

EDA

Load Libraries

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

Load Data

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

```
head(df)
```

	X	city	area	rooms	bathroom	parking-spaces	floor	animal	furniture
1	0	1	240	3	3	4	-	accept	furnished
2	1	0	64	2	1	1	10	accept	not furnished
3	2	1	443	5	5	4	3	accept	furnished
4	3	1	73	2	2	1	12	accept	not furnished
5	4	1	19	1	1	0	-	not accept	not furnished
6	5	1	13	1	1	0	2	accept	not furnished

```
df<- df %>% select(-c(floor,animal,hoa,total,X,city,property.tax))
```

Type of Data

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

```
## fire.insurance    rent.amount
## Length:6080      Length:6080
## Class :character  Class :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("[/\\$,%]", "",x) # eliminate special charaters
  x<-as.numeric(x) # transform to numeric data
  return(x) # return data
}
attach(df)
```

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

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.

```
histogram<-function(x,...){
  df %>%
    ggplot(aes(x=x,y=-.density..)) +
    geom_histogram(color="black",fill="#FFD0C0") +
    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$rent.amount,lab(x="Rent Amount",title = "Rent Amount"))
area_histogram=histogram(df$area,lab(x="Area",title = "Area"))
fire_insurance_histogram=histogram(df$fire.insurance,lab(x="Fire Insurance",title = "Fire Insurance"))
rooms_histogram=histogram(df$rooms,lab(x="Rooms",title = "Rooms"))
bathroom_histogram=histogram(df$bathroom,lab(x="Bathroom",title = "Bathroom"))
parking_spaces_histogram=histogram(df$parking.spaces,lab(x="Parking Spaces",title="Parking Spaces"))
```

```
scatter_plot<-function(x,feature,...){
  ggplot(data=df,aes(x=x,feature,y=rent.amount)) +
  geom_point(color="#77D077",alpha=0.5) +
  theme_light() +
  geom_smooth(method = "lm",color="red") +
  ...
}
```

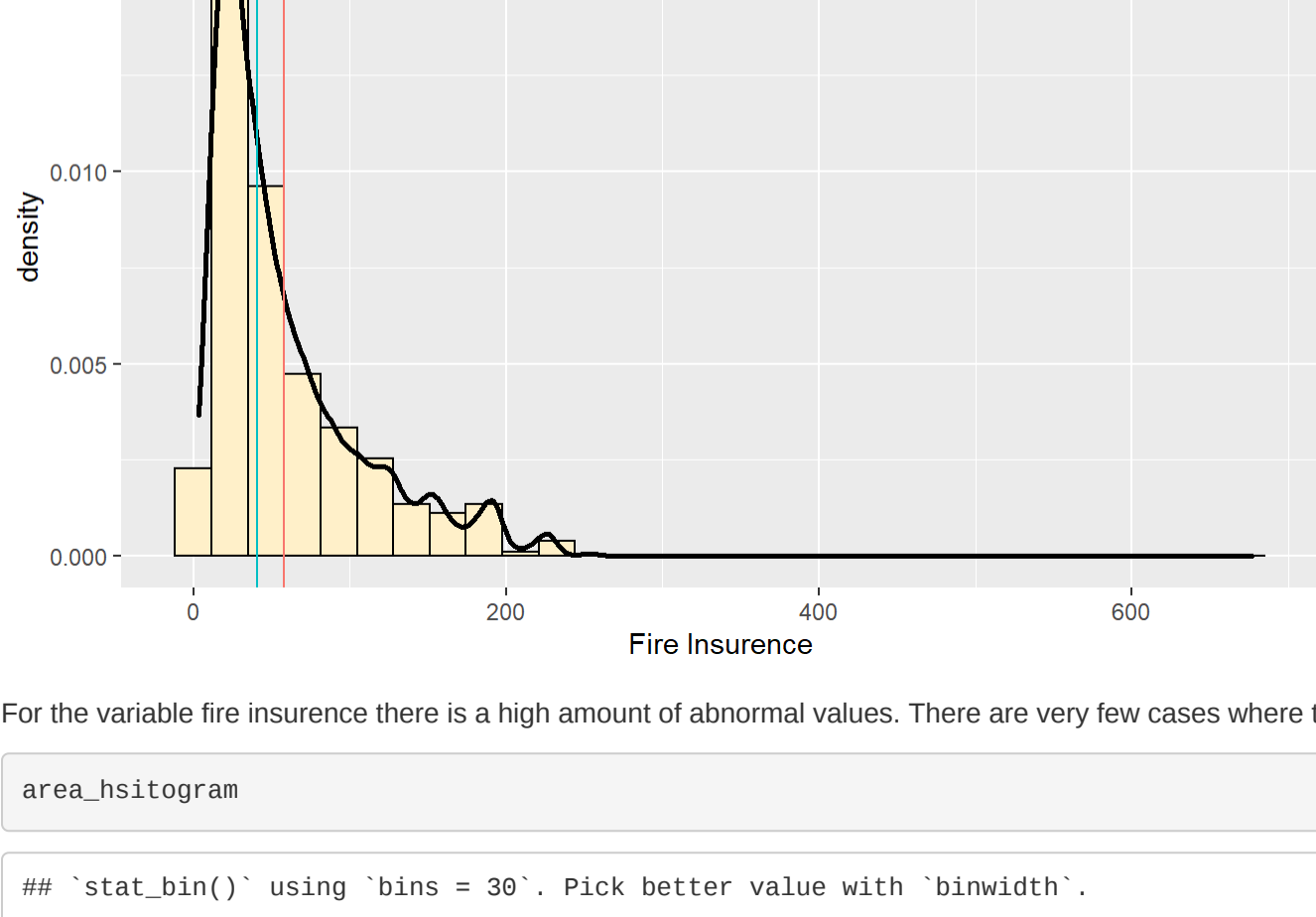
```
area_scatter<-scatter_plot(df$area,lab(x="Area",y="Rent Amount"))
fire_scatter<-scatter_plot(df$fire.insurance,lab(x="Fire Insurance",y="Rent Amount"))
rooms_scatter<-scatter_plot(df$rooms,lab(x="Rooms",y="Rent Amount"))
bathroom_scatter<-scatter_plot(df$bathroom,lab(x="Bathroom",y="Rent Amount"))
parking_scatter<-scatter_plot(df$parking.spaces,lab(x="Parking Spaces",y="Charges"))
```

```
library(gridExtra)
```

