

Modeling How Living Area Impacts Home Prices Across Neighborhoods

Katelyn McDonald

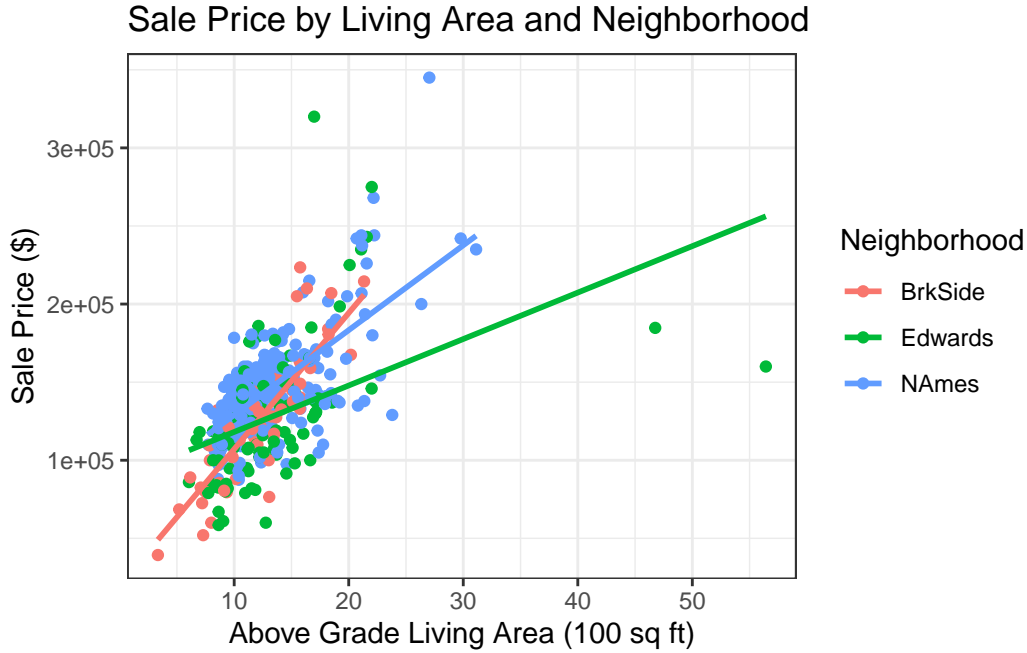
Table of contents

Problem Statement	1
Model Selection	2
Addressing Diagnostics	3
Influential Point Analysis	5
Final Model and Interpretations	7
Conclusion	8
Take Home Message	8
Appendix	9

Problem Statement

Century 21 Ames is interested in understanding how the sale price of homes in three neighborhoods, Brookside, Edwards, and North Ames, is related to the size of the home's above grade living area. Because realtors usually discuss square footage in increments of 100 square feet, the analysis uses living area measured in 100 square feet units. The goal of this project is to quantify how much additional living area contributes to higher sale prices and to determine whether this relationship differs across neighborhoods. The output will be a model for the agency of sale price based on neighborhood and living area.

A scatterplot of sale price versus living area, colored by neighborhood, displays the differing price patterns across the three neighborhoods. The fitted regression lines illustrate that Brookside line has the steepest slope, while both Edwards and North Ames show shallower increases with living area.



Model Selection

The first model fit was where the neighborhoods differ only by intercept which implies that the added value of extra living area is the same across all three locations. The second model included an interaction between neighborhood and Gr100, allowing each neighborhood to have its own slope. Comparing the two models, the model with the interaction term has all predictors are significant at the .001 level excluding the intercept, while the base model has neighborhood = Edwards as not statistically significant. The R^2 also drastically improved from 0.397 to 0.447 with the interaction model. This interaction model is written as:

$$\text{SalePrice}_i = \beta_0 + \beta_1 \cdot \text{Edwards}_i + \beta_2 \cdot \text{NAmes}_i + \beta_3 \cdot \text{Gr100}_i + \beta_4 \cdot (\text{Edwards}_i \times \text{Gr100}_i) + \beta_5 \cdot (\text{NAmes}_i \times \text{Gr100}_i) + \varepsilon_i$$

