

R Notebook

leo datos

```
data<-read.csv("entrena.csv",stringsAsFactors = F)
```

Creo las 3 bases

```
knitr::opts_chunk$set( message = FALSE, results='hide')  
idx <- sample(seq(1, 3), size = nrow(data), replace = TRUE, prob = c(.7, 2, .1))  
train <- data[idx == 1,]  
test <- data[idx == 2,]  
cal <- data[idx == 3,]
```

summary

```
knitr::opts_chunk$set(echo = TRUE, message = FALSE, results = 'show')  
summary(train)
```

```

##      hotel      is_canceled      lead_time      arrival_date_year
## Length:22972      Length:22972      Min.   : 0.00      Min.   :2015
## Class :character      Class :character      1st Qu.: 14.00      1st Qu.:2016
## Mode  :character      Mode  :character      Median  : 57.00      Median :2016
##                                     Mean   : 96.12      Mean   :2016
##                                     3rd Qu.:144.00      3rd Qu.:2016
##                                     Max.   :737.00      Max.   :2017
##
## arrival_date_month arrival_date_week_number arrival_date_day_of_month
## Length:22972      Min.   : 1.00      Min.   : 1.00
## Class :character      1st Qu.:13.00      1st Qu.: 8.00
## Mode  :character      Median :30.00      Median :16.00
##                                     Mean   :28.13      Mean   :15.79
##                                     3rd Qu.:41.00      3rd Qu.:23.00
##                                     Max.   :53.00      Max.   :31.00
##
## stays_in_weekend_nights stays_in_week_nights      adults
## Min.   : 0.0000      Min.   : 0.000      Min.   : 0.000
## 1st Qu.: 0.0000      1st Qu.: 1.000      1st Qu.: 2.000
## Median : 1.0000      Median : 2.000      Median : 2.000
## Mean    : 0.8931      Mean    : 2.459      Mean    : 1.835
## 3rd Qu.: 2.0000      3rd Qu.: 3.000      3rd Qu.: 2.000
## Max.    :19.0000      Max.    :50.000      Max.    :55.000
##
##      children      babies      meal      country
## Min.   :0.00000      Min.   :0.000000      Length:22972      Length:22972
## 1st Qu.:0.00000      1st Qu.:0.000000      Class :character      Class :character
## Median :0.00000      Median :0.000000      Mode  :character      Mode  :character
## Mean    :0.08755      Mean    :0.007487
## 3rd Qu.:0.00000      3rd Qu.:0.000000
## Max.    :3.00000      Max.    :2.000000
## NA's    :1
## market_segment      distribution_channel is_repeated_guest
## Length:22972      Length:22972      Min.   :0.00000
## Class :character      Class :character      1st Qu.:0.00000
## Mode  :character      Mode  :character      Median :0.00000
##                                     Mean    :0.03334
##                                     3rd Qu.:0.00000
##                                     Max.    :1.00000
##
## previous_cancellations previous_bookings_not_canceled reserved_room_type
## Min.   : 0.0000      Min.   : 0.0000      Length:22972
## 1st Qu.: 0.0000      1st Qu.: 0.0000      Class :character
## Median : 0.0000      Median : 0.0000      Mode  :character
## Mean    : 0.1133      Mean    : 0.1371
## 3rd Qu.: 0.0000      3rd Qu.: 0.0000
## Max.    :26.0000      Max.    :61.0000
##
## assigned_room_type booking_changes deposit_type      agent
## Length:22972      Min.   : 0.000      Length:22972      Length:22972
## Class :character      1st Qu.: 0.000      Class :character      Class :character
## Mode  :character      Median : 0.000      Mode  :character      Mode  :character
##                                     Mean    : 0.214
##                                     3rd Qu.: 0.000
##                                     Max.    :21.000
##
##      company      days_in_waiting_list customer_type      adr

```

```
## Length:22972      Min.   : 0.000      Length:22972      Min.   : 0.00
## Class :character  1st Qu.: 0.000      Class :character  1st Qu.: 64.00
## Mode  :character  Median : 0.000      Mode  :character  Median : 87.00
##                               Mean  : 3.072      Mean  : 92.89
##                               3rd Qu.: 0.000      3rd Qu.:114.75
##                               Max.   :391.000      Max.   :508.00
##
## required_car_parking_spaces total_of_special_requests
## Min.   :0.00000      Min.   :0.0000
## 1st Qu.:0.00000      1st Qu.:0.0000
## Median :0.00000      Median :0.0000
## Mean   :0.06599      Mean   :0.5257
## 3rd Qu.:0.00000      3rd Qu.:1.0000
## Max.   :2.00000      Max.   :5.0000
##
```

proporcion de cancelados

```
a <- table(train$is_canceled)
prop <- prop.table(a)
prop
```

```
##
##   cancelado no_cancelado
## 0.3618753   0.6381247
```

```
sub_cancelados <- subset(train, is_canceled == "cancelado")
```

#funcion para graficar las variables originales

```
analisis_cancelaciones <- function(datos,nombre) {
  cancel <- table(datos)
  par(cex=0.5) #control size of labels
  g_1 <- barplot(cancel, main = nombre)
  return(g_1)
}
```

#funcion para graficar las variables ordenadas

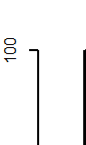
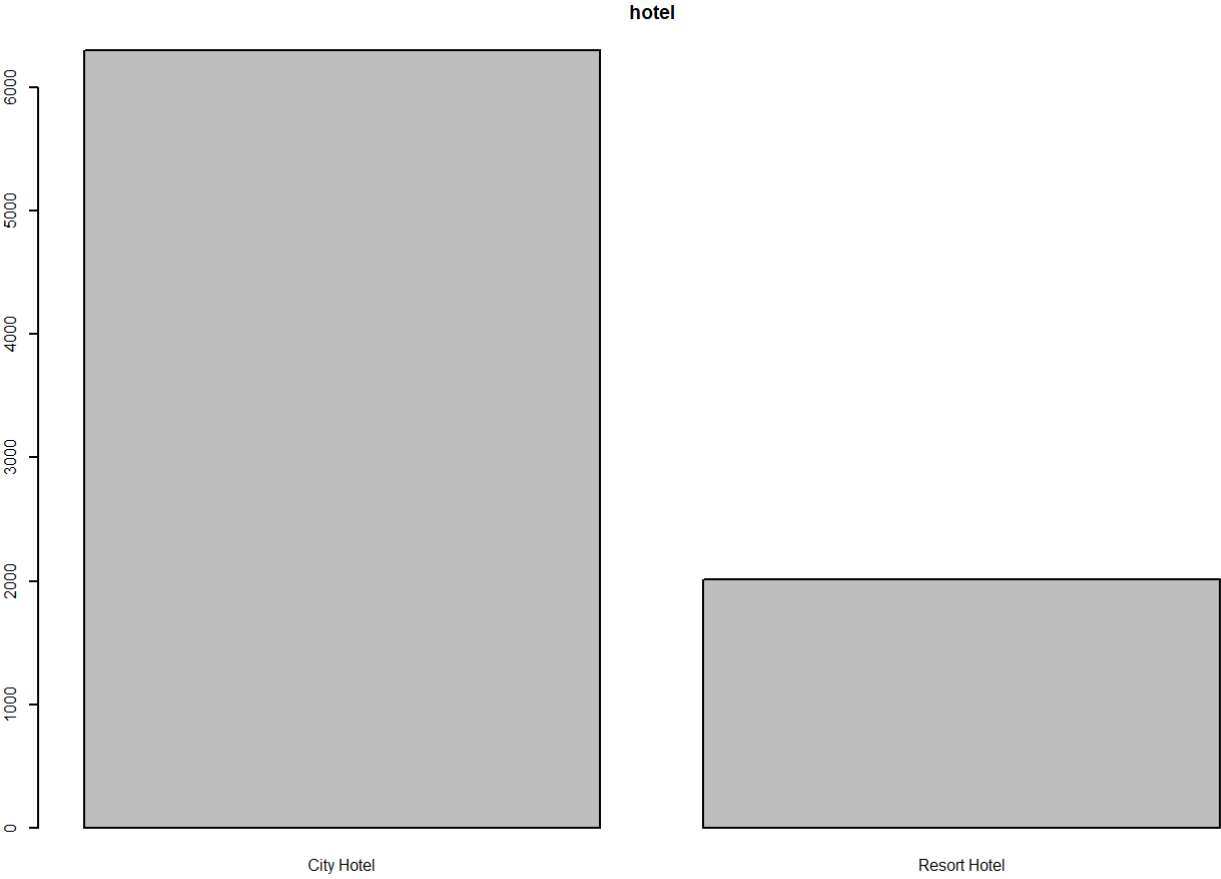
```
analisis_cancelaciones_ord <- function(datos,nombre) {
  x <- table(datos)
  cancel <- x[order(x,decreasing = TRUE)]
  par(cex=0.5) #control size of labels
  g_1 <- barplot(cancel, main = nombre)
  return(g_1)
}
```

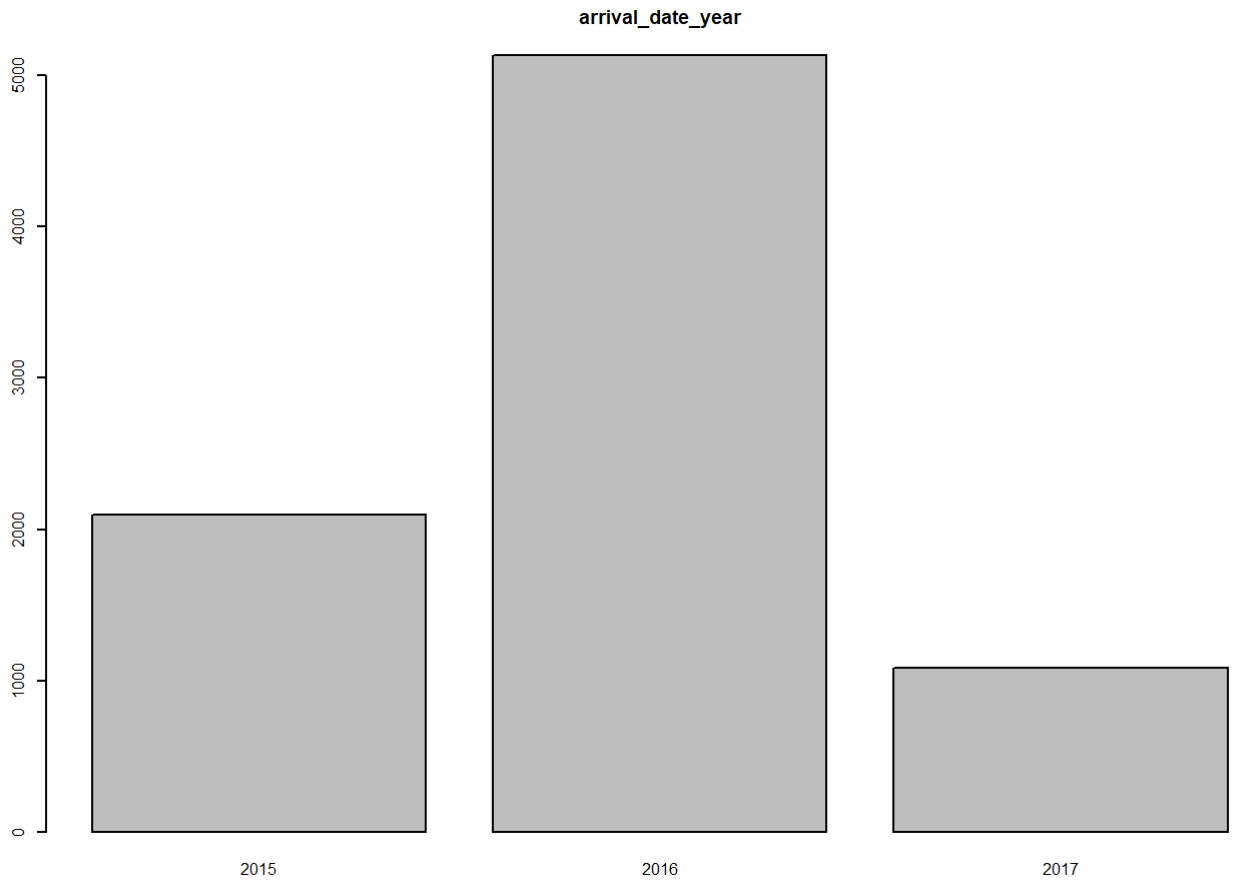
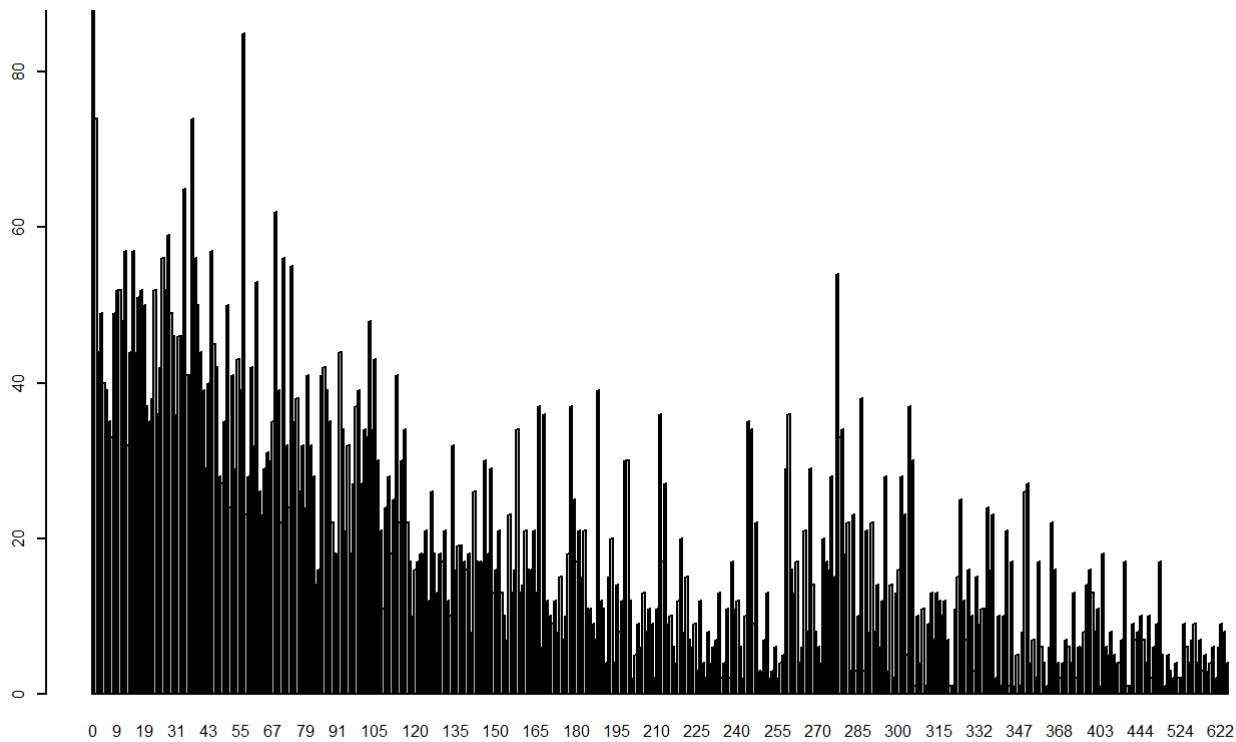
analisis de todas las variables con el subconjunto de canceladas

```
ncol(sub_cancelados)
```

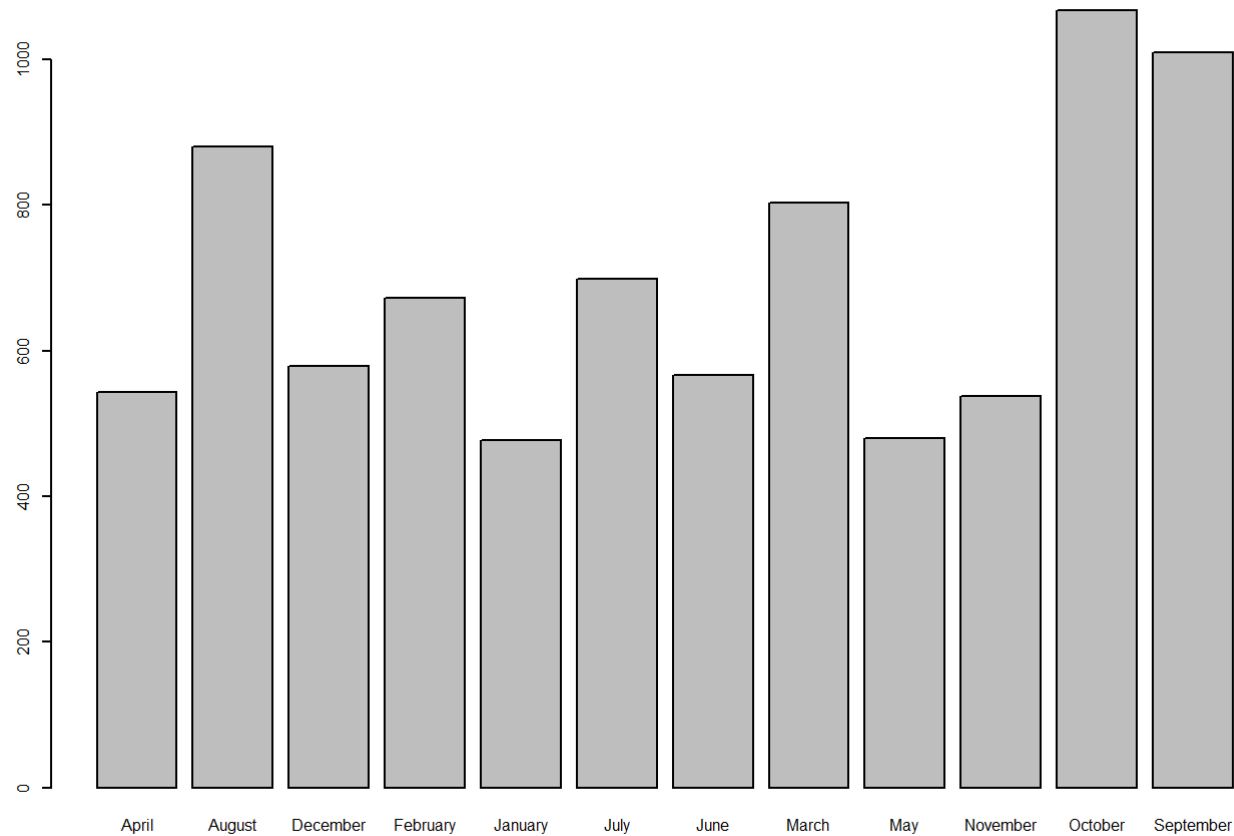
```
## [1] 30
```

```
nombres <- names(sub_cancelados)
for (i in c(1:ncol(sub_cancelados))){
  analisis_cancelaciones(sub_cancelados[,i],nombres[i])
}
```