Histogram Plots

```
fire_insurance_histogram
```

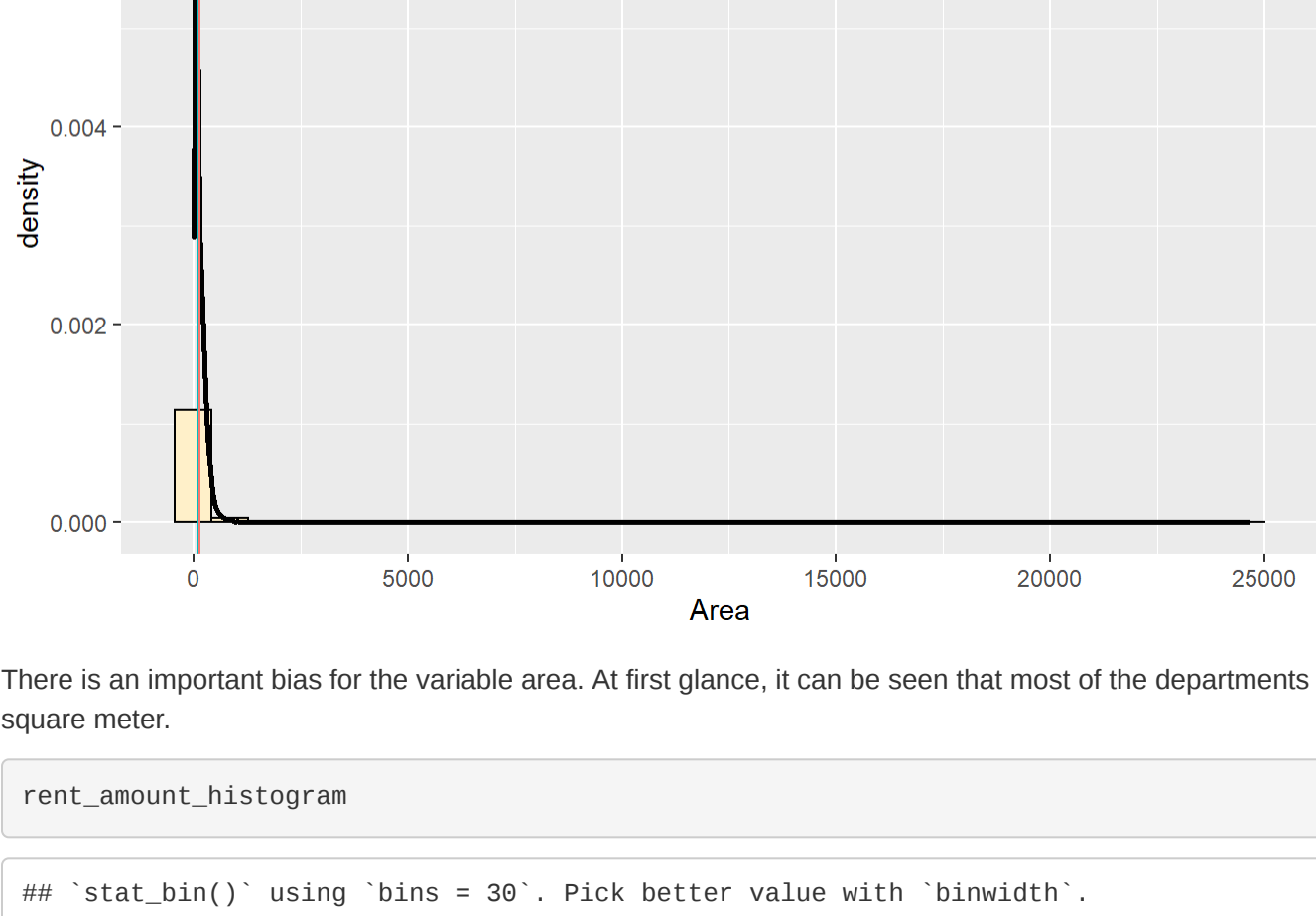
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



