as a slope for the regressor in question is being skewed away from the bulk of the data.

Unusual Observations

Once the necessary assumptions regarding the errors and the structure of the model have been satisfied, it is important to check if some unusual observations are dramatically affecting the model fit. There are three important types of unusual observations that we are looking for.

1. Outliers,
2. Leverage values, and
3. Influential observations.

Outliers

Outliers are observations whose response don't fit the pattern of the data. Outliers have the potential to affect how the model fits the data, so it is important to determine whether they are present in the data.

Leverage Values

Leverage values are unusual observations in the regressor space. As with outliers, these can potentially affect how well the model fits the data.

Influential Observations

Influential observations are observations that cause an even greater impact on model fit than outliers and leverage values. Influential observations don't have to be observations that are outliers or ones with large leverage, but they usually have one of these two characteristics. To test for influential observations, I will use an influence plot. The influence plot compares studentized residuals versus the leverage values in the model. Studentized residuals are the difference between the observed value and the fitted value of a model excluding the one observation, divided by the estimated standard error. By doing this, the influence plot effectively helps to find how outliers and leverage values are influencing the overall model fit. In the influence plot, I will be looking for observations that have a high studentized residual and/or high leverage value.

Results

In the prior section I gave an overview of the model assumptions that must be satisfied in order for an OLS model to be valid. In this section, I will report some results of the model, I will go

through and conduct each of the procedures I described in the last section, and I will show that the model satisfies these assumptions.

Model Results

Table 1 gives some key characteristics of the model's overall fit. These include the estimated coefficients, $\hat{\beta}$, their estimated standard errors, and their associated p-values.

Table 1: Model Coefficient Estimates

Coefficient	Estimate	Estimated Standard Error	P-value
Baths	.0940	.007503	< 2e-16
Garage	.0119	.002364	.00023
Log(property_size)	.7533	.017010	< 2e-16
Beds	-.0579	.006440	< 2e-16

The residual standard error of the model is .2573, and the R^2 is .8336. All estimated coefficients make sense, except the `beds` coefficient. Within context of the problem a negative estimated coefficient value implies that as the number of beds increases, the listing price of the property decreases by 5.9%.

Assumption of the Errors

Constant Variance

In order to show this assumption is satisfied, I will plot the residuals, $\hat{\epsilon}$, versus the fitted values, \hat{y} .

