

Data Models and Linear Regression

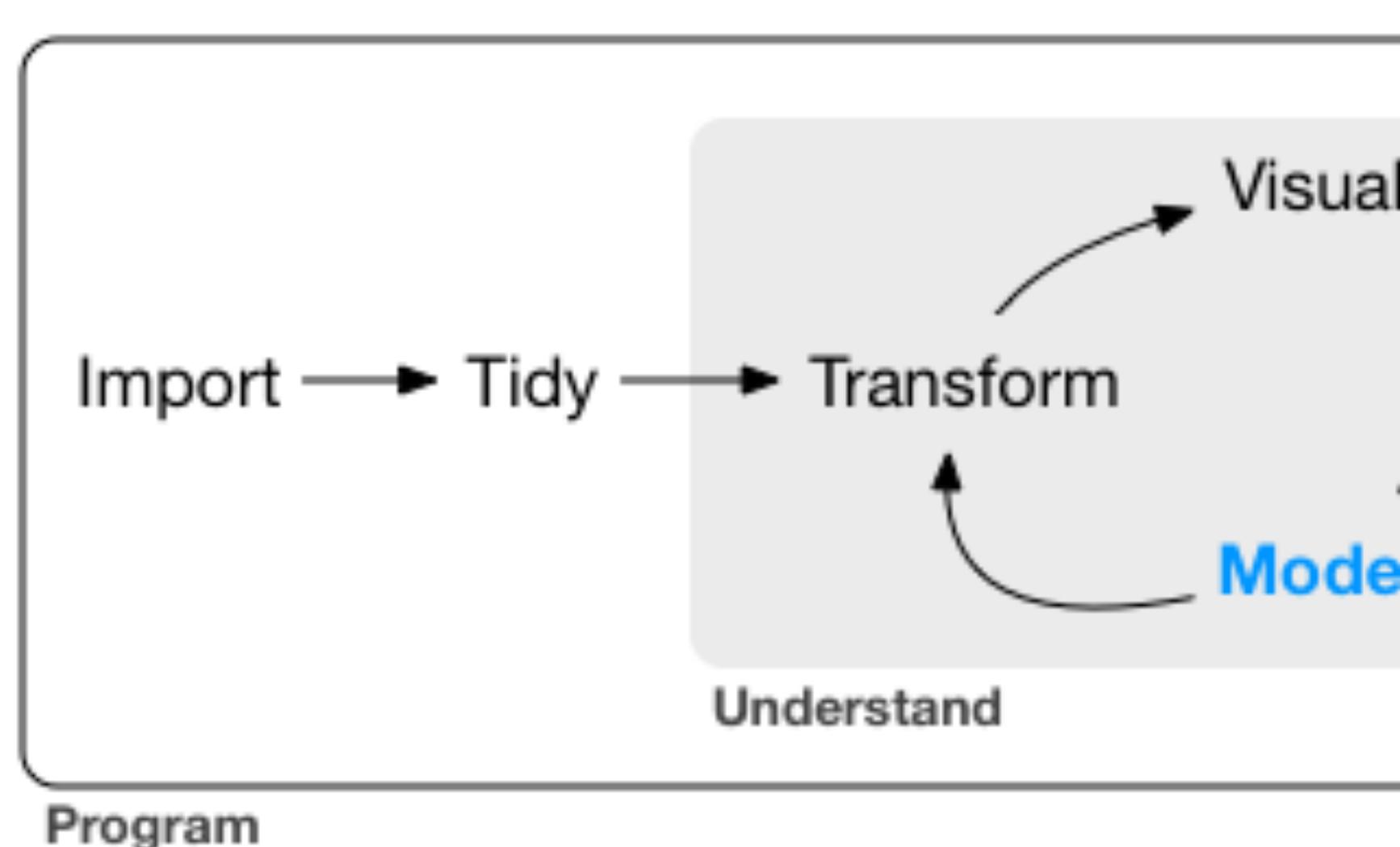
Kylie A. Bem

Northeastern University
Khoury College of Computer Sciences



Northeastern University

Modeling



R for Data Science, Wickham an

Learning goals

- What are the goals of regression?
- Why linear regression?
- Fitting linear models
- Model diagnostics

DATA MOD

Modeling goals

- Predict new values
 - ◆ Regression — continuous values
 - ◆ Classification — categorical values
- Discover hidden structures
 - ◆ Cluster into groups based on similarities
- Interpret and understand
 - ◆ Statistical inference

Regression

- Given predictors $x_1, x_2, x_3 \dots$
- Predict continuous-valued response variable y
- Supervised: requires known response values

Classification

- Given predictors $x_1, x_2, x_3 \dots$
- Predict categorical-valued response variable y
- Supervised: requires known category labels

Clustering

- Given features $x_1, x_2, x_3 \dots$
- Group similar data points together into clusters
- Unsupervised: expects unlabeled data

Statistical inference

- Given sample data $x_1, x_2, x_3 \dots$
- Infer properties of the underlying population
- Confidence intervals
- Hypothesis tests

Machine learning &

Statistics and machine learning share many methods, and some models (but not all) can be used for both.

- Supervised learning
 - ◆ Regression
 - ◆ Classification
- Unsupervised learning
 - ◆ Clustering
 - ◆ Dimension reduction
- Statistical inference

“All models are wrong, but some are useful.”
— George E. P. Box

- No mathematical model is a perfect representation of reality
- But some models are more illuminating than others
- Choice of model should consider:
 - ◆ What are the overall goals of the model?
 - ◆ How useful is the model for the purpose?

LINEAR REGR

Linear regression

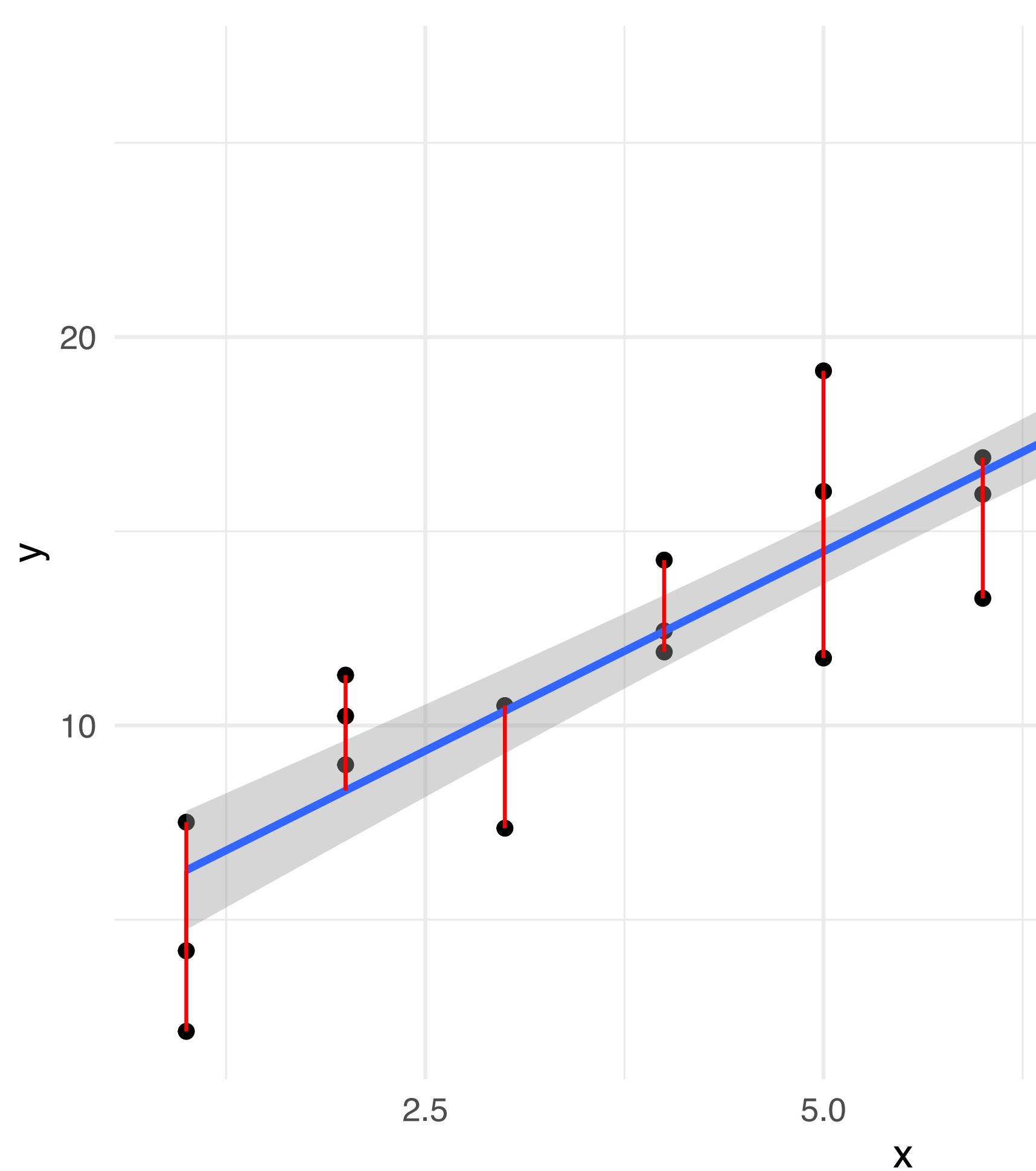
- Predicts a scalar continuous variable (dependent variable)
- Linear function of one or more variables (explanatory/independent variables)
- Flexible model with many applications
 - ◆ Linear constraint is rarely a problem
 - ◆ Useful for statistical inference: Goodness-of-fit tests, confidence intervals