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
```

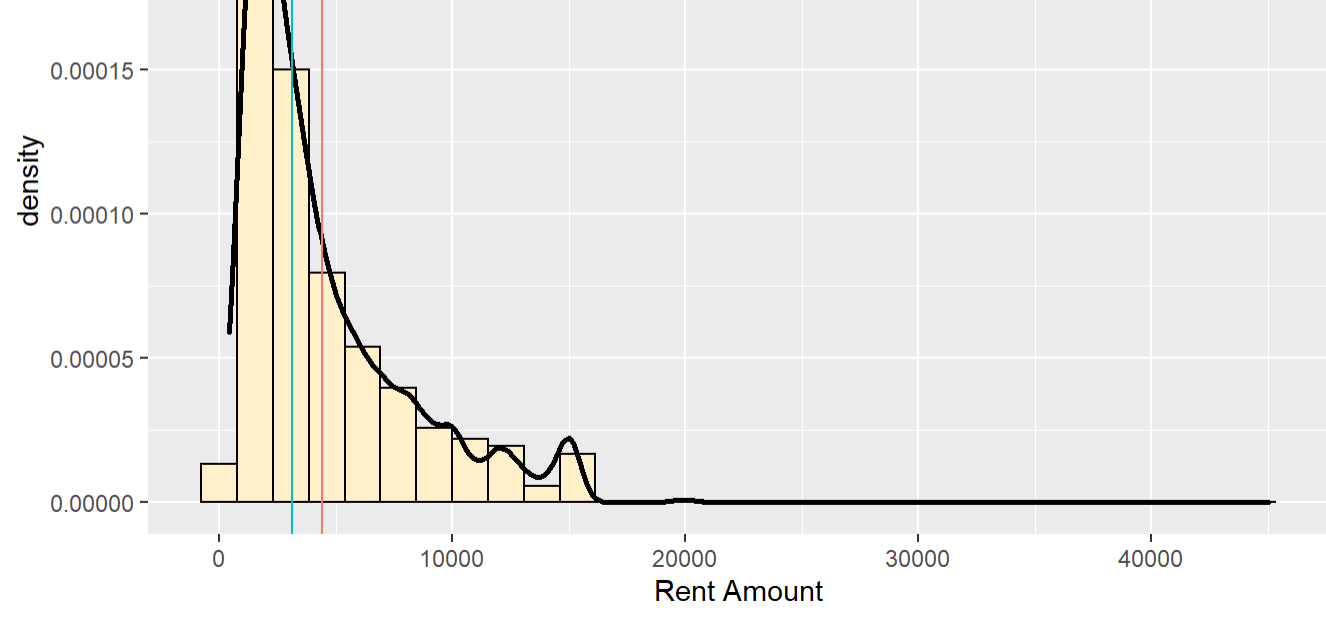
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



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
```