Figure 7: Residuals versus Fitted Values

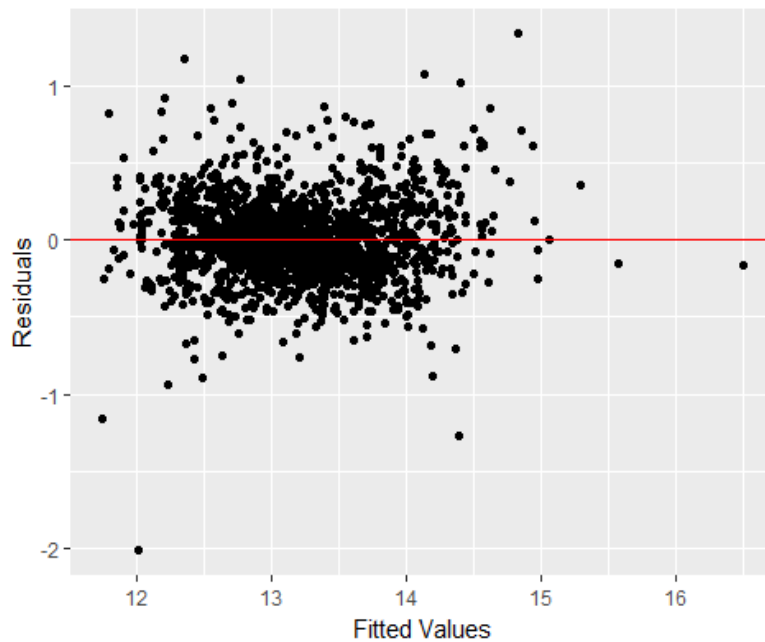


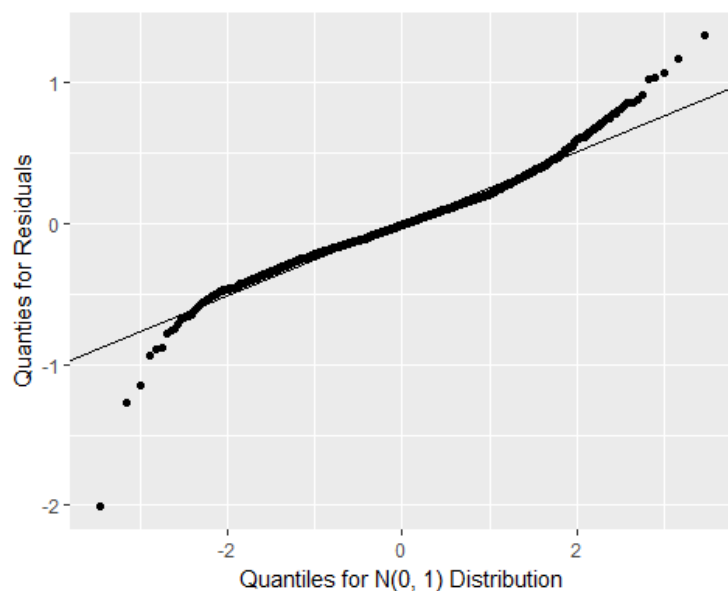
Figure 7 is the residuals versus the fitted values for my model. First, it appears that there is a smooth scattering of points across the line $y = 0$. This is a good indicator that the variance is constant. The second thing to note is that there is no overarching structure to the errors. To me, this plot demonstrates that the constant variance of the errors has been shown.

Normality

Next, I will verify that the errors come from a normal distribution. To do this I will present a Q-Q plot of the residuals using Figure 8.

From this plot, it is evident that the errors do not come from a perfectly normal distribution, because our points begin to depart from the line on the graph, which represents what truly normal data would look like. This is something that raises concern. As mentioned before, it is important that our errors come from a normal distribution, because this ensures any inference made using the model is valid. Figure 10 is a density plot of the residuals. This will be an additional aid in determining if the residuals are normal or not.

Figure 8: Q-Q Plot of Residuals



Based on Figure 9, I believe that the errors are close to a normal distribution, however, they depart slightly. This shouldn't be a problem, but I want to dig a little deeper to determine if this is acceptable or not.

In order to test whether this is truly a problem in my model, I will fit a robust linear model, and compare the estimated coefficients, $\hat{\beta}$, to the ones in my OLS model. If these estimates differ by a substantial amount for each regressor, then I would consider using the robust method over the OLS model.

Before fitting the robust model, I want to give an overview of it and why I would use it in this case.

Robust regression is a technique used to combat the fact that an OLS model does not satisfy all the necessary assumptions. In this case, the residuals from my model do not correspond to a perfectly normal distribution. Robust regression attempts to alleviate this by assigning weights to the observations in the data set, where outliers are given smaller weights. The outliers receive smaller weights so they aren't as influential on the model's fit.

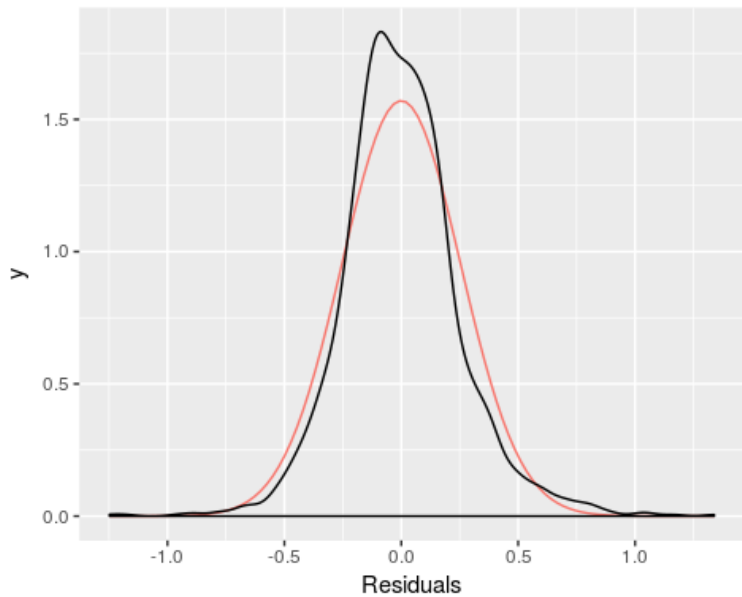
Table 2 shows the output of the coefficients from the two models.

