

Feature engineering

Load Libraries

```
library(tidyverse) # common libraries
```

Load Data

```
df<-read.csv("rent-amount-brazil.csv")

head(df)

  area  rooms  bathroom  parking_spaces  furniture  rent.amount  fire.insurance
1   240     3         3             4    furnished      8000             121
2    64     2         1             1    not furnished      820             11
3   443     5         5             4    furnished      7000             89
4    73     2         2             1    not furnished     1250             16
5    19     1         1             0    not furnished     1200             16
6    13     1         1             0    not furnished     2200             28
6 rows
```

In EDA we conclude that the data set has too many outliers, which is why they require a type of processing.

I can think of 3 ways to deal with the situation.

- Remove outliers.
- Replace abnormal values by statistical measures.
- Perform a logarithmic transformation
- Calculate an upper range and based on the range replace outliers with normal values.

The first and second options are ruled out, since the data set has very few observations. And if we replace the outliers with statistical measures such as the mean or the median, we would be drying the data set, we would be affecting the distribution of the data.

We can combine the 3 and 4 option, we can establish an **upper range**. Based on the interval we create a criterion, we make a sample of random values that oscillate in those ranges. To be able to **replace them** and then perform a **logarithmic transformation** to smooth the data, improving their distribution.

```
attach(df)

upper_limit<-function(x, limit){
  mean<-mean(x) # calculate mean
  sd<-sd(x) # calculate standard division
  mean+limit*sd # calculate upper range
}

calculate_upper_limit<-function(feature){
  limit<-c(2,2.5,3,3.5,4)
  for(index in 1:5){
    print(paste("Upper Limit", limit[index]))
    print(upper_limit(feature, limit[index]))
    print("=====")
  }
}
```

Upper Limits

```
calculate_upper_limit(feature = rent.amount)

## [1] "Upper Limit 2"
## [1] 11549.38
## [1] "=====
## [1] "Upper Limit 2.5"
## [1] 13397.52
## [1] "=====
## [1] "Upper Limit 3"
## [1] 15325.86
## [1] "=====
## [1] "Upper Limit 3.5"
## [1] 16954.19
## [1] "=====
## [1] "Upper Limit 4"
## [1] 18702.52
## [1] "=====
```

The best upper limit is 3 which we can round up to R \$15,000. Since there are more rental houses that are around those prices as we can see in the histogram.

```
Fire Insurance

calculate_upper_limit(feature = fire.insurance)

## [1] "Upper Limit 2"
## [1] 156.8332
## [1] "=====
## [1] "Upper Limit 2.5"
## [1] 181.4879
## [1] "=====
## [1] "Upper Limit 3"
## [1] 286.1446
## [1] "=====
## [1] "Upper Limit 3.5"
## [1] 236.8632
## [1] "=====
## [1] "Upper Limit 4"
## [1] 255.458
## [1] "=====
```

With an upper limit of 4 it is good. With an upper limit of 4 it is good. We can round it up to R \$250 since there are very few cases where the price of fire insurance exceeds that amount.

```
Area

calculate_upper_limit(feature = area)

## [1] "Upper Limit 2"
## [1] 992.2629
## [1] "=====
## [1] "Upper Limit 2.5"
## [1] 1096.043
## [1] "=====
## [1] "Upper Limit 3"
## [1] 1277.822
## [1] "=====
## [1] "Upper Limit 3.5"
## [1] 1465.692
## [1] "=====
## [1] "Upper Limit 4"
## [1] 1653.382
## [1] "=====
```

With an upper limit of 2.5 it is a good point. Since most of the departments do not exceed 1000 square meters.

```
replace_outliers<-function(data, normal_sample, upper_limit){
  return(ifelse(data>upper_limit, normal_sample, data))
}
```

Normal Values Sample for Upper Limit

```
set.seed(2018) # we define random seed, so that the numbers generated by the sample do not vary.
sample_rent<-sample(14500:15800, size=25, replace = T)
sample_fire_insurance<-sample(240:250, size=25, replace = T)
sample_area<-sample(980:1000, size=25, replace = T)
```

```
normal_values_rent<-replace_outliers(data=df$rent.amount, sample_rent, upper_limit = 15000)
normal_fire_insurance<-replace_outliers(data=df$fire.insurance, sample_fire_insurance, upper_limit = 250)
normal_area<-replace_outliers(data=area, sample_area, upper_limit = 1000)
```

Replace Outliers Values

```
df<- df %>%
  mutate(rent.amount=normal_values_rent) %>%
  mutate(fire.insurance=normal_fire_insurance) %>%
  mutate(area=normal_area)
```

With the mutate function we make modifications. We replace outliers with normal values, using a random sample.