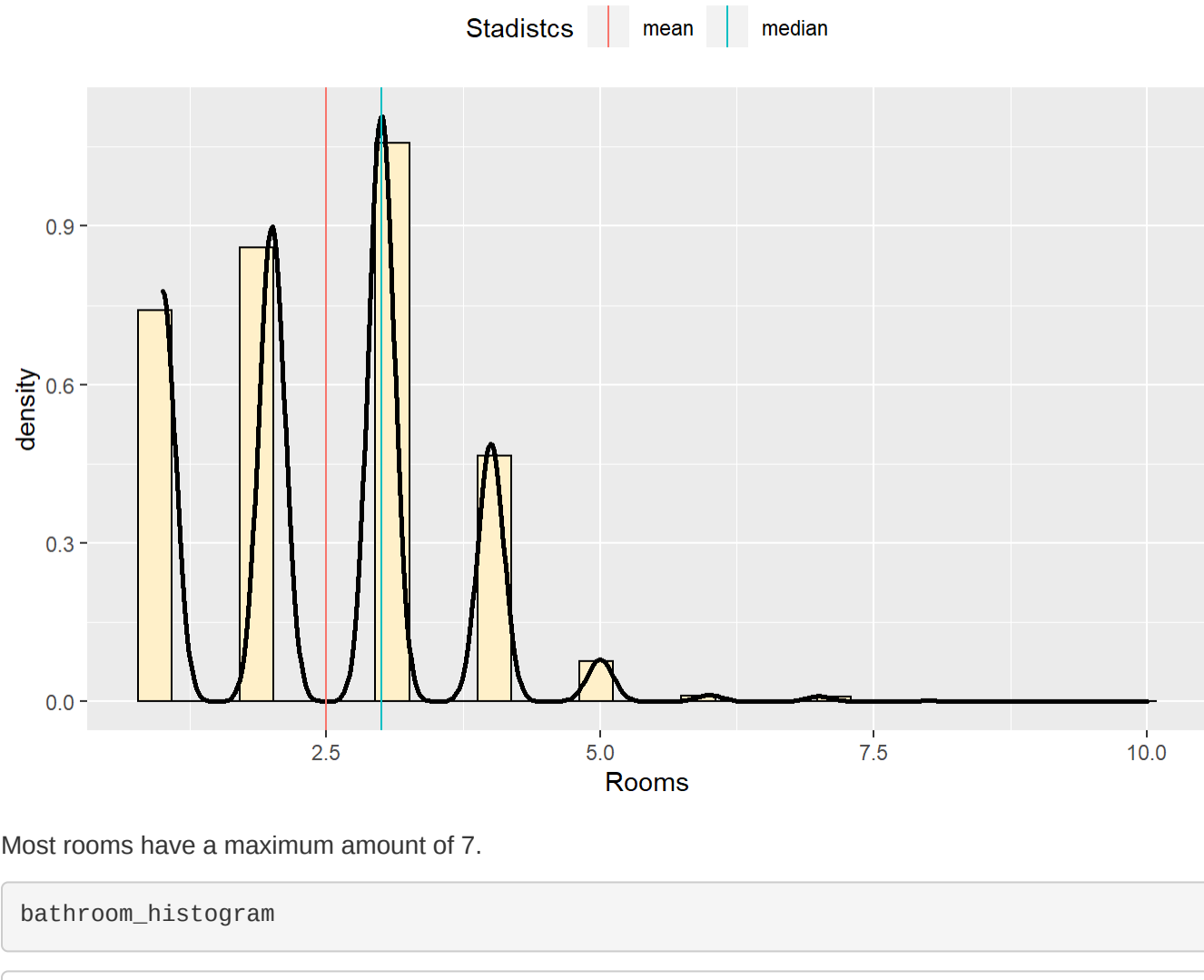
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



