

EDA

Load Libraries

```
library(tidyverse) # data manipulation, create plots
library(DataExplorer) # correlation matrix
```

Load Data

```
df<-read.csv("rent_amount.csv")
```

head(df)								
city	area	rooms	bathroom	parking.spaces	floor	animal	furniture	rentAmount
1 yes	240	3	3	4	·	accept	furnished	R\$8,000
2 no	64	2	1	1	10	accept	not furnished	R\$820
3 yes	443	5	5	4	3	accept	furnished	R\$7,000
4 yes	73	2	2	1	12	accept	not furnished	R\$1,250
5 yes	19	1	1	0	·	not accept	not furnished	R\$1,200
6 yes	13	1	1	0	2	accept	not furnished	R\$2,200

```
df %>% select(fire.insurance, rent.amount, floor) %>% summary()
```

fire.insurance	rent.amount	floor
Length:6080	Length:6080	Length:6080
Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character

We note that the price of fire insurance, floor and the price of rent. They are character type values, due to the fact that you have special characters, R standardizes and converts them to character type.

Convert variables to numeric

```
conv_to_numeric<-function(x){
  x<-gsub("[/R$]", "", x) # eliminate special characters
  x<-as.numeric(x) # transform to numeric data
  return(x) # return data
}
```

```
attach(df)
```

The following objects are masked from df (pos = 4):

animal, area, bathroom, city, fire.insurance, floor, furniture, parking.spaces, rent.amount, rooms, total

```
fire.insurance<-apply(fire.insurance, conv_to_numeric)
rent.amount<-apply(rent.amount, conv_to_numeric)
```

```
df<-df %>%
  mutate(fire.insurance=fire.insurance,
         rent.amount=rent.amount)
# With the mutate function we make modifications to the data frame.
```

```
unique(floor)
```

[1] "·" "18" "3" "12" "2" "16" "6" "4" "1" "7" "13" "9" "14" "5" "8" "15" "11" "19" "20"
[20] "24" "23" "17" "22" "27" "85" "28" "29" "35" "21" "31" "99" "26" "68" "32" "51"

We note that for unique values of the floor variable an underscore appears. Which indicates that the rental house is independent, that is, it is not established in a building.

```
floor_modify<-ifelse(floor=="·", 0, floor)
# Where the hyphen appears, we transform it to 0 for the reasons explained above.
floor_modify<-as.numeric(floor_modify) # transform to numeric data
```

```
df<- df %>% mutate(floor=floor_modify)
```

Do furnished houses have a higher rental price compared to those that are not?

```
df %>% group_by(furniture) %>% summarise(rent.amount.mean=mean(rent.amount))
```

furniture	rent.amount.mean
furnished	5387.092
not furnished	4047.211

The average price of furnished houses is higher than those that are not.

Accepting animals increases the price?

```
df %>% group_by(animal) %>% summarise(animal.mean=mean(rent.amount))
```

animal	animal.mean
accept	4585.849
not accept	3768.856

The average price of apartments where animals are accepted is higher. Since they are allowing you the luxury of putting your pets in the apartment.

```
histogram<-function(x,...){
  df %>%
    ggplot(aes(x=x, y=-.density..)) +
    geom_histogram(color="black", fill="white", width=1) +
    geom_density(color="black", lwd=1) +
    geom_vline(aes(x=intercept=mean(x), color="mean")) +
    geom_vline(aes(x=intercept=median(x), color="median")) +
    labs(col="Statistics") +
    theme(legend.position = "top") +
    ...
}
```

```
rent.amount_histogram=histogram(df$rent.amount, labs(x="Rent Amount", title = "Rent Amount"))
area_histogram=histogram(df$area, labs(x="Area", title = "Area"))
fire.insurance_histogram=histogram(df$fire.insurance, labs(x="Fire Insurance", title = "Fire Insurance"))
rooms_histogram=histogram(df$rooms, labs(x="Rooms", title = "Rooms"))
bathroom_histogram=histogram(df$bathroom, labs(x="Bathroom", title = "Bathroom"))
floor_histogram=histogram(df$floor, labs(x="Floor", y="Rent Amount", title = "Floor"))
```

```
scatter_plot<-function(x, feature,...){
  ggplot(data=df, aes(x=x, y=rent.amount)) +
  geom_point(color="red", alpha=0.5) +
  theme_light() +
  geom_smooth(method = "lm", color="red") +
  ...
}
```

