

Homework 4

Problem 1

- a. Give the names of the team members for your course project. You can work alone, or with one more person in the class.
- b. Provide descriptions or suggestions of the project that you will be interested in. This is for my information only - it will help me decide on the topics that we will cover for the rest of the semester. You can fine-tune or change the scope of the project later on. You are welcome to use a project from your research, provided that you make a new specific effort for the class.

Problem 2

'KNNL' refers to the book by Kutner, Nachtsheim, neter and Li.

- a. KNNL 1.3
- b. KNNL 1.5
- c. KNNL 1.16
- d. KNNL 1.18
- e. KNNL 1.22
- f. KNNL 1.26

Problem 3

First load the packages.

Load packages

```
library(tidyr, quietly = TRUE)
library(dplyr, quietly = TRUE)
```

We will use the Sacramento dataset available at

<http://samplecsvs.s3.amazonaws.com/Sacramentorealestatetransactions.csv>

(<http://samplecsvs.s3.amazonaws.com/Sacramentorealestatetransactions.csv>)

```
sestates <- read.csv("http://samplecsvs.s3.amazonaws.com/Sacramentorealestatetransactions.csv", header = TRUE, stringsAsFactors = FALSE)
head(sestates)
```

```
##           street      city  zip state beds baths sq__ft      type
## 1      3526 HIGH ST SACRAMENTO 95838    CA    2    1    836 Residential
## 2        51 OMAHA CT SACRAMENTO 95823    CA    3    1   1167 Residential
## 3     2796 BRANCH ST SACRAMENTO 95815    CA    2    1    796 Residential
## 4    2805 JANETTE WAY SACRAMENTO 95815    CA    2    1    852 Residential
## 5     6001 MCMAHON DR SACRAMENTO 95824    CA    2    1    797 Residential
## 6 5828 PEPPERMILL CT SACRAMENTO 95841    CA    3    1   1122         Condo
##           sale_date price latitude longitude
## 1 Wed May 21 00:00:00 EDT 2008 59222 38.63191 -121.4349
## 2 Wed May 21 00:00:00 EDT 2008 68212 38.47890 -121.4310
## 3 Wed May 21 00:00:00 EDT 2008 68880 38.61830 -121.4438
## 4 Wed May 21 00:00:00 EDT 2008 69307 38.61684 -121.4391
## 5 Wed May 21 00:00:00 EDT 2008 81900 38.51947 -121.4358
## 6 Wed May 21 00:00:00 EDT 2008 89921 38.66260 -121.3278
```

Use functions in dplyr and tidyr to answer the following questions. Provide code and output to justify your answers.

a. Replace the column “sale_date” with three coluns “day_week”, “month”, “day_month”

```
##           street      city  zip state beds baths sq__ft      type
## 1      3526 HIGH ST SACRAMENTO 95838    CA    2    1    836 Residential
## 2        51 OMAHA CT SACRAMENTO 95823    CA    3    1   1167 Residential
## 3     2796 BRANCH ST SACRAMENTO 95815    CA    2    1    796 Residential
## 4    2805 JANETTE WAY SACRAMENTO 95815    CA    2    1    852 Residential
## 5     6001 MCMAHON DR SACRAMENTO 95824    CA    2    1    797 Residential
## 6 5828 PEPPERMILL CT SACRAMENTO 95841    CA    3    1   1122         Condo
##   day_week month day_month price latitude longitude
## 1    Wed    May        21 59222 38.63191 -121.4349
## 2    Wed    May        21 68212 38.47890 -121.4310
## 3    Wed    May        21 68880 38.61830 -121.4438
## 4    Wed    May        21 69307 38.61684 -121.4391
## 5    Wed    May        21 81900 38.51947 -121.4358
## 6    Wed    May        21 89921 38.66260 -121.3278
```