Table 2: Models' Coefficient Estimates

Coefficients	Ordinary Least Squares	Robust Regression
Garage	.0119	.0111
Log(property_size)	.7533	.7371
Baths	.0940	.0814
Beds	-.0579	-.0474

In order to interpret the estimates given from each model, I will interpret them within context of the listing price. In the OLS model, the `garage` regressor increases the listing price of a property by 1.11% for one vehicle increase in garage size or parking availability. Compared to

Figure 9: Density Plot of Residuals (Black) Compared to $N(0, .06)$ Distribution Density (Red)

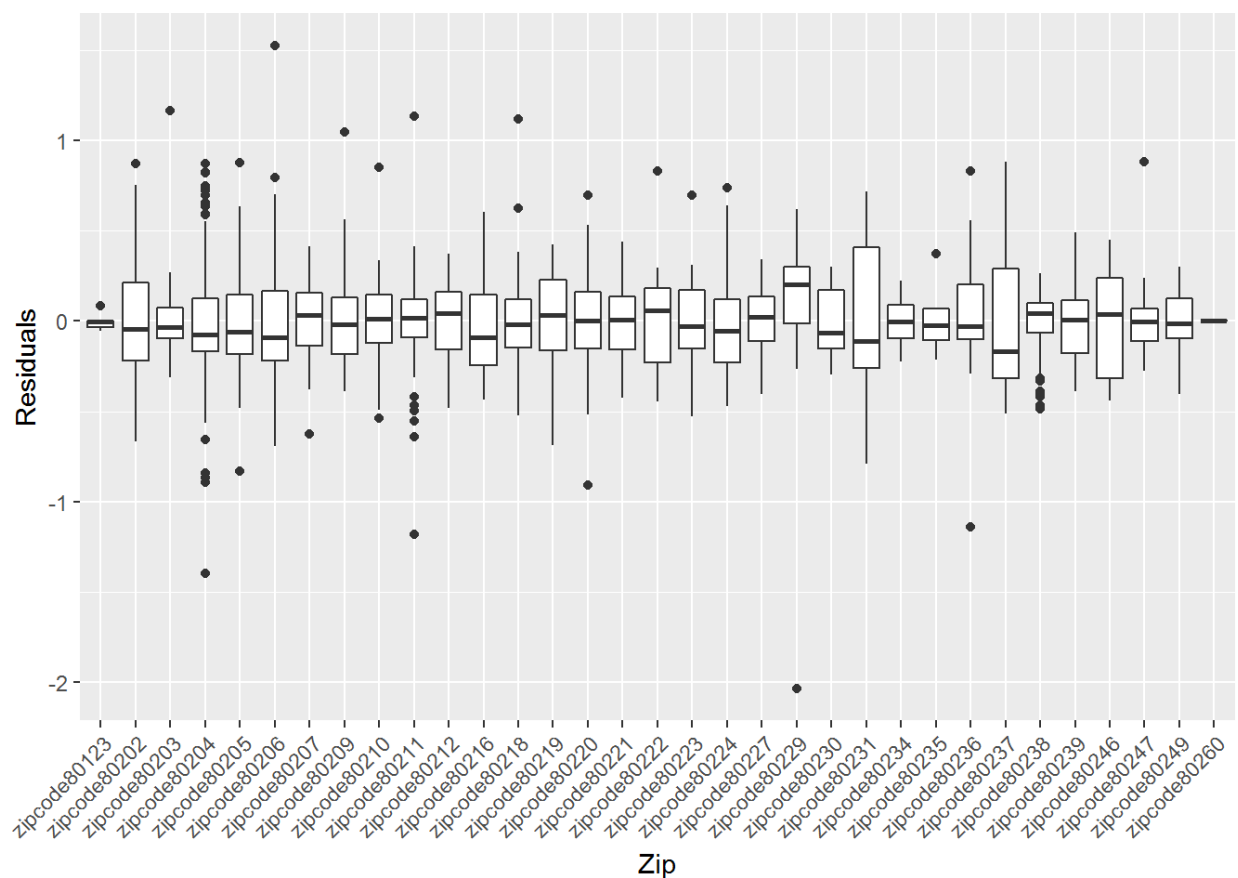


the robust model, where a one vehicle increase in size increases the listing price of a property by 1.19%. For the $\log(\text{property_size})$ regressor, each 10% increase in square footage a property has, the property listing price will increase by 7.44% in the OLS model, and by 7.27% in the robust model. For one additional bathroom in a property, the listing price will increase by 9.8% in the OLS model, and by 8.50% in the robust model. Finally, one additional bedroom in a property decreases the value of the listing price by 5.62% in the OLS model, and by 4.63% in the robust model. To me, this is not a substantial change, and I will therefore move forward with the OLS model.

Geospatial Correlation

Below is a visualization representing a boxplot of the residuals for each zipcode in my model.