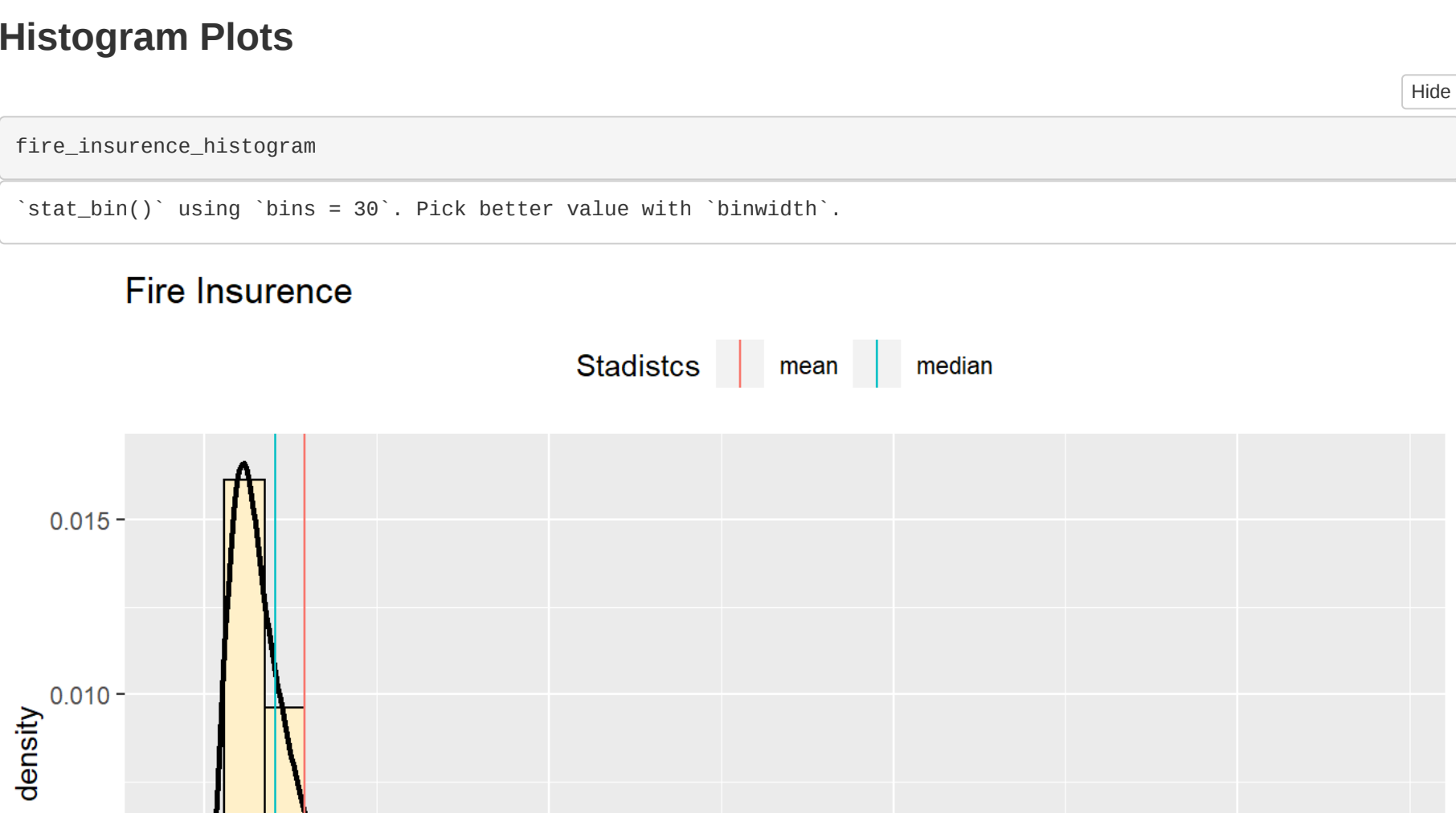
```
area_scatter<-scatter_plot(df$area, labs(x="Area", y="Rent Amount"))
fire_scatter<-scatter_plot(df$fire.insurance, labs(x="Fire Insurance", y="Rent Amount"))
rooms_scatter<-scatter_plot(df$rooms, labs(x="Rooms", y="Rent Amount"))
bathroom_scatter<-scatter_plot(df$bathroom, labs(x="Bathroom", y="Rent Amount"))
floor_scatter<-scatter_plot(df$floor, labs(x="Floor", y="Rent Amount"))
```

```
library(gridExtra)
```