where

$$\text{Edwards}_i = 1 \quad \text{if the home is in Edwards, 0 otherwise}$$

$$\text{NAmes}_i = 1 \quad \text{if the home is in North Ames, 0 otherwise}$$

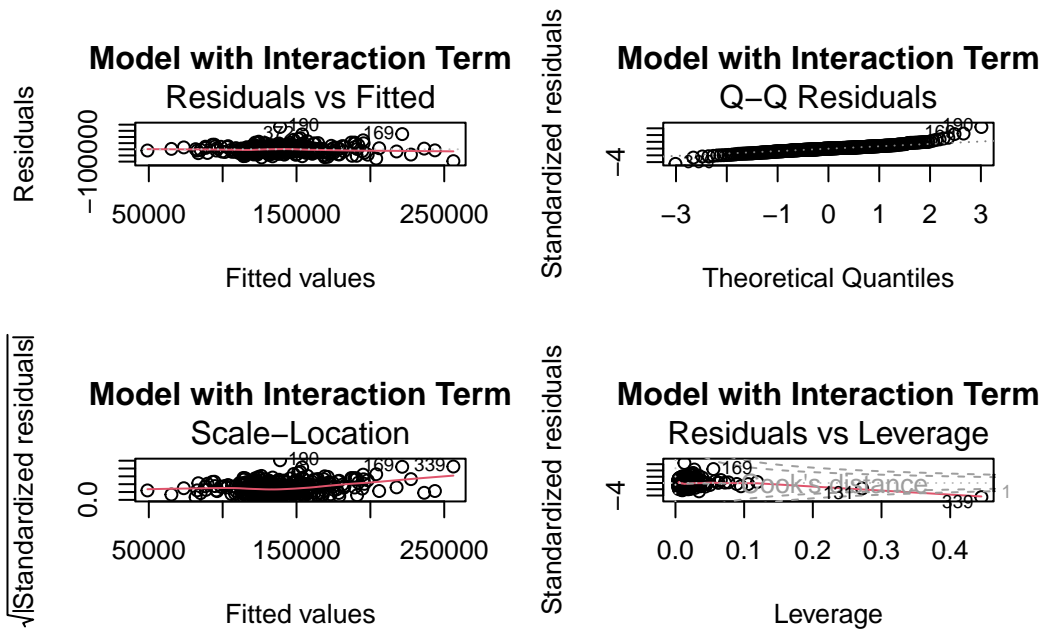
$$\text{Gr100}_i = \text{GrLivArea}_i / 100$$

Addressing Diagnostics

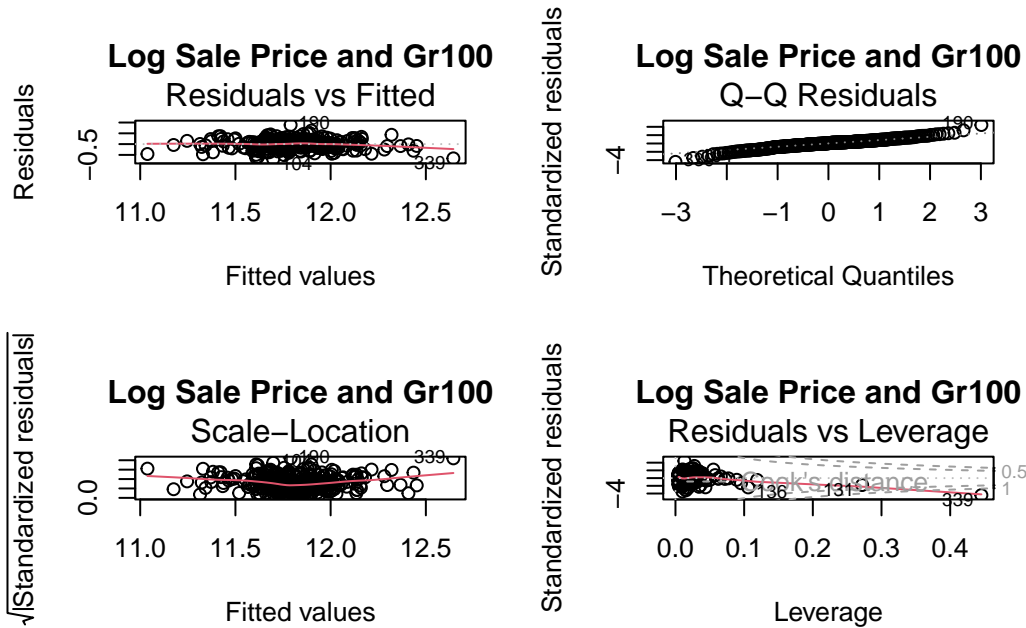
The model with the interaction terms using SalePrice on the original scale, a log-linear model using $\log(\text{SalePrice})$, and a $\log(\text{SalePrice})$ – $\log(\text{GrLivArea})$ model were all compared for diagnostics. Although the log–log model produced the highest adjusted R^2 , this is less suitable for the task at hand because this improvement came at the cost of interpretability. The coefficient on living area represented an elasticity rather than a dollar or percent based effect tied directly to 100 square foot increments. Since the client specifically prefers interpretation in units of 100 square feet, the log–log model was not selected.