Figure 10: Residuals by Zip Code



Based on this plot, I believe that no spatial correlation is occurring. If there were spatial correlation, there would be a clear trend in how the boxplots moved in the y-direction. Since all of the boxplots are reasonably centered near 0, I think this is sufficient to show there is no spatial correlation.

Model Structure

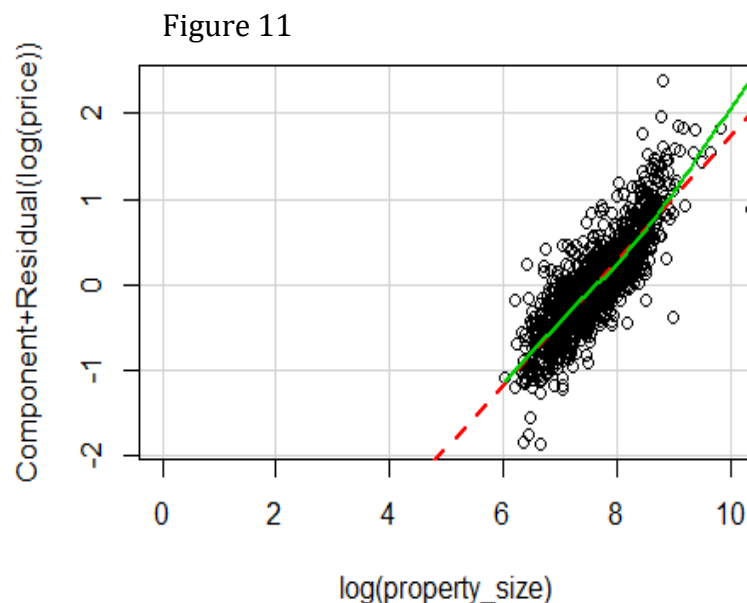
In order to test the structural assumption of the model, $E(y) = X\beta$, I will use Component Plus Residual plots as well as an Added Variable plots for the garage, baths, beds, and

`log(property_size)` regressors in the model. I am including plots for these regressors and not the `zip code` regressor for two reasons. First, because these are the significant regressors in the model at the 5% level. Second, the `zipcode` variable is a dummy variable, and it is expanded out to over 20 separate regressors. Obviously, this would take a lot of plotting to verify the structural assumptions for each of these.

Component Plus Residual Plot

Below is the Component Plus Residual plot for the `log(property_size)` regressor. The Component Plus Residual and Added Variable plots for the other regressors are in Appendix B with a description of whether or not they give sufficient information to determine the structural assumption has been satisfied.

Component Plus Residual plot for the `log(property_size)` regressor:



The red line in Figure 11 signifies the estimated coefficient value $\hat{\beta}_{\log(\text{property size})}$ and the green line is a line generated internally to try to fit the data in a smooth fashion. When the red and green lines deviate from one another dramatically, this indicates there is nonlinear structure present in the data. I think this plot shows that the `log(property_size)` regressor is fitting the response well, and there isn't any odd or nonlinear structure present.