Histogram Plots

```
fire.insurance_histogram
```

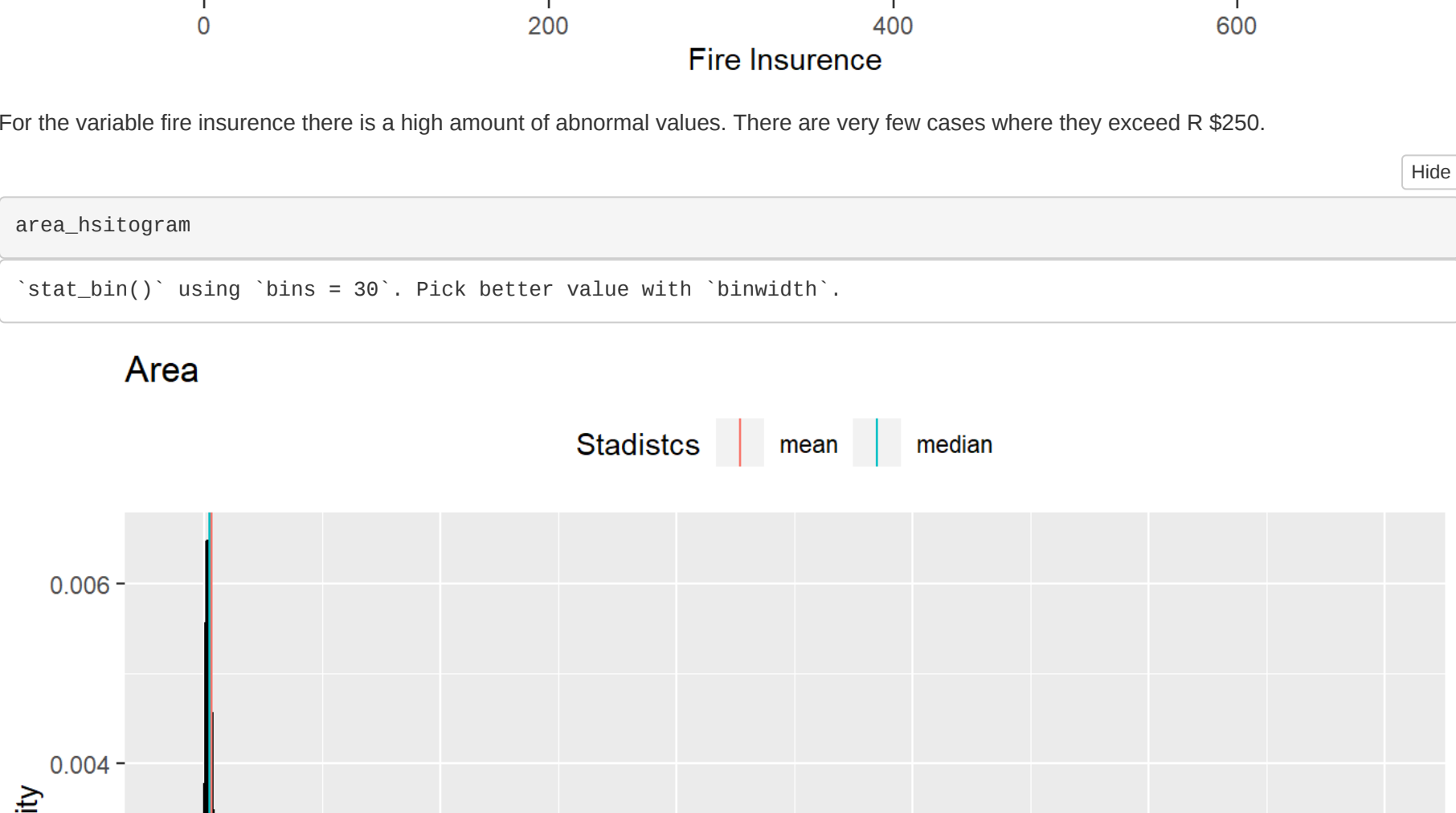
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



For the variable fire insurance there is a high amount of abnormal values. There are very few cases where they exceed R \$250.

```
area_histogram
```