We check the changes

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

##      rent.amount      area      fire.insurance
##  Min.   : 420  Min.   : 18.0  Min.   : 3.00
##  1st Qu.: 1000  1st Qu.: 50.0  1st Qu.: 23.00
##  Median : 3313  Median : 188.0  Median : 41.00
##  Mean   : 4383  Mean   : 145.2  Mean   : 58.11
##  3rd Qu.: 5952  3rd Qu.: 248.0  3rd Qu.: 77.00
##  Max.   : 15000  Max.   : 1088.0  Max.   : 250.00

df<- df %>% mutate(fire.insurance=ifelse(fire.insurance>3, 5, fire.insurance))
```

We replaced the minimum price value of fire insurance by five real Brazilian.

```
df %>% select(fire.insurance) %>% summary()

##      fire.insurance
##  Min.   : 4.00
##  1st Qu.: 23.00
##  Median : 41.00
##  Mean   : 58.11
##  3rd Qu.: 77.00
##  Max.   : 250.00
```

logarithmic transformation

```
df<- df %>%
  mutate(area_log=log(area)) %>%
  mutate(fire.insurance_log=log(fire.insurance)) %>%
  mutate(rent.amount_log=log(rent.amount))
```

We perform logarithmic transformation to improve the distribution of continuous data.

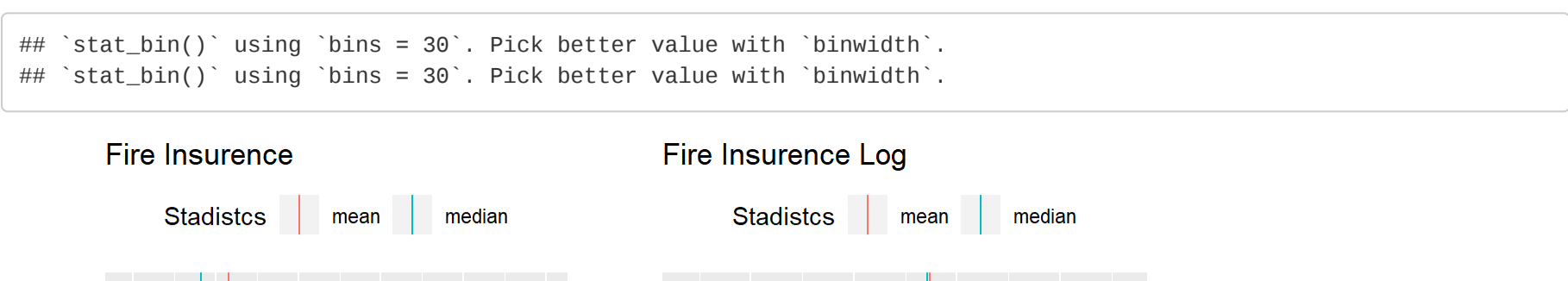
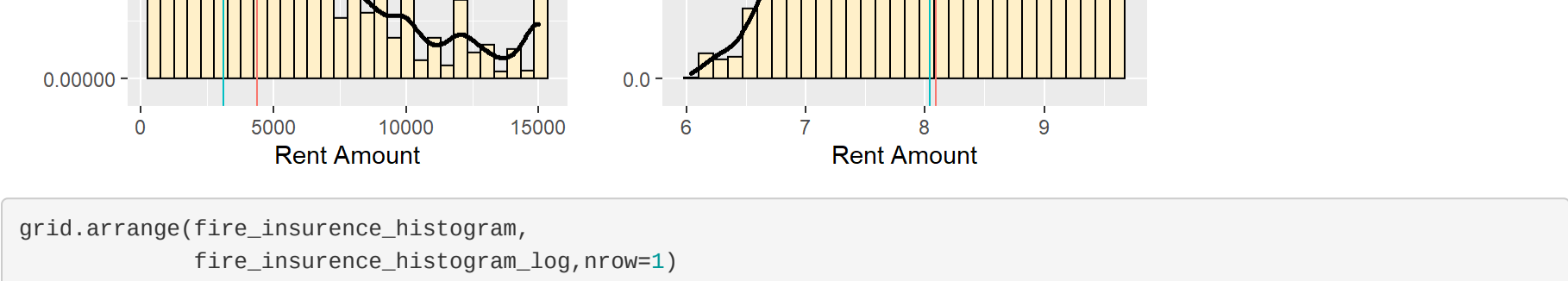
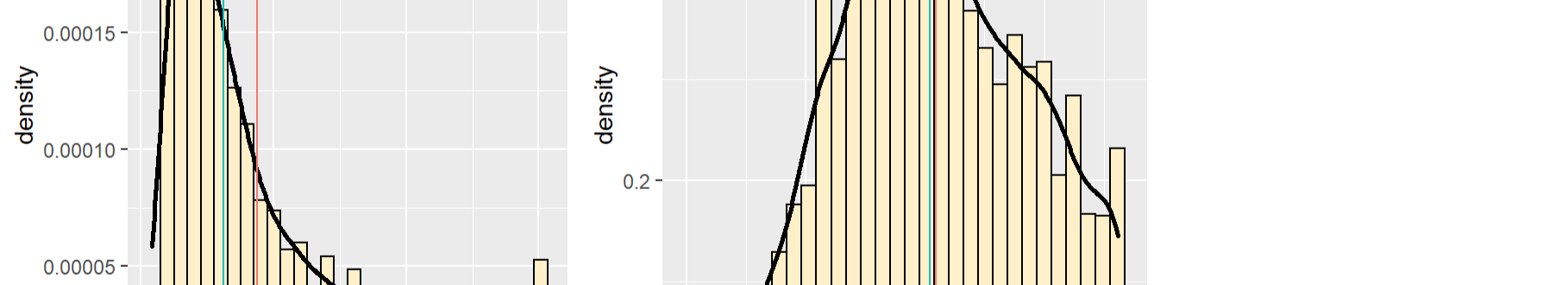
```
histogram<-function(x,...){
  df %>%
    ggplot(aes(x=x,y=-.density...)) +
    geom_histogram(color="black", fill="transparent") +
    geom_density(color="black", lwd=1) +
    geom_vline(aes(x=intercept=mean(x), color="mean")) +
    geom_vline(aes(x=intercept=median(x), color="median")) +
    labs(coi="Statistics")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

