

Regression

Meinhard Capucan, Khang Thai
This data set contains information about the price and various attributes of about 54,000 diamonds.

[Here is a link to the dataset.](#)

Linear regression is used for regression. This has its strengths and weaknesses describe below.

Linear regression consists of predictor values (x) and target values (y), where the goal is to find the relationship between x and y and be able to predict future values from this relationship. Simple linear regression has only one predictor variable. Adding more predictors makes it multiple linear regression. This algorithm works with quantitative data. Linear regression strength lies in that it is useful when the relationship between predictor and target values indicate a linear relationship, although many algorithms are better, when we know the data is linear linear regression excels. Linear regression has low variance as well. However, most data will not be linear, causing linear regression to be less favorable in most cases. Linear regression also tends to have high bias.

Data reading and Installing Packages

First, install the diamond.csv data set. Then, we look at the columns and their specific data types.

```
install.packages("tidyverse", repos = "http://cran.us.r-project.org")

# Installing package into 'C:/Users/meinc/AppData/Local/R/win-library/4.2'
# (as 'lib' is unspecified)

# package 'tidyverse' successfully unpacked and MD5 sums checked
# The downloaded binary packages are in
# C:/Users/meinc/AppData/Local/Temp/Rtmp6kxW/downloaded_packages

library(tidyverse)

# --- Attaching packages
# ---
# tidyverse 1.3.2 ---

# ✔ ggplot2 3.3.6      ✔ purrr  0.3.4
# ✔ tibble 3.1.8       ✔ dplyr  1.0.10
# ✔ tidyr  1.2.1       ✔ stringr 1.4.1
# ✔ readr  2.1.2       ✔ forcats 0.5.2

# --- Conflicts --- tidyverse_conflicts() ---
# ✔ dplyr::filter() masks stats::filter()
# ✔ dplyr::lag()    masks stats::lag()

diamonds <- read_csv("diamonds.csv")

# New names:
# Row: 53840 Columns: 11
# --- Column specification
# ---
#   (1) cut, color, clarity and (2) ... 3, carat, depth, table, price, x, y, z
# I use 'spec()' to retrieve the full column specification for this data, I
# Specify the column types or set 'show_col_types = FALSE' to quiet this message.
# ---
#   > #>   1

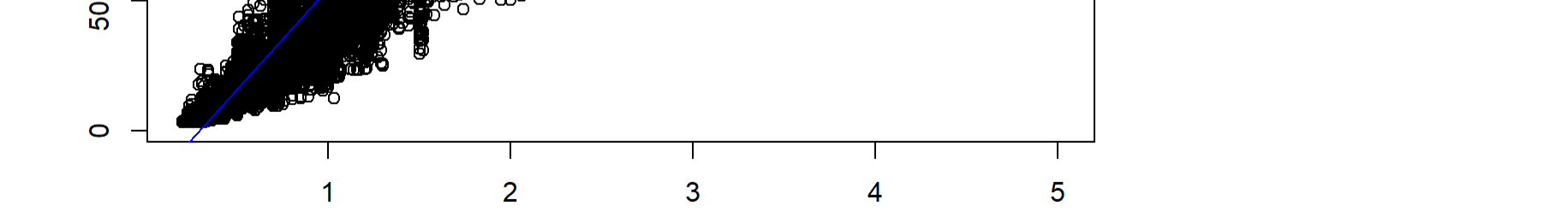
str(diamonds)

#> spec_tbl_df [53,840 x 11] (53: spec, tbl_df/tbl_info/tibble.frame)
#>   ...1 :   num [1:53840] 1 2 3 4 5 6 7 8 9 10 ...
#>   $ carat :   num [1:53840] 0.23 0.21 0.23 0.29 0.31 0.24 0.24 0.26 0.22 0.23 ...
#>   $ cut   :   chr [1:53840] "Ideal" "Premium" "Good" "Premium" ...
#>   $ color :   chr [1:53840] "E" "E" "I" "I" ...
#>   $ clarity:   chr [1:53840] "VS2" "VS1" "VS2" "VS2" ...