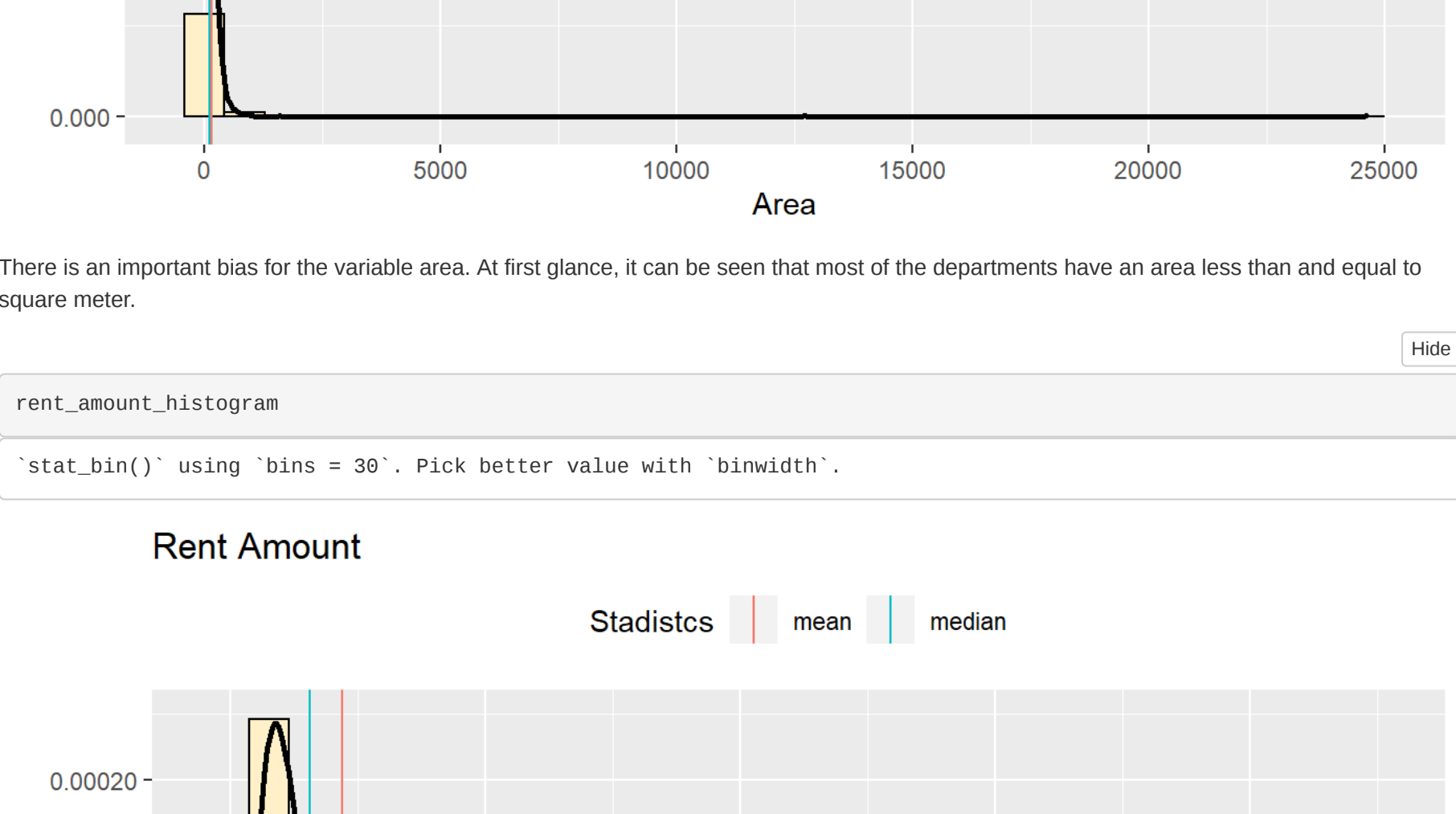
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



There is an important bias for the variable area. At first glance, it can be seen that most of the departments have an area less than and equal to square meter.

```
rent.amount_histogram
```