The primary comparison was then between the interaction model and the log-linear model. Both models had statistically significant coefficients for neighborhood, living area, and the interaction terms. However, the diagnostic plots had significant improvement from the log-linear model.

Diagnostic plots will be reproduced to demonstrate how well each model satisfied the regression assumptions. With the interaction model, the residuals vs fitted plot displayed a clear funnel pattern, indicating strong heteroscedasticity or unequal variance as fitted values increased. The scale–location plot also showed this expansion of variance, and the Q–Q plot showed slight deviations from normality. Several observations displayed both high leverage and elevated Cook’s distances, suggesting possible influential data points. Diagnostic plots for the interaction model is shown below.



By contrast, with the log-linear model, the residuals vs fitted plot showed a more even scatter around zero. The scale-location plot was considerably flatter, indicating a more equal spread of variance across the range of fitted values. Normality also improved as the Q-Q plot was much closer to a straight line. There were still a view values displaying both high leverage and elevated Cook's distances, so these will be investigated in the next section. These improvements show that the log-linear model satisfies the assumptions of linearity, constant variance, and normality for regression assumptions and for these reasons this model will be used for the analysis. Diagnostic plots for the log-linear model is shown below.



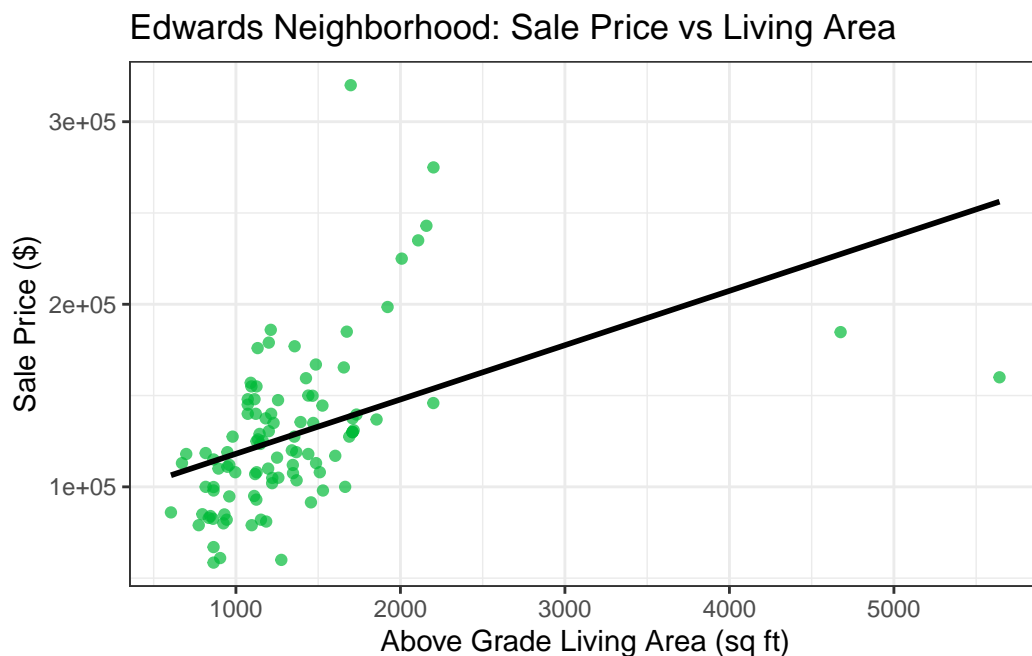
Influential Point Analysis

There is one home in the Edwards neighborhood that stands out in every diagnostic measure. This home has a Cook's distance of 2.61, a studentized residual of -4.53 , and a leverage value of 0.445 . With 5,642 square feet of living area, it is the largest home in Edwards. The next largest home is 4,676 square feet, and the majority of Edwards homes fall well below this range. Because this property is a legitimate part of the Edwards housing market, and because the goal of the analysis is to model all homes in the neighborhood, the point is retained despite its large influence on the fitted line.