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

```
rooms_histogram
```

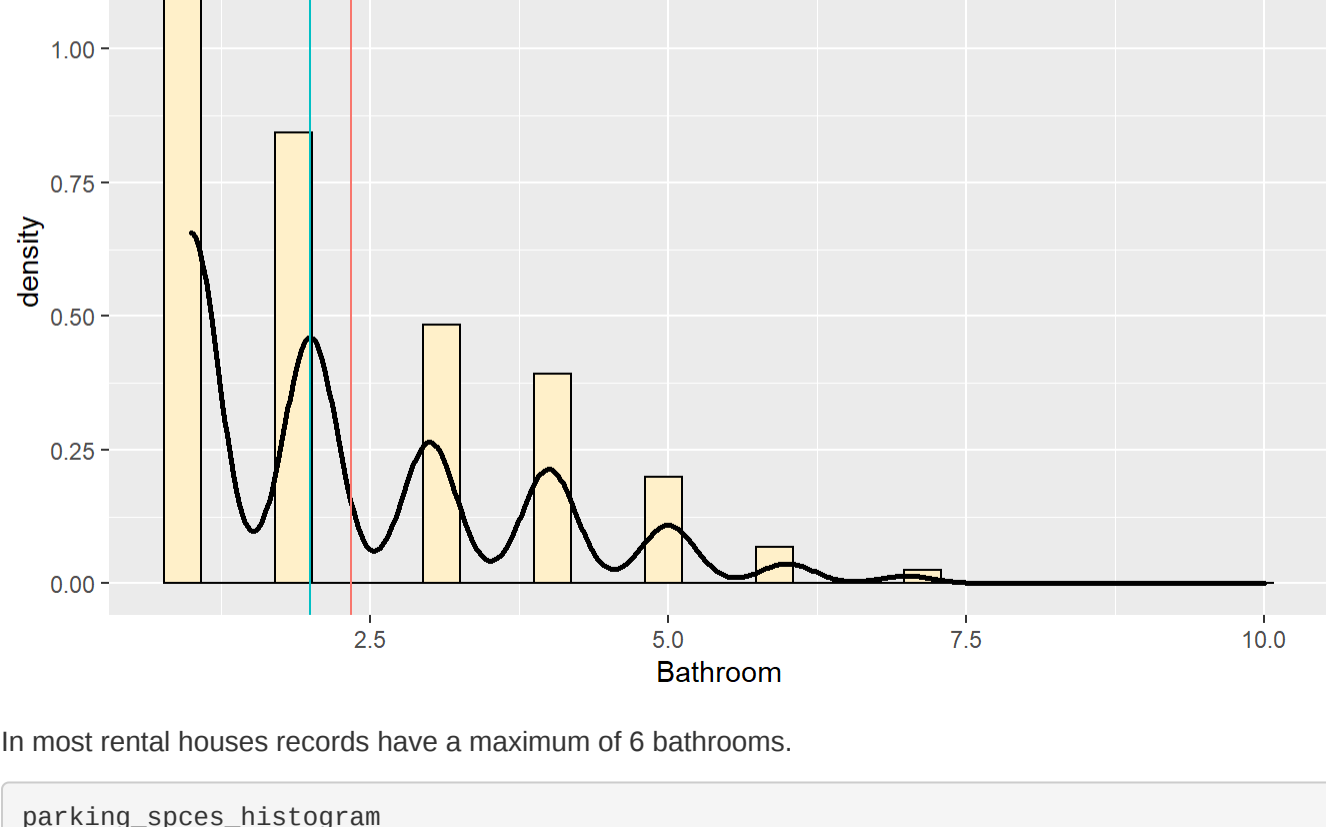
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