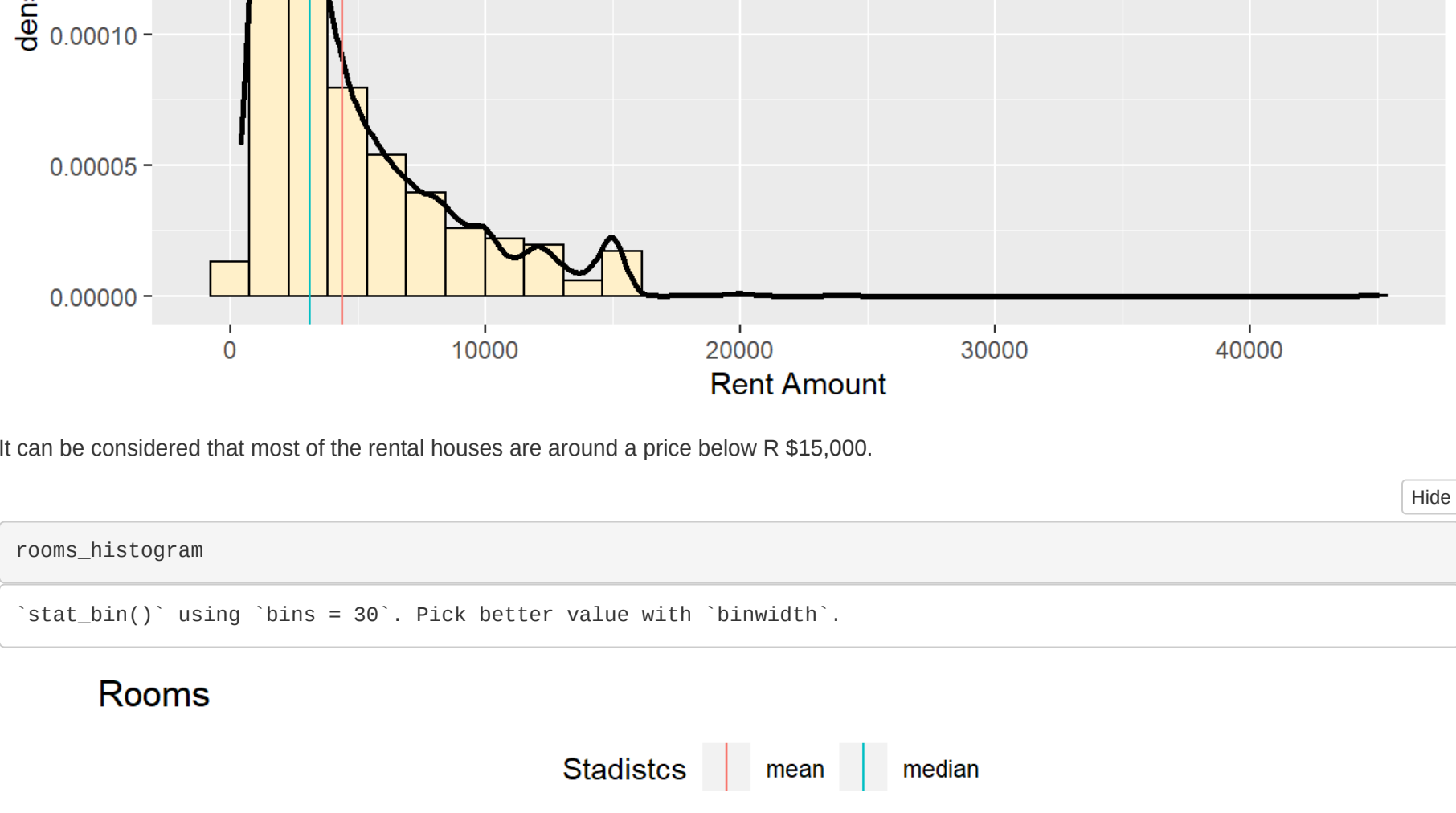
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



It can be considered that most of the rental houses are around a price below R \$15,000.

```
rooms_histogram
```