Linear regression

- Parametric model
 - ◆ Create a low-dimensional summary
 - ◆ Forms a constraint on possible
- Linear model form:
$$y = \beta_0 + \beta_1 x + \dots + \beta_p x_p + \epsilon$$

error $\sim N(0, \sigma^2)$
- Fit the model through data
 - ◆ Maximize log-likelihood (statistic)
 - ◆ Minimize squared errors (loss function)

Fitting the regression



Find parameters (β 's) that minimize

Assumptions of linear regression

- Linearity
 - ◆ The relationship must be linear
 - ◆ Transformations can help make it linear
- Constant variance (homoscedasticity)
 - ◆ Especially important for data collection
- Independent errors
 - ◆ Especially important for data collection
- Lack of perfect multicollinearity
 - ◆ Predictors can't be (too) correlated

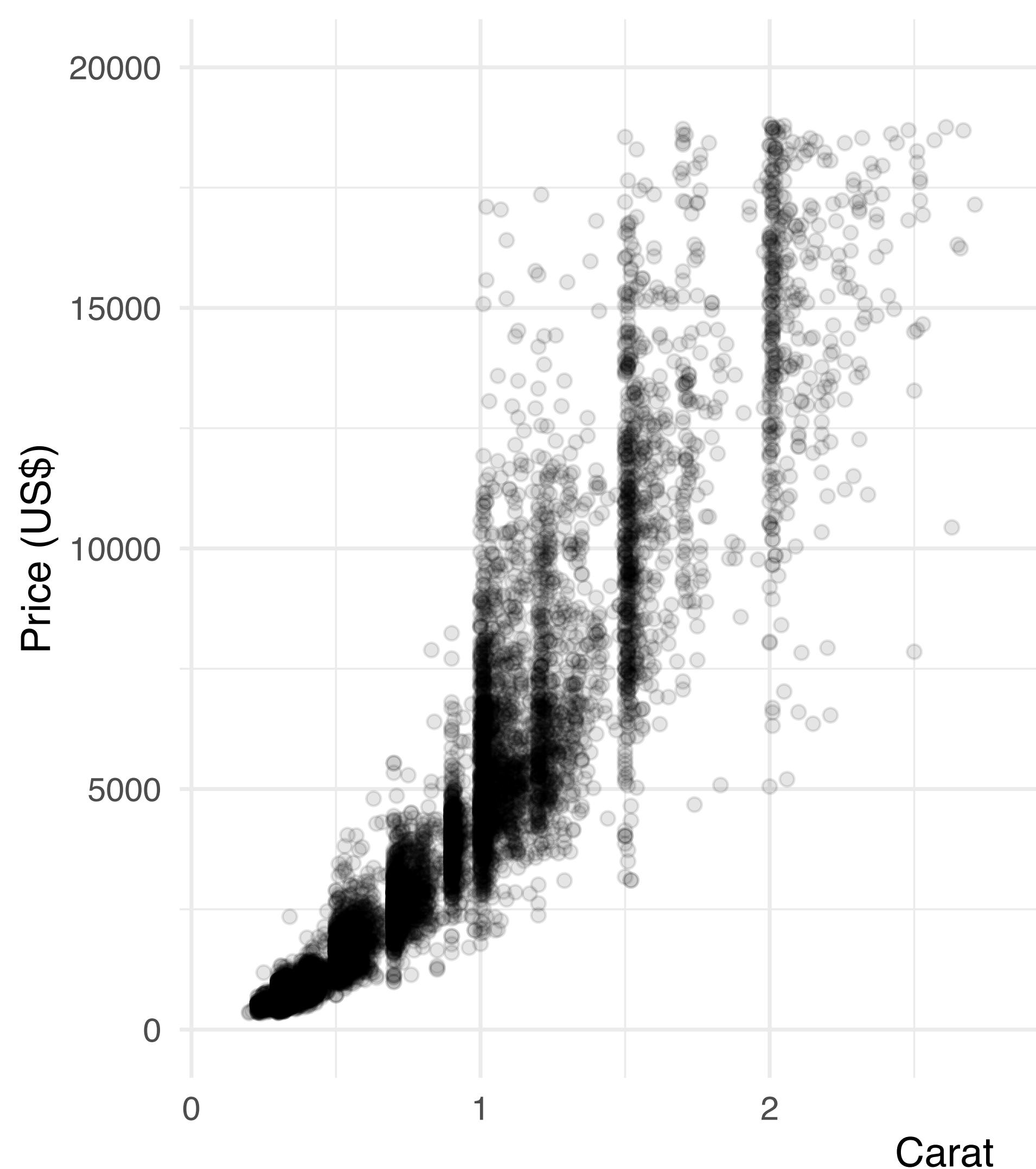
Modeling real

Prices and measurements of diamonds

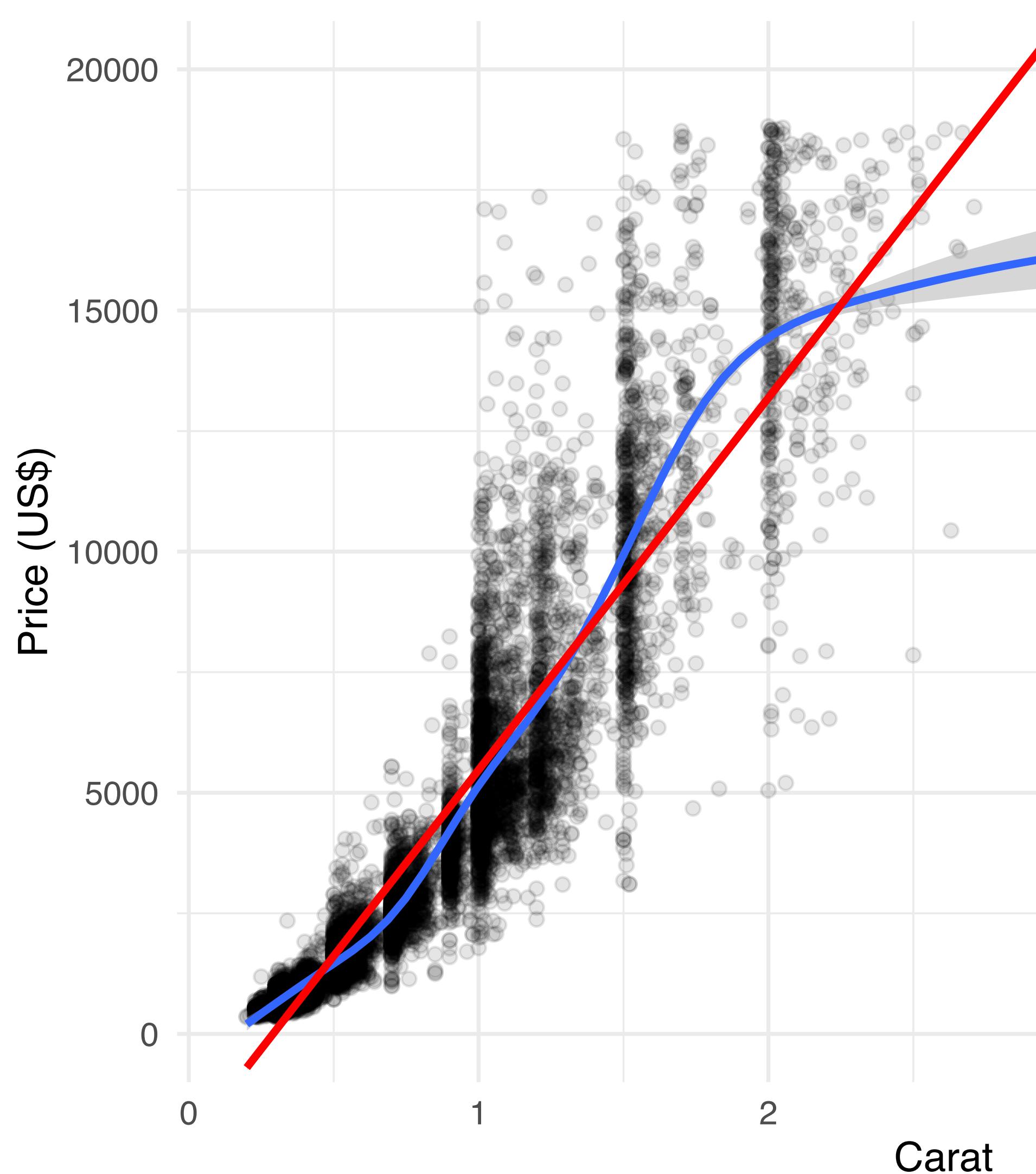
```
## #     tibble: 53,940 x 10
##       carat   cut      color clarity depth t
##       <dbl> <ord>    <ord> <ord>   <dbl> <
## 1 0.23   Ideal     E      SI2     61.5
## 2 0.21   Premium   E      SI1     59.8
## 3 0.23   Good      E      VS1     56.9
## 4 0.290  Premium   I      VS2     62.4
## 5 0.31   Good      J      SI2     63.3
## 6 0.24   Very Good J      VVS2    62.8
## 7 0.24   Very Good I      VVS1    62.3
## 8 0.26   Very Good H      SI1     61.9
## 9 0.22   Fair       E      VS2     65.1
## 10 0.23  Very Good H      VS1     59.4
## # ... with 53,930 more rows
```