b. What’s the top 10 cities with the most transactions?

```
## Source: local data frame [10 x 2]
##
##           city n_trans
##      (chr)   (int)
## 1  SACRAMENTO    439
## 2   ELK GROVE   114
## 3    LINCOLN    72
## 4   ROSEVILLE   48
## 5 CITRUS HEIGHTS   35
## 6    ANTELOPE    33
## 7 RANCHO CORDOVA   28
## 8 EL DORADO HILLS   23
## 9         GALT    21
## 10 NORTH HIGHLANDS 21
```

c. What's the accumulated number of transactions from May 15 to May 21 in city ELK GROVE?

```
## Source: local data frame [5 x 4]
## Groups: month [1]
##
##   month day_month n_trans n_cumtrans
##   (chr)   (chr)   (int)   (int)
## 1  May      15      13        13
## 2  May      16      28        41
## 3  May      19      19        60
## 4  May      20      22        82
## 5  May      21      32       114
```

d. For each type of house (Condo, Multi-Family, and Residential), what's the highest 3 transaction prices? In which cities?

```
## Source: local data frame [9 x 3]
## Groups: type [3]
##
##           type           city price
##      (chr)         (chr)   (int)
## 1      Condo      SACRAMENTO 360000
## 2      Condo      ROSEVILLE 350000
## 3      Condo      GOLD RIVER 300000
## 4 Multi-Family      SACRAMENTO 416767
## 5 Multi-Family      SACRAMENTO 297000
## 6 Multi-Family      AUBURN 285000
## 7 Residential      WILTON 884790
## 8 Residential EL DORADO HILLS 879000
## 9 Residential      LOOMIS 839000
```

e. Are the values in column sq__ft looking okay? If not, what changes need to be made? For each type of house (Condo, Multi-Family, and Residential), what's the top 3 cities with the highest price per square foot?

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##           0     952    1304    1315    1718    5822
```

Not possible to have a house of square foot zero. We should set the zero values to NA (missing value). Do so and return summary output.

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
##      484    1144    1418    1591    1851    5822    171
```

```
## Source: local data frame [9 x 3]
## Groups: type [3]
##
##           type      city price_sq
##           (chr)    (chr)    (dbl)
## 1      Condo SACRAMENTO 304.9233
## 2      Condo SACRAMENTO 214.3740
## 3      Condo   AUBURN 207.1713
## 4 Multi-Family   AUBURN 296.8750
## 5 Multi-Family SACRAMENTO 238.9719
## 6 Multi-Family SACRAMENTO 134.2070
## 7 Residential SACRAMENTO 619.6660
## 8 Residential   LOOMIS 516.6256
## 9 Residential SACRAMENTO 459.2652
```

f. For each type of house (Condo, Multi-Family, and Residential) in city SACRAMENTO, what's the number of transactions, average price, and average price per square foot?

```
## Source: local data frame [3 x 4]
##
##           type n_trans ave_price ave_price_sq
##           (chr)  (int)    (dbl)    (dbl)
## 1      Condo      27  137690.7    132.8999
## 2 Multi-Family     10  214189.7    104.5987
## 3 Residential    402  201359.6    137.5226
```

Problem 4

In this assignment, we'll use ggplot2 to reverse-engineer a plot constructed using another package. You'll see the output of the ggplot2 code, but you'll have to provide that code yourself.

We'll example the dataset *diamonds*, a dataset containing the prices and other attributes of almost 54,000 diamonds. Run `?diamonds` to learn about the data and each variable it contains.