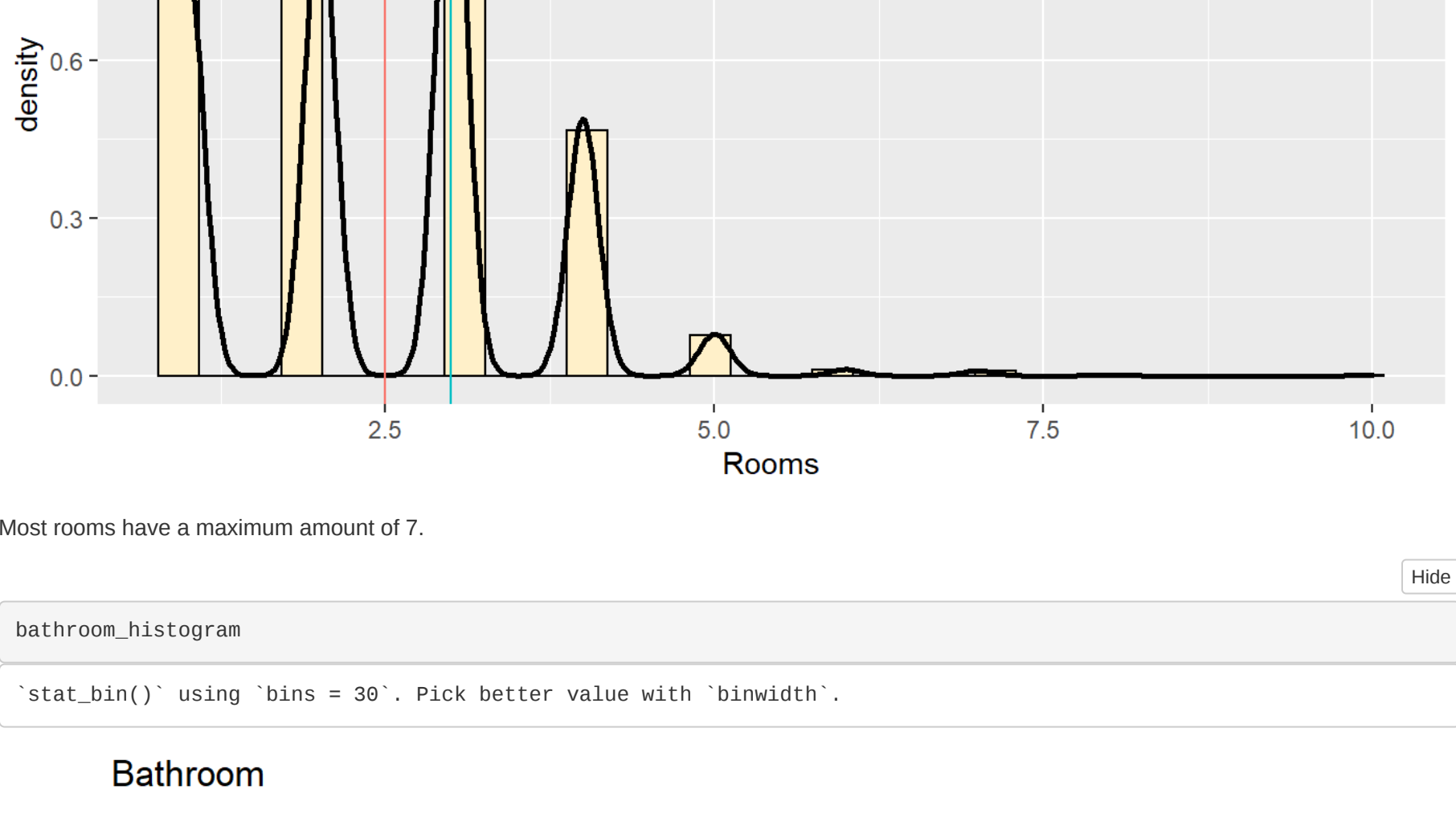
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