All points with large studentized residual are within the Edwards neighborhood. These homes have unusual sale prices than what the model expects. There are two homes that sold much more than predicted, with 1698 sqft and studentized residual of 4.55 and 2201 sqft with studentized residual of 3.19. There were three homes that sold for less than what was expected by the model. One of those homes was the home discussed earlier with 5642 sqft and a studentized residual of -4.53 which sold for far less than expected given its size. Two other homes with 864 sqft (-3.21 studentized residual) and 1276 sqft (-3.53 studentized residual) also sold for less than what we expected. While these points are rare, they are not necessarily impracticable. These large residuals could have resulted from remodeling or a lack of remodeling that affected the prices, condition of the home, or even location within the neighborhood. None of the points present themselves to be an error or significant deviation from what is expected so they will all be retained in the final model.

All remaining observations with large studentized residuals also come from the Edwards neighborhood. These homes have sale prices that differ more from the model's predictions than expected. Two homes sold for noticeably more than predicted. One with 1,698 square feet (studentized residual 4.55) and one with 2,201 square feet (studentized residual 3.19). Three others sold for less than expected, including the 5,642 square foot home discussed above (-4.53), as well as homes with 864 square feet (-3.21) and 1,276 square feet (-3.53). While these residuals are large, they are not unreasonable. Differences in remodeling, home condition, or location within the neighborhood could easily lead to homes selling far above or below the general trend. None of these observations appear to be data errors or implausible homes, so all are retained in the final model.

A scatterplot restricted to the Edwards neighborhood (shown below) helps illustrate the distribution of living area and sale prices among these homes. The influential points stand out, but they fall within the realistic range of values for Edwards.



Several additional homes in the full dataset have high leverage values simply because they have unusually large or unusually small living areas. High leverage alone does not make an observation problematic. It only indicates that the point has the potential to affect the slope. For most of these homes, the studentized residuals are small and Cook's distances fall well below the standard cutoff. While the leverage is high, they still align well with their neighborhood's trend. Because the neighborhoods genuinely differ in typical home size, these extreme values naturally occur and represent meaningful parts of the housing market. As a result, these points were also retained.

Final Model and Interpretations

Before fitting the final model, living area was centered at its sample mean. Centering reduces the correlation between living area and the interaction terms which improves interpretability. Before centering, neighborhood coefficients corresponded to the expected price at 0 square feet, an unrealistic and extrapolated value that increased their statistical significance. After centering, neighborhood differences are evaluated at the average home size, which provides a more meaningful comparison and explains why the neighborhood difference for Edwards was no longer significant, while all slope and interaction terms retained their significance. With the new variable $Gr100.c$ equalling $Gr100 - \text{mean}(Gr100)$ and Brookside as the reference neighborhood ($Edwards = 0$, $NAmes = 0$), the final model is the following:

$$\begin{aligned} \log(SalePrice) = & 11.752637 - 0.048604(Edwards) + 0.112660(NAmes) \\ & + 0.073822(Gr100.c) - 0.052153(Edwards \times Gr100.c) - 0.041410(NAmes \times Gr100.c) \end{aligned}$$

Brookside (reference group)

For homes in Brookside, each additional 100 sq ft above the average increases the median sale price by about 7.7%. The median sale price for an average sized Brookside home is approximately \$120,452 to \$134,090.

$$\log(SalePrice) = 11.7526 + 0.07382(Gr100.c)$$

Edwards

For an average sized home in Edwards, the median sale price is about 4.7% lower than Brookside, but this difference is not statistically significant because the interval includes 0. Each additional 100 sq ft increases the median price by about 2.2%, significantly less than Brookside's 7.7% due to the negative interaction term.

$$\log(SalePrice) = 11.7040 + 0.02167(Gr100.c)$$

North Ames

For an average sized home in North Ames, the median sale price is about 11.9% higher than Brookside, a statistically significant difference. Each additional 100 sq ft increases the median price by about 3.3%, also smaller than Brookside's slope, reflecting the negative interaction term.

$$\log(SalePrice) = 11.8653 + 0.03241(Gr100.c)$$

Overall, extra living area increases sale price in all neighborhoods, but the strength of that increase differs meaningfully across neighborhoods. The interactions highlight that Edwards and North Ames have flatter slopes, meaning additional square footage adds less value there compared to Brookside.