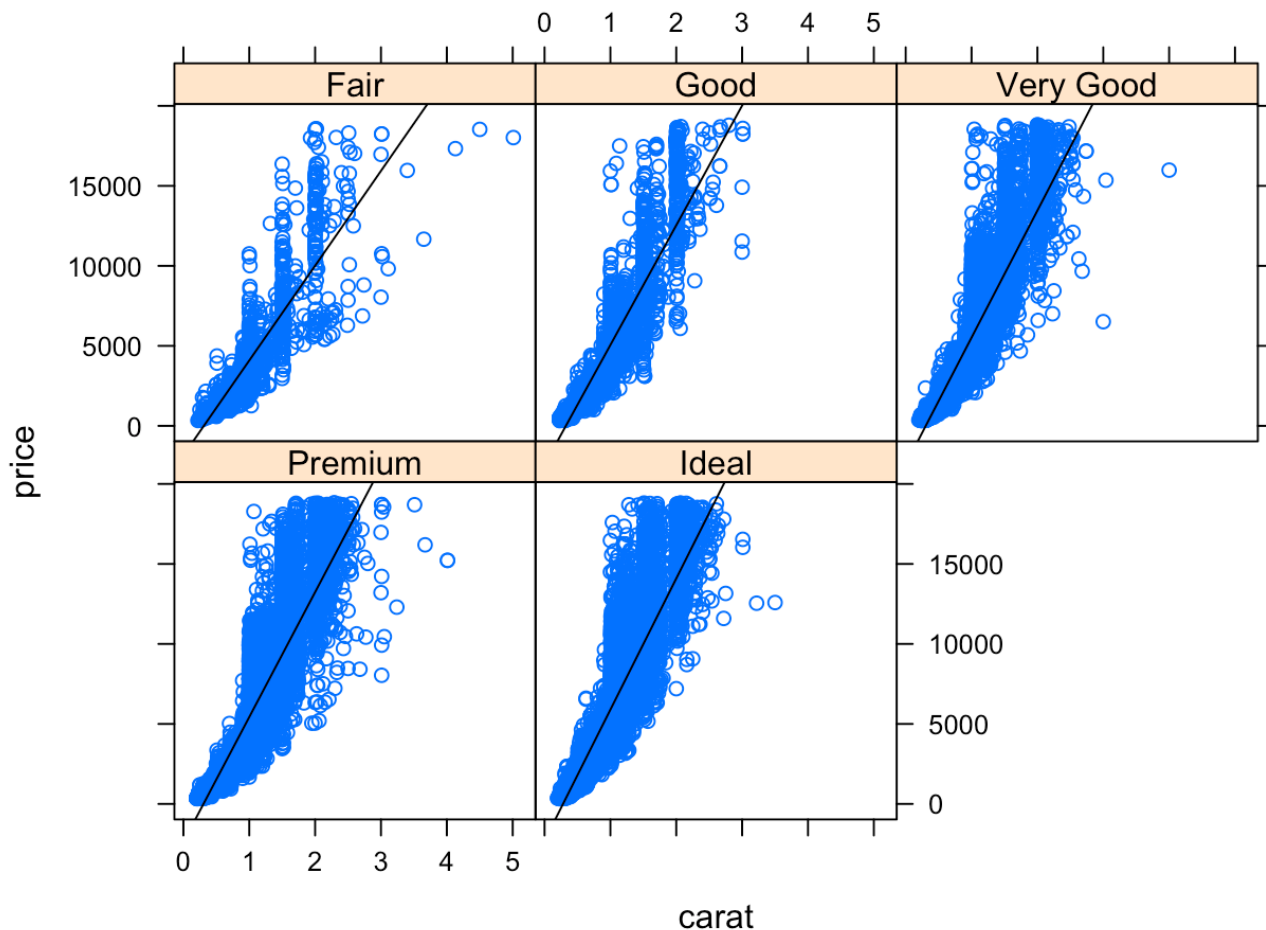
```
library(lattice, quietly = TRUE)
library(ggplot2, quietly = TRUE)
```

```
## Warning: package 'ggplot2' was built under R version 3.2.3
```

```
data(diamonds)
```

Run the following code using functions from the lattice package.

```
xyplot(price ~ carat | cut,  
        data = diamonds,  
        panel = function(x, y, ...) {  
          panel.xyplot(x, y, ...)  
          lm1 <- lm(y ~ x)  
          panel.abline(a = lm1$coefficients[1],  
                      b = lm1$coefficients[2])  
        },  
        as.table = TRUE)
```