Added Variable Plot

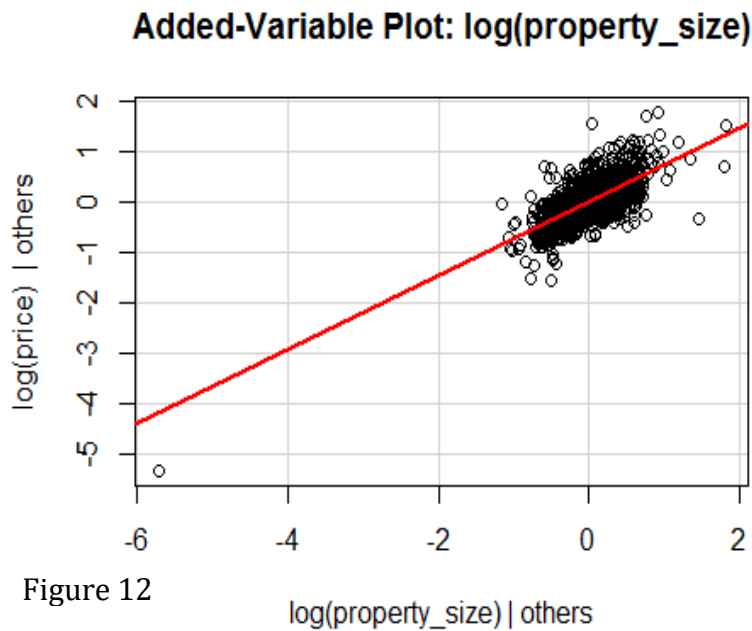


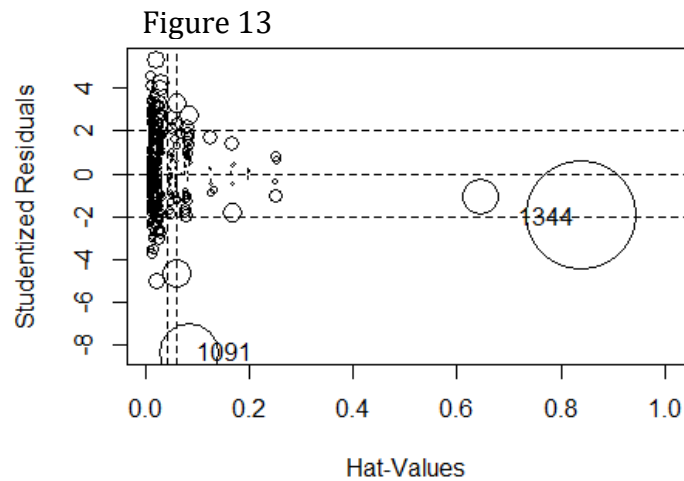
Figure 12

This is the Added Variable plot for the $\log(\text{property_size})$ regressor. From this plot, I think that everything looks good with the $\log(\text{property_size})$ regressor. I want to note that there is one observation that seems to be nonconforming (lower left of plot). However, I don't believe this is affecting how the model is fitting the data.

This concludes the structural diagnostics. Now I want to look at potential influential observations.

Influential Observations

In order to check for influential observations, I will use an influence plot. The influence plot compares the studentized residuals versus the hat-value for each observation. The area of the circle for each observation is associated with each observation's Cook's distance. Observations that are potentially influential will have either a large circle area, or will be towards the top right of the plot. If any observations come up as influential, I will then fit a model with them removed and determine if the estimated coefficients change dramatically.



Based on the influence plot, it seems that observation 1344 and 1091 are the most influential, as they correspond to the observations with the largest circle areas. Now I will fit the same model without these points and compare the estimated coefficients, $\hat{\beta}$, for the garage, $\log(\text{property_size})$, baths, and beds regressors. If these estimates differ substantially, then I will consider this to be an influential observation.

Table 3 shows the estimated coefficients for the original model, one without observation 1091, and one without observation 1344.

Table 3: Comparison of Models without Potential Influential Points

Coefficients	Original Model	Without Obs. 1091	Without Obs. 1344
Garage	.0119	.0119	.0167
Log(property_size)	.7533	.7533	.7540
Baths	.0940	.0940	.0933
Beds	-.0579	-.0579	-.0601

From Table 3, it is clear that the estimated coefficient values between the original model and the model with observation 1091 excluded don't differ at all. However, there are some differences that occur between the original model and the model with observation 1344 excluded. For the garage regressor, one additional vehicle for a property equates to a 1.11% increase in the listing price in the original model, and a 1.67% increase in the model without observation 1344. For 10% increase in the size of a property in square feet, the listing price will increase 7.44% in the original model, and 7.45% in the model with observation 1344 removed. For one additional bathroom in a property, the listing price of the property increases by 9.8% in the original model, and by 9.7% in the model with observation 1344 excluded. Finally, one additional bed in a property leads to a decrease in the listing price by 5.62% in the original model, and by 5.8% in the model with observation 1344 excluded. After further digging, it appears that observation

1344 was the observation with a value of 102 for the `garage` regressor. Moving forward, I would remove this from my data set, as it represents an incorrect value for the garage regressor. This property is actually an apartment, and the listing gives the number of parking spaces in the entire complex as the garage value.

Conclusions

The model I have chosen to keep is the OLS model using `baths`, `log(property_size)`, `garage`, `beds`, and `zipcode` as regressors, and `log(price)` as the response. The model satisfied nearly all assumptions, and only failed to exhibit a normal distribution in the errors. I chose the OLS model with these regressors, because using a Robust Linear Model didn't have a substantial impact how the regression coefficients affected the response.

Almost all estimated model coefficients made logical sense. However, the `beds` coefficient was negative, which was somewhat surprising, since it implies a property loses listing value for increases in the number of bedrooms.

Figure 16 shows interpretations of changes in the expected value of the property listing price for each of the regressors, given all other regressor values remain the same. These numbers reflect the proper transformation back to the original scale.