Most rooms have a maximum amount of 7.

```
bathroom_histogram
```

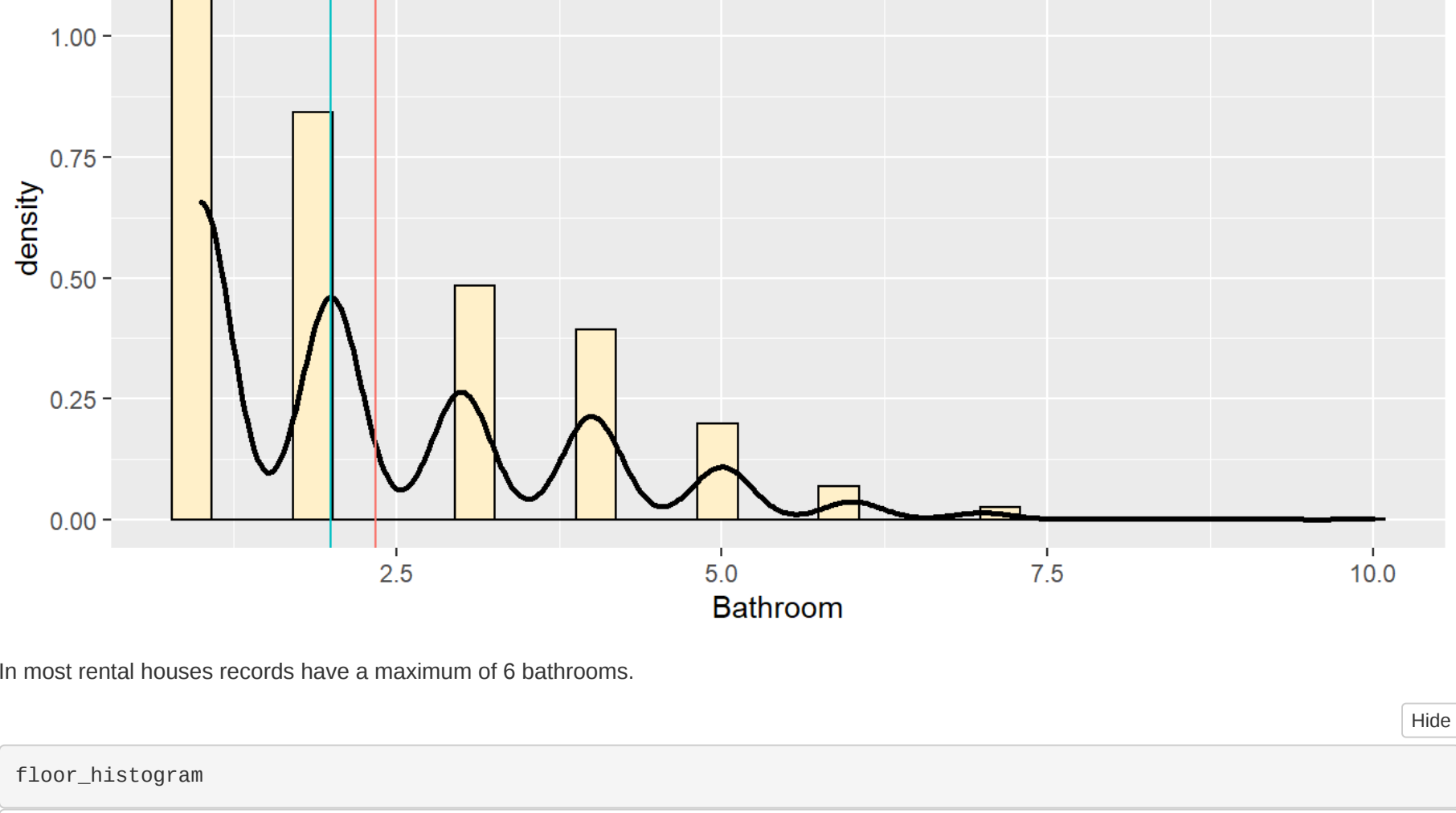
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



In most rental houses records have a maximum of 6 bathrooms.

```
floor_histogram
```

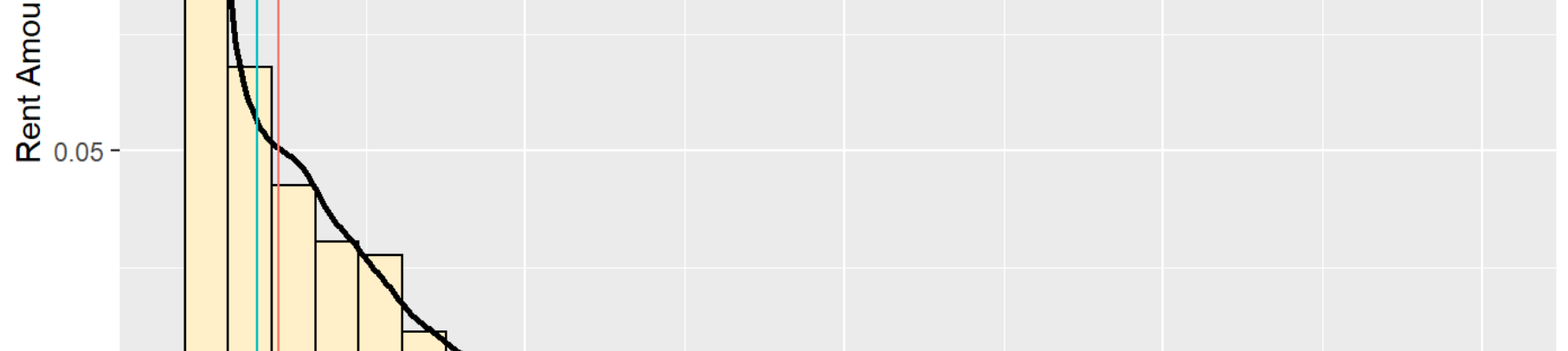
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



Most of the buildings where the apartments are rented are less than 25 stories.

Scatter Plots

```
grid.arrange(area_scatter,
              fire_scatter,
              rooms_scatter,
              bathroom_scatter,
              floor_scatter)
```



The area variable apparently does not show any possible correlation with the rental price. Since it contains a good number of outliers, which causes the data to be skewed.