#>   $ depth :   num [1:53840] 61.5 59.8 56.9 62.4 63.3 62.8 62.3 61.9 65.1 59.4 ...
#>   $ table :   num [1:53840] 55 65 58 58 57 57 55 62 61 ...
#>   $ price :   num [1:53840] 326 326 327 334 335 336 336 337 337 338 ...
#>   $ x     :   num [1:53840] 3.95 3.89 4.05 4.2 4.34 3.84 3.95 4.07 3.87 4 ...
#>   $ y     :   num [1:53840] 3.98 3.88 4.07 4.21 4.35 3.98 3.98 4.11 3.78 4.05 ...
#>   $ z     :   num [1:53840] 2.43 2.31 2.31 2.63 2.75 2.48 2.47 2.53 2.49 2.39 ...
#>   $      :   chr [1:53840] "VS2" ...
#>   ... col:
#>   ...   cut = col_double(),
#>   ...   color = col_double(),
#>   ...   clarity = col_character(),
#>   ...   depth = col_double(),
#>   ...   table = col_double(),
#>   ...   price = col_double(),
#>   ...   x = col_double(),
#>   ...   y = col_double(),
#>   ...   z = col_double(),
#>   ...
#>   attr(,"problems")=externalptr
```

Graphs and Plotting

Here, we plotted cut as a function of carat to observe any relationships or trends. The abline function plots a general line through the graph.

```
plot(diamonds$price~diamonds$carat, xlab="Carat", ylab="Price")
abline(lm(diamonds$price~diamonds$carat), col="blue")
```



Function of carat is split into colors and made easier to read, where any carat values above 2.0 are colored blue. Anything below that is colored red.

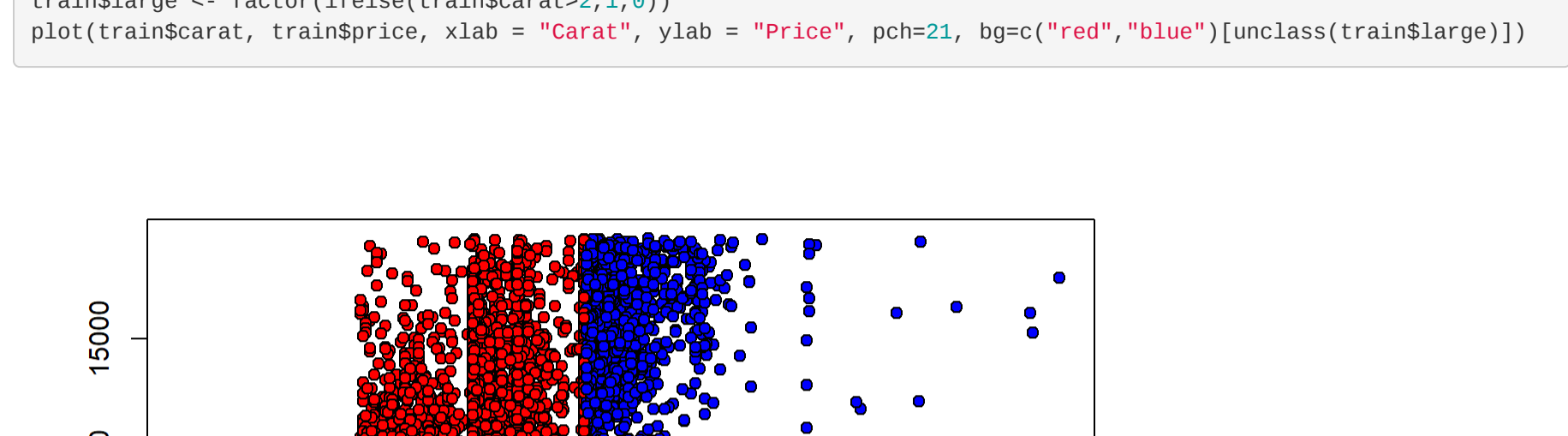
Divide into Train and Test Data

We will go ahead and divide the test data into training and test data, with 80% being training and 20% being test.