Figure 16: Interpretations of Changes in Single Regressors

Regressor	Change in Regressor	Change in Expected Value of Property Listing
Beds	1 additional bedroom	Decrease of 5.9% in property listing price
Baths	1 additional bedroom	Increase of 9.8% in property listing price
Log(property_size)	10% increase in property size	Increase of 7.6% in property listing price
Garage	1 additional car fit in garage	Increase of 1.3% in property listing price

It is important to clarify a few things when interpreting this table. The changes in the expected value of the response are given by comparing two identical properties, other than a change in the regressor in question. For example, two properties located in the same zip code, with the same

garage size, same number of bedrooms, and the same property size, could be compared to determine how an additional bathroom would affect the expected value of the listing price. In this case, if one property had an additional bathroom, we would expect this property to have a listing price that is 9.8% greater than the other property on average.

Moving forward, I would like to build upon this model in a few ways. First, I want to give a more detailed geographical view of these properties by classifying each property into a neighborhood of Denver – Capitol Hill, River North, etc. Second, I want to obtain some publicly available data on how old each property is and include that in the model as a regressor.

Appendix A

Figure 17: Variance Inflation Factor (VIF) and Generalized Variance Inflation Factor (GVIF) for Regressors

Regressors	GVIF/VIF
Baths	3.572
Zip code	1.502
Garage	1.151
Log(property_size)	2.605
Beds	2.974

Figure 17 shows either the Variance Inflation Factors (VIF) or the Generalized Variance Inflation Factors (GVIF) for each regressor. GVIF is calculated instead of VIF for dummy regressors or polynomial term regressors. Thus, the `zipcode` regressor is the only one that has a GVIF value.

The Variance Inflation Factor (VIF) is given by:

$$VIF_k = \frac{1}{1 - R_k^2}$$

Where the R_k^2 is the squared multiple correlation of the k th regressor on all other regressors.

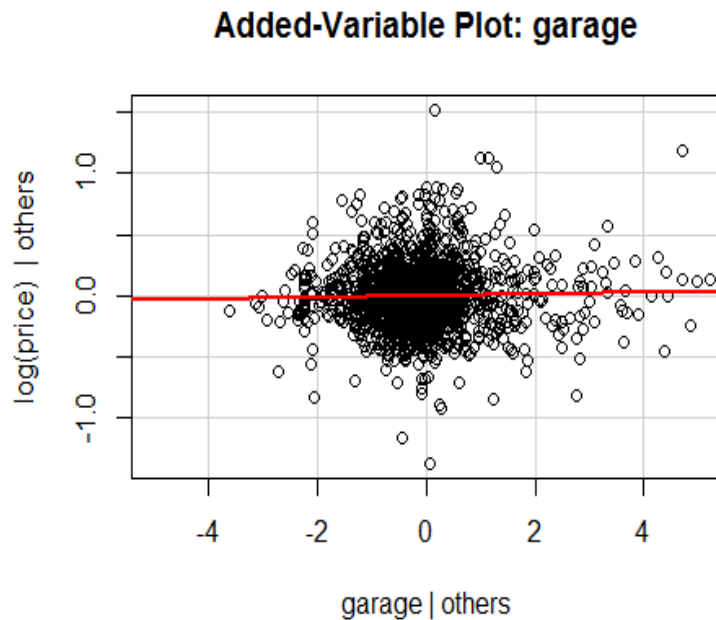
GVIF values are defined as

$$GVIF = \frac{\det \mathbf{R}_{11} \det \mathbf{R}_{22}}{\det \mathbf{R}}$$

Where \mathbf{R}_{11} is the correlation matrix of the columns created by the `zipcode` dummy variable, \mathbf{R}_{22} is the correlation matrix of all the other columns, and \mathbf{R} is the correlation matrix of all columns. The GVIF represent how much larger confidence intervals or confidence regions are compared to if there were truly orthogonal data in the model. Since larger regions imply less certainty with estimates, we prefer smaller GVIF values.