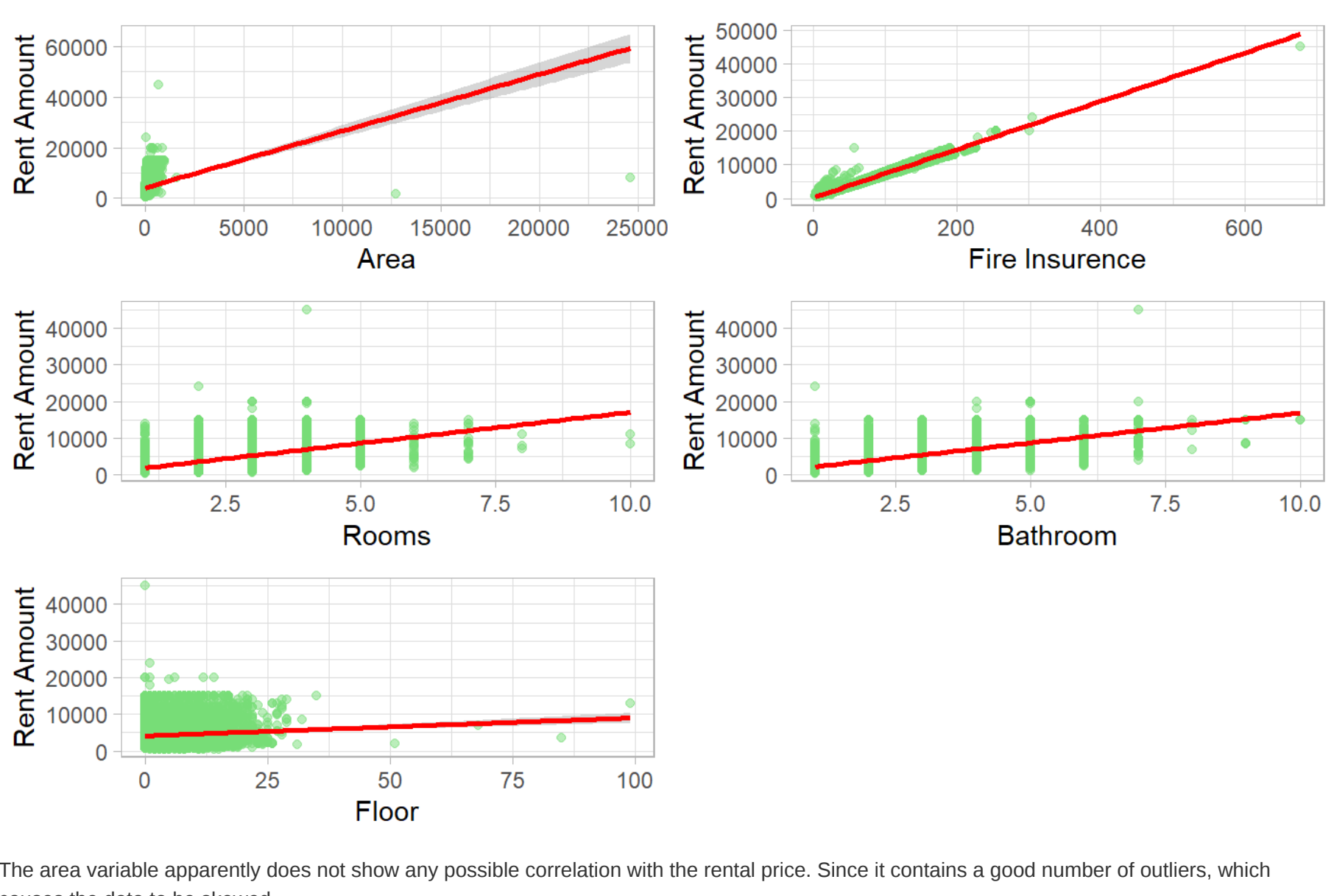
There is a clear correlation between the price of fire insurance with respect to the price of the rental house, since the more expensive the price of the apartment, the more you will have to pay to protect it against fire. These variables have a linear relationship, that is, they increase proportionally with another.

The variable floor, rooms and bathrooms have to raise the price of rent rent something that is logical.

Correlation Matrix

Shows the degree of relationship of the variables. They are measured from 0 to 1 if it is a positive correlation, otherwise it is measured from 1 to -1.

```
plot_correlation(df, title = "Correlation Matrix")
```



Save Dataset

```
write.csv(df, "rent-amount.csv", row.names = FALSE)
```

Conclusion

There is a strong presence of outliers in the data set. Especially for the area variable. For continuous variables, that is, those values with decimals, we can perform a logarithmic transformation to transform the outliers.