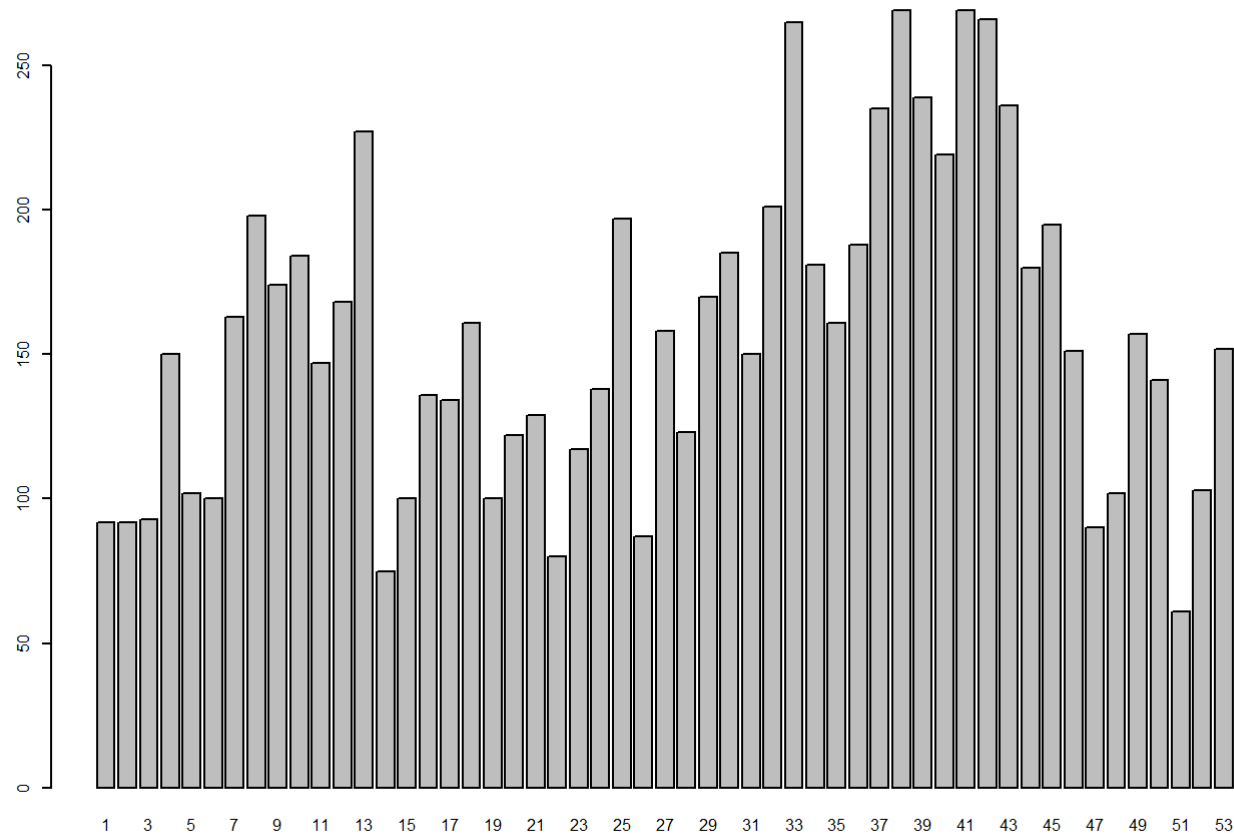




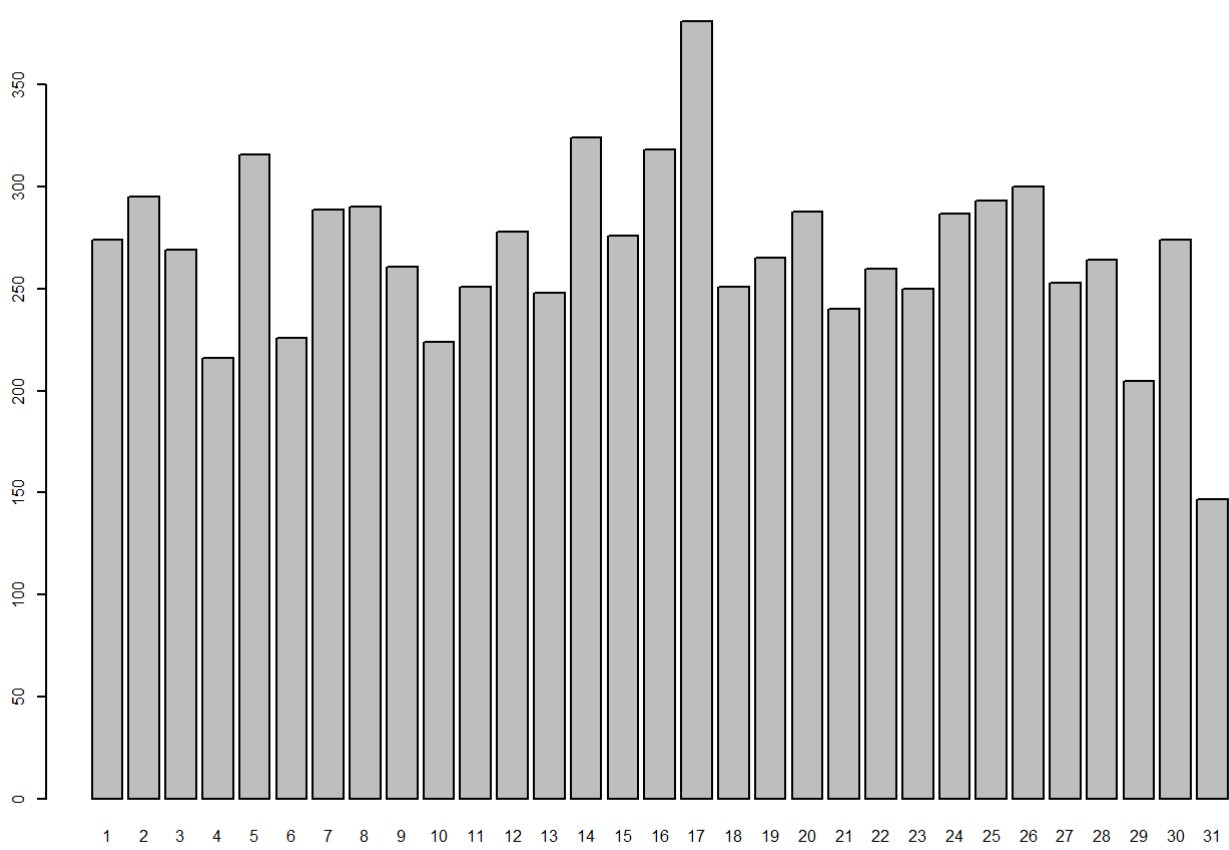
arrival_date_month



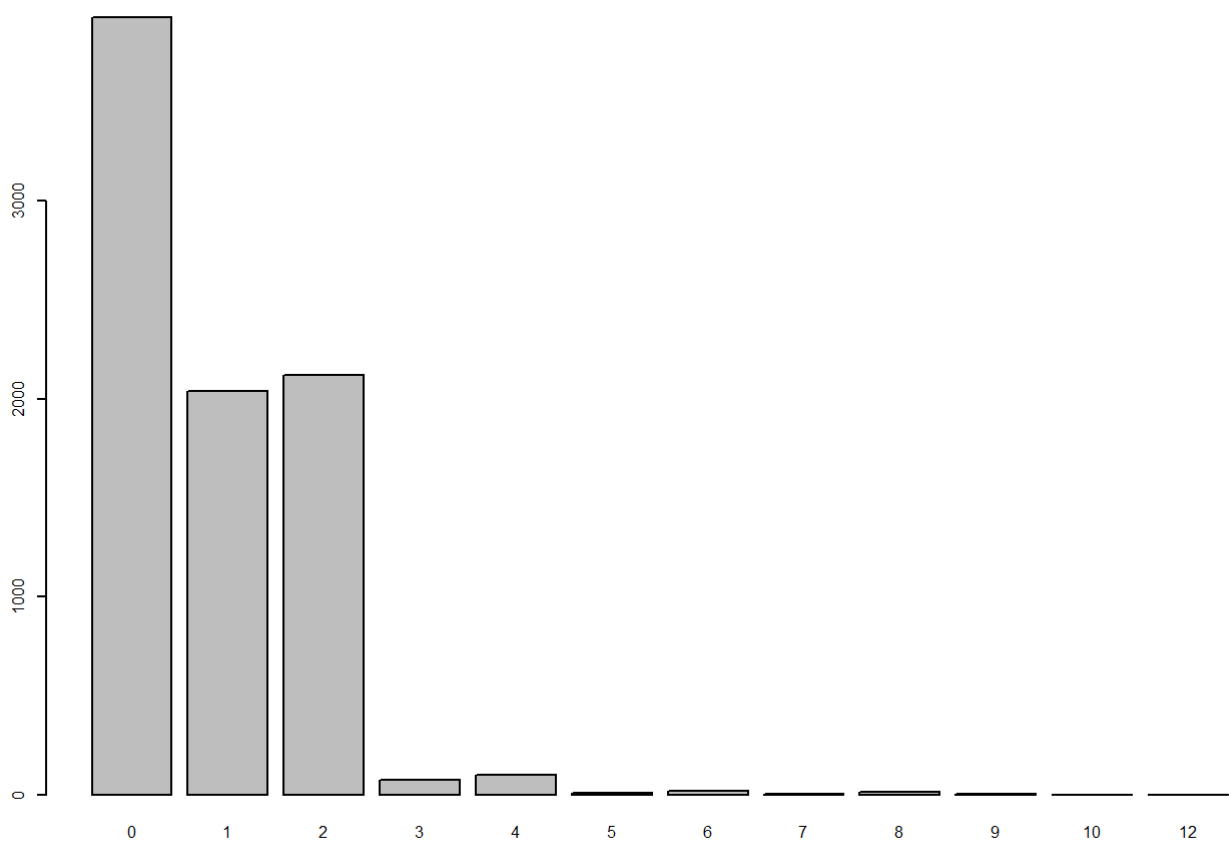
arrival_date_week_number



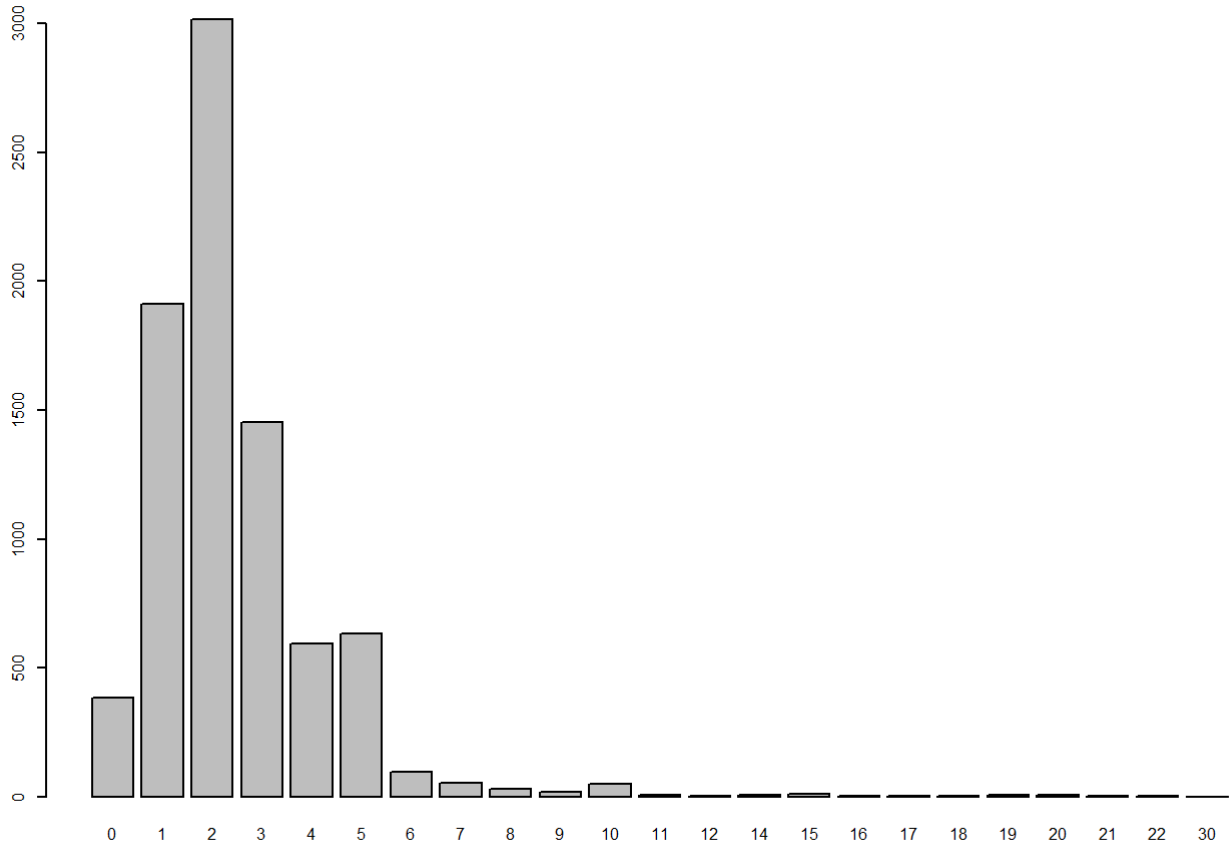
arrival_date_day_of_month



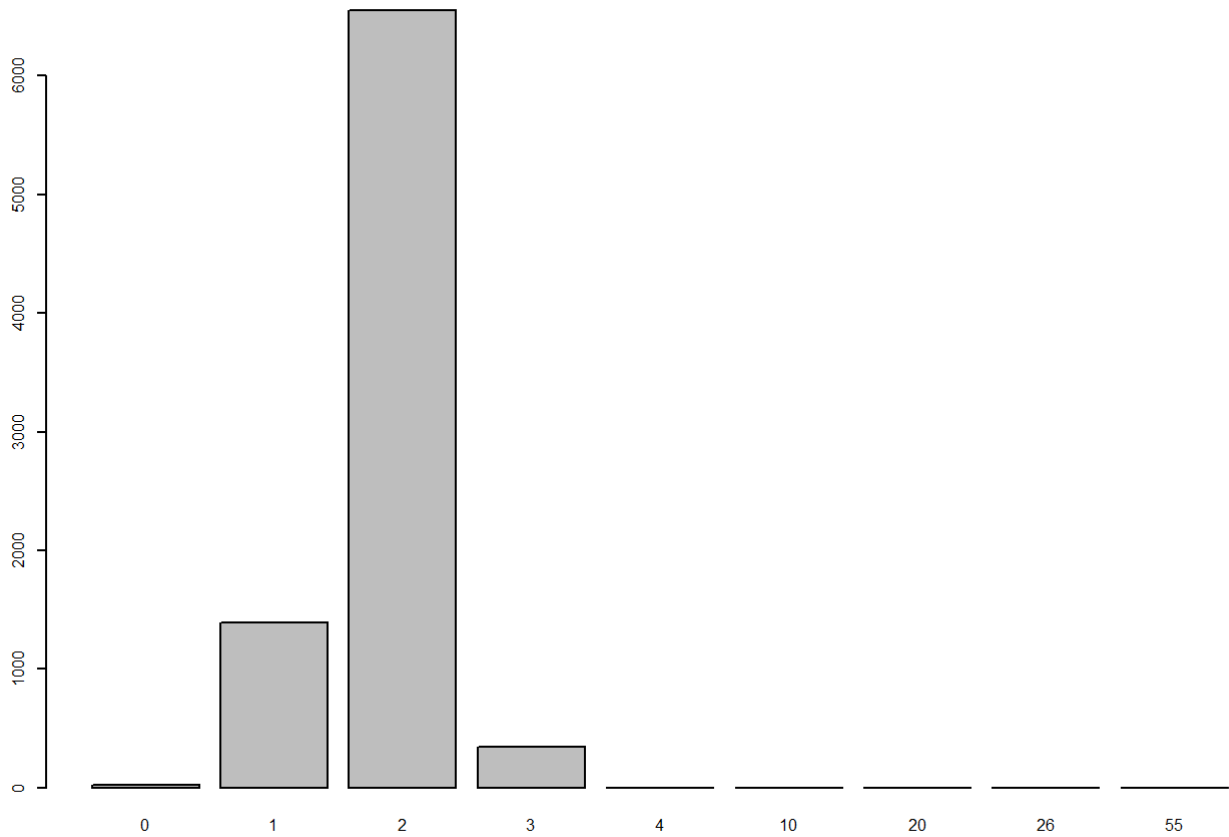
stays_in_weekend_nights



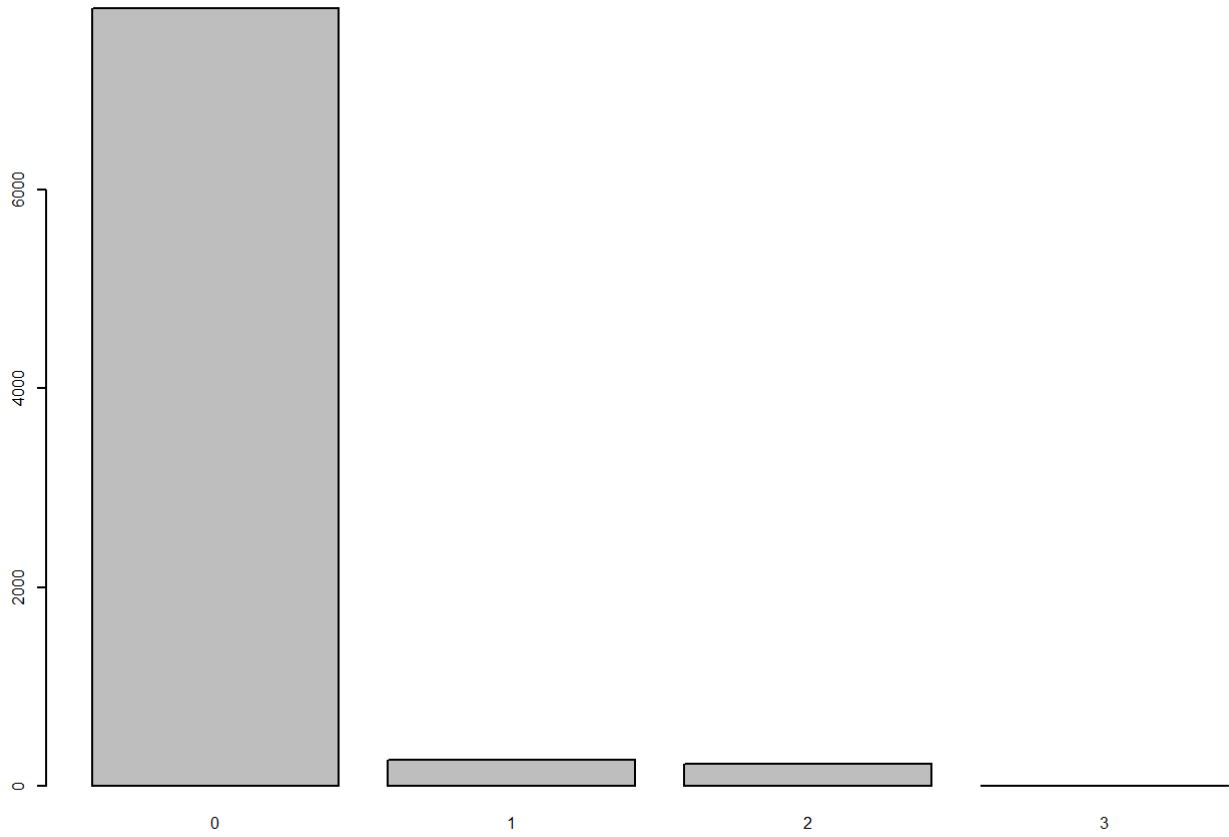
stays_in_week_nights