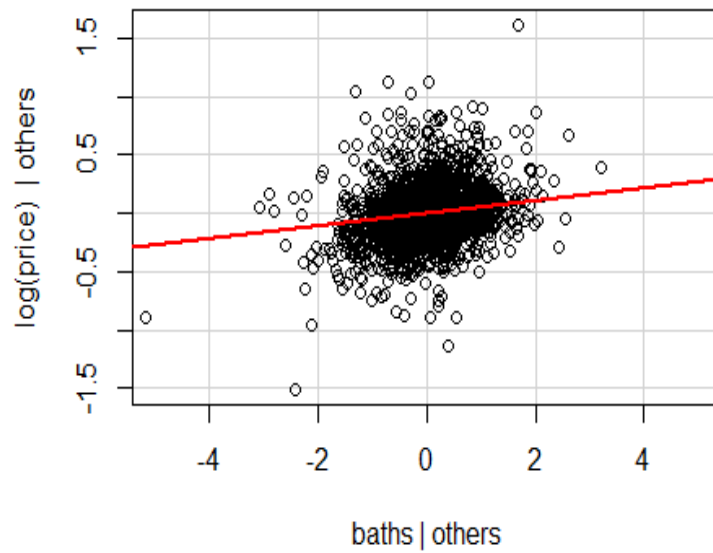
Appendix B

This appendix contains the Component Plus Residual and Added Variable plots for the `garage` and `baths` regressors.



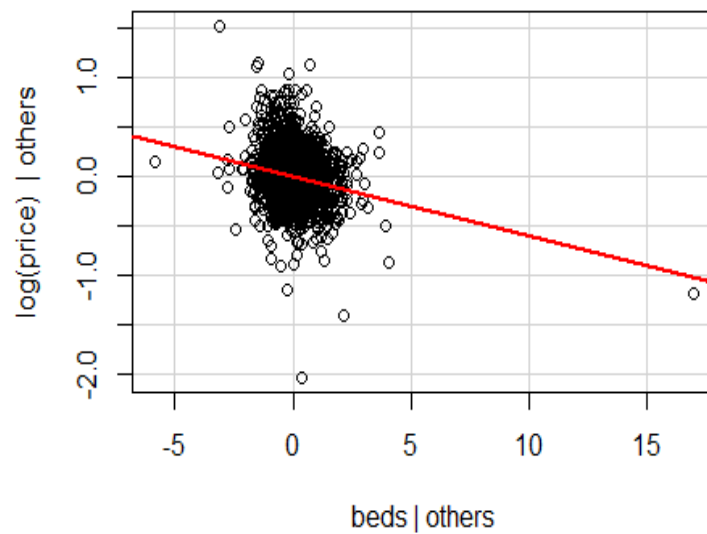
This is the Added Variable Plot for the `garage` regressor. From this plot, it is evident that there aren't any observations affecting how the model fits the data.

Added-Variable Plot: baths

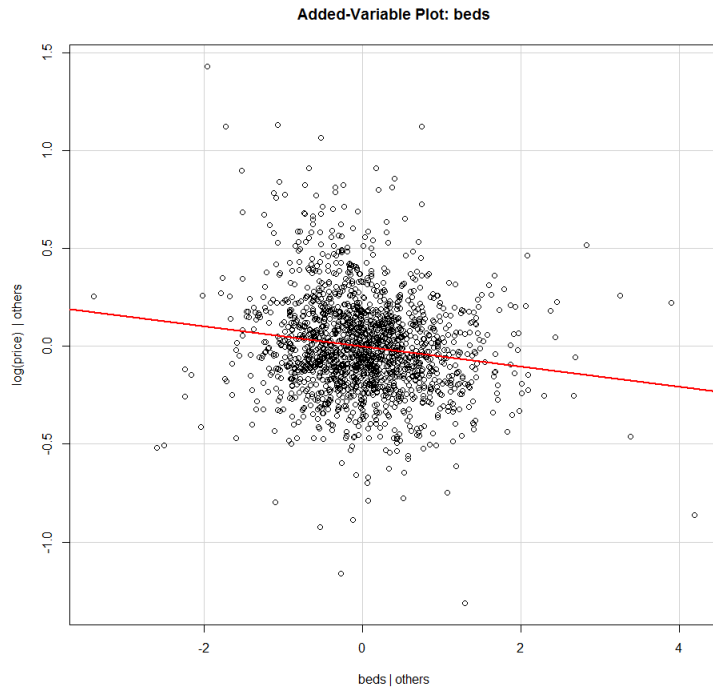


This is the Added Variable plot for the `baths` regressor. Based on this plot, there doesn't appear to be any observations affecting how the model fits the data.