```
set.seed(1234)
i <- sample(nrow(diamonds), nrow(diamonds)*0.8, replace = FALSE)
train <- diamonds[i,]
test <- diamonds[-i,]
```

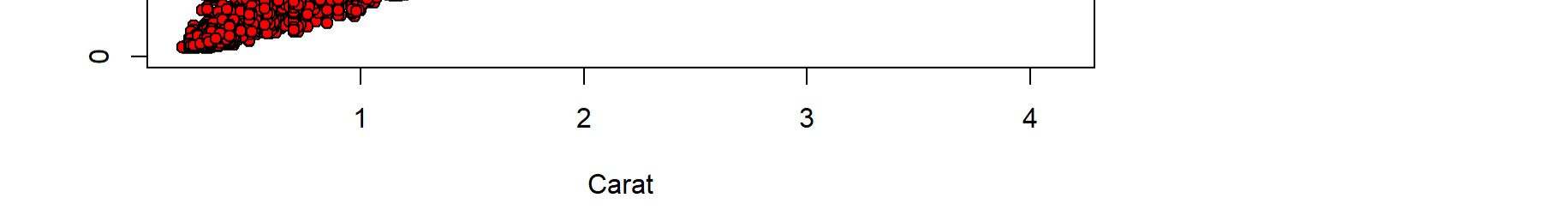
Creating Graphs

```
par(mfrow=c(1,2))
trainlarge <- factor(ifelse(train$carat>2.1,0))
plot(train$carat, train$price, xlab="Carat", ylab="Price", pch=21, bty="n", col=c("red", "blue"))[unclass(trainlarge)]
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



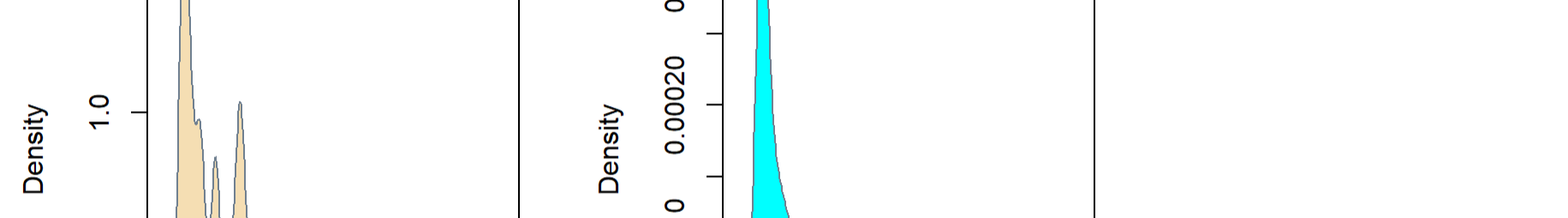