Available in ggplot2 p

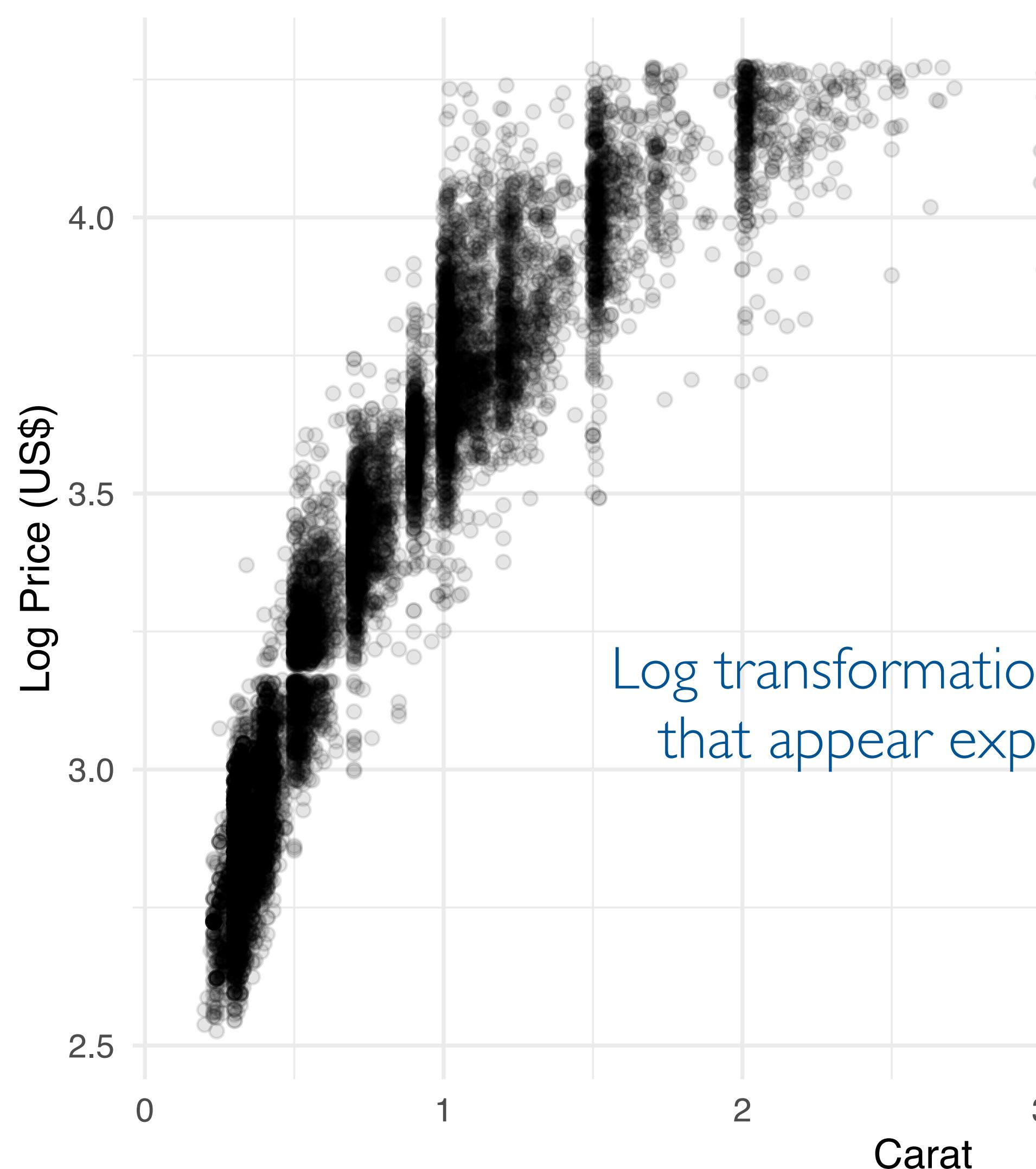
Does diamond size effect price?