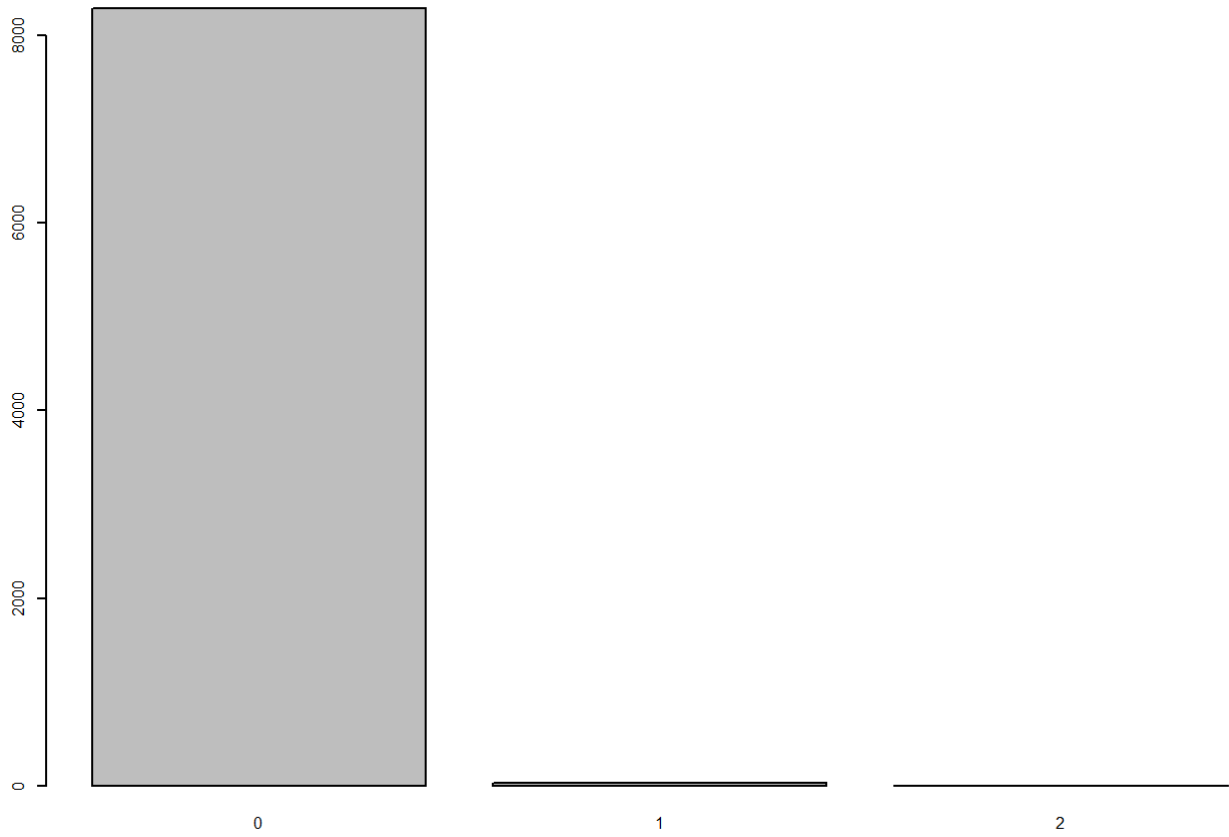
adults



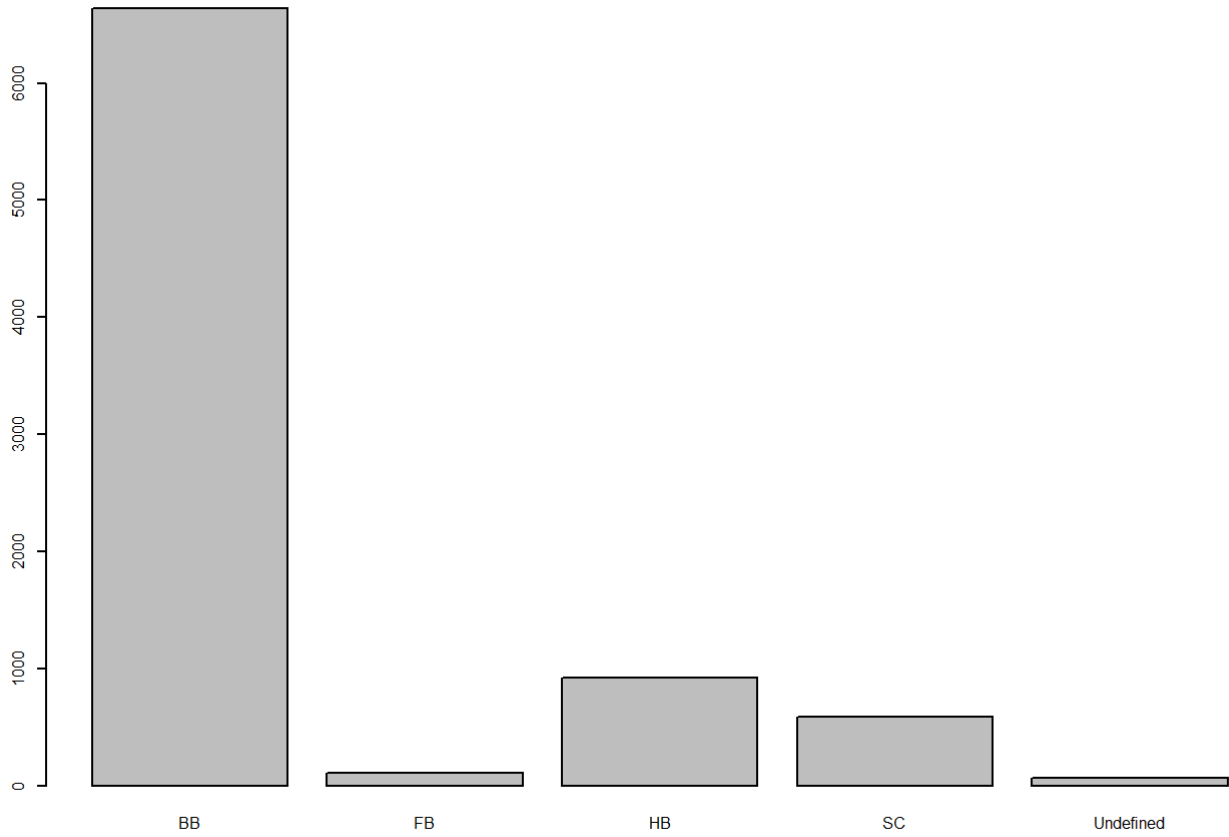
children



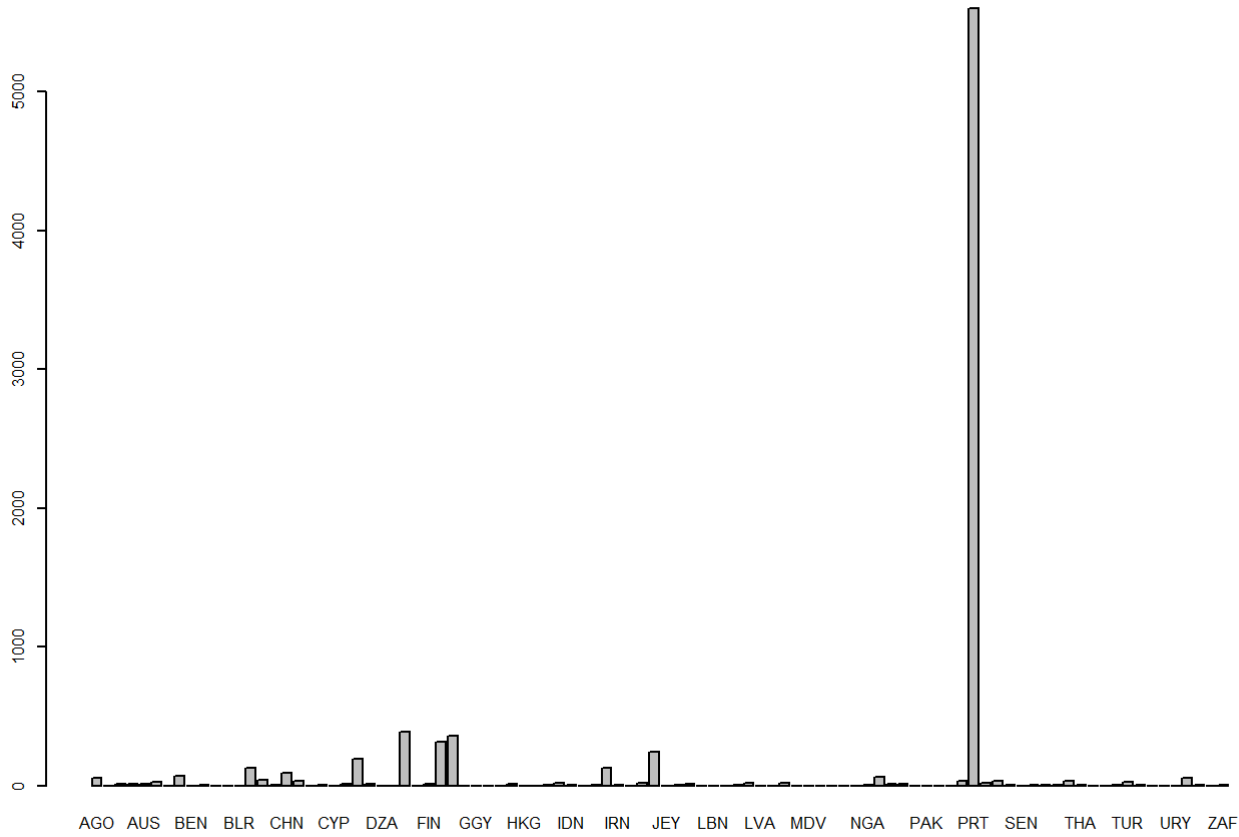
babies



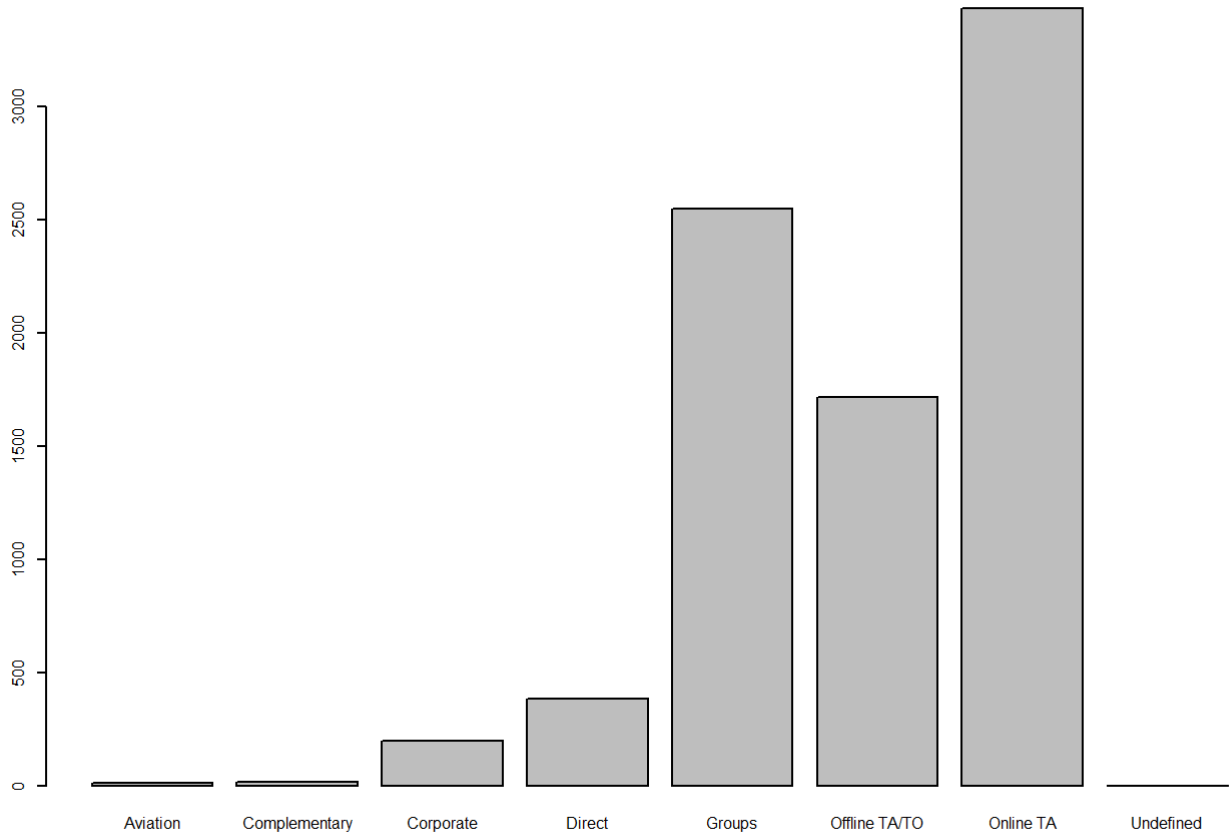
meal



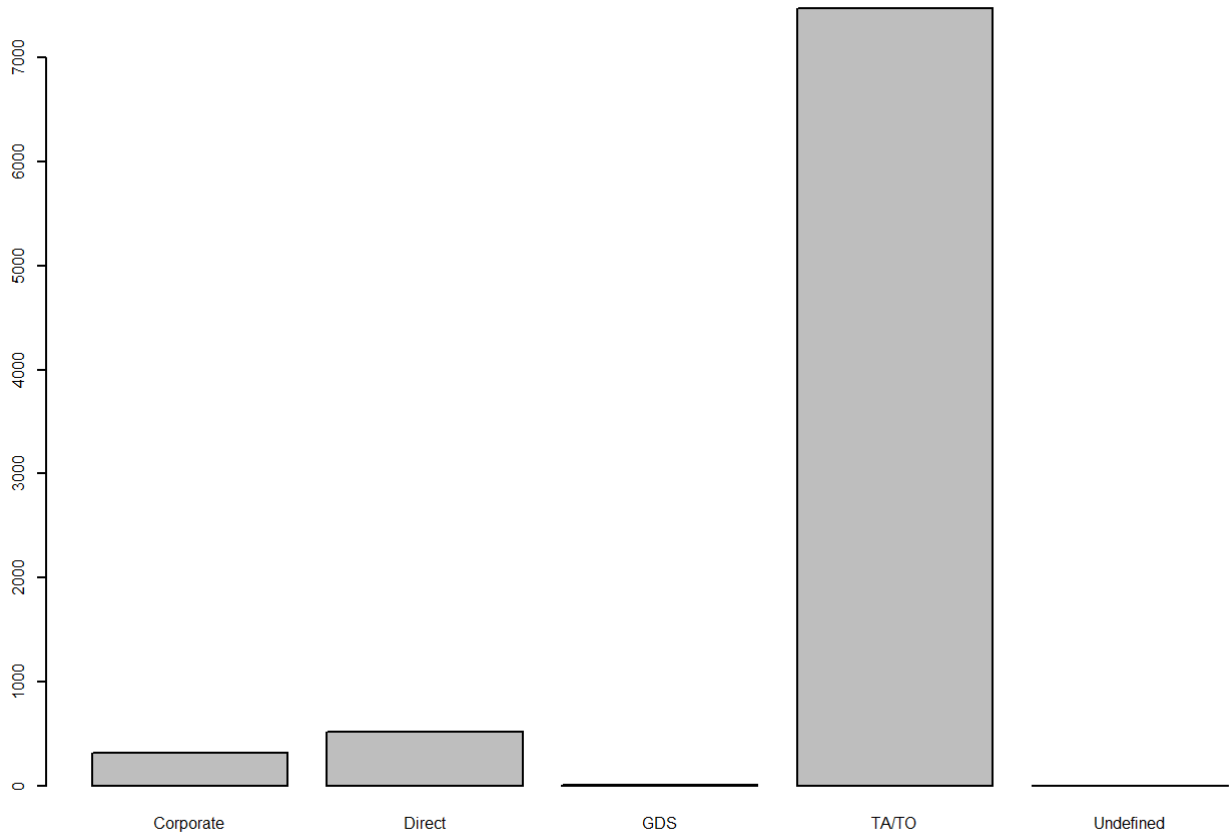
country



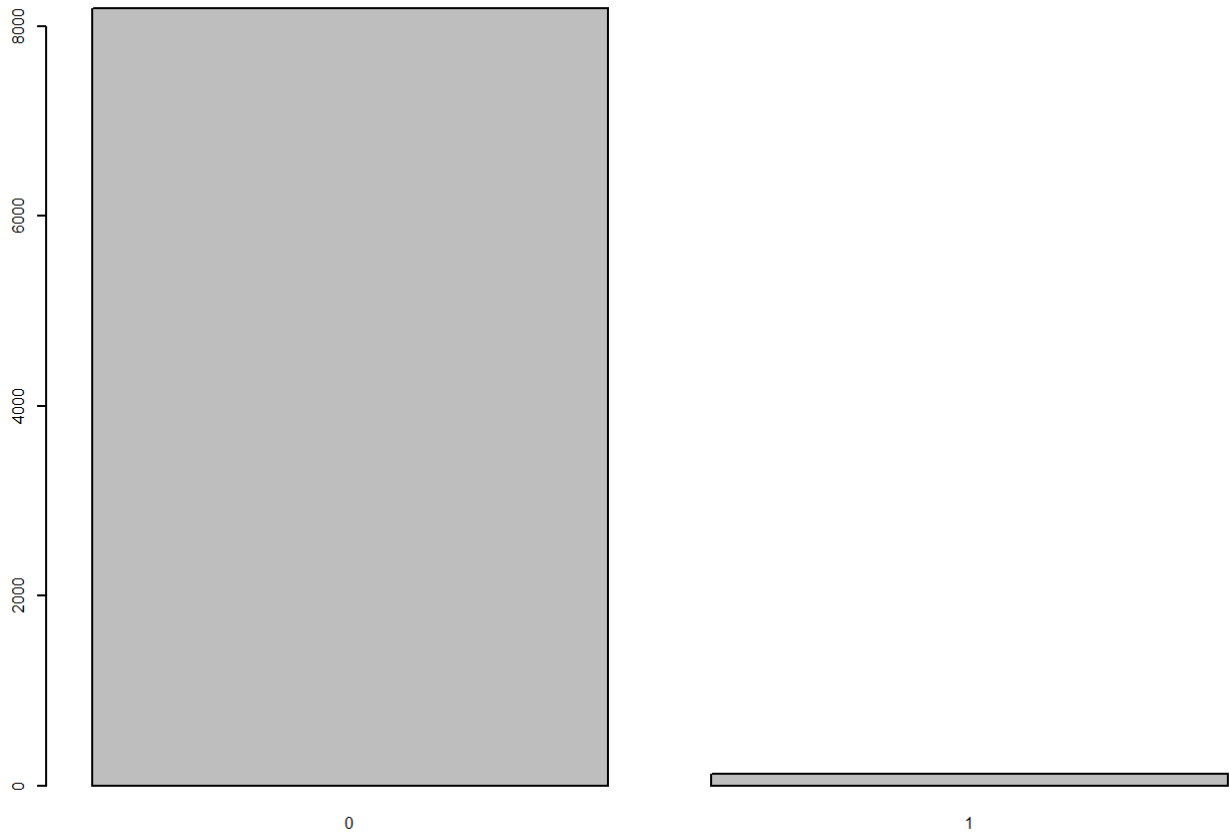
market_segment



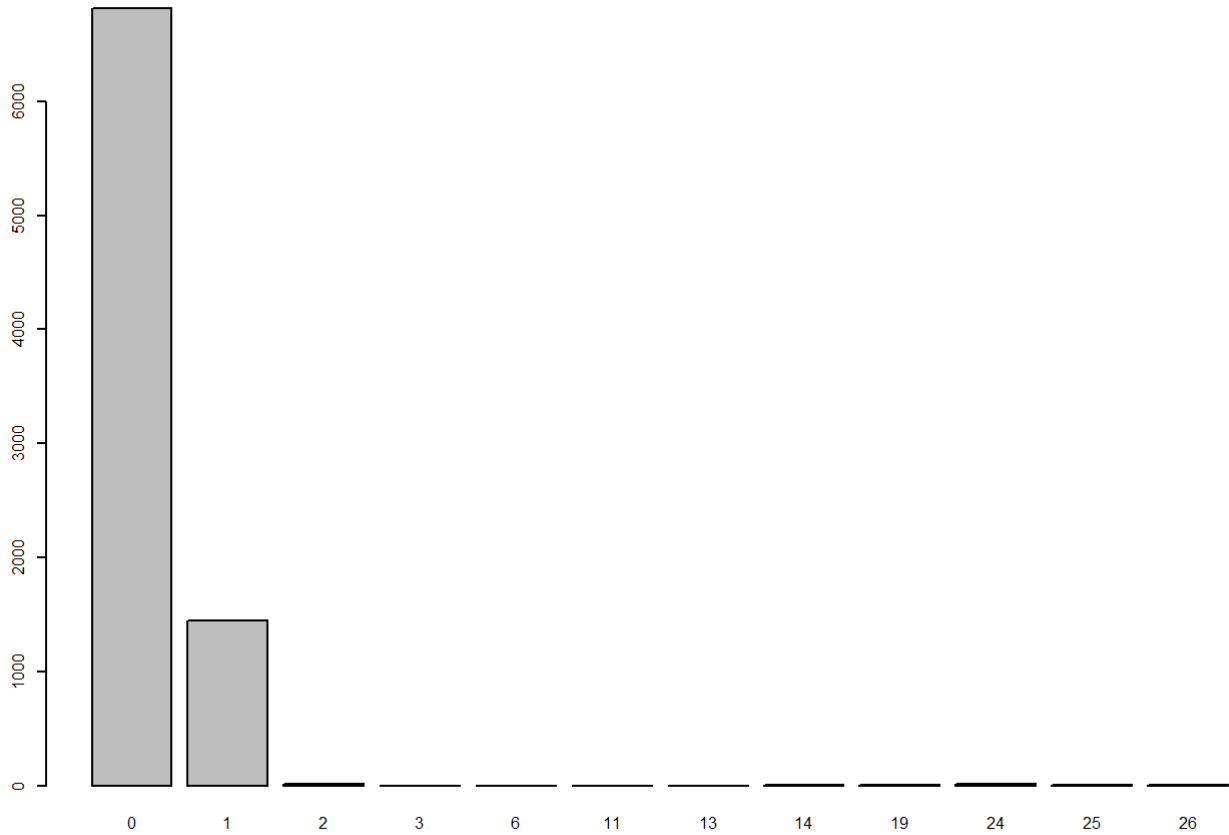
distribution_channel



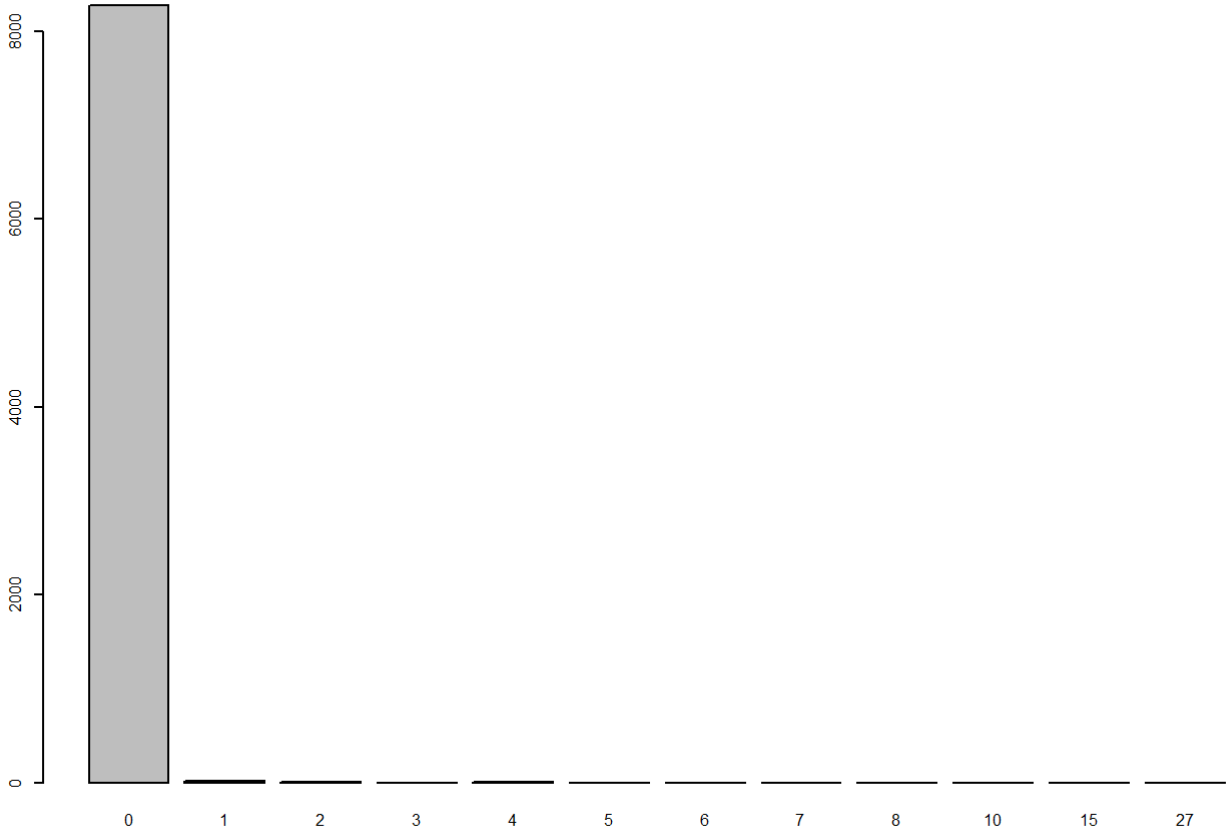