library(gridExtra)

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$fire.insurance, labs(x="Fire Insurance", title = "Fire Insurance"))
rent_amount_histogram_log<-histogram(df$rent.amount_log, labs(x="Rent Amount", title = "Rent Amount Log"))
area_histogram_log<-histogram(df$area_log, labs(x="Area", title = "Area Log"))
fire_insurance_histogram_log<-histogram(df$fire.insurance_log, labs(x="Fire Insurance", title = "Fire Insurance Log"))

grid.arrange(rent_amount_histogram,
              rent_amount_histogram_log, nrow=1)

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

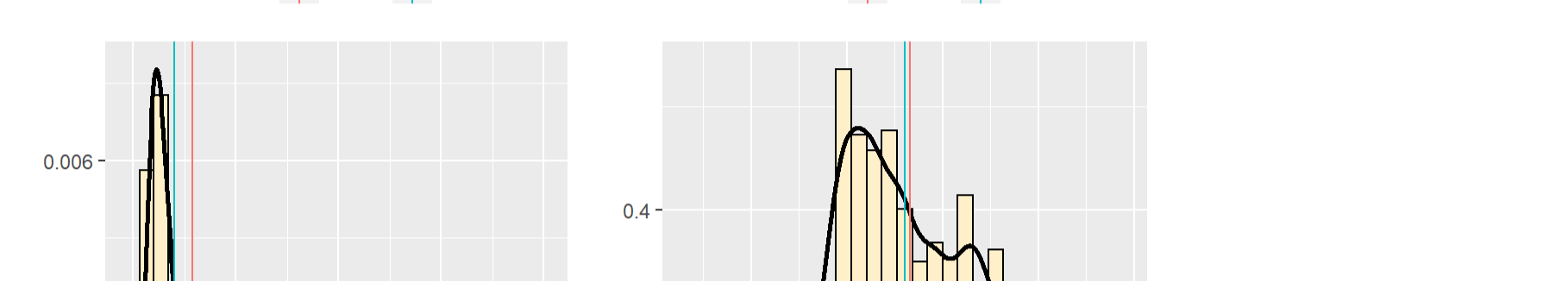
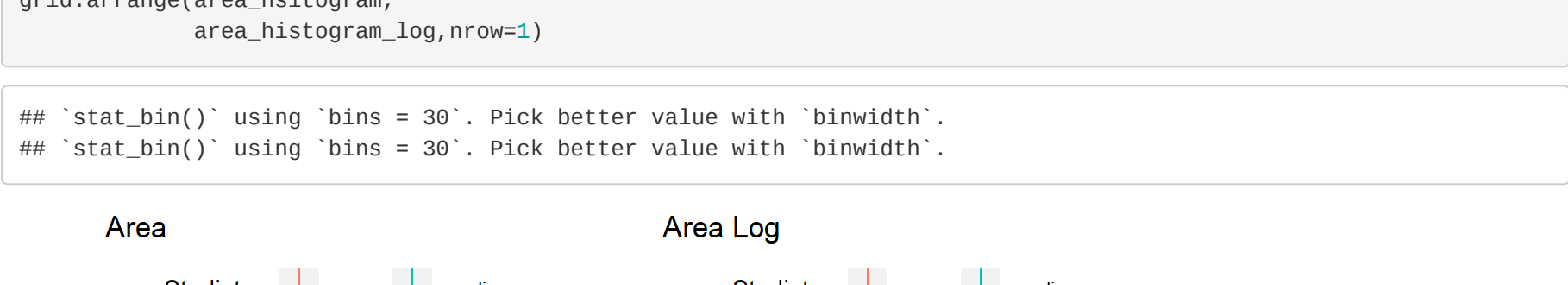
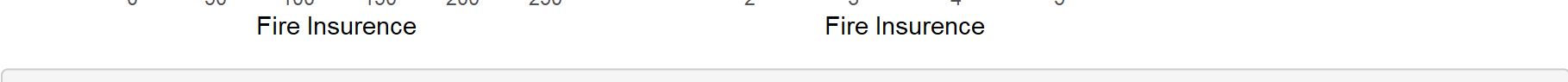


We observed a great improvement in the distribution of the data.

```
rooms_histogram<-histogram(df$rooms, labs(x="Rooms"))
bathroom_histogram<-histogram(df$bathroom, labs(x="Bathroom"))
parking_spaces_histogram<-histogram(df$parking_spaces, labs(x="Parking Spaces"))

rooms_histogram

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



We can replace both variables, the values that are out of a normal range by the amount of 5, since it is a good confidence interval.