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



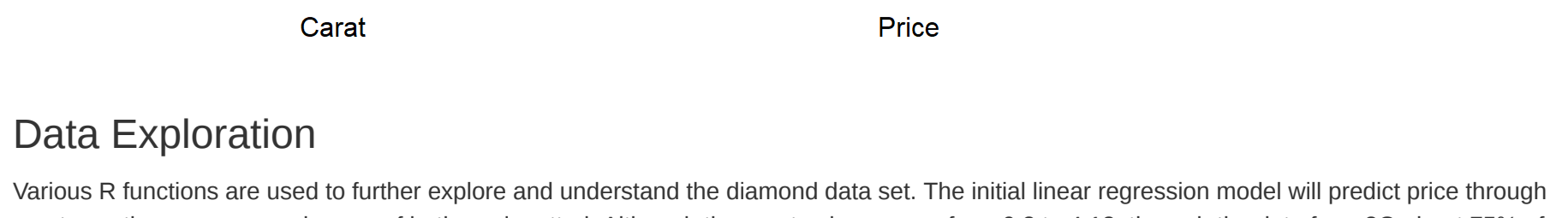
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



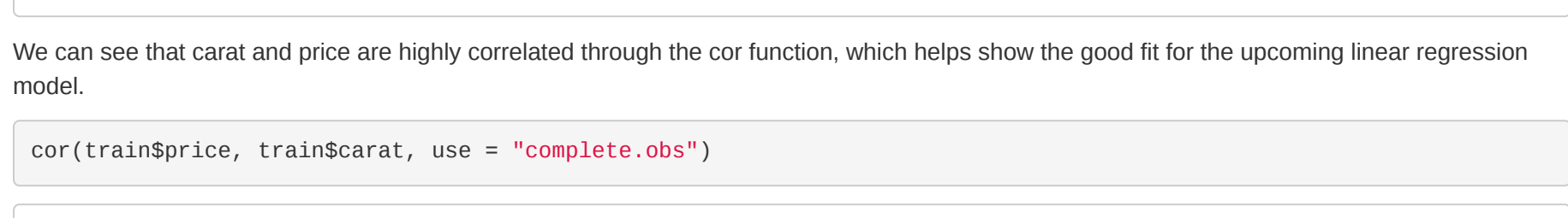
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



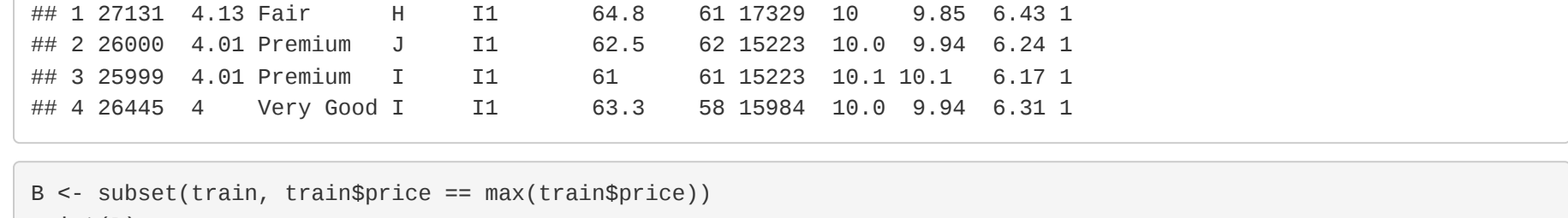
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



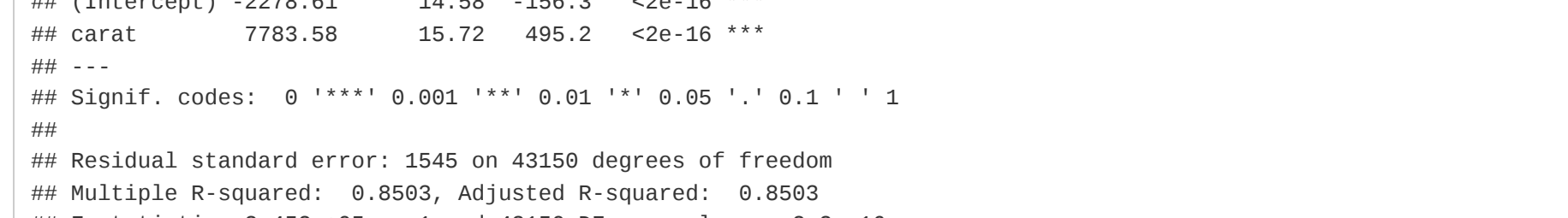
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



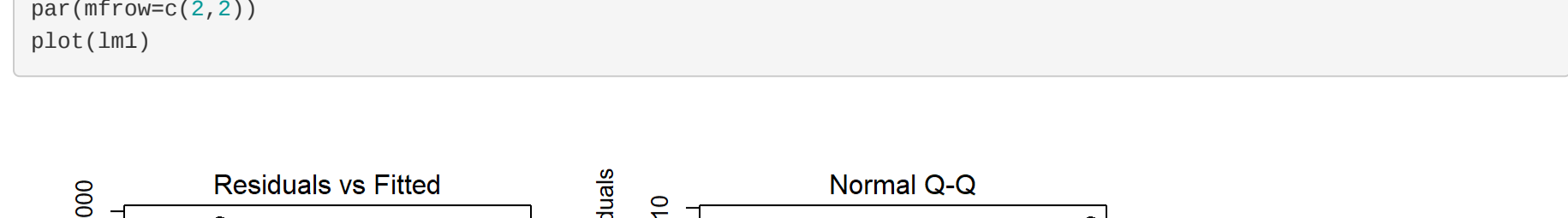