is_repeated_guest



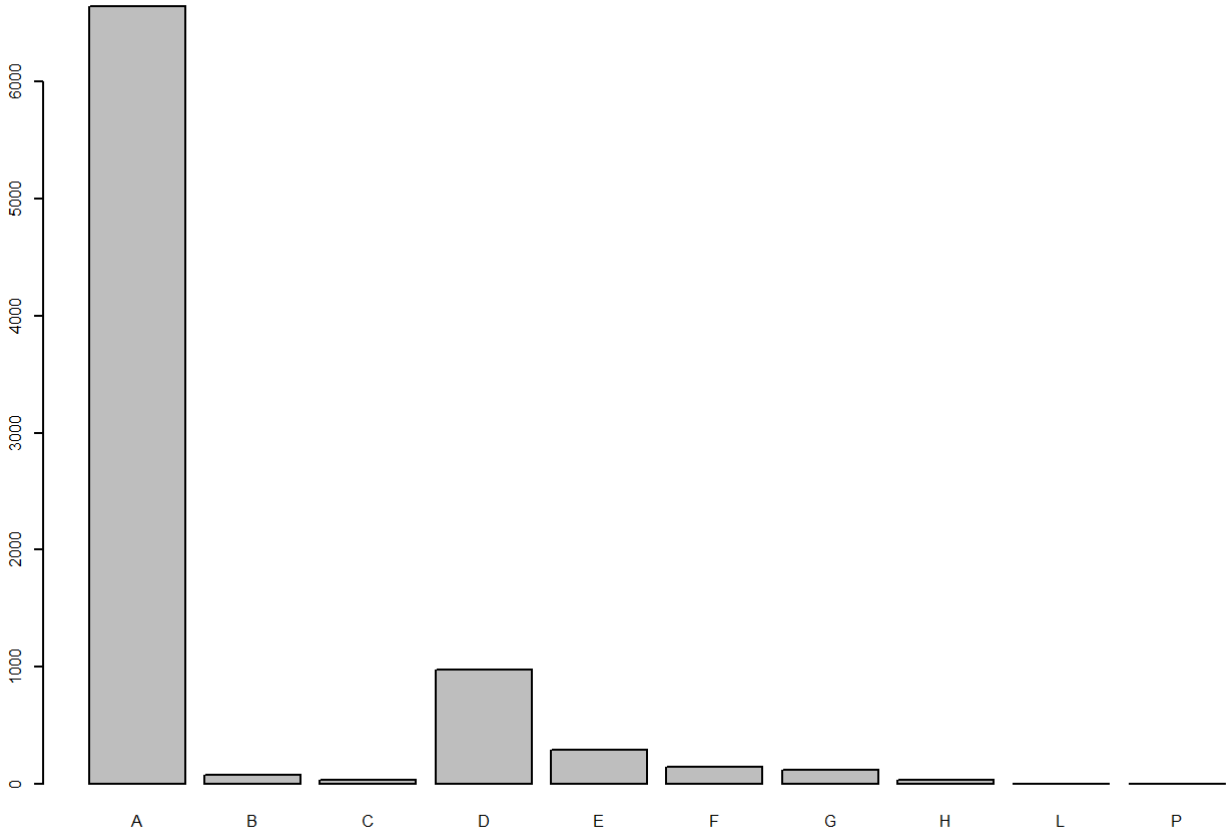
previous_cancellations



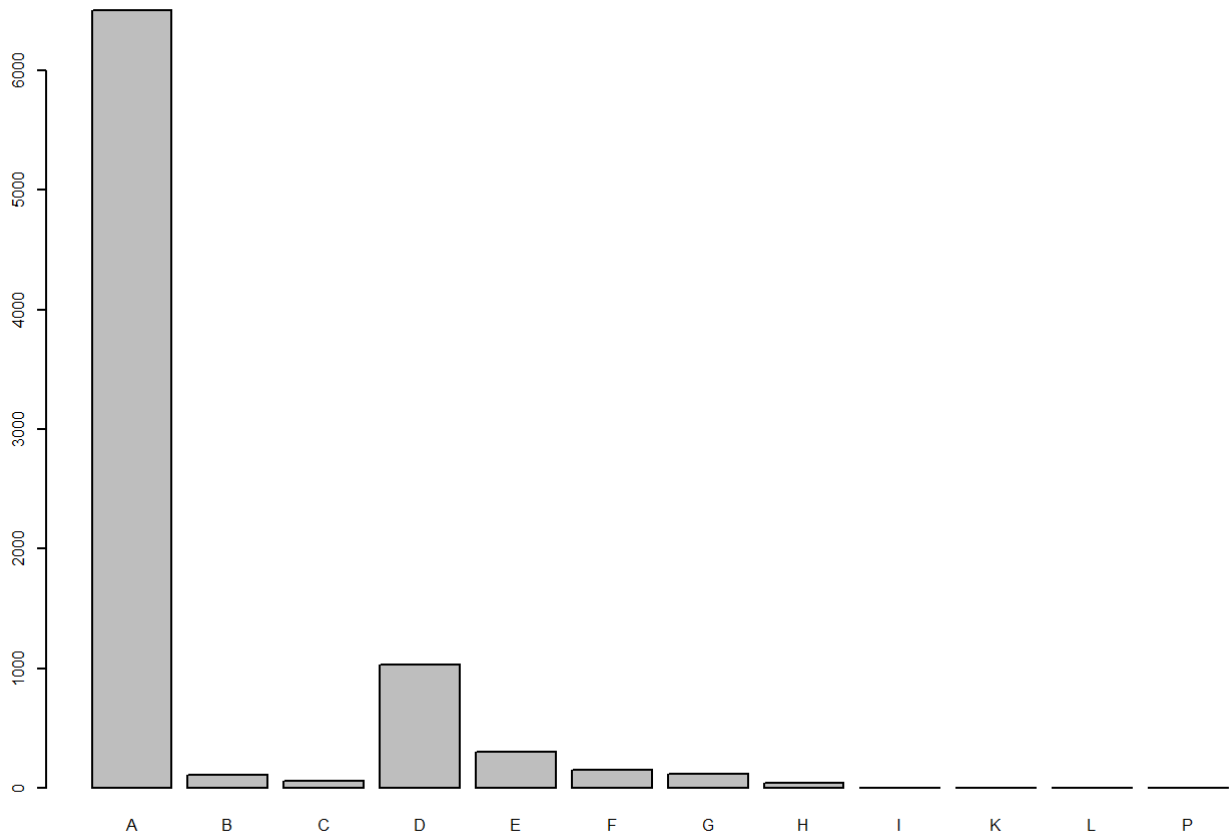
previous_bookings_not_canceled



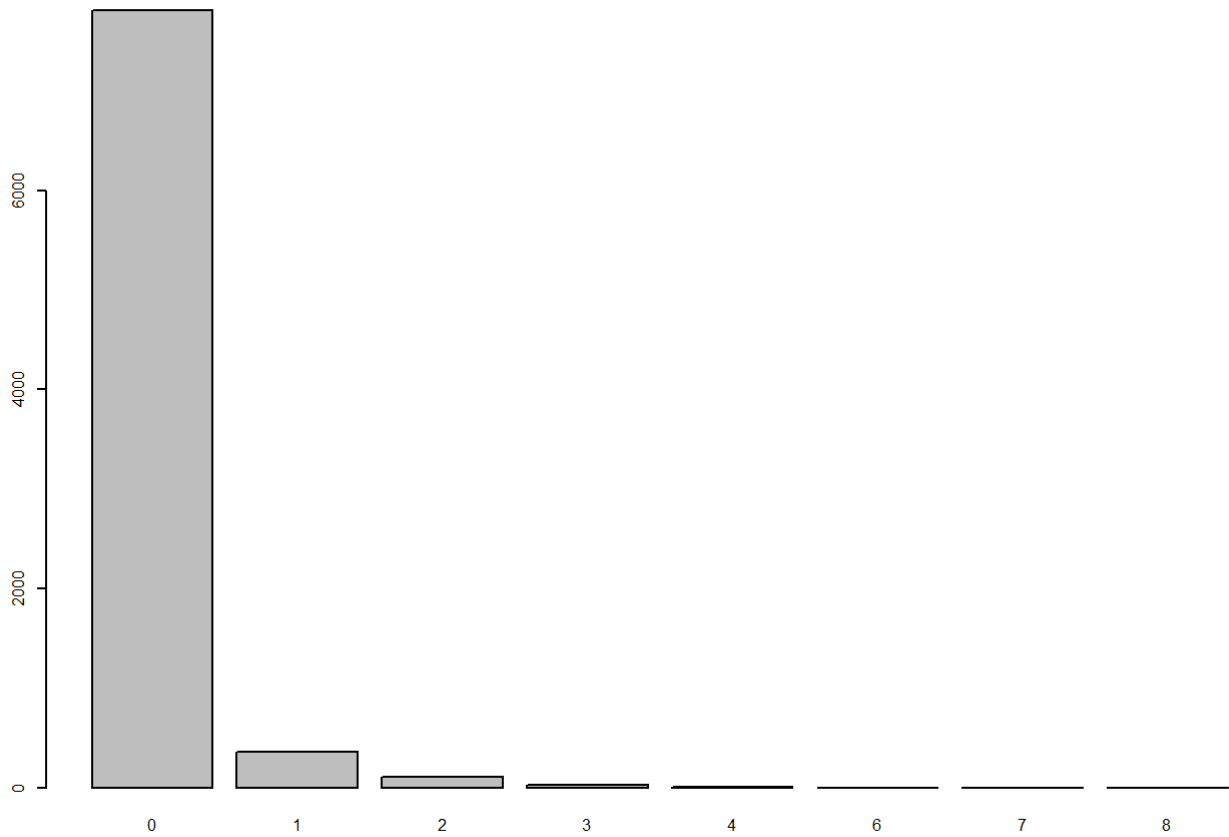
reserved_room_type



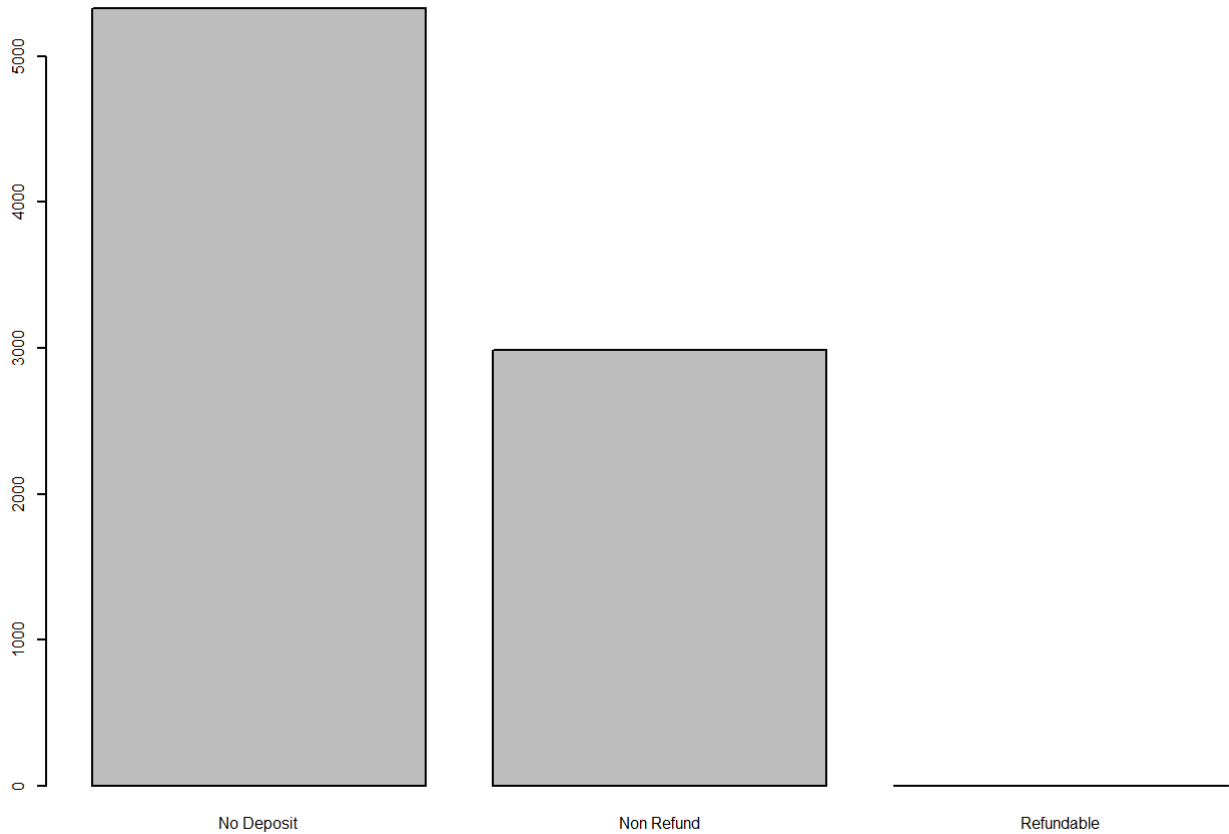
assigned_room_type



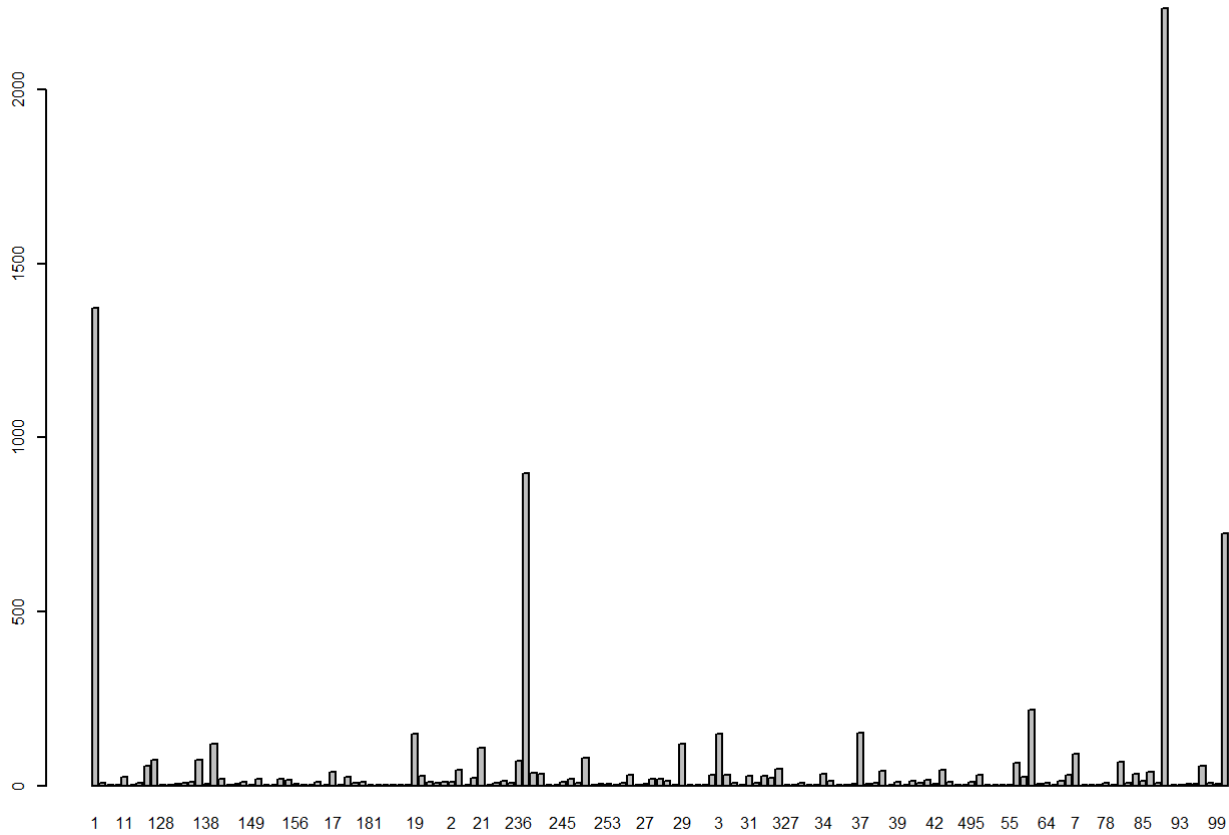
booking_changes



deposit_type