```
replace_values<-function(x){
  ifelse(x>5, 5, x)
}

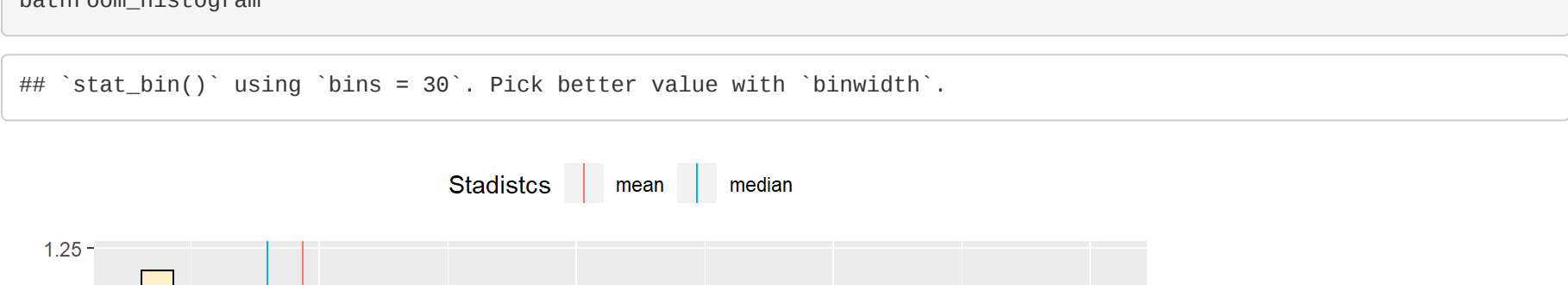
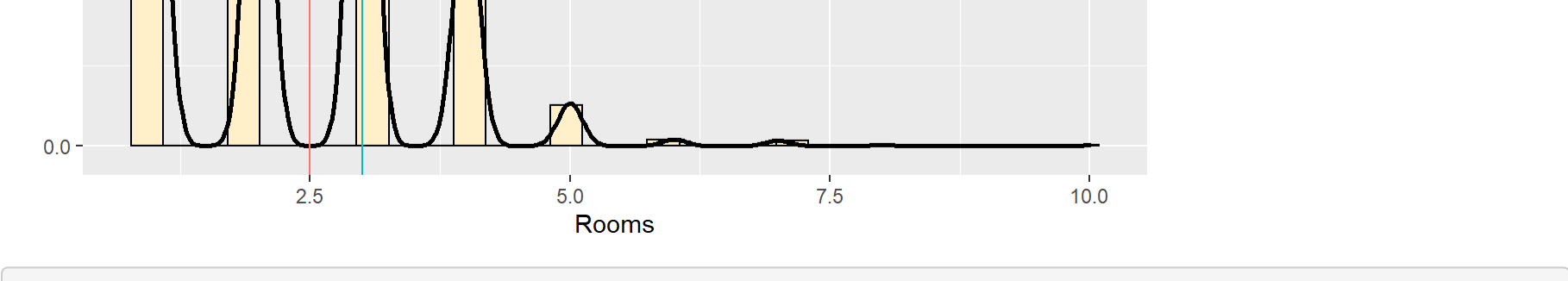
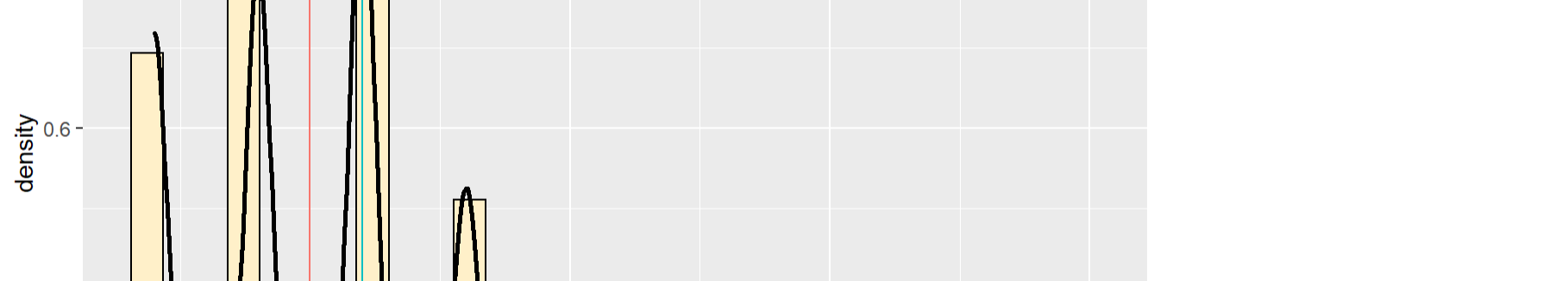
replace_rooms<-mapply(replace_values, rooms)
replace_bathrooms<-mapply(replace_values, bathroom)
replace_spaces<-mapply(replace_values, parking_spaces)
```

```
df<- df %>%
  mutate(replace_rooms=replace_rooms) %>%
  mutate(replace_bathrooms=replace_bathrooms) %>%
  mutate(replace_parking_spaces=replace_parking_spaces)
```

```
rooms_hist_replace<-histogram(df$replace_rooms, labs(x="Rooms", title = "Replace Values"))
bathroom_hist_replace<-histogram(df$replace_bathrooms, labs(x="Bathroom", title = "Replace Values"))
parking_hist_replace<-histogram(df$replace_parking_spaces, labs(x="Parking Spaces", title = "Replace Values"))

grid.arrange(rooms_histogram, rooms_hist_replace, nrow=1)

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



We finish with the cleaning of our data.