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



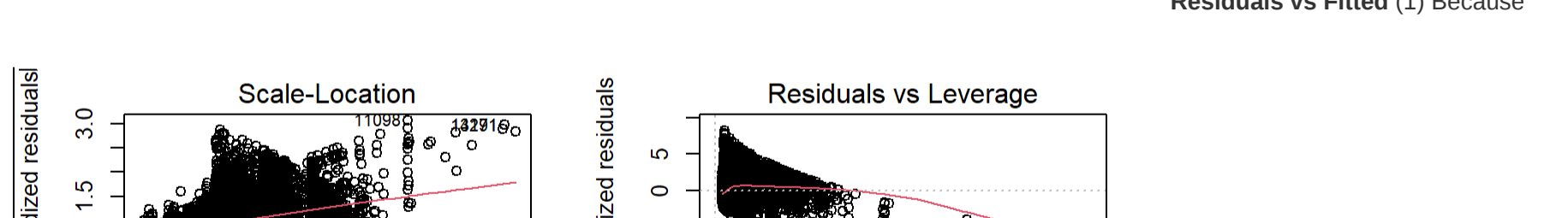
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



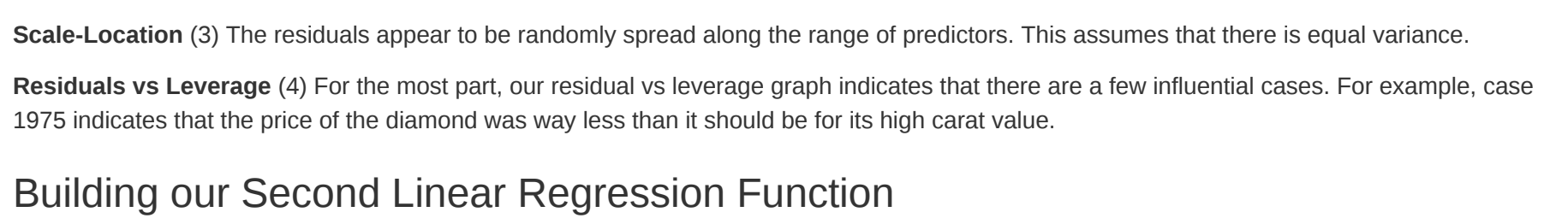
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



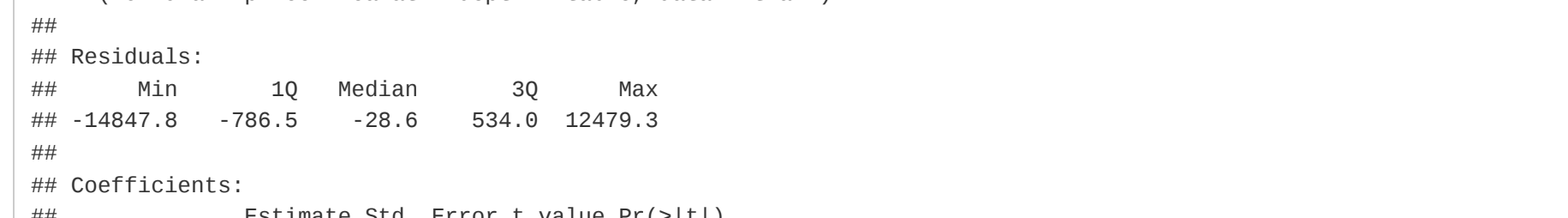
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



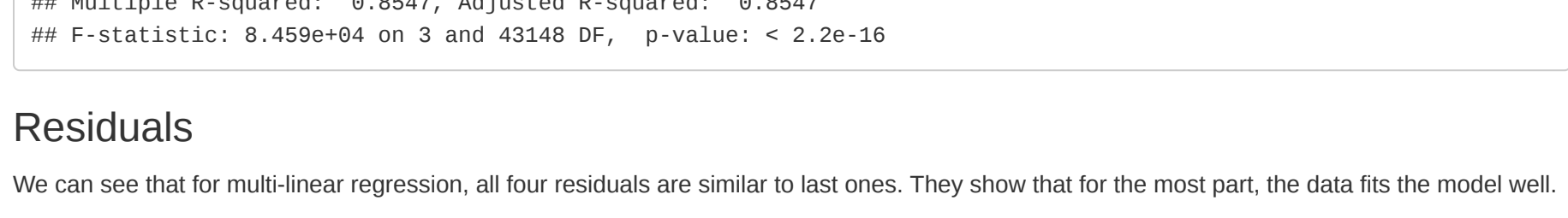
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



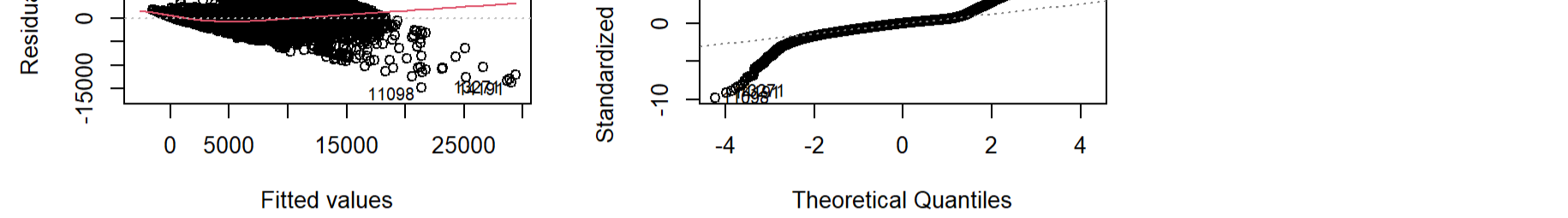