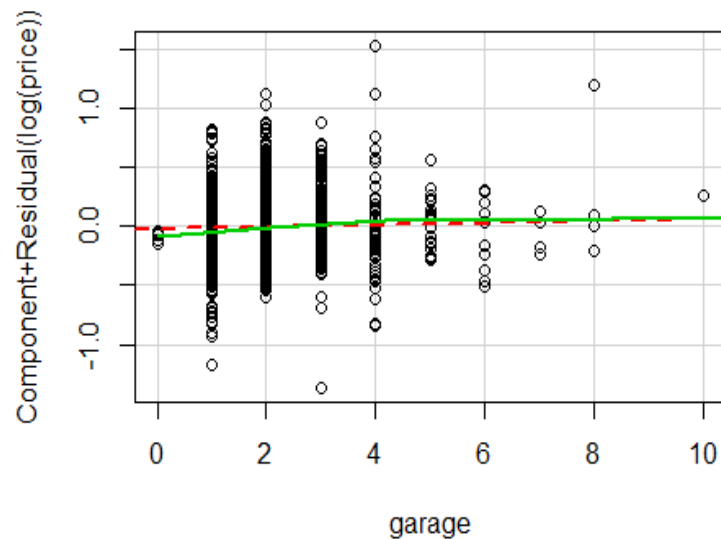
Added-Variable Plot: beds



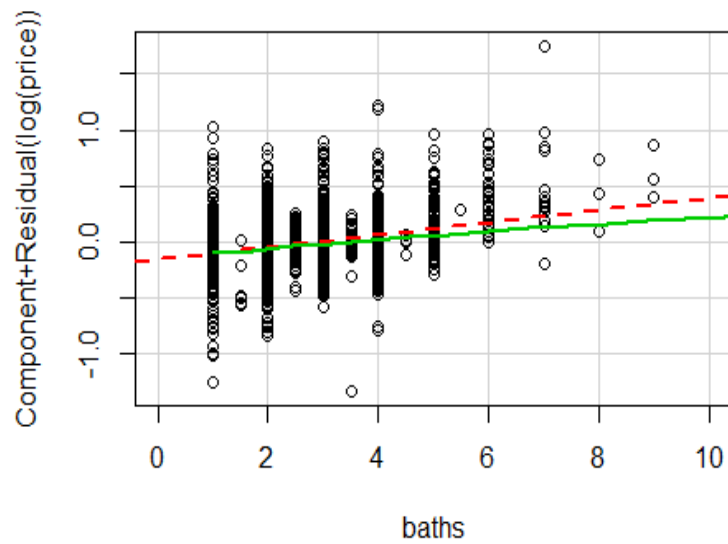
This is the Added-Variable Plot for the `beds` regressor. We can see from this plot that there is one observation that is not conforming to the rest of the data. However, it doesn't appear that this observation is effecting the fit of the line through the bulk of the data.



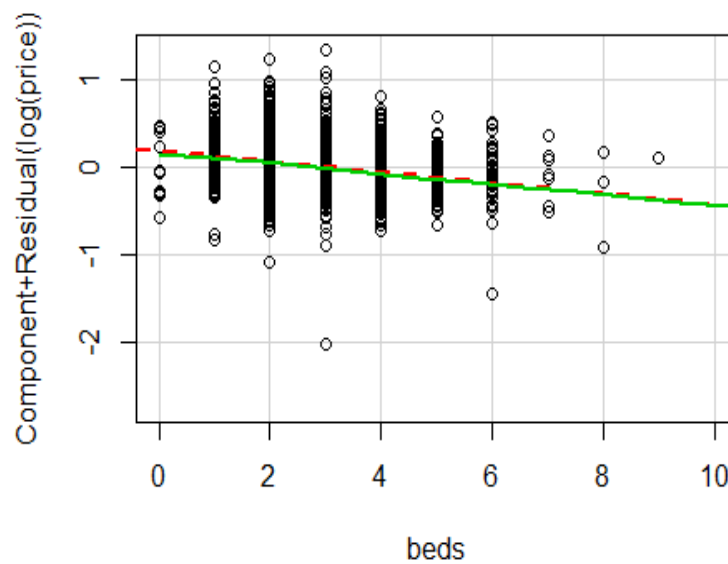
This is the Added-Variable Plot for the beds regressor without outlier above. As we can see, there is still a negative relationship between the beds regressor and the response.



This is the Component Plus Residual Plot for the garage regressor. There are vertical lines formed by the points, because the garage is a discrete variable. Based on this plot, I believe that the linear structure of the data is accurately defined by the model.



This is the Component Plus Residual Plot for the `baths` regressor. From this plot, I believe that the linear relationship between `baths` regressor and the response is confirmed.



This is the Component Plus Residual Plot for the `beds` regressor. From this plot, I believe that the linear relationship between the `beds` regressor and the response is confirmed.