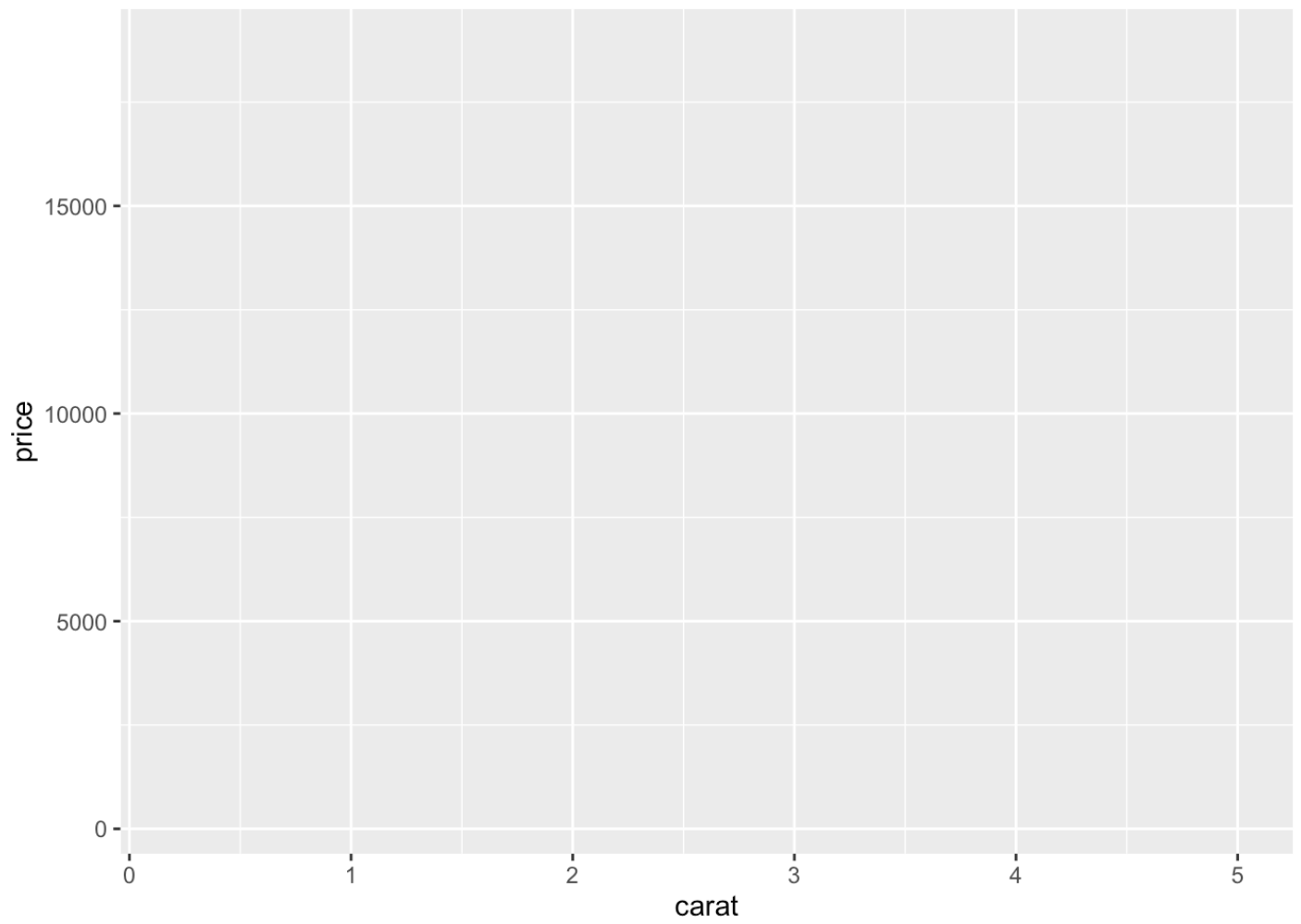


This plots the number of carats on the x axis and the price on the y axis. The panels group the data by the diamond quality. A linear regression line is fit to the data in each panel, and is plotted.