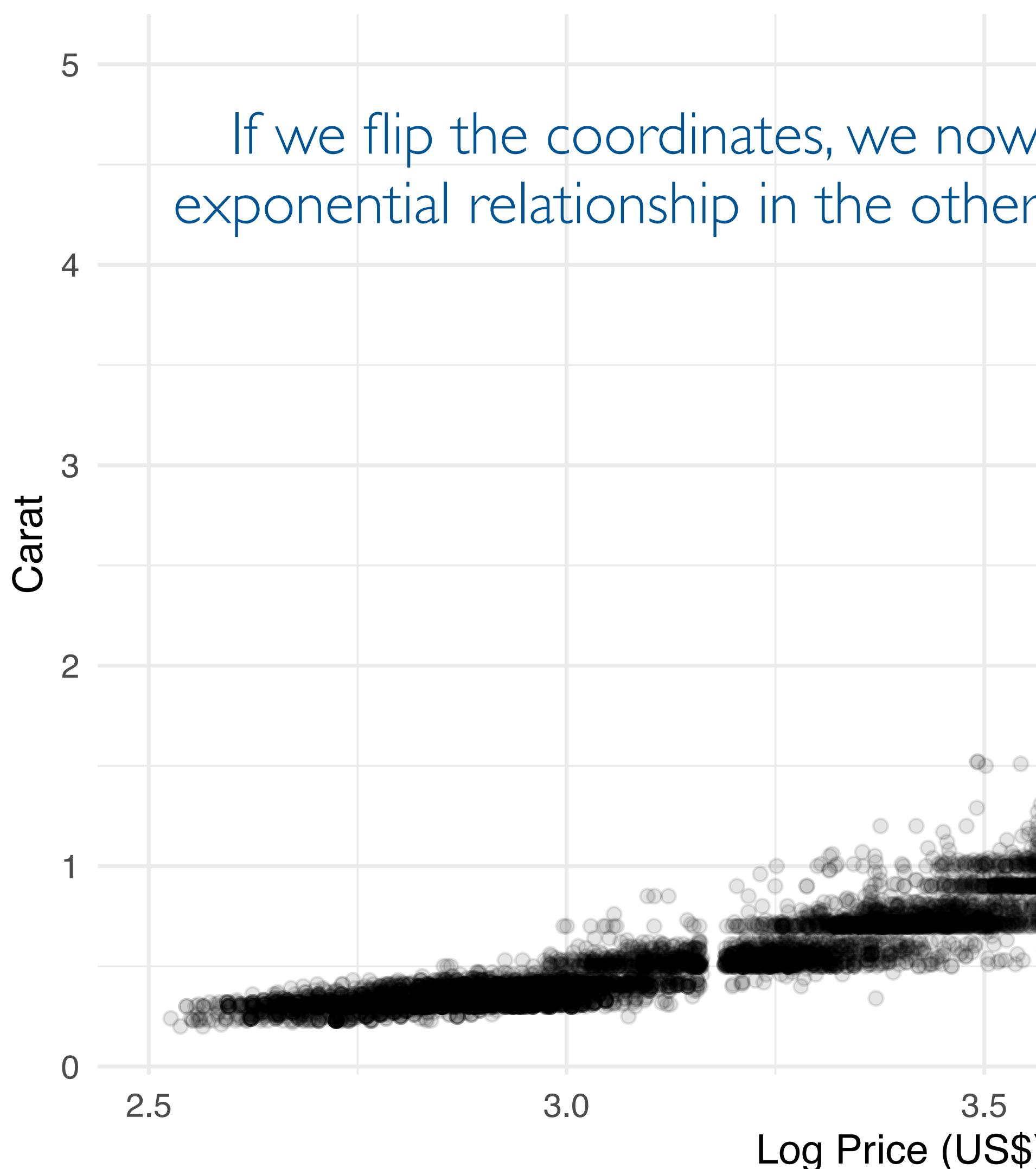
Is it linear



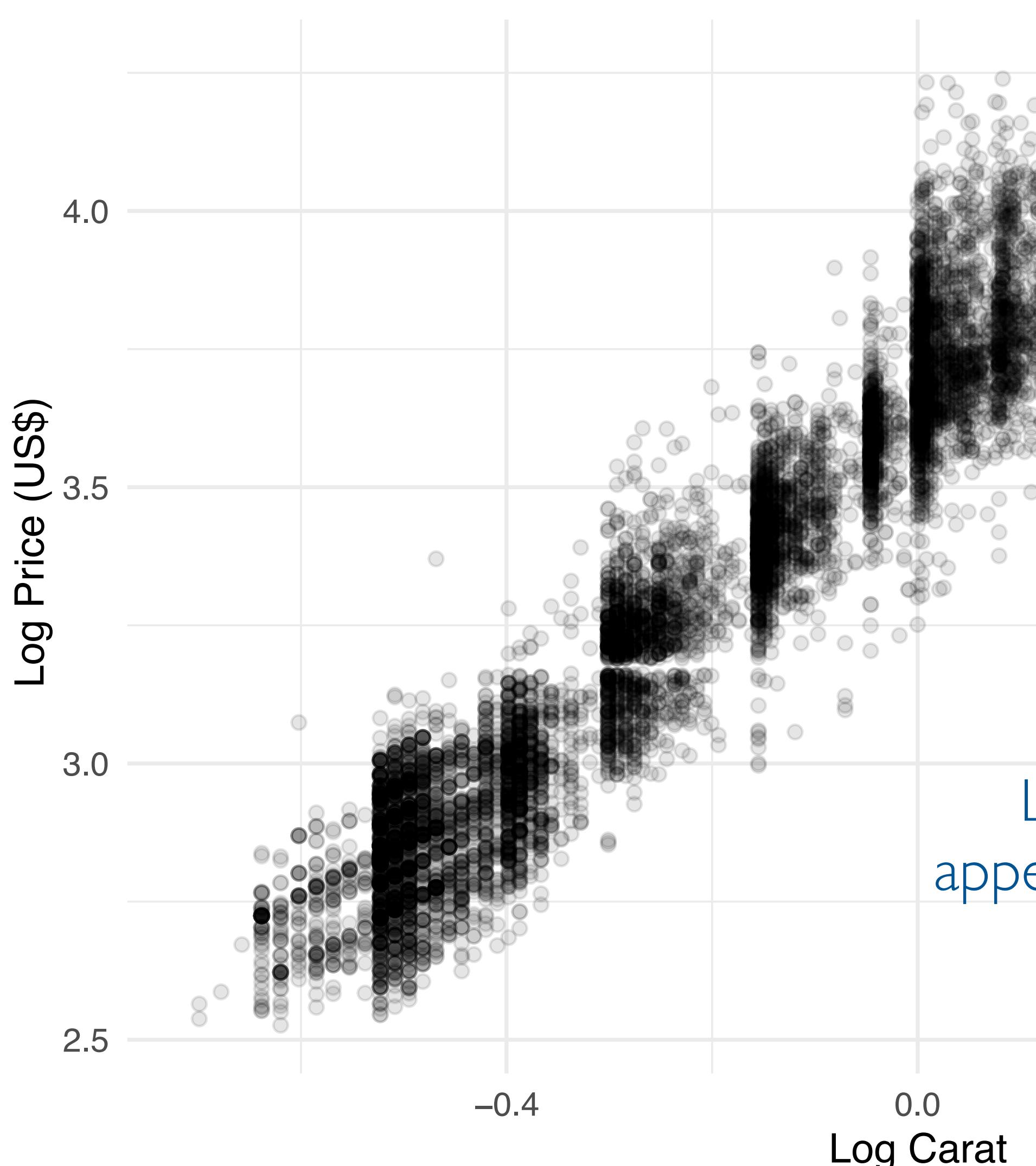
Log-transform



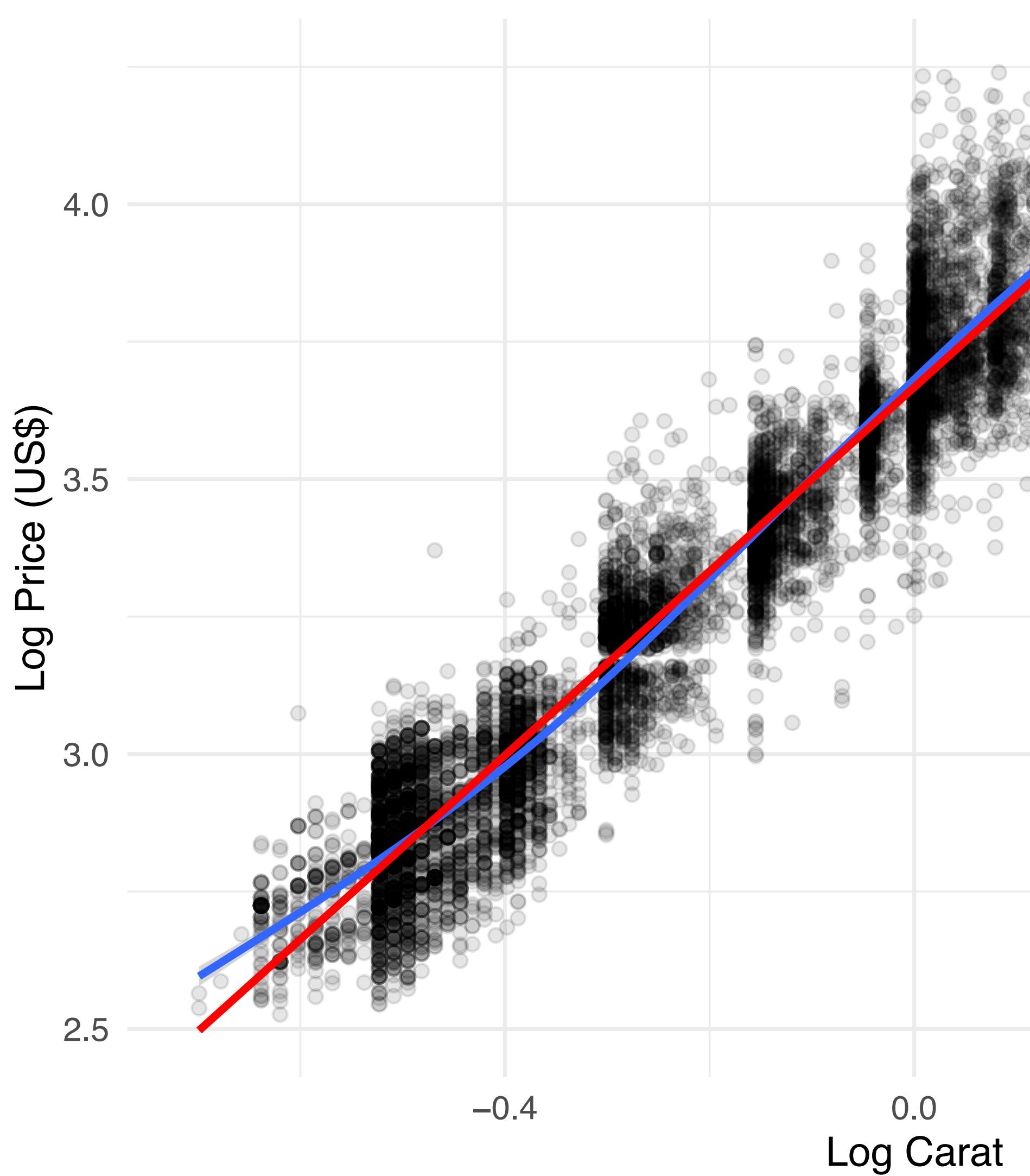
A different pers



Log-log transformation



Is it linear now?