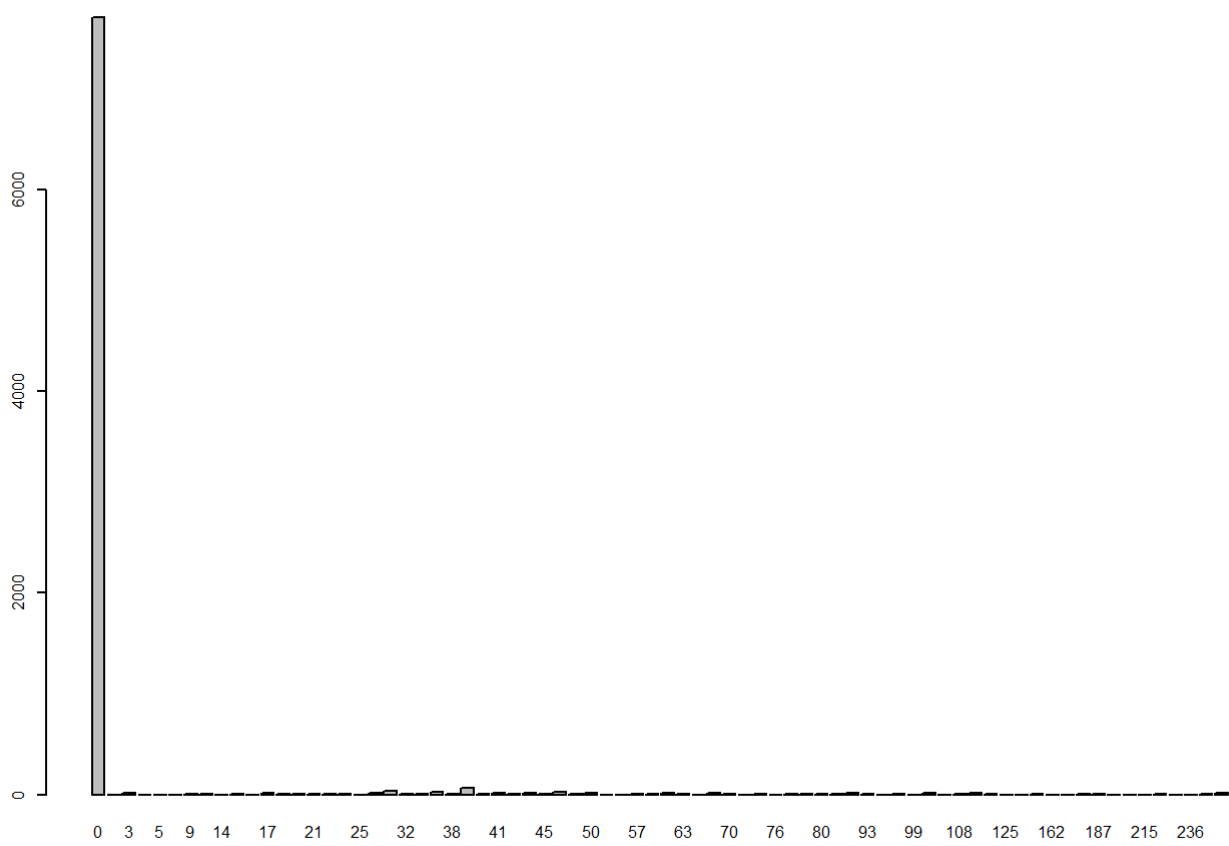
agent



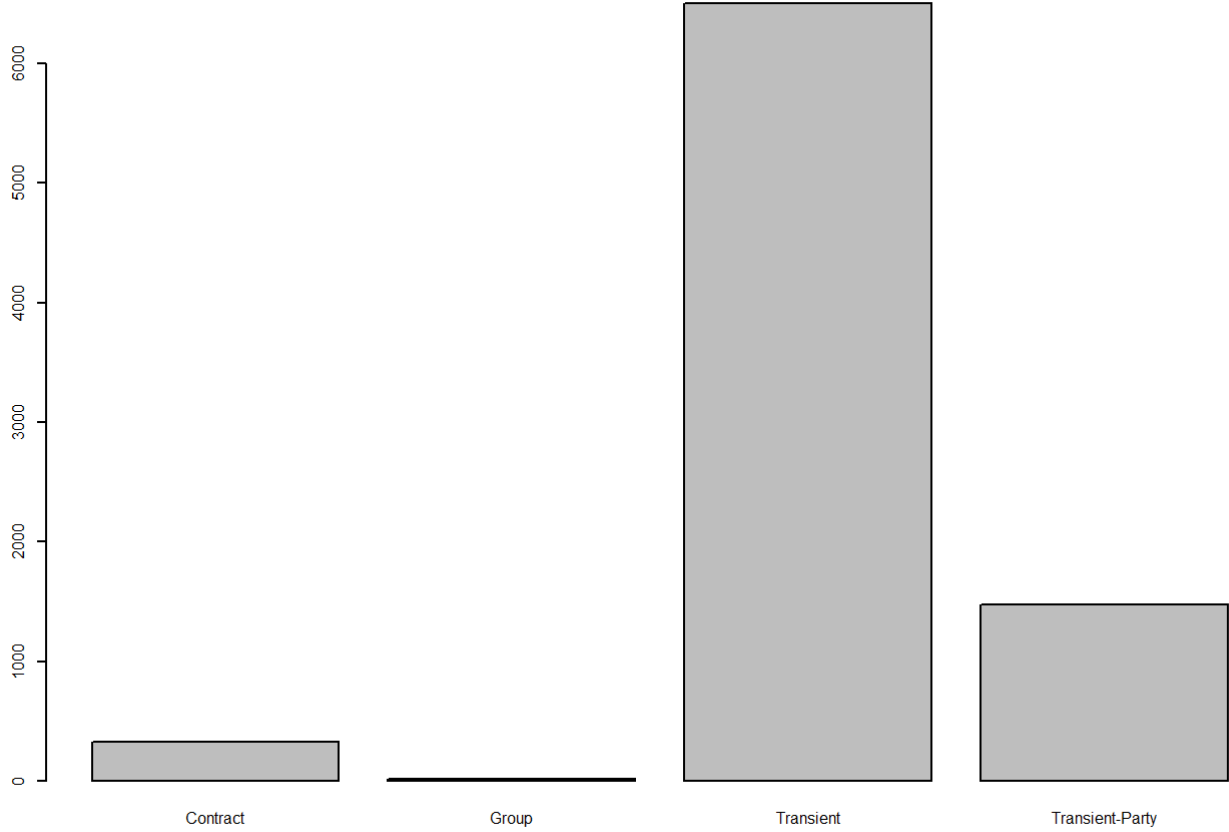
company



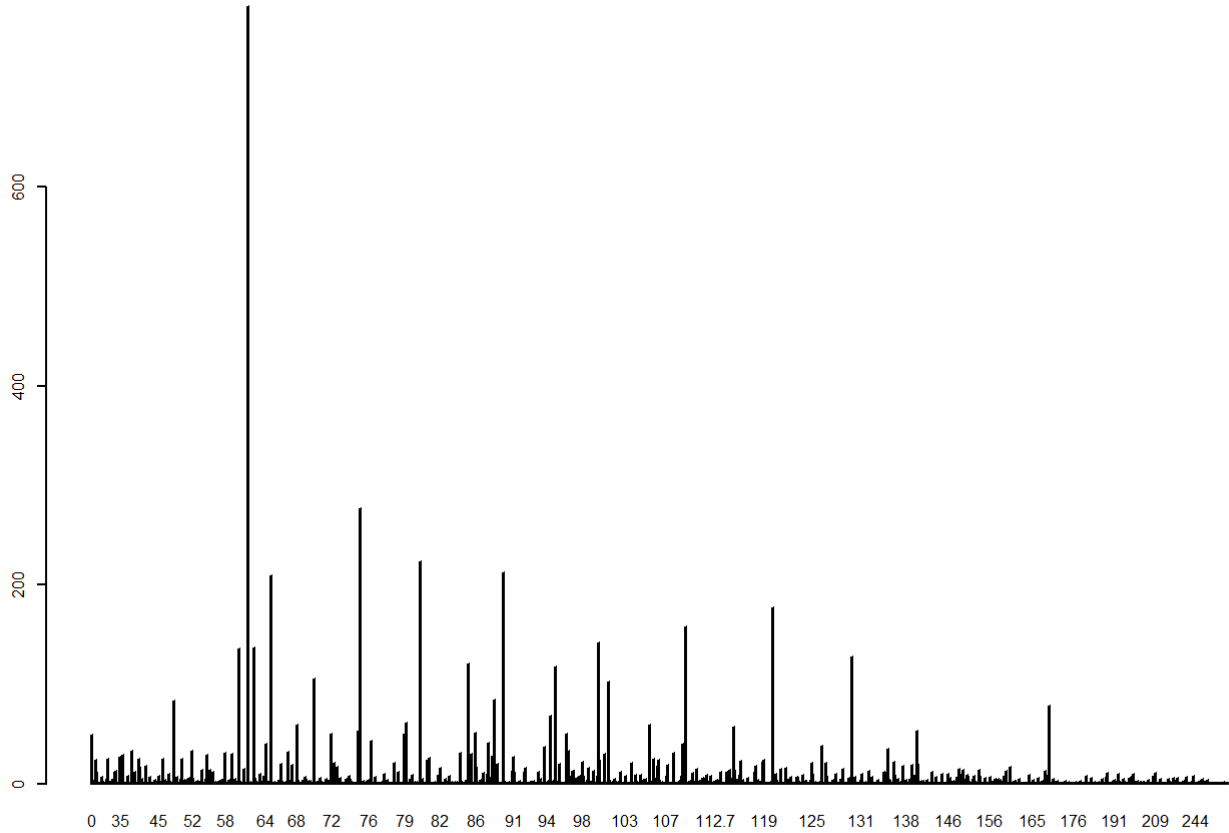
days_in_waiting_list



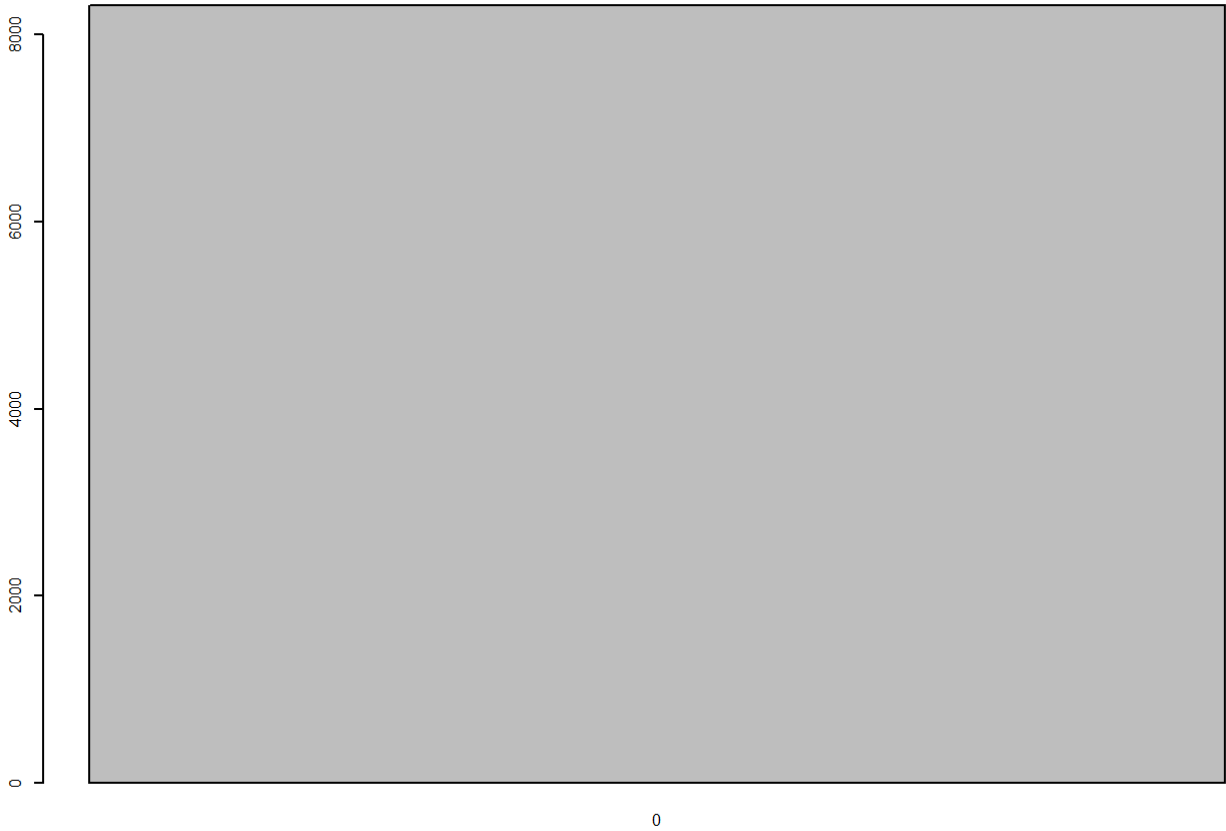
customer_type



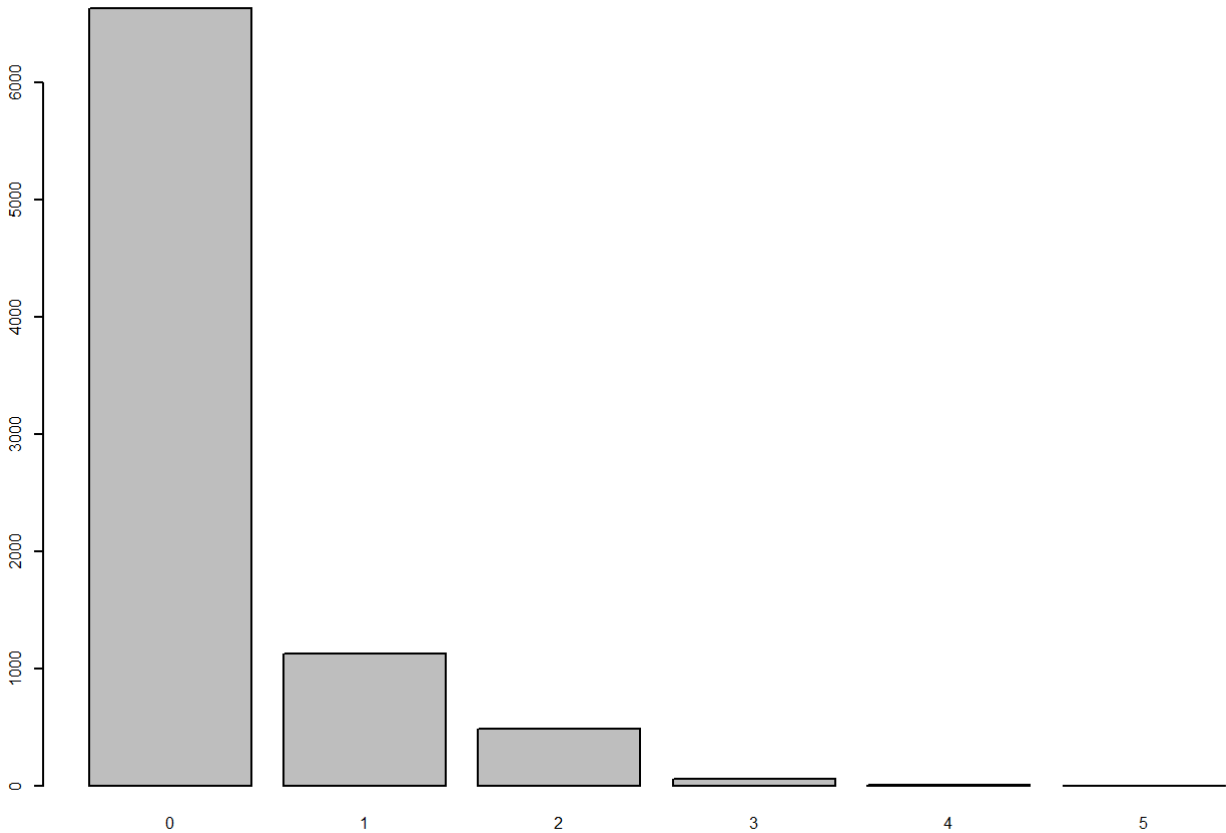
adr



required_car_parking_spaces



total_of_special_requests

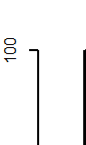
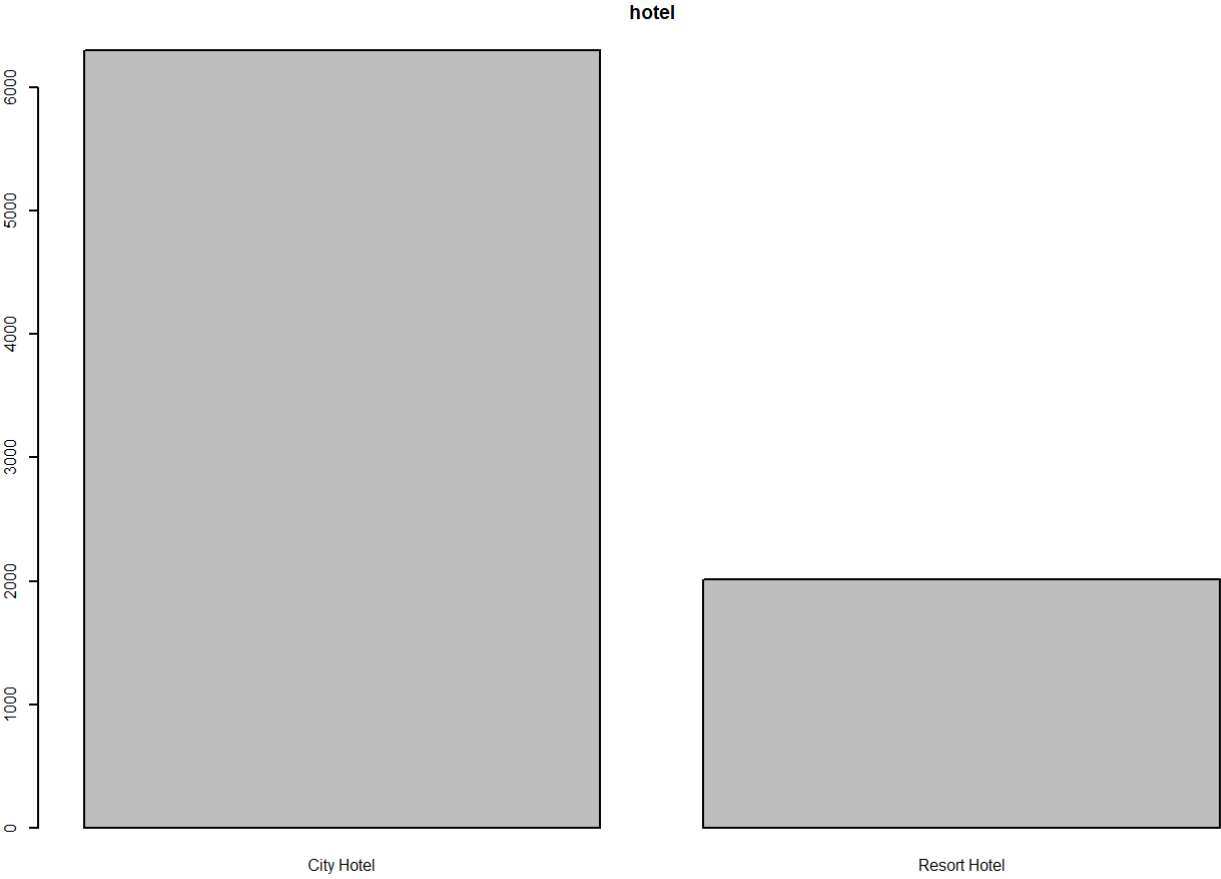