Most rooms have a maximum amount of 7.

```
bathroom_histogram
```

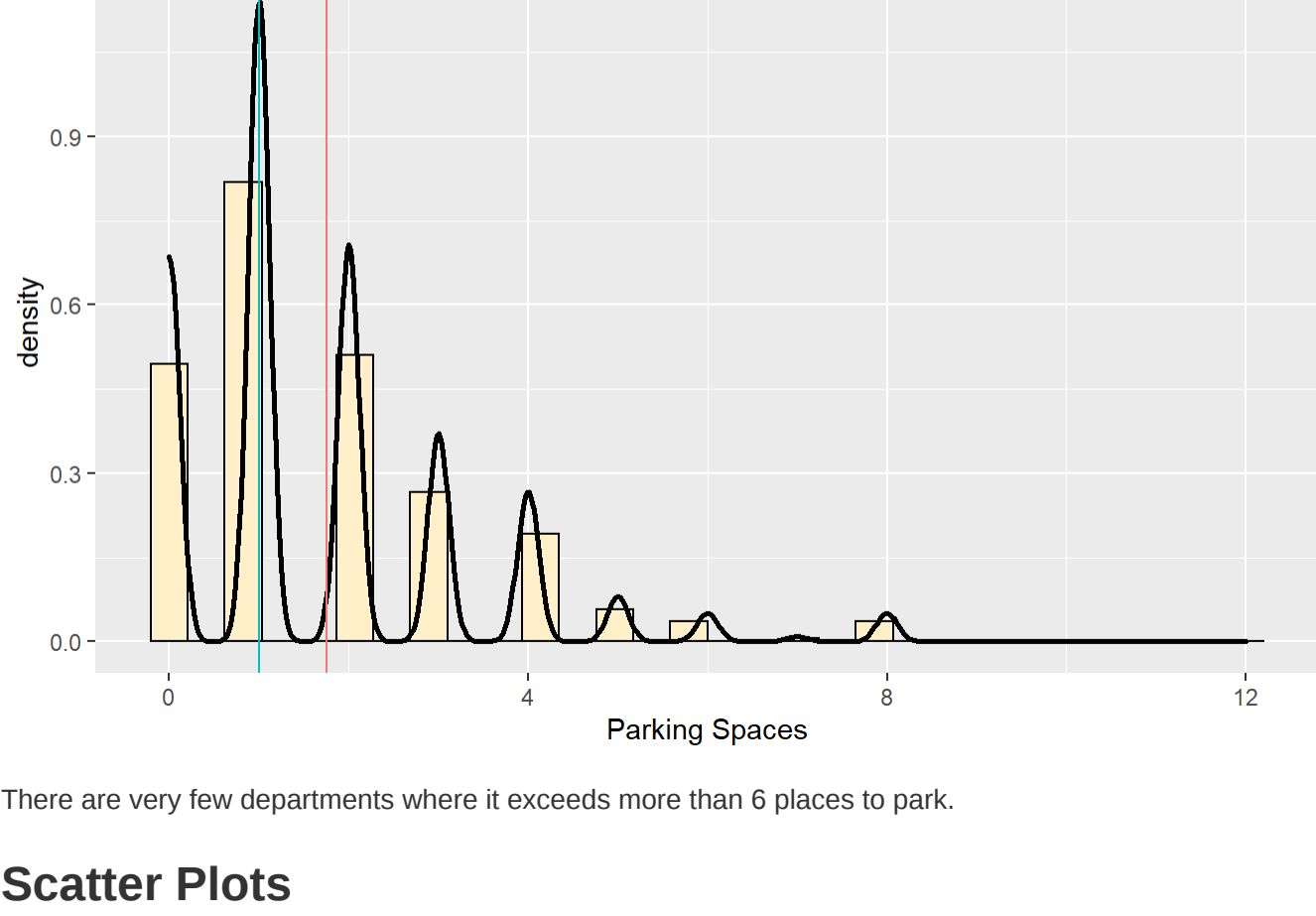
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



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

```
parking_spaces_histogram
```

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



There are very few departments where it exceeds more than 6 places to park.

Scatter Plots

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

```
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
```



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.

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")
```

	area	rooms	bathroom	parking-spaces	rent.amount	fire.insurance	furniture_furnished	furniture_not_furnished
area	0.02	0.09	-0.02	0	-0.16	-0.13	-1	1
rooms	-0.02	-0.09	0.02	0	0.16	0.13	1	-1
bathroom	0.25	0.56	0.66	0.6	0.99	1	0.13	-0.13
parking-spaces	0.24	0.53	0.66	0.58	1	0.99	0.16	-0.16
rent.amount	0.25	0.62	0.69	1	0.58	0.6	0	0
fire.insurance	0.28	0.75	1	0.69	0.66	0.66	0.02	-0.02
furniture_furnished	0.27	1	0.75	0.62	0.53	0.56	-0.09	0.09
furniture_not_furnished	1	0.27	0.28	0.25	0.24	0.25	-0.02	0.02

```
write.csv(df,"rent-amount-brazil.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.

Variables such as the size of the apartment area, the number of bedrooms, number of bathrooms and fire insurance. It makes all the sense in the world for it to increase prices, since these qualities increase the size of houses. The higher the cost of the apartment, the higher the price of fire insurance, since you will have to cover more costs due to the proportion of the apartment.