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



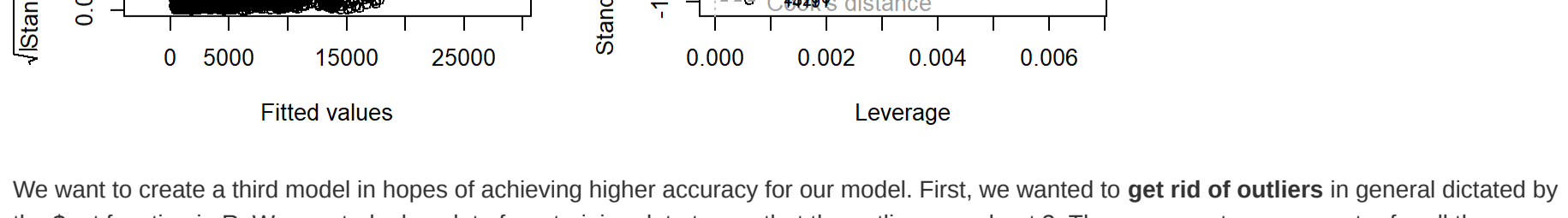
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



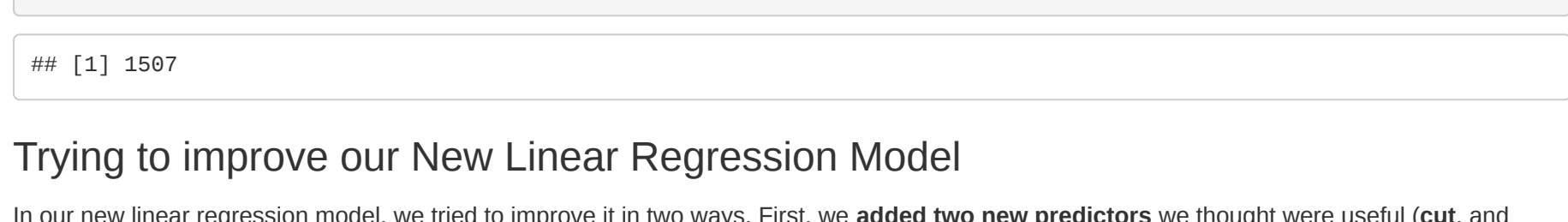
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



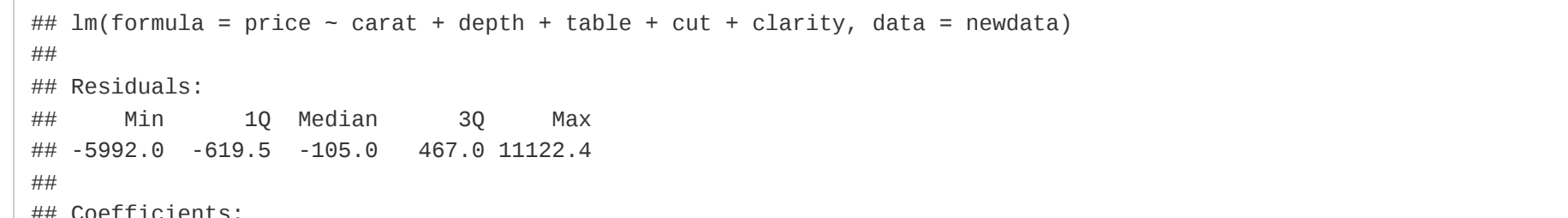
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



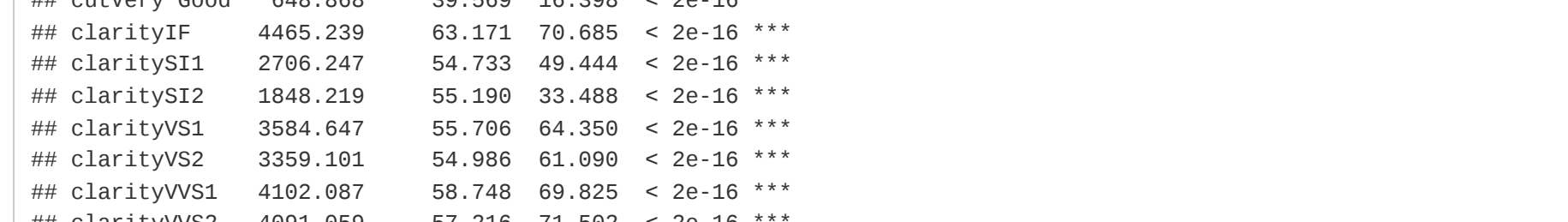
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



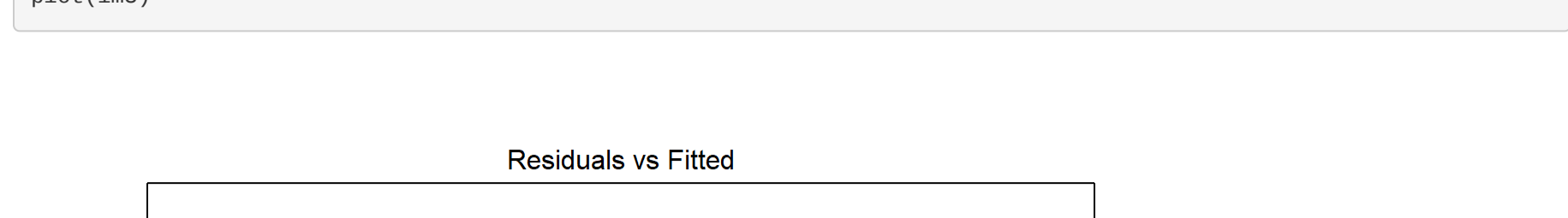