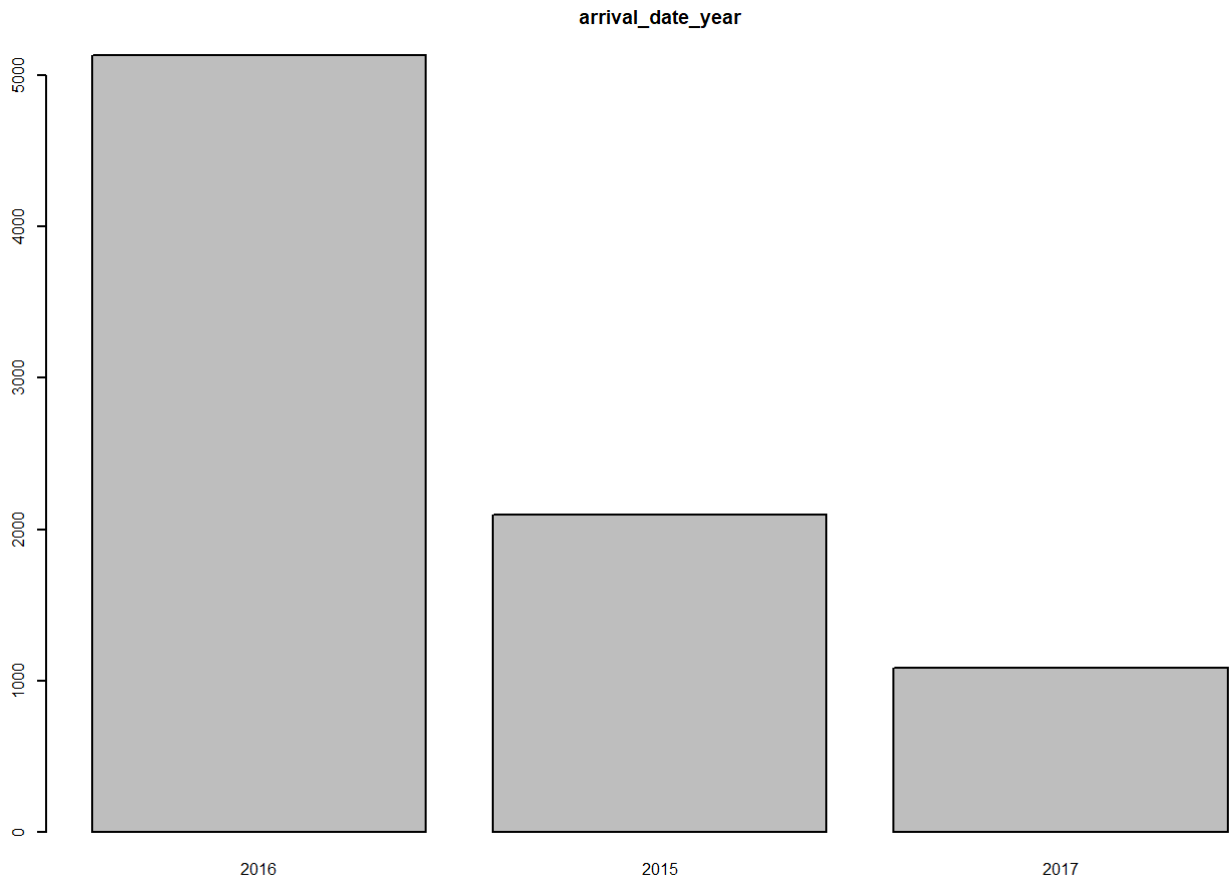
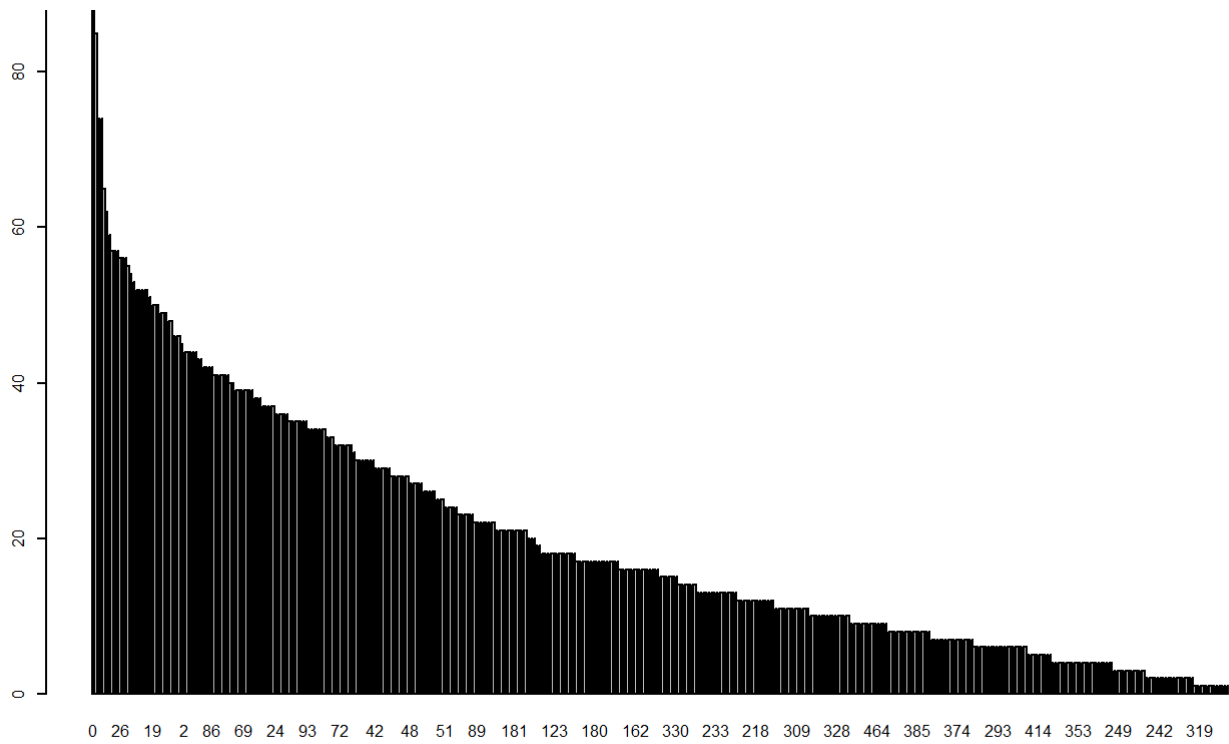
analisis de todas las variables con el subconjunto de canceladas ordenas

```
ncol(sub_cancelados)
```

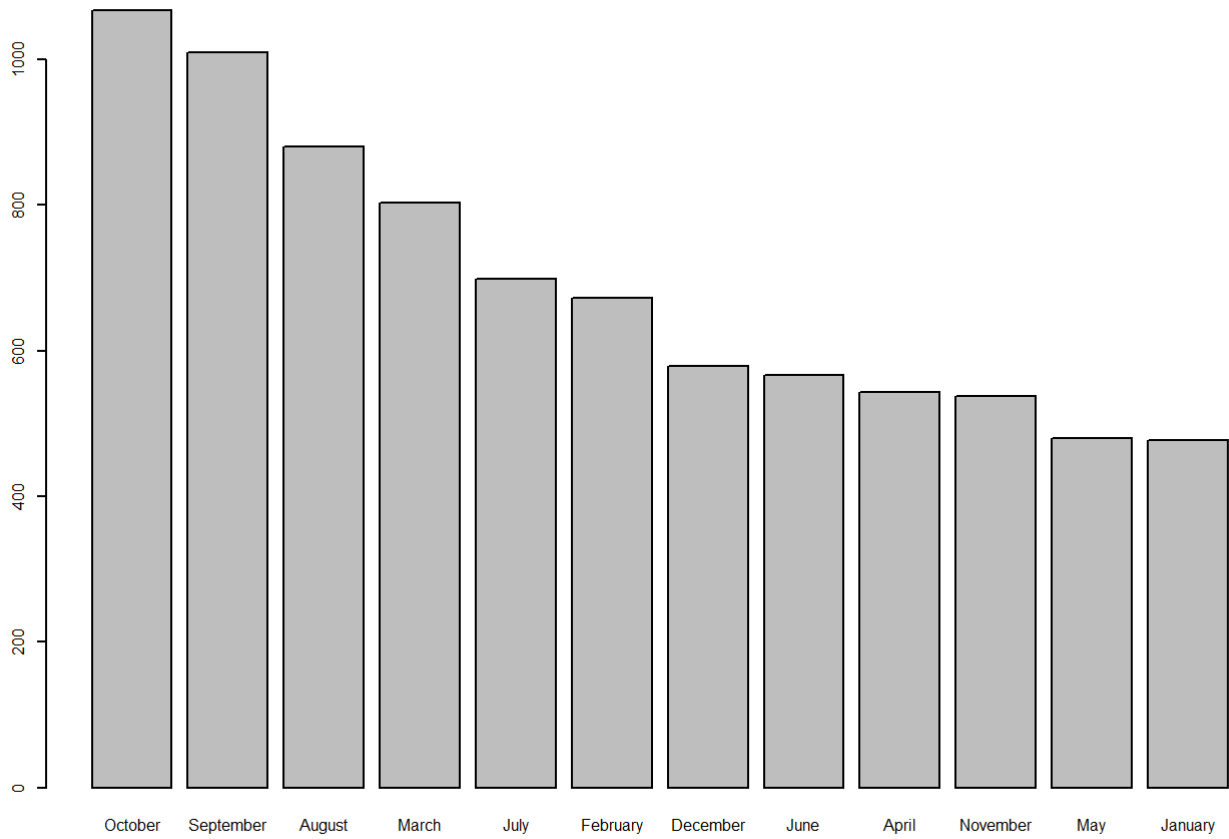
```
## [1] 30
```

```
nombres <- names(sub_cancelados)
for (i in c(1:ncol(sub_cancelados))){
  analisis_cancelaciones_ord(sub_cancelados[,i],nombres[i])
}
```