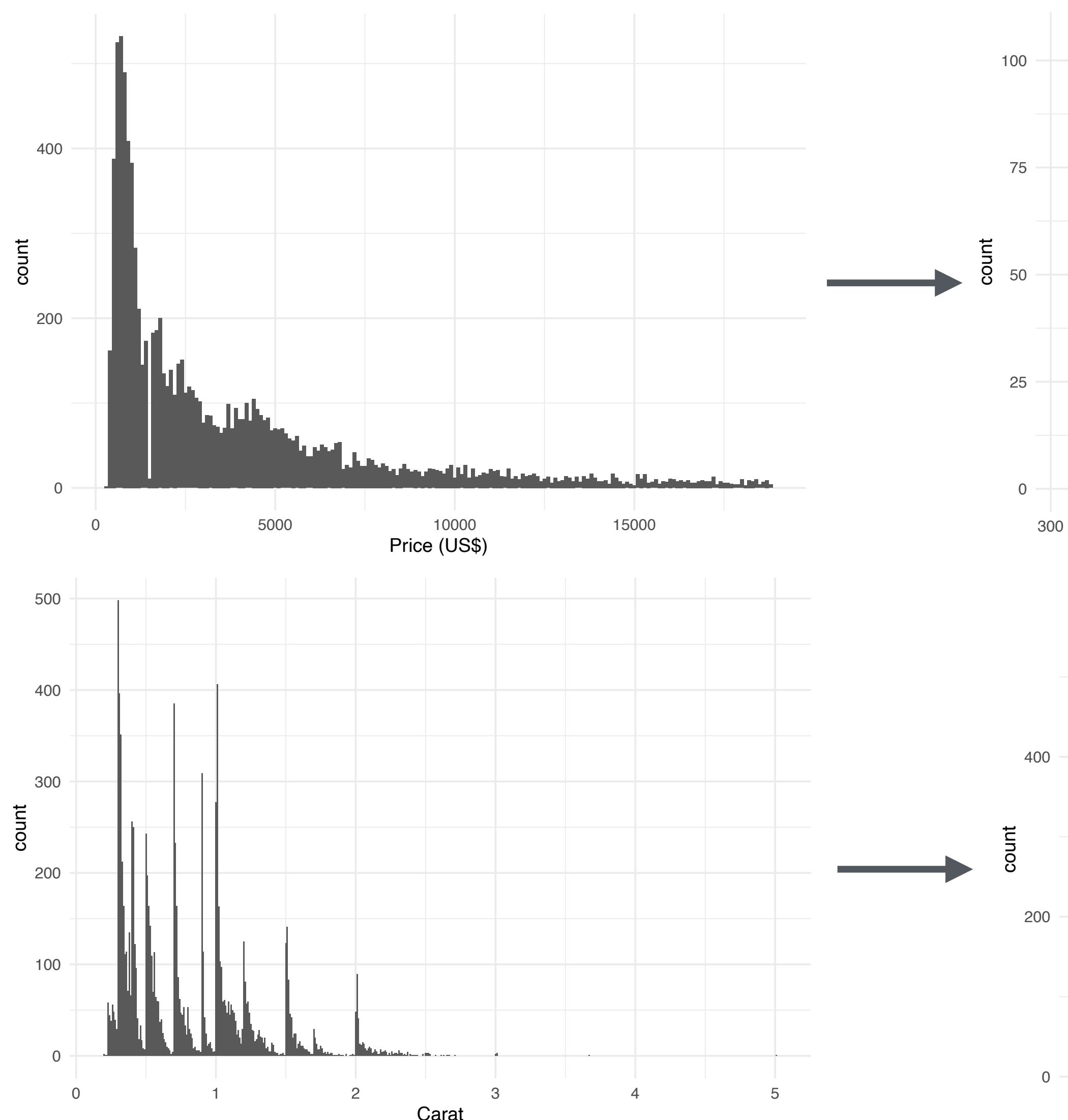
A note on log transformations

- Log transformations are useful for:
 - ◆ Right-skewed data → more symmetric
 - ◆ Exponential relationships → more linear
 - ◆ Log-normal distributions → more normal
- Has effect of “subduing” extreme values
- Good for visualization –
 - Makes linear models more appropriate

Log transformat

- Square root has similar
- Log transformations have:
 - ◆ Monotonic increasing \rightarrow preserve order
 - ◆ Constant ratios \rightarrow constant differences
 - ◆ Easy to differentiate \rightarrow useful for regression
- Transformations for left-skewed data:
 - ◆ Power functions
 - ◆ Exponentiation

Log transformat



Fitting the model: pr

```
fit1 <- lm(log10(price) ~ log10(carat),  
summ ry(fit1)  
  
##  
## Call:  
## lm(formula = log10(price) ~ log10(carat))  
##  
## Residuals:  
##      Min       1Q   Median       3Q  
## -0.58268 -0.07321 -0.00186  0.07172  
##  
## Coefficients:  
##                 Estimate Std. Error t value  
## (Intercept) 3.667107  0.001393  2663.  
## log10(carat) 1.673639  0.004536  370.5  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01  
##  
## Residual standard error: 0.115 on 9998 degrees of freedom  
## Multiple R-squared:  0.9316, Adjusted R-squared:  0.9316  
## F-statistic: 1.361e+05 on 1 and 9998 DF, p-value: < 2.2e-16
```

Interpretation

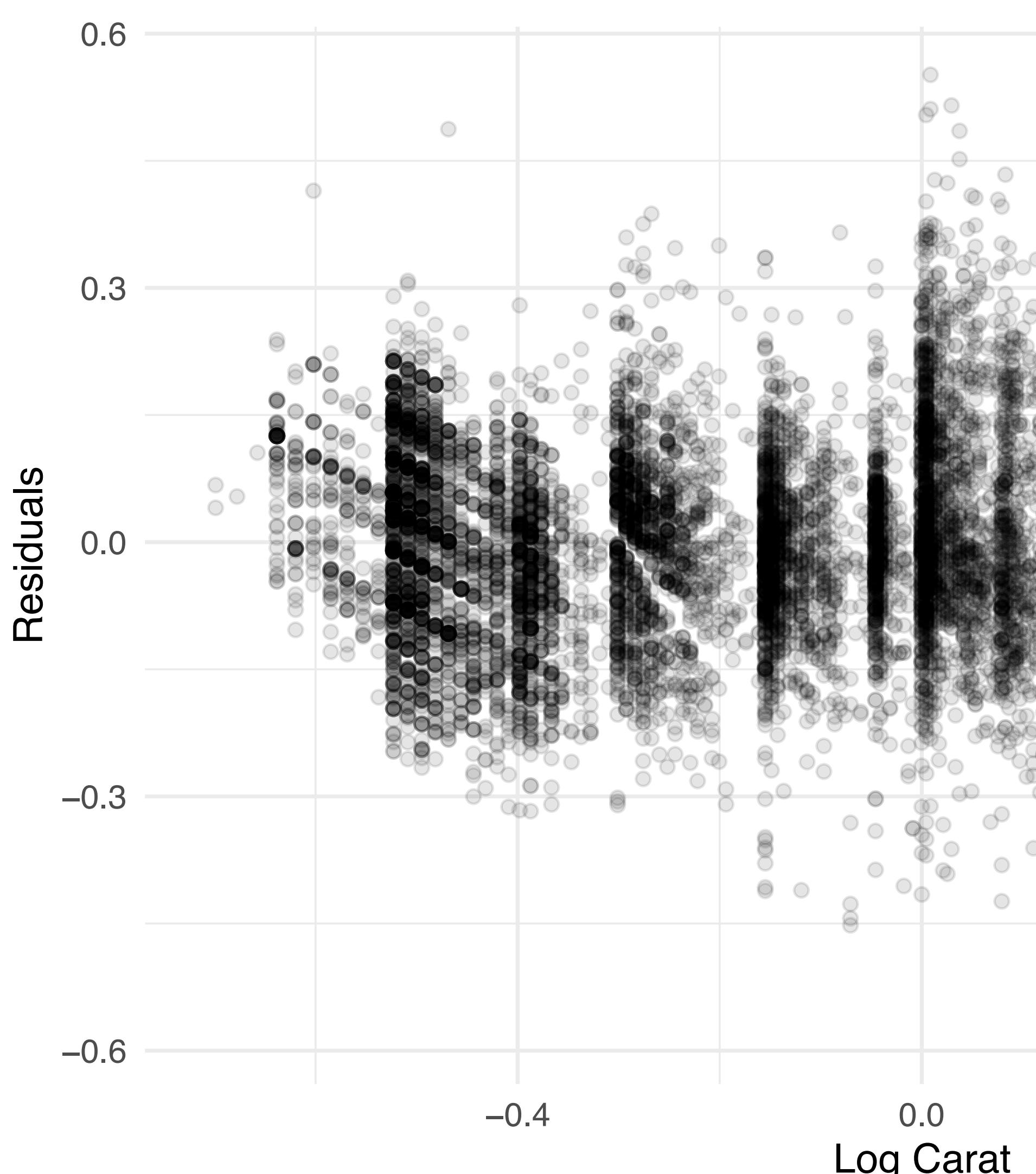
- Model coefficients describe the relationship between predictor and outcome
- Assume all other predictors are held constant
- For a one-unit change in a predictor, the resulting change in \mathbf{y} is the coefficient value