We remove old variables

```
df<-df %>%
  select(-rent.amount, -area, -fire.insurance) %>%
  select(-rooms, -bathroom, -parking_spaces)
```

Rename Features

```
df_clear<-df %>%
  rename(fire.insurance=fire.insurance_log,
         area=area_log,
         rent.amount=rent.amount_log) %>%
  rename(bathrooms=replace_bathrooms,
         parking_spaces=replace_parking_spaces,
         rooms=replace_rooms)
```

```
quality_department<-cut(df_clear$rent.amount, breaks = 5, labels = c(1,2,3,4,5))
```

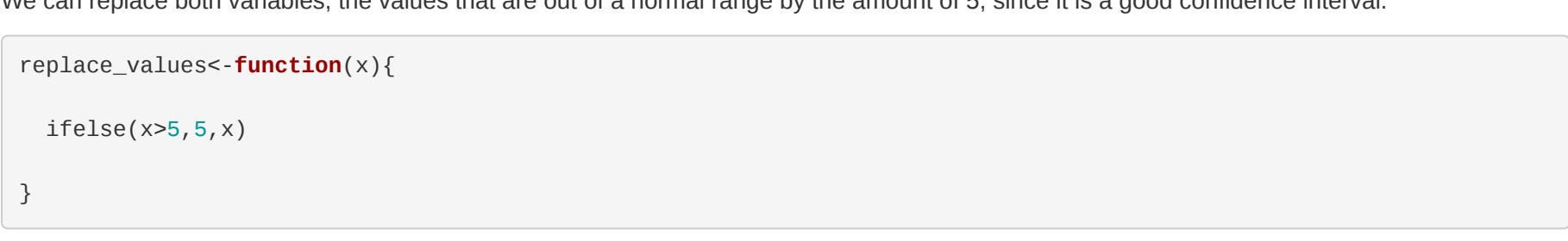
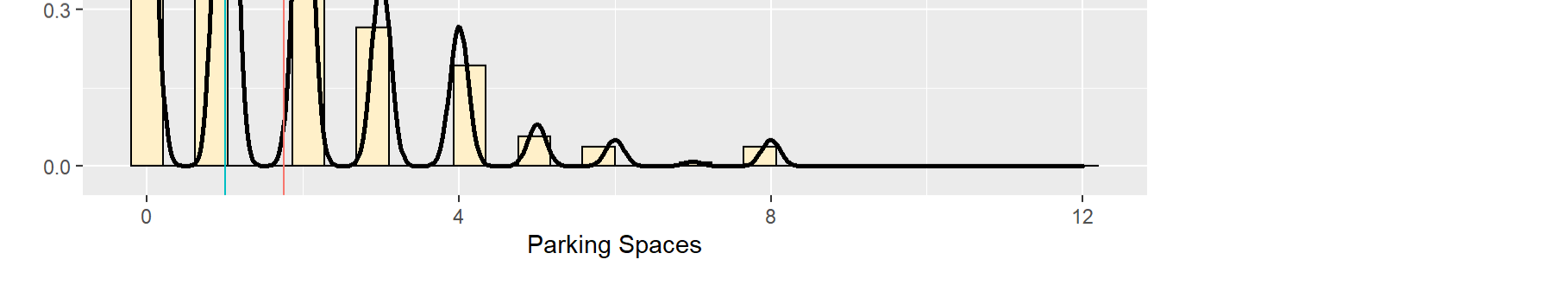
```
df_clear<-df_clear %>% mutate(quality_department=quality_department)
```

```
df_clear<-df_clear %>% mutate(quality_department=as.numeric(quality_department))
```

We can classify the degree of quality of the rental house. Generally the higher the price, the higher the quality of the house. We can group it from a range of 1 to 5, the higher the number, the greater the impact of the department.

```
library(bateExplorer) # correlation matrix

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



Save Dataframe clear

```
write.csv(df_clear, "rent-amount-brazil-clear.csv", row.names = FALSE)
```