Conclusion

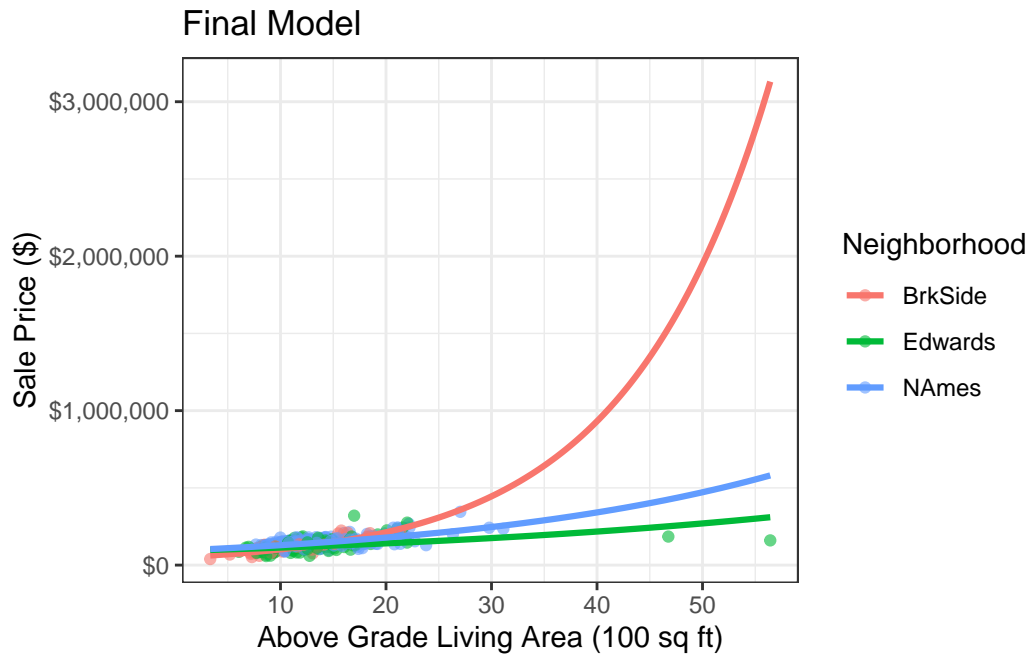
The centered log-linear interaction model shows that living area is a strong and statistically significant predictor of sale price across all three neighborhoods. However, the effect of added square footage differs considerably by location. Brookside has the steepest slope, meaning buyers in this neighborhood place the highest value on extra living space. Both Edwards and North Ames have flatter slopes, indicating that prices raise more moderately with additional square footage compared to Brookside. These differences in slope do not reflect whether a neighborhood is more or less expensive overall but reflect how each neighborhood's market responds to increased home size. A possible limitation of the model is that Brookside does not have extremely large homes in the dataset, unlike Edwards, which means the slope estimate for Brookside may be less precise for larger homes. Despite this, the model accurately reflects the general patterns of home prices in the neighborhood.

Because living area was centered, the intercepts correspond to the expected price of an average sized home in each neighborhood, and the slopes measure the percent change in median sale price for each additional 100 square feet. The model satisfies key regression assumptions, and provides interpretable estimates.

Overall, the results show that all three neighborhoods value increased living area, but Brookside values additional space most strongly. Edwards and North Ames have more moderate price growth as homes increase in size. Incorporating the interaction between neighborhood and living area is important for accurately describing the housing market and for providing pricing guidance.

Take Home Message

Home prices rise as houses get larger in all three neighborhoods, but the amount of value added by extra space isn't the same everywhere. In Brookside, buyers place the highest premium on additional square footage meaning bigger homes tend to see the strongest jump in price. In Edwards and North Ames, the starting prices of homes differ from Brookside, but extra space adds value more gradually. This means the same increase in size can lead to very different increases in price depending on the neighborhood. When advising clients, setting asking prices, or comparing homes, it is important to consider not just how big a house is, but where it is located. Neighborhood differences shape how much extra space is worth in the local market, and this model gives a clear, data-based view of those patterns. A scatterplot with the final model is produced below.