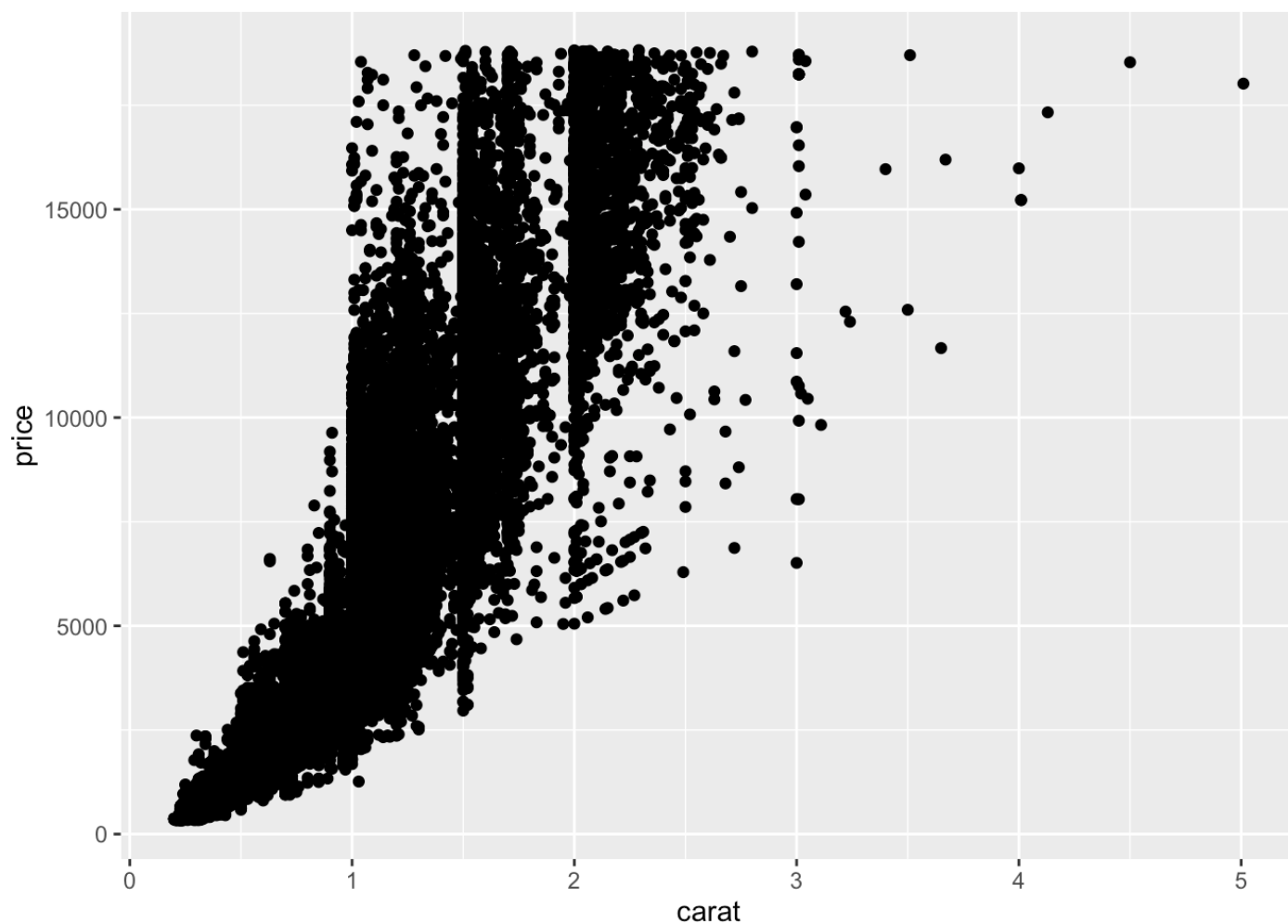
Your task is to reproduce this plot in ggplot2.