MODEL DIAGN

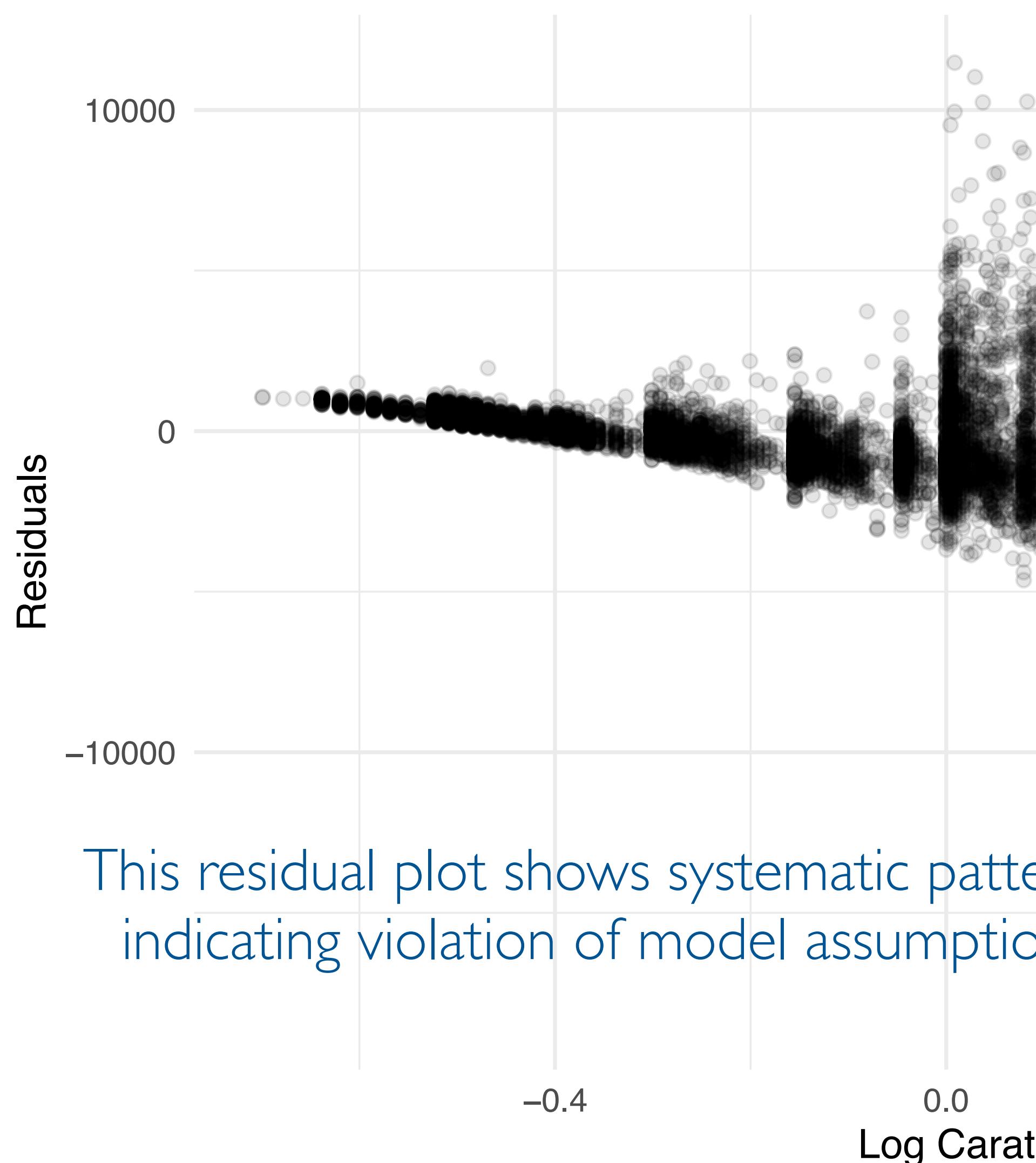
Model diagnosis

- Visualize the residuals (residual plots)
- Useful diagnostic step for detecting violations of model assumptions
- Can be used to identify relationships to model variables