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



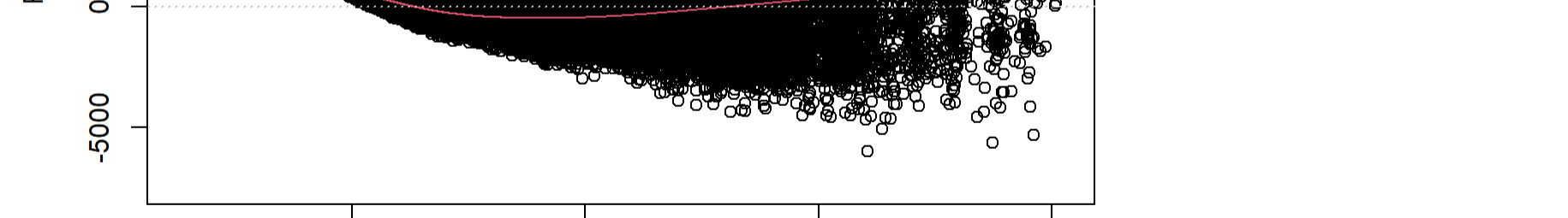
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



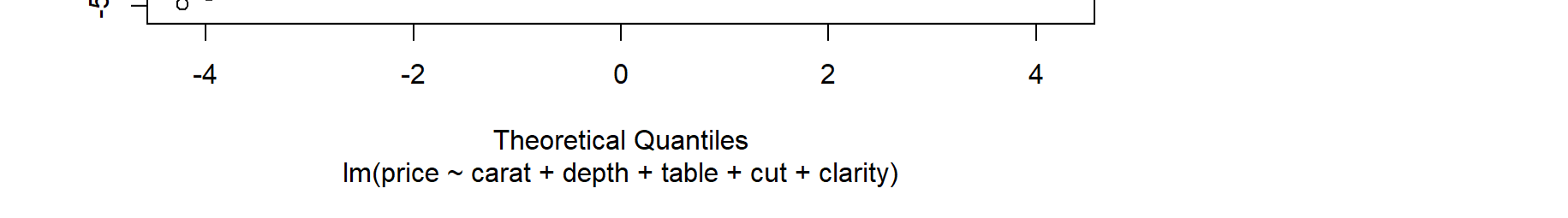
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



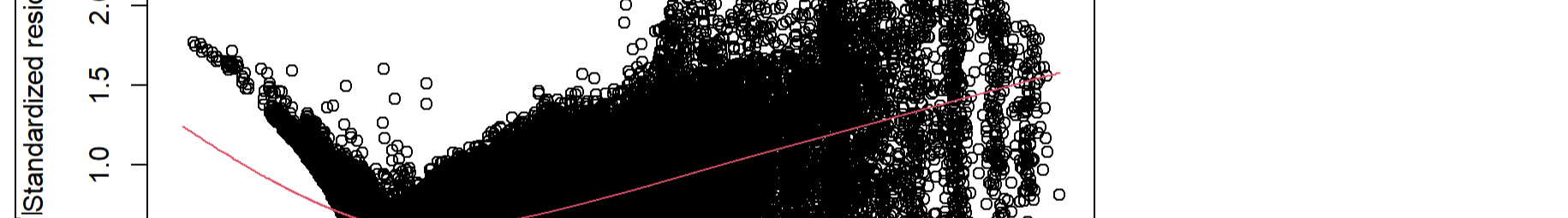
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



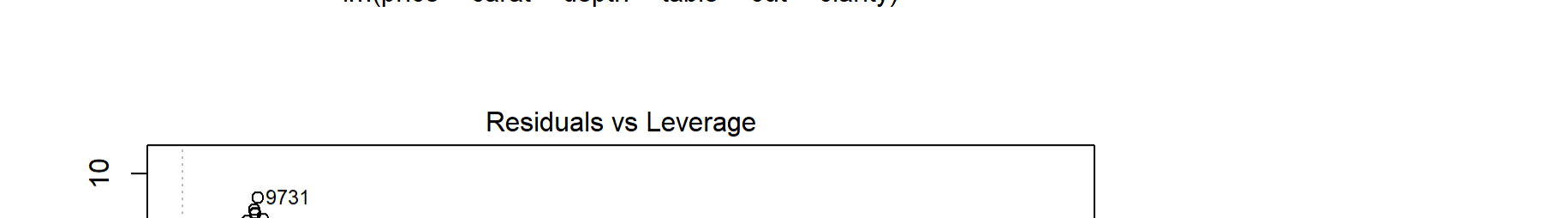
price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```



price and density. This reveals how diamond carats concentrate around values from 0.0 to 2.0. For the price, a similar trend is seen as well where they concentrate towards the left side.

```
par(mfrow=c(1,2))
d <- density(train$carat, na.rm = TRUE)
plot(d, main = "Density Plot for Carat", xlab = "Carat")
polygon(d, col="yellow", border="slategrey")
d <- density(train$price, na.rm = TRUE)
plot(d, main = "Density Plot for Price", xlab = "Price")
polygon(d, col="cyan", border="slategrey")
```