- Initialize a ggplot object, map carat to the x axis and price to the y axis. You should get the

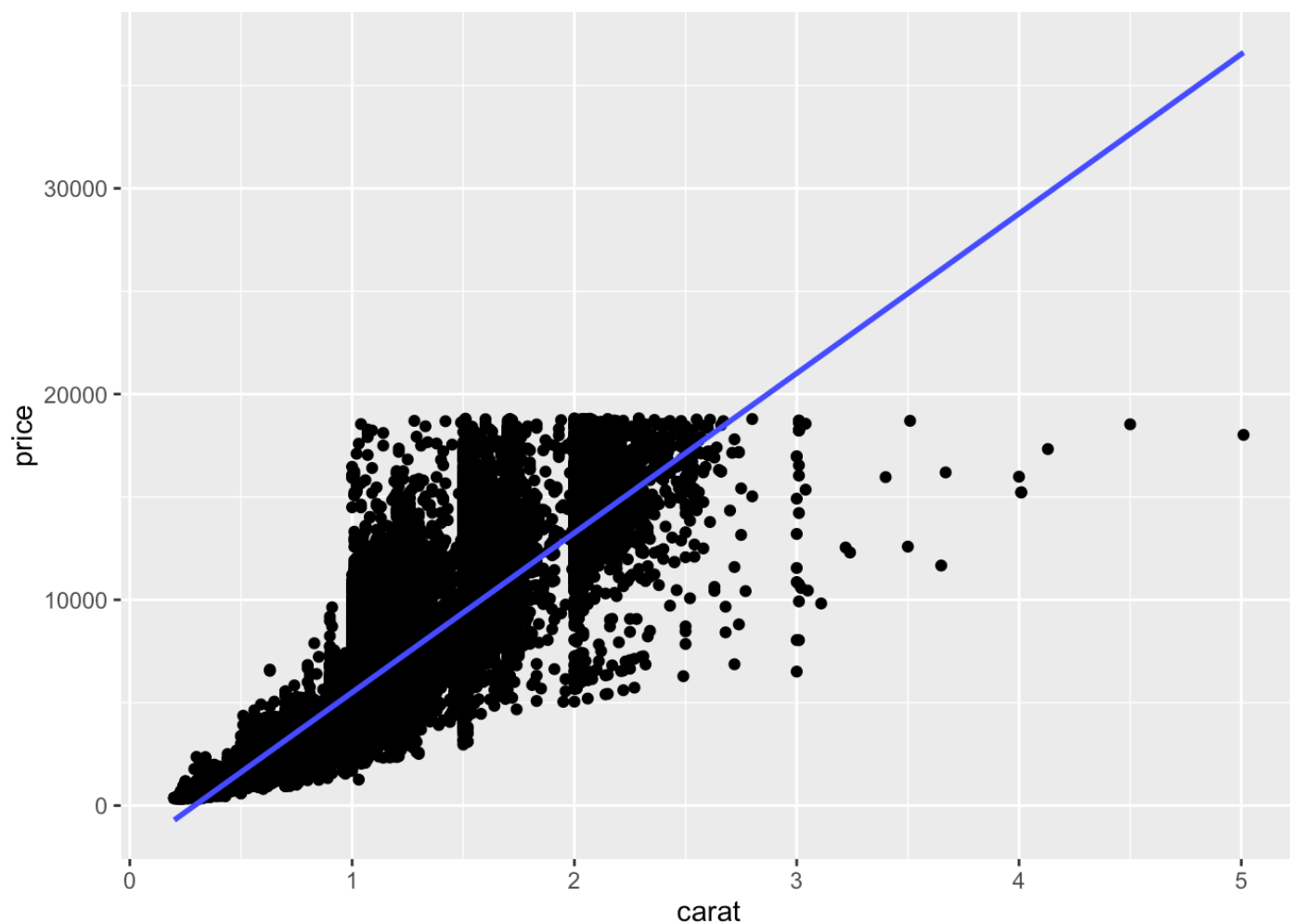
following output.



b. Next, plot the points. (hint: “geom_points”)



c. Next, we need to add a linear regression line. (Hint: The “geom_smooth” function fits various lines to data. Run `?geom_smooth` to check the help documentation and make sure you fit a linear regression line.)



d. Finally, separate out into panels according to diamond quality (indicated by the “cut” variable). (hint: You need “facet_grid”, and you will have to pass a formula object to the formula argument. Look at the help docs and play around.)