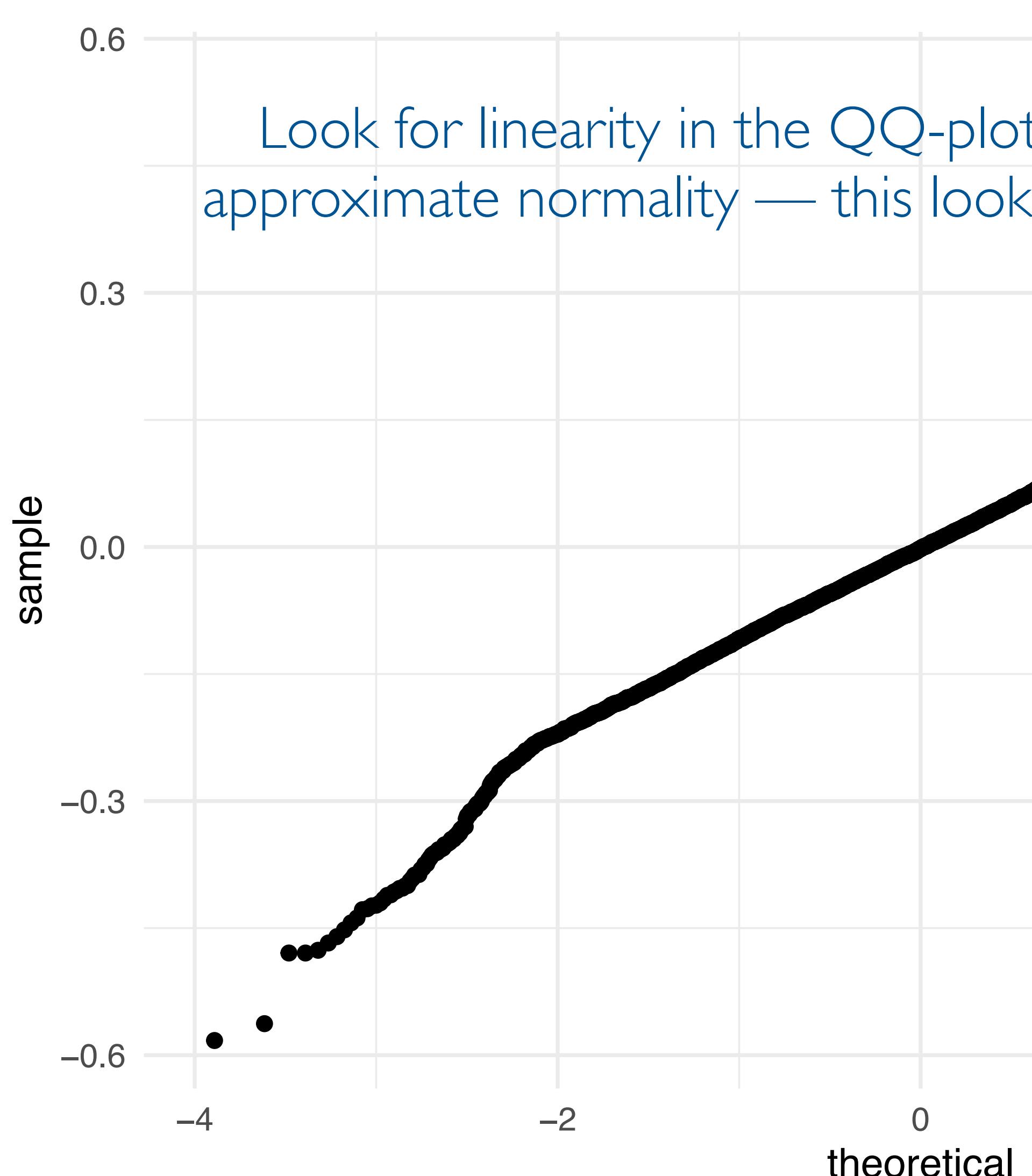
Plotting the residuals



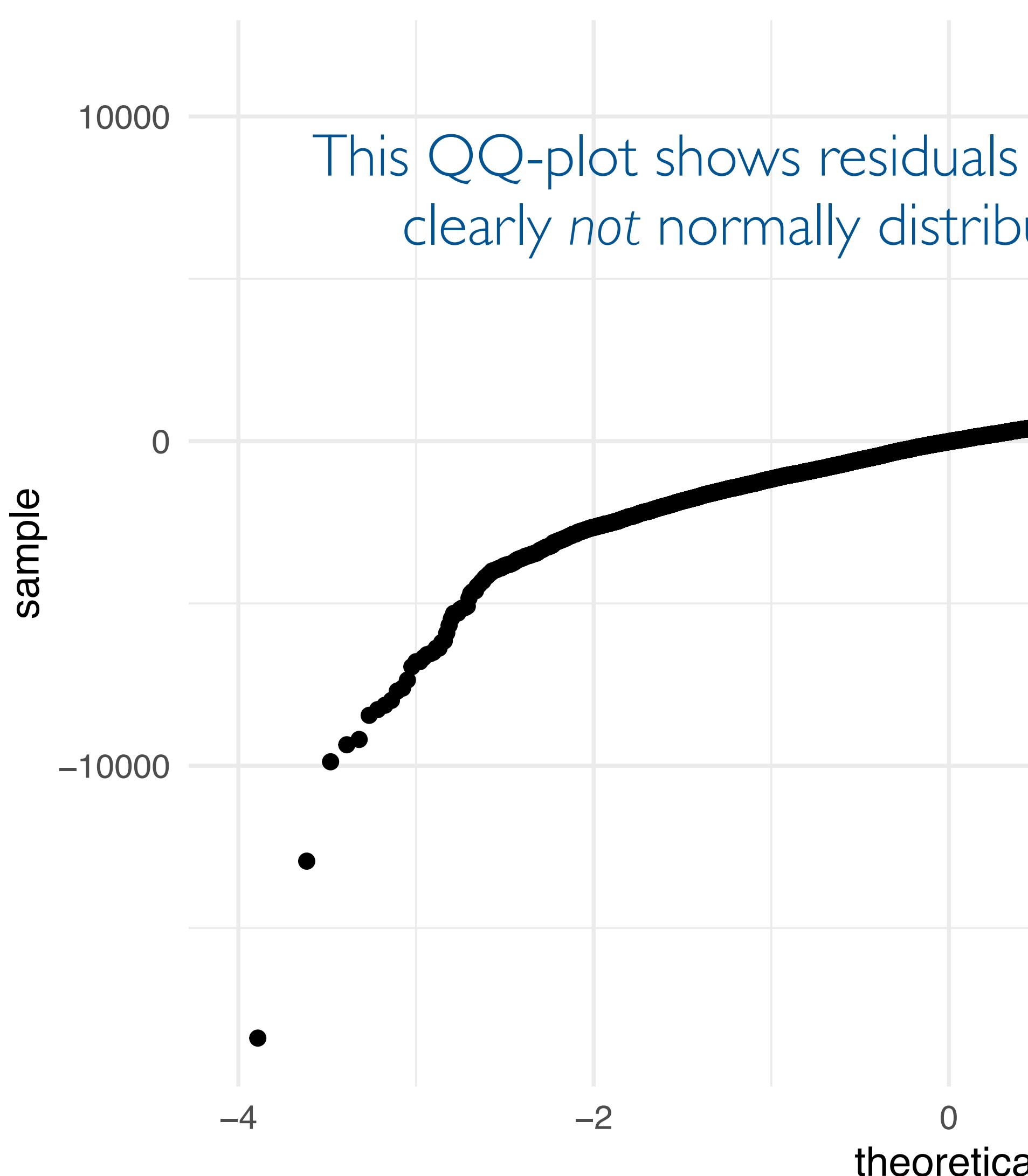
Bad residuals



Normality of re



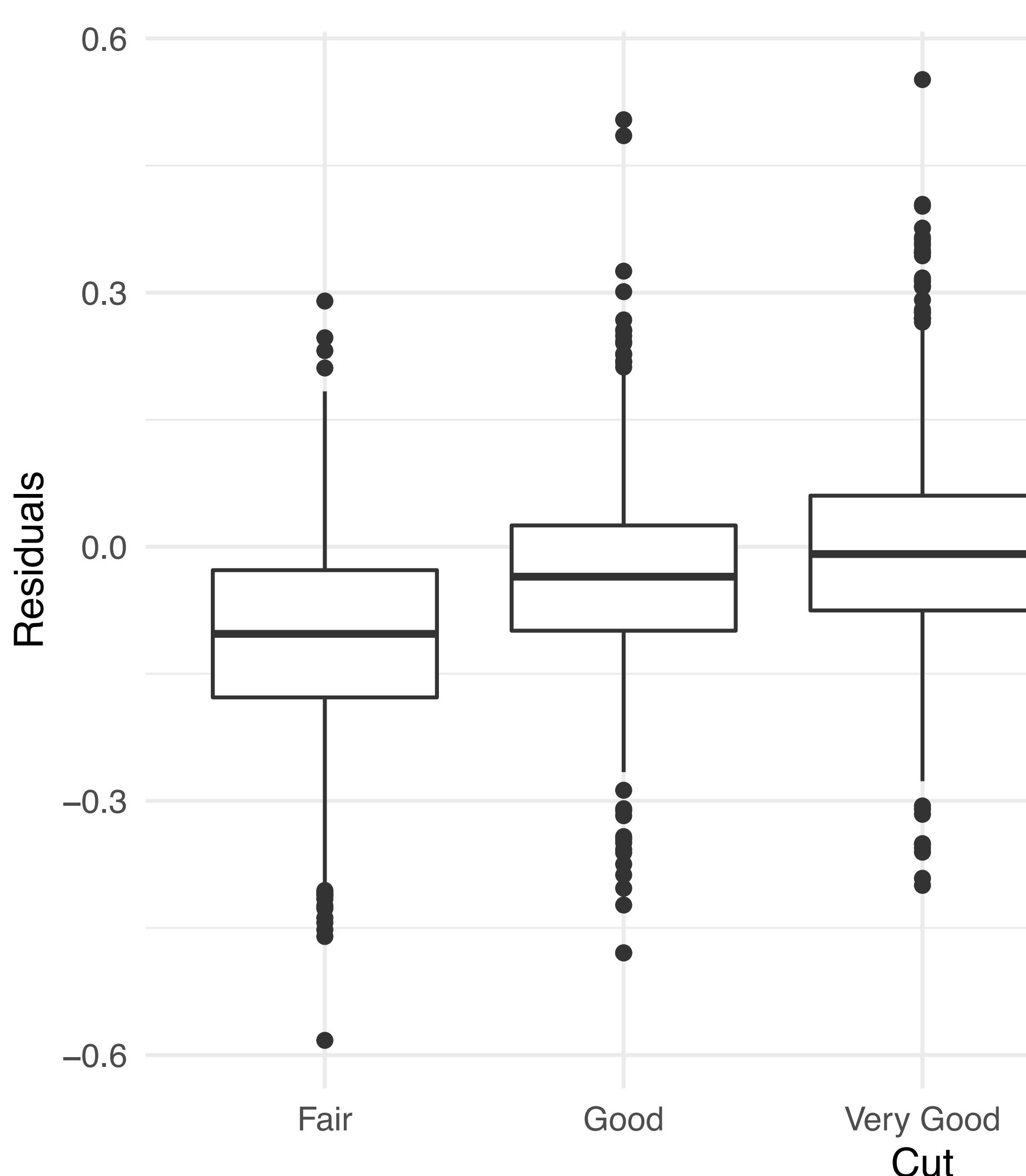
Non-normal residuals



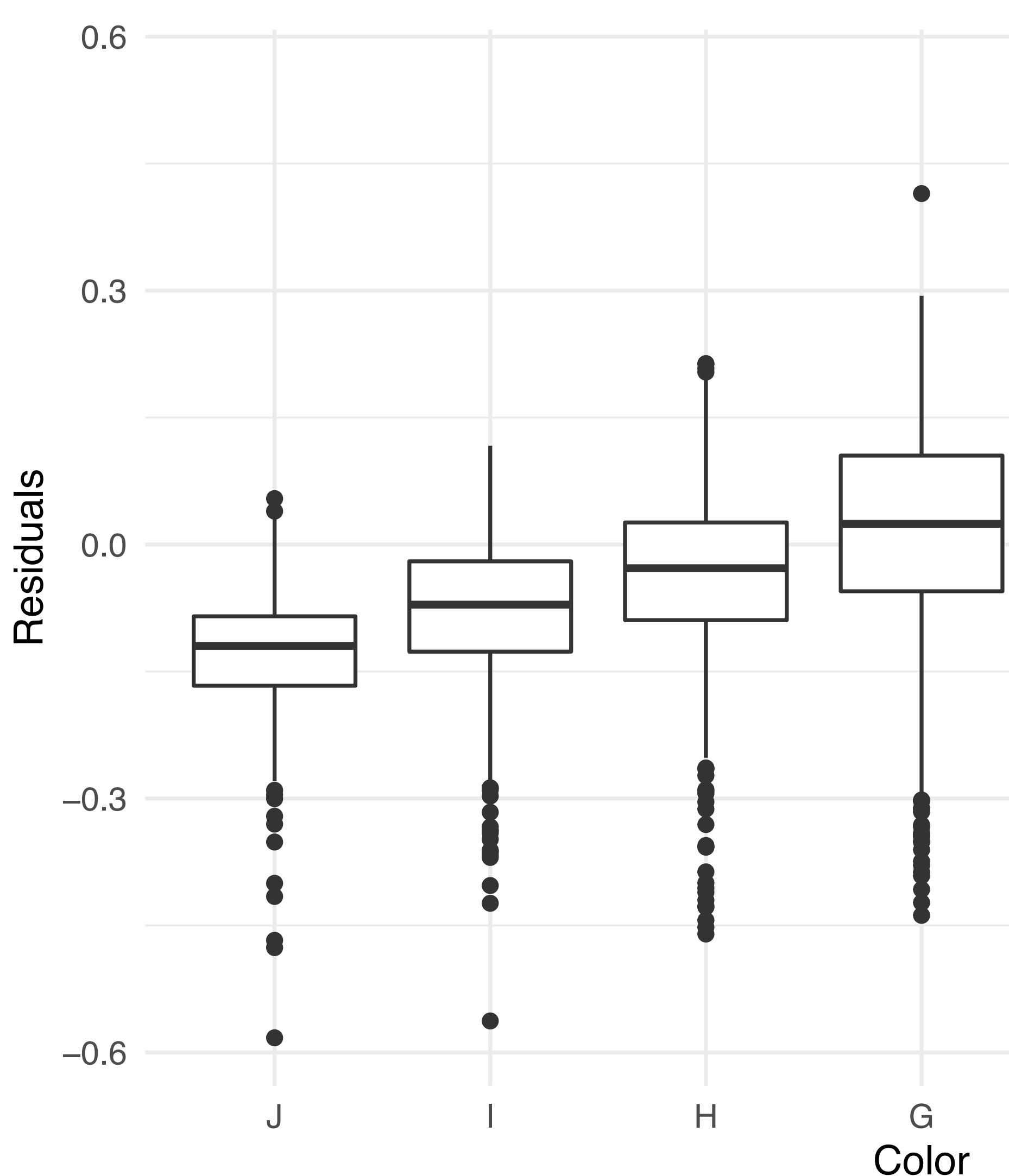
Checking other variables

- The residuals can act as the “leftover variable” for the response variable
- Modeling “removes” variance from the response variable that is not explained by the predictors in the model
- Residuals contain the “leftover” variance from the model
- Look for relationship between residuals and unused variables

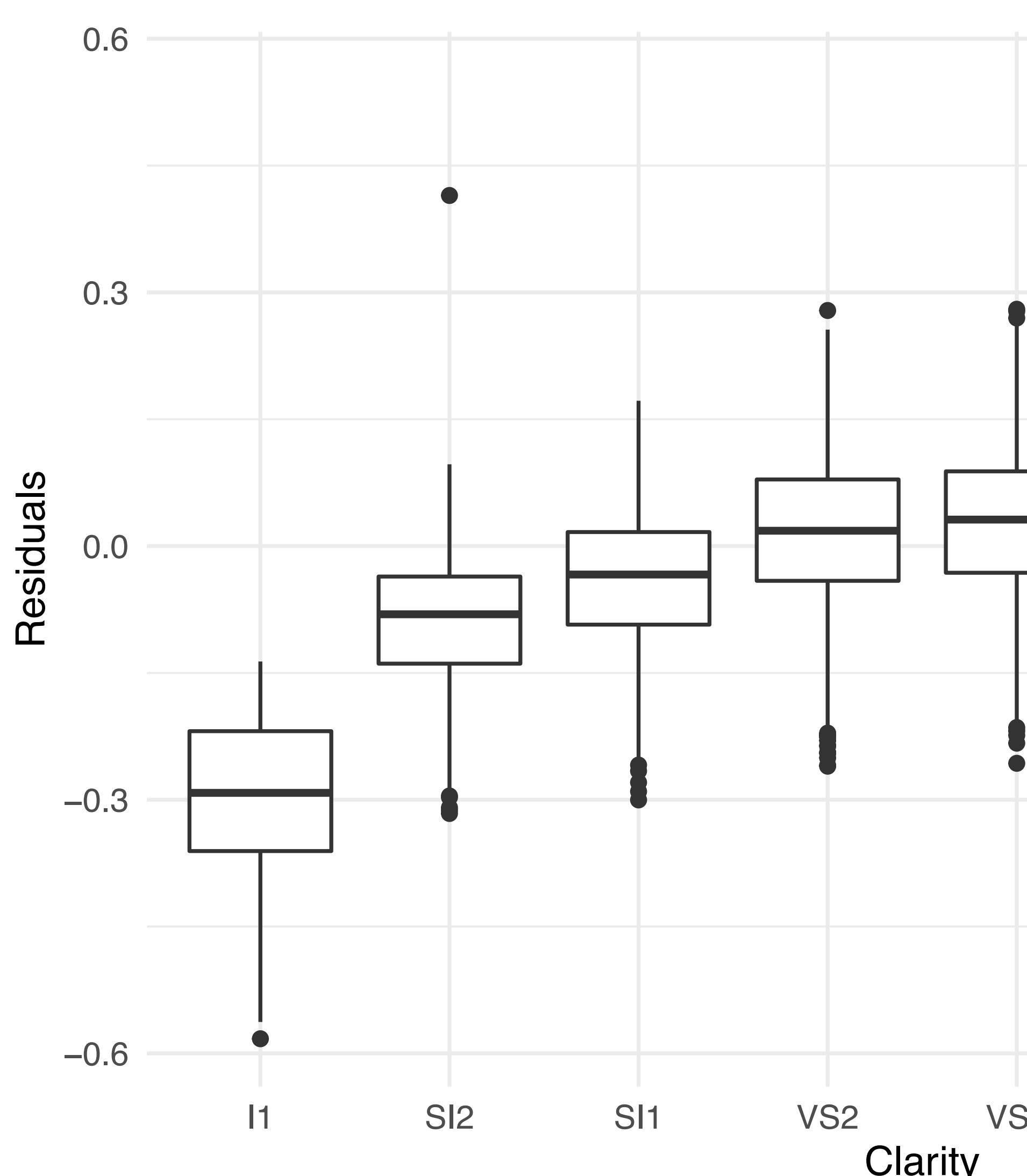
Better cuts → high



Better color → high



Better clarity → hi



Modeling with categorical variables

- Categorical variables are modeled using indicator variables
- Each indicator variable is 0 or 1
- For a categorical variable with m categories, model it with **$m-1$** indicator variables
- Otherwise, model is overfitted

A better model for diamonds

```
fit2 <- lm(log10(price) ~ log10(carat) + cut + color + clarity)
summary(fit2)

##
## Call:
## lm(formula = log10(price) ~ log10(carat) + cut + color + clarity,
##      data = diamonds10k)
##
## Residuals:
##    Min     1Q Median     3Q    Max 
## -0.30130 -0.03785  0.00030  0.03564  0.61635 
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 3.406711  0.005657 602.249 <2e-16 ***
## log10(carat) 1.885043  0.002622 719.022 <2e-16 ***
## cutGood      0.032996  0.003916   8.425 <2e-16 ***
## cutIdeal     0.068733  0.003574  19.233 <2e-16 ***
## cutPremium   0.061333  0.003600  17.039 <2e-16 ***
## cutVery Good 0.051455  0.003637  14.147 <2e-16 ***
## colorE       -0.025128  0.002143 -11.725 <2e-16 ***
## colorF       -0.042578  0.002166 -19.662 <2e-16 ***
## ...
## clarityVVS2  0.420432  0.005195  80.931 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 
##
## Residual standard error: 0.05791 on 9981 degrees of freedom
## Multiple R-squared:  0.9827, Adjusted R-squared:  0.9827 
## F-statistic: 3.148e+04 on 18 and 9981 DF,  p-value: < 2.2e-16
```