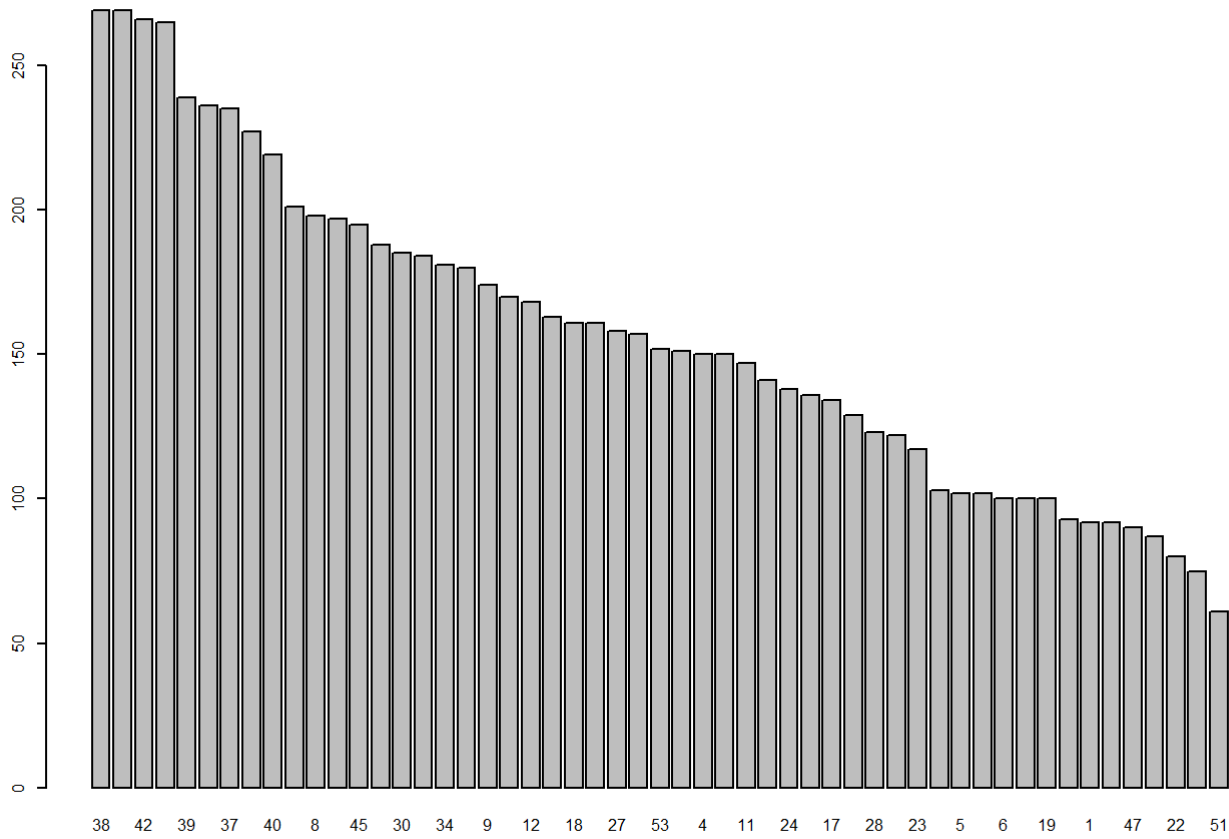




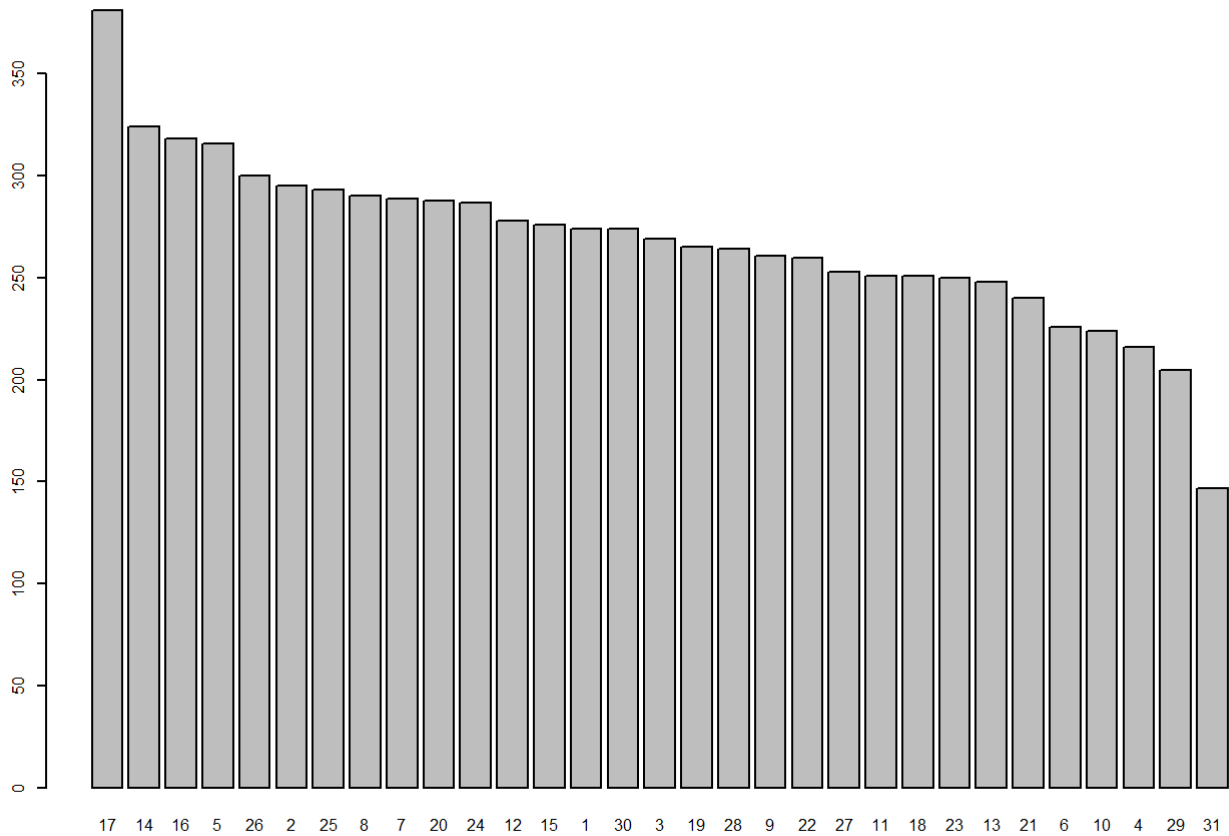
arrival_date_month



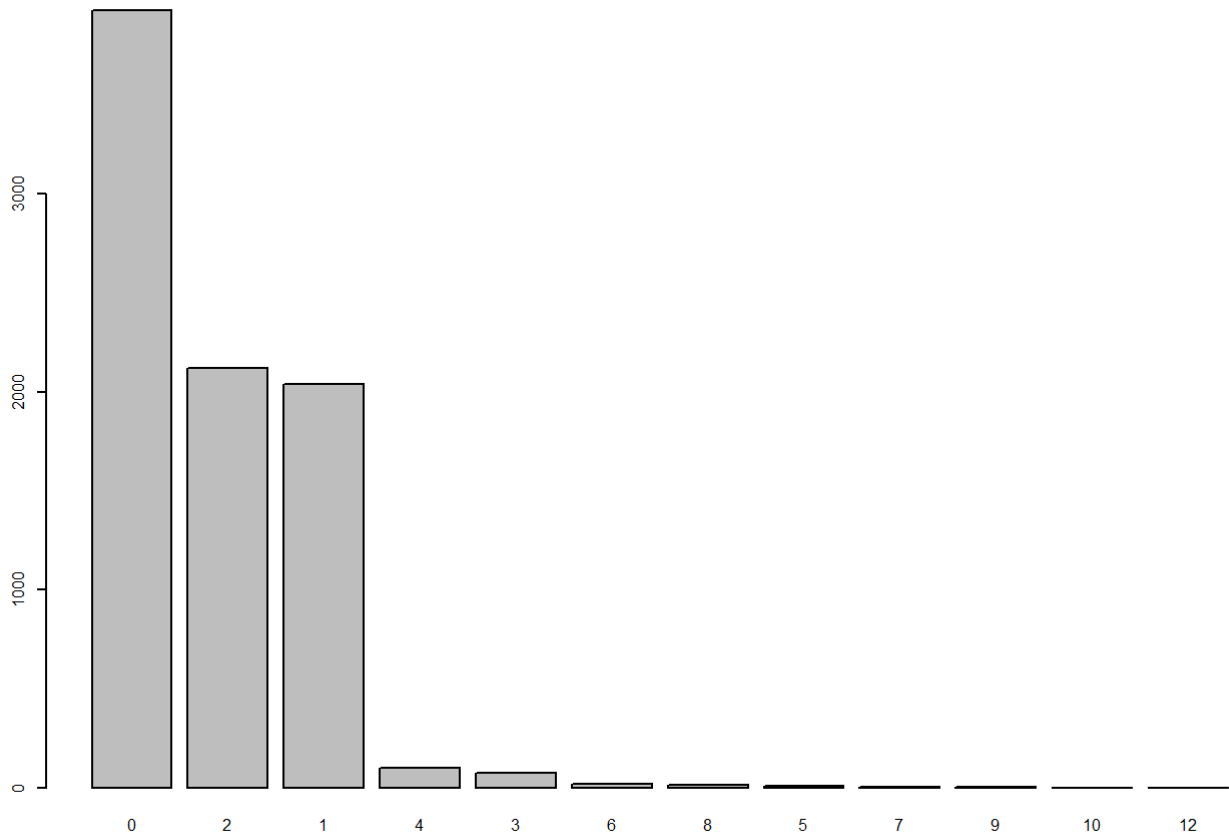
arrival_date_week_number



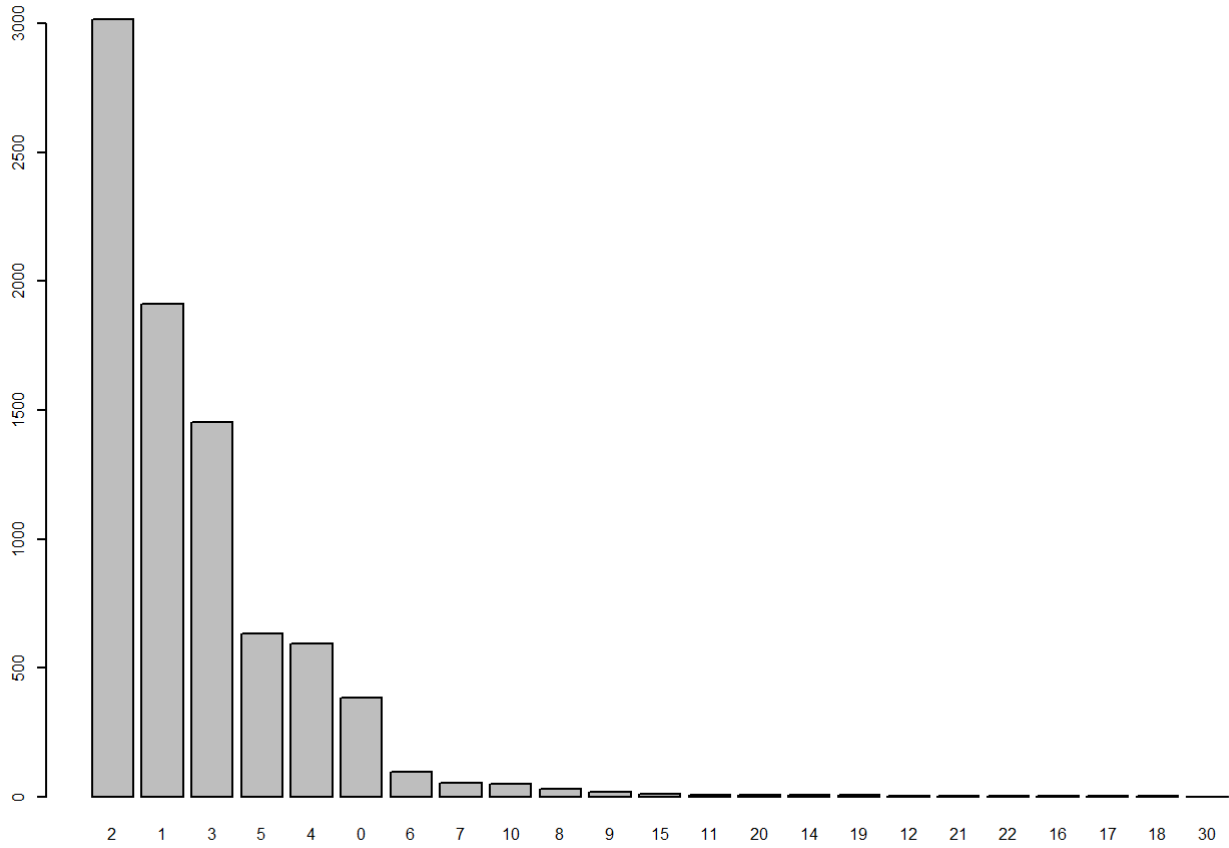
arrival_date_day_of_month



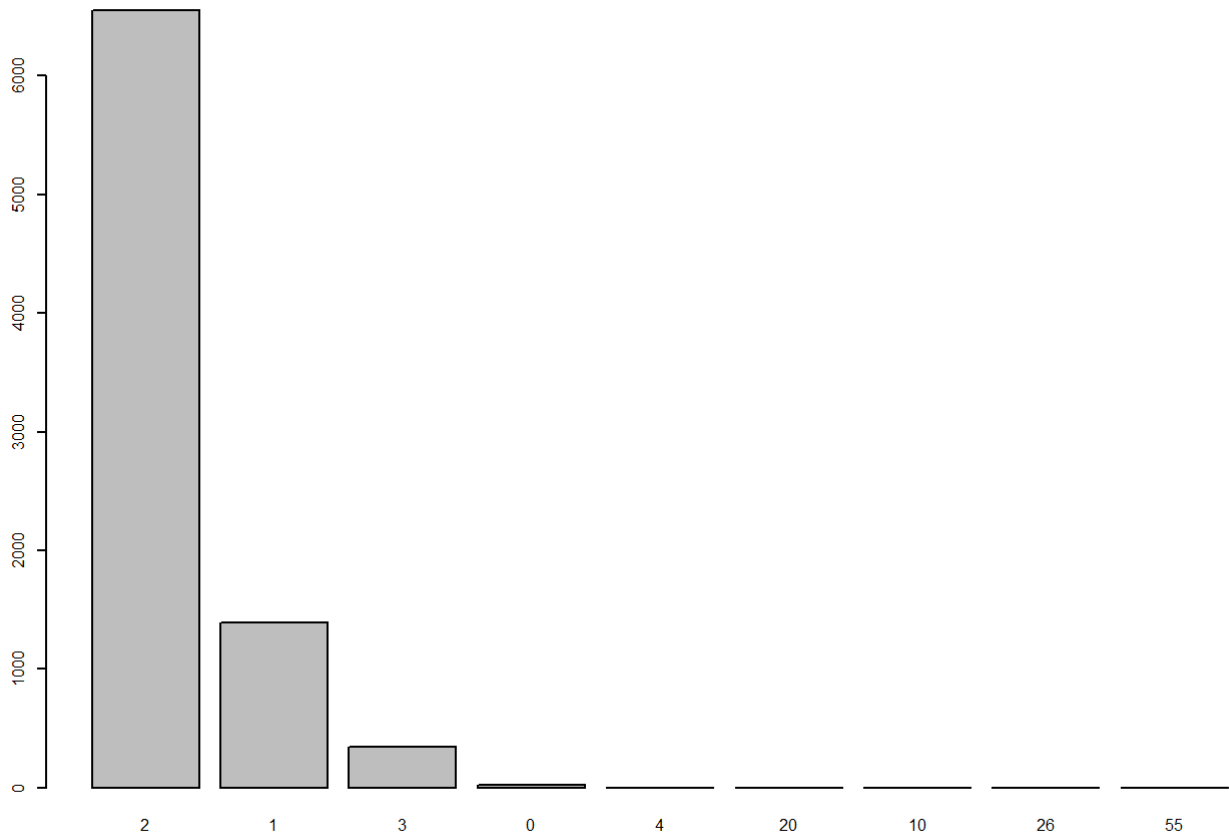
stays_in_weekend_nights



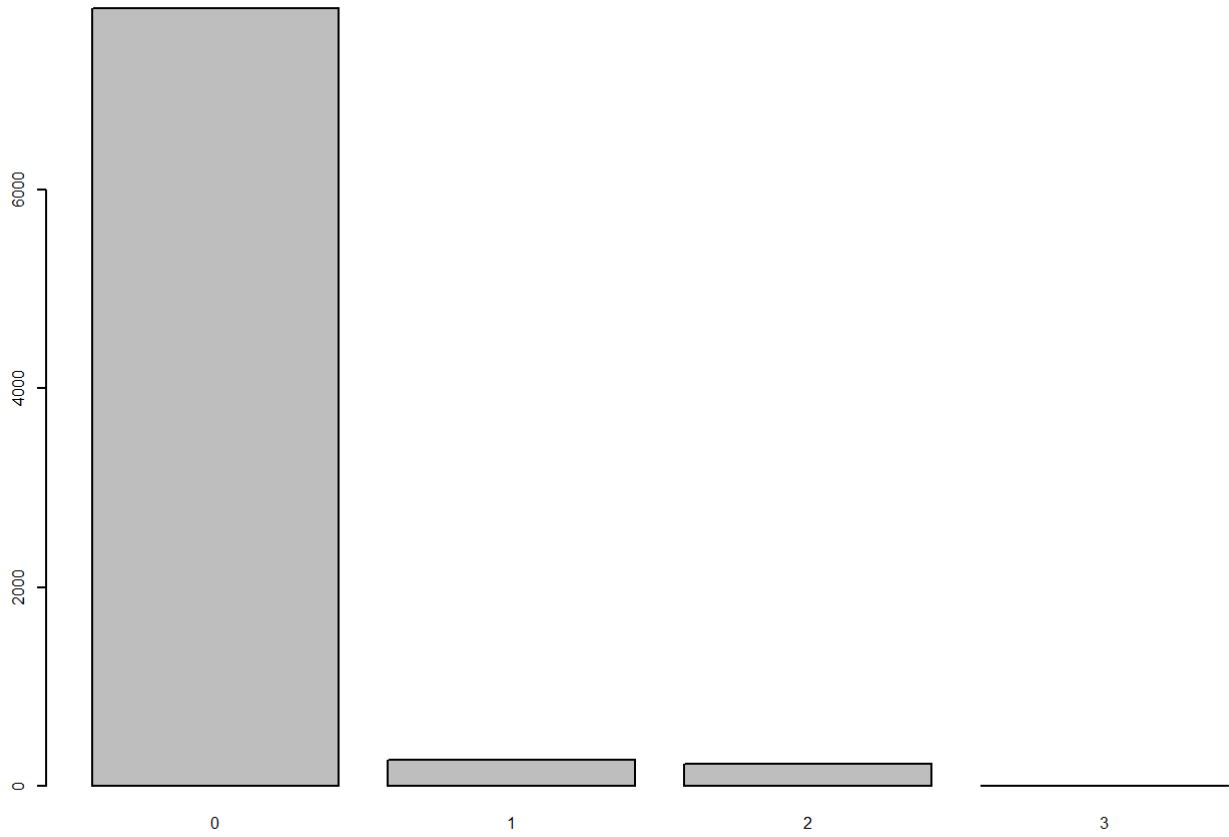
stays_in_week_nights



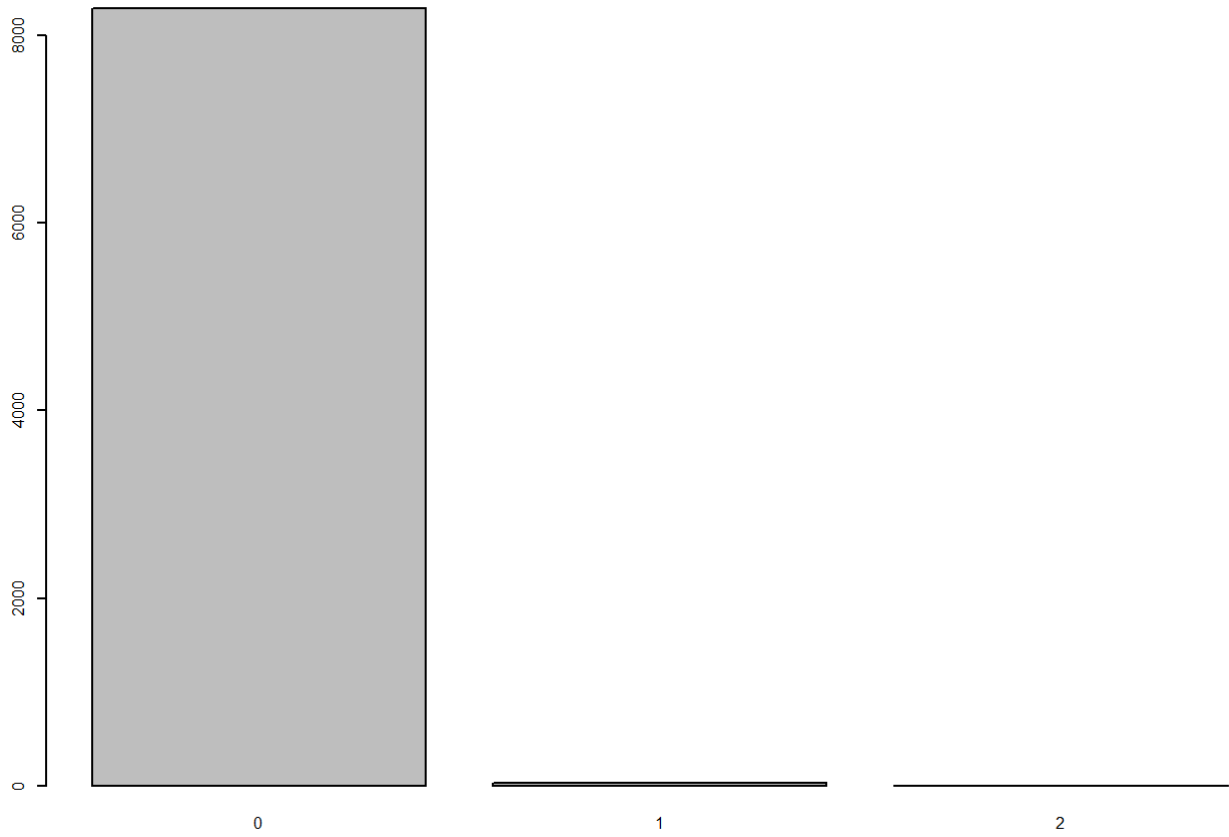
adults



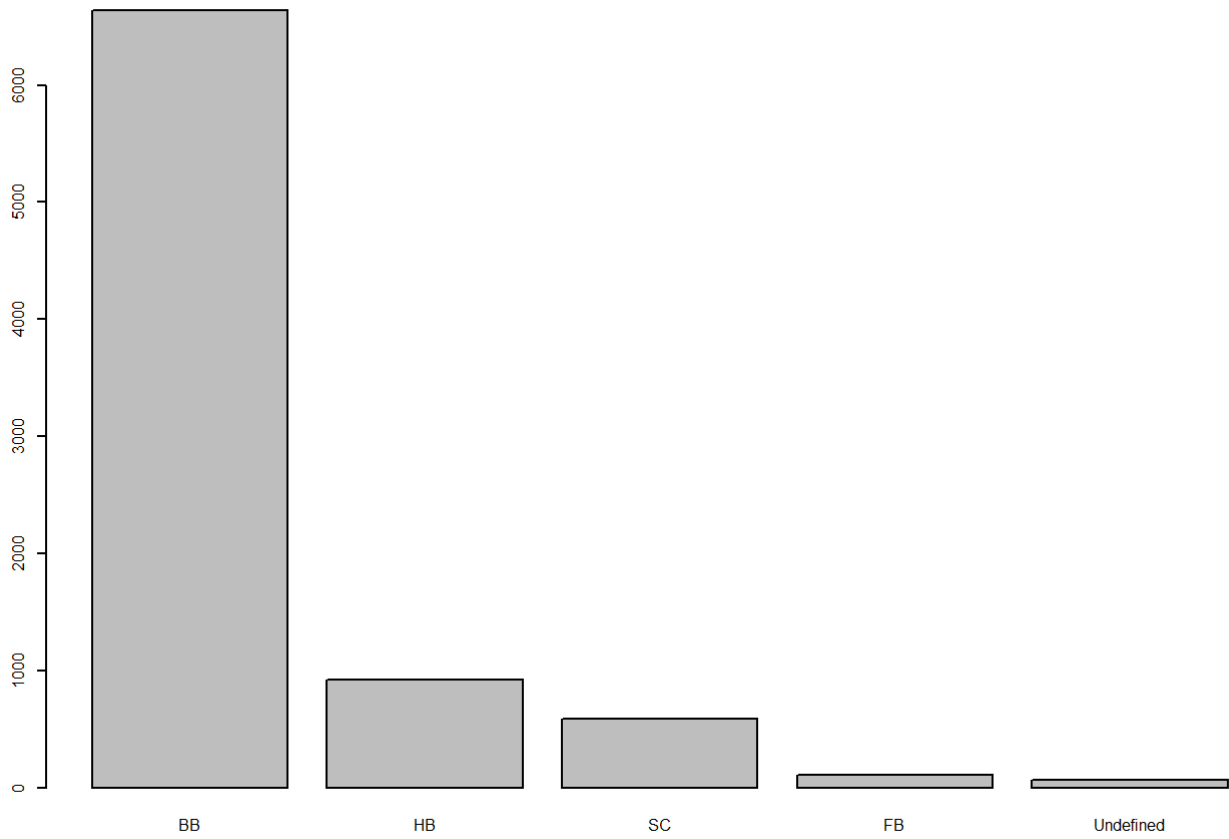
children