Appendix

```

1 # load libraries & read in data
2 library(tidyverse)
3 library(janitor)
4 library(naniar)
5 library(ggplot2)
6 library(dplyr)
7 library(car)
8 library(scales)
9 library(MASS)
10
11 train = read.csv("/Users/katelyn/Desktop/Projects/House Price Project/data/train.csv")
12 glimpse(train)

```

```

1 # subset data
2 data = train %>%
3   select(Neighborhood, GrLivArea, SalePrice) %>%
4   filter(Neighborhood %in% c("BrkSide", "Edwards", "NAmes")) %>%
5   mutate(Neighborhood = factor(Neighborhood),
6          Gr100 = GrLivArea / 100)

```

```

1 # base model
2 model.1 = lm(SalePrice ~ Neighborhood + Gr100, data = data)
3 summary(model.1)
4 confint(model.1)
5
6 # interaction term
7 model.2 = lm(SalePrice ~ Neighborhood * Gr100, data = data)
8 summary(model.2)
9 confint(model.2)
10
11 AIC(model.1,model.2) # model.2 has the lower aic

```

```

1 # log y
2 model.logy = lm(log(SalePrice) ~ Neighborhood * Gr100, data = data)
3 summary(model.logy)
4 par(mfrow=c(2,2))
5 plot(model.logy, , main = "Log Sale Price")
6
7 # log both
8 model.logboth = lm(log(SalePrice) ~ Neighborhood * log(Gr100*100), data = data)
9 summary(model.logboth)
10 par(mfrow=c(2,2))
11 plot(model.logboth , main = "Log Sale Price and Gr100")

```

```

1 # diagnostics of best two models
2 par(mfrow = c(2,2))
3 plot(model.2, main = "Model with Interaction term")
4
5 par(mfrow = c(2,2))
6 plot(model.logy, main = "Log Sale Price and Gr100")

```

```

1 # studentized residuals, leverage, cook's D
2 diagnostics = data %>%
3   mutate(
4     row.id = row_number(),
5     student.resid = rstudent(model.logy),
6     leverage = hatvalues(model.logy),
7     cooksD = cooks.distance(model.logy) )
8
9 # cutoffs
10 n = nrow(data)
11 p = length(coef(model.logy))

```

```

12 lev.cutoff = (2 * p) / n
13 cooks.cutoff = 4 / (n - p)
14
15 # diagnostics table
16 diagnostics = diagnostics %>%
17   mutate(
18     high.leverage = leverage > lev.cutoff,
19     large.student = abs(student.resid) > 3,
20     large.cooks = cooksD > cooks.cutoff,
21     all = high.leverage | large.student | large.cooks )
22
23 diagnostics %>%
24   filter(high.leverage) %>%
25   arrange(desc(high.leverage)) %>%
26   dplyr::select(1:8)

```

```

1 # vif and centering
2 library(car)
3 vif(model.logy)
4 vif(model.logy, type = "predictor")
5
6 # centering
7 data = data %>% mutate(Gr100.c = Gr100 - mean(Gr100))
8
9 model.logy.c = lm(log(SalePrice) ~ Neighborhood * Gr100.c, data = data)
10 vif(model.logy.c)
11
12 summary(model.logy.c)
13 confint(model.logy.c)
14
15 diagnostics2 = data %>%
16   mutate(
17     row.id = row_number(),
18     student.resid = rstudent(model.logy.c),
19     leverage = hatvalues(model.logy.c),
20     cooksD = cooks.distance(model.logy.c) )
21
22 # cutoffs
23 n = nrow(data)
24 p = length(coef(model.logy.c))
25 lev.cutoff = (2 * p) / n
26 cooks.cutoff = 4 / (n - p)

```

```

27
28 # influential points didn't change
29 diagnostics2 = diagnostics2 %>%
30   mutate(
31     high.leverage = leverage > lev.cutoff,
32     large.student = abs(student.resid) > 3,
33     large.cooks = cooksD > cooks.cutoff,
34     all = high.leverage | large.student | large.cooks )
35
36 diagnostics2 %>%
37   filter(high.leverage) %>%
38   arrange(desc(high.leverage)) %>%
39   dplyr::select(1:8)
40
41
42 # vif after centering
43 vif(model.logy.c)

```

```

1 # looking at ranges
2 range(data$Gr100[data$Neighborhood == "BrkSide"])
3 range(data$Gr100[data$Neighborhood == "Edwards"])
4 range(data$Gr100[data$Neighborhood == "NAmes"])
5
6 range(data$SalePrice[data$Neighborhood == "BrkSide"])
7 range(data$SalePrice[data$Neighborhood == "Edwards"])
8 range(data$SalePrice[data$Neighborhood == "NAmes"])

```