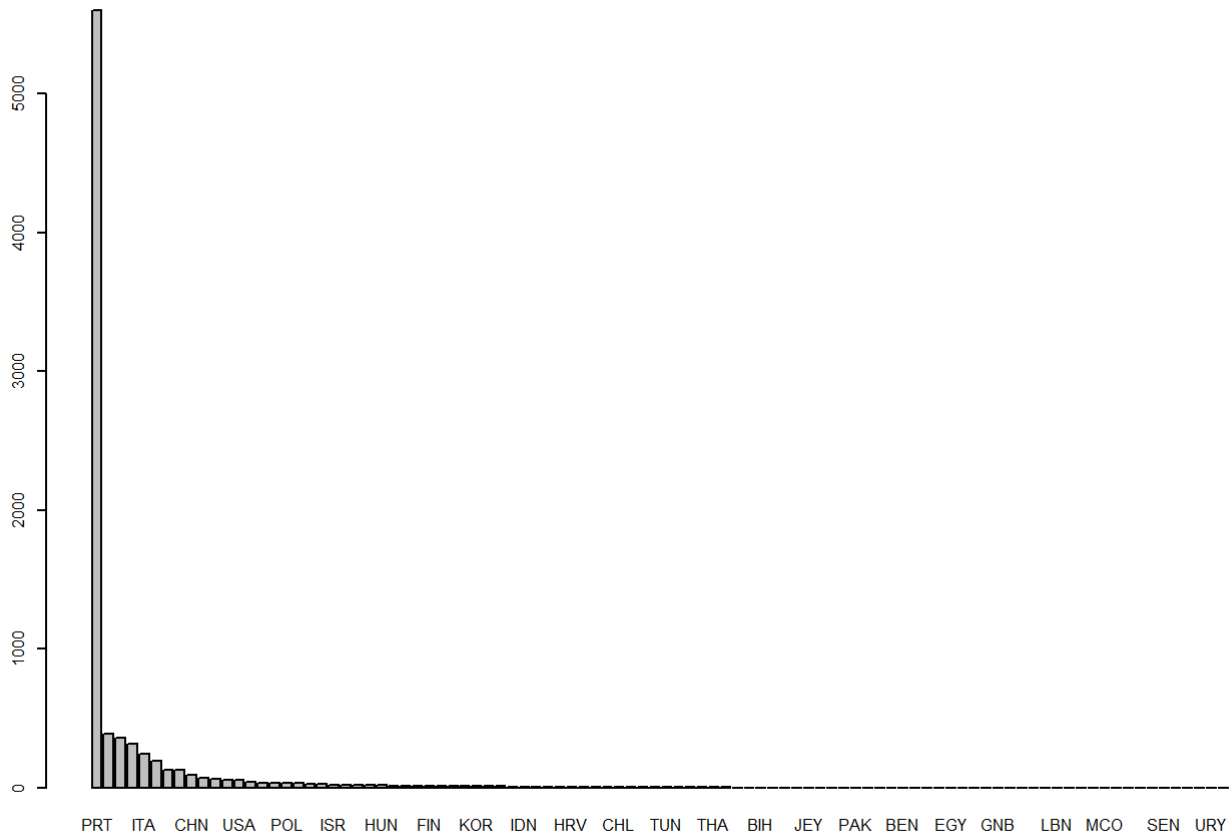
babies



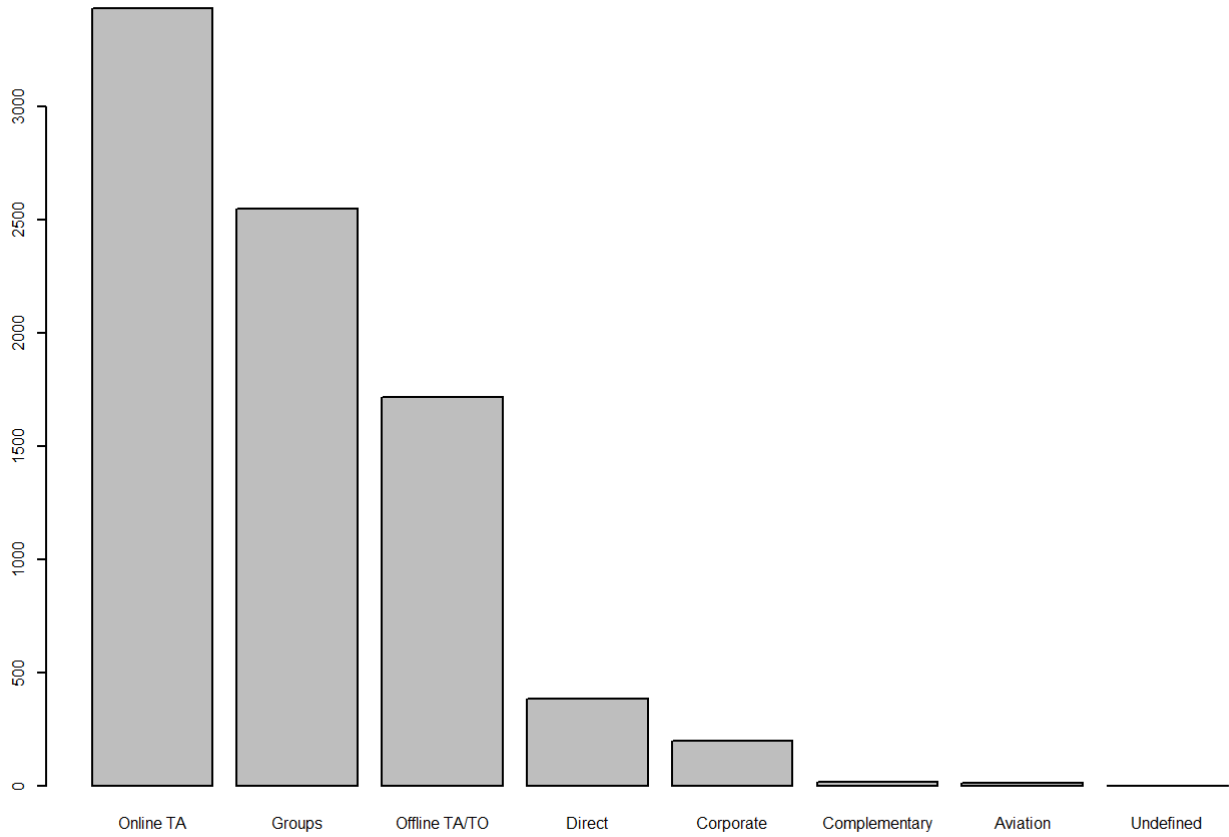
meal



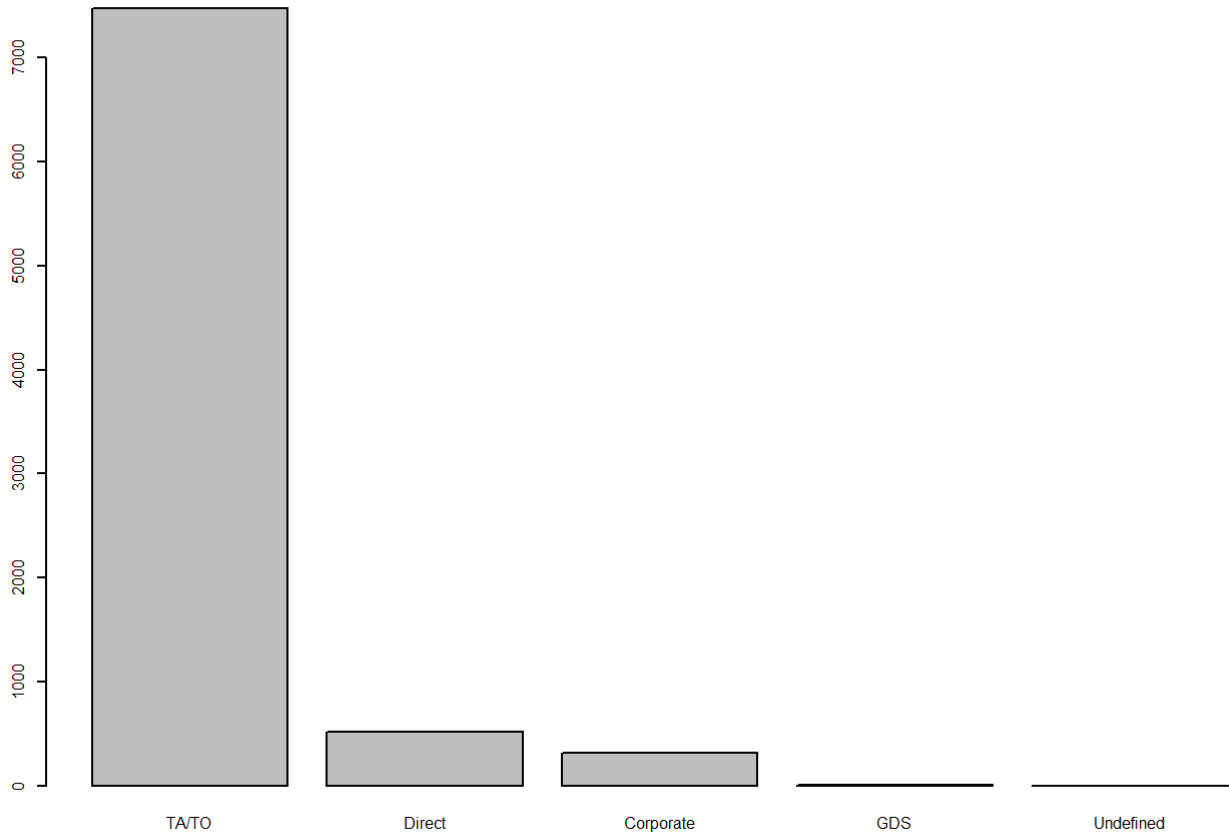
country



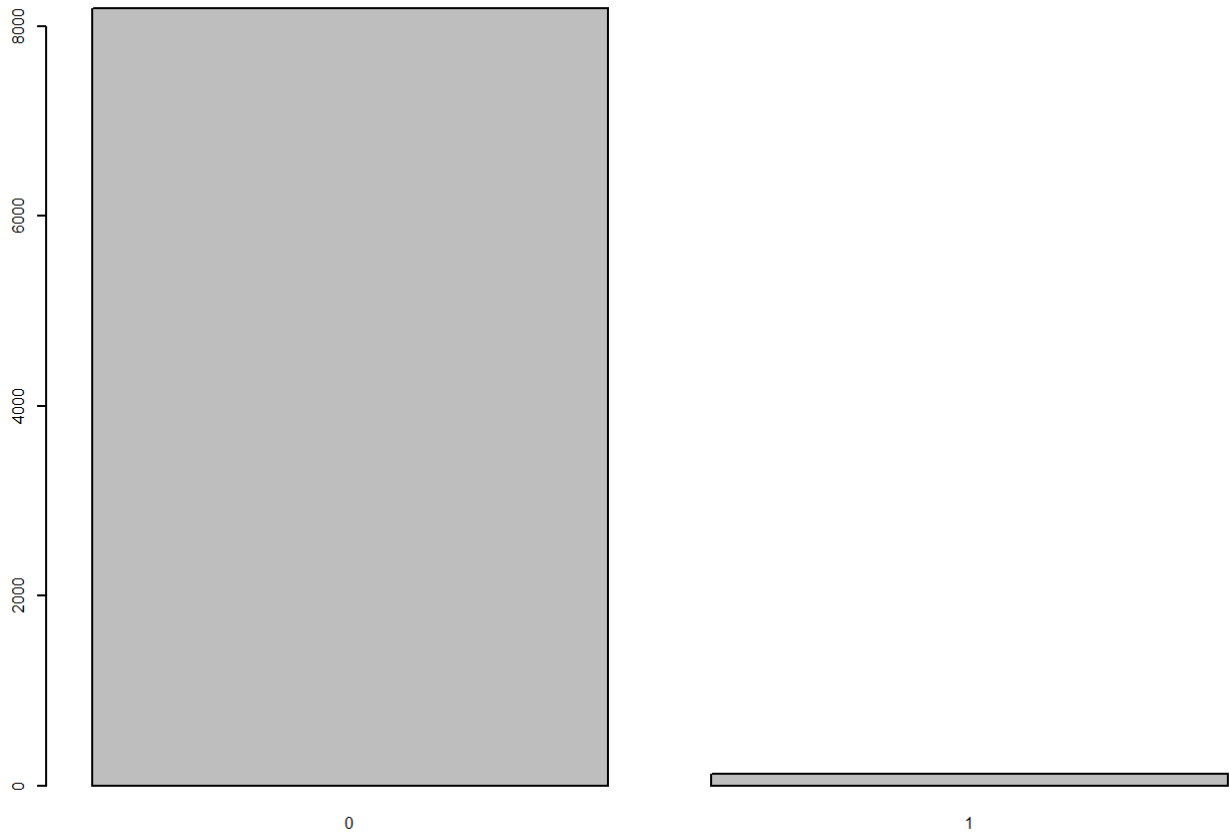
market_segment



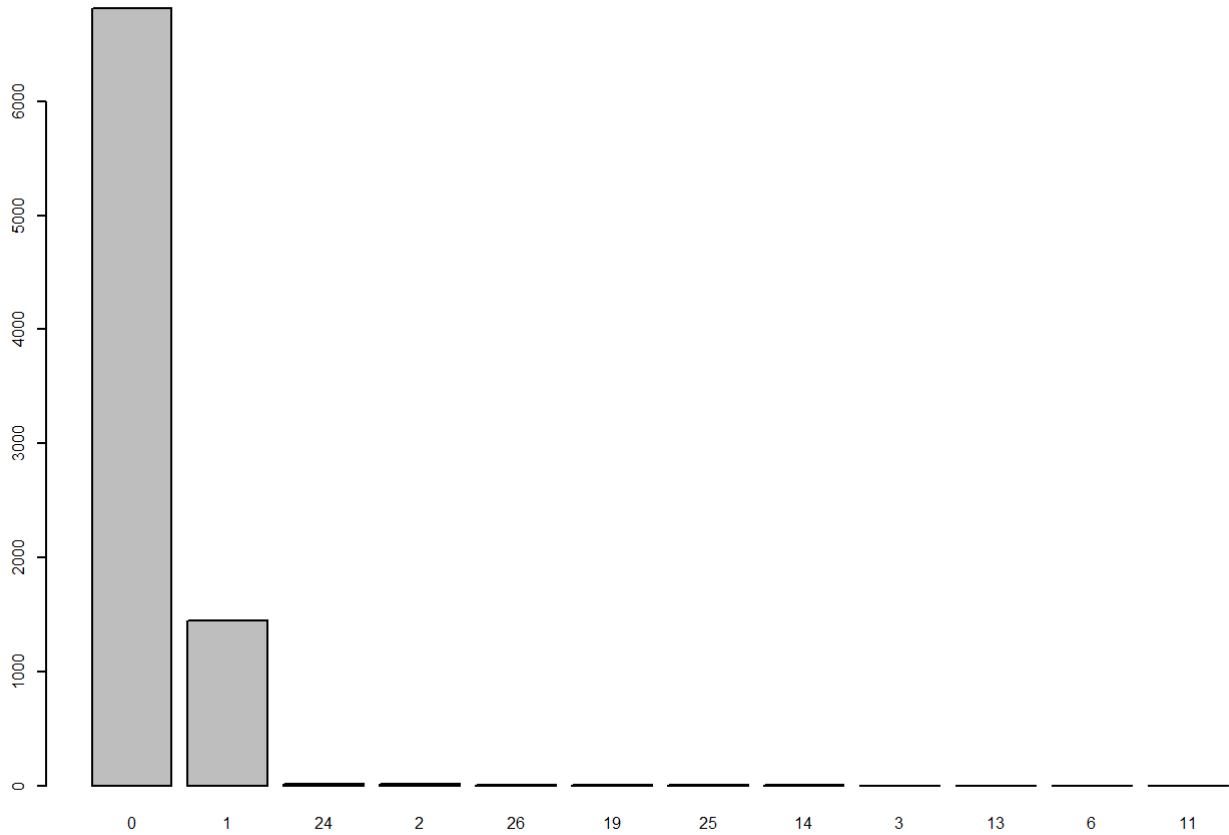
distribution_channel



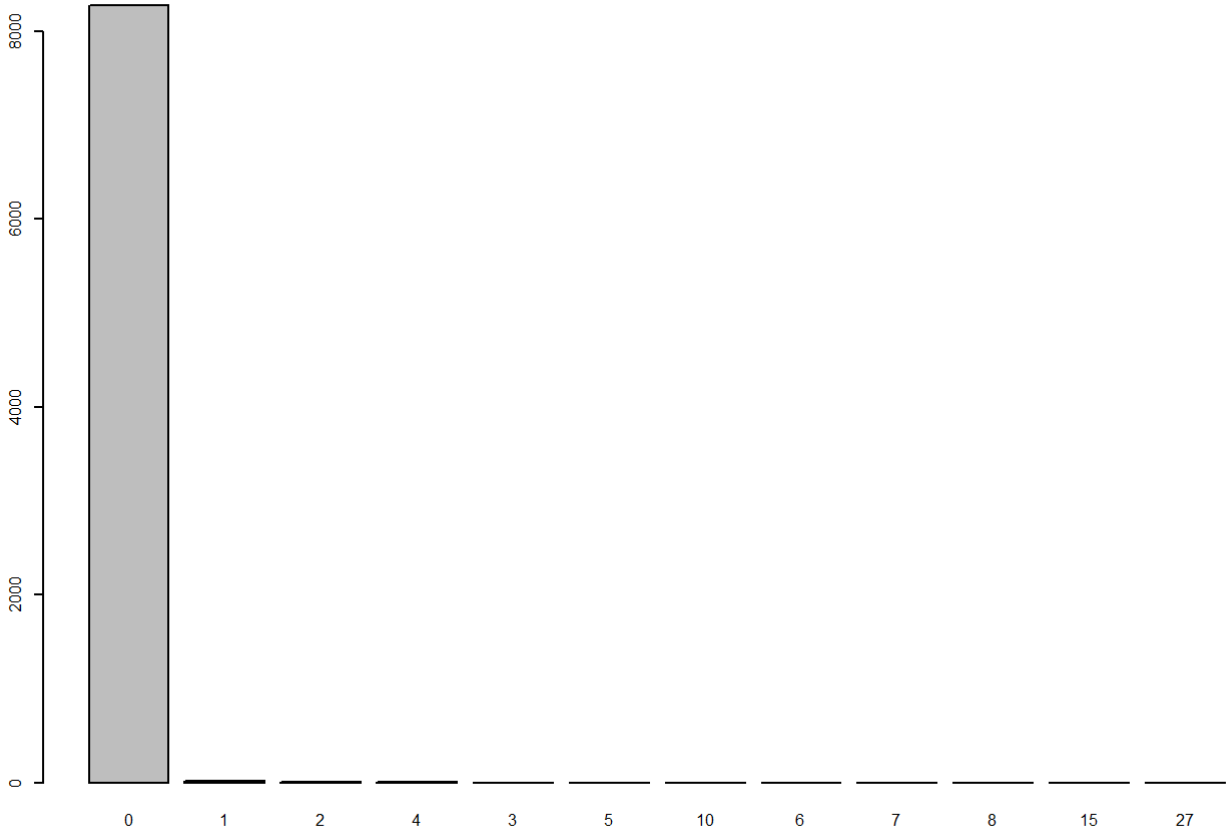
is_repeated_guest



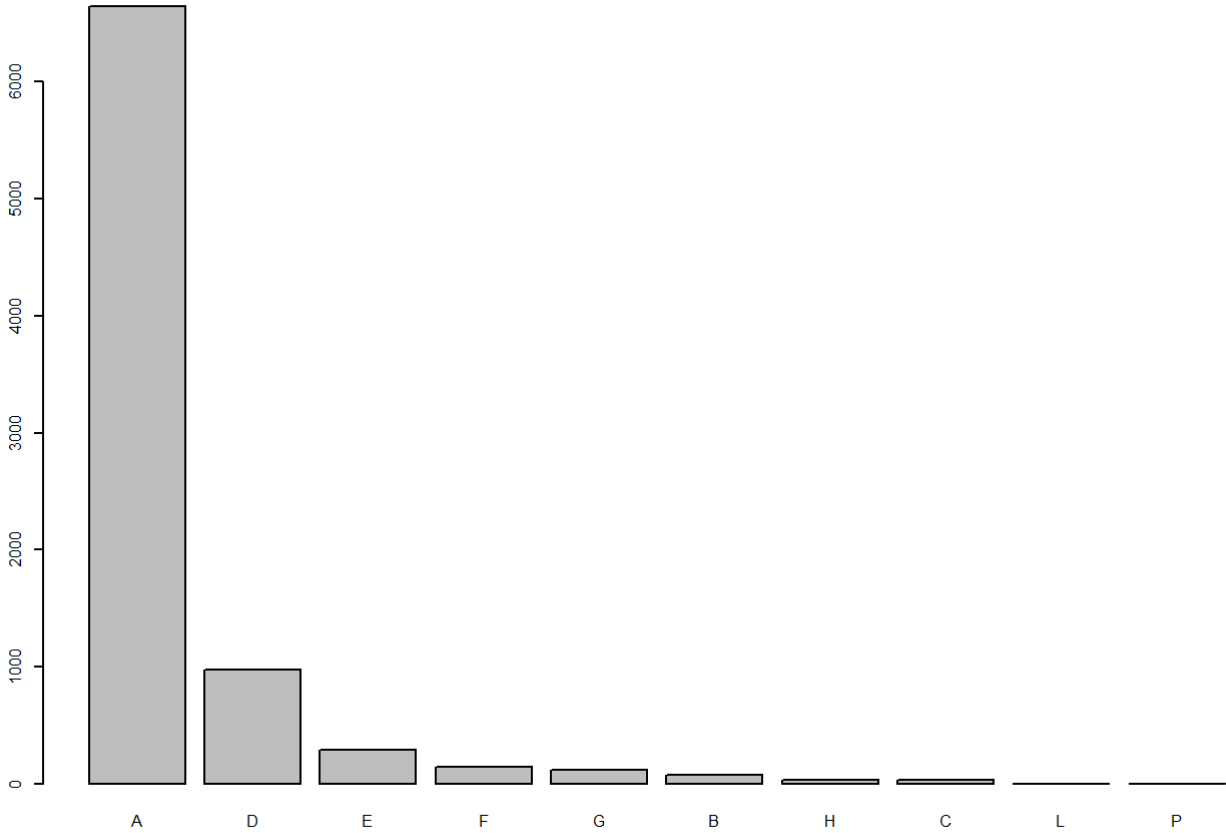
previous_cancellations



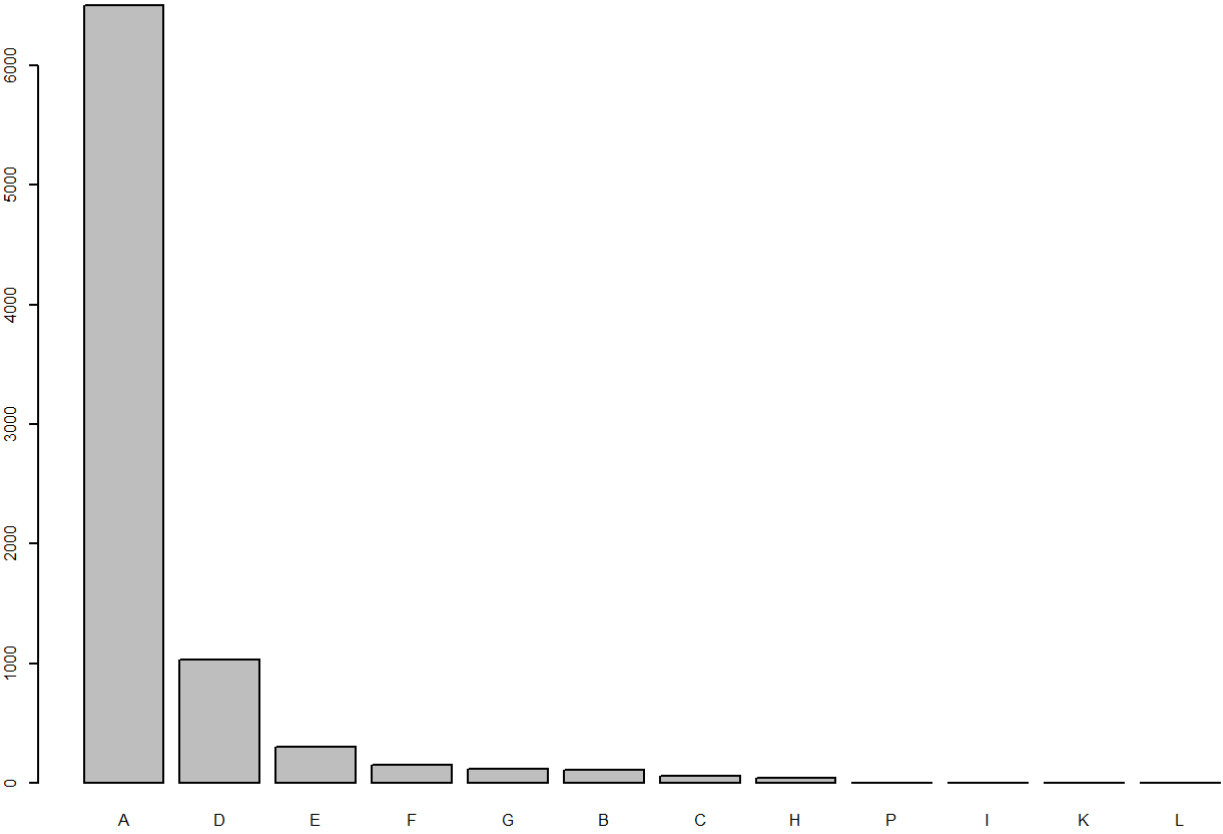
previous_bookings_not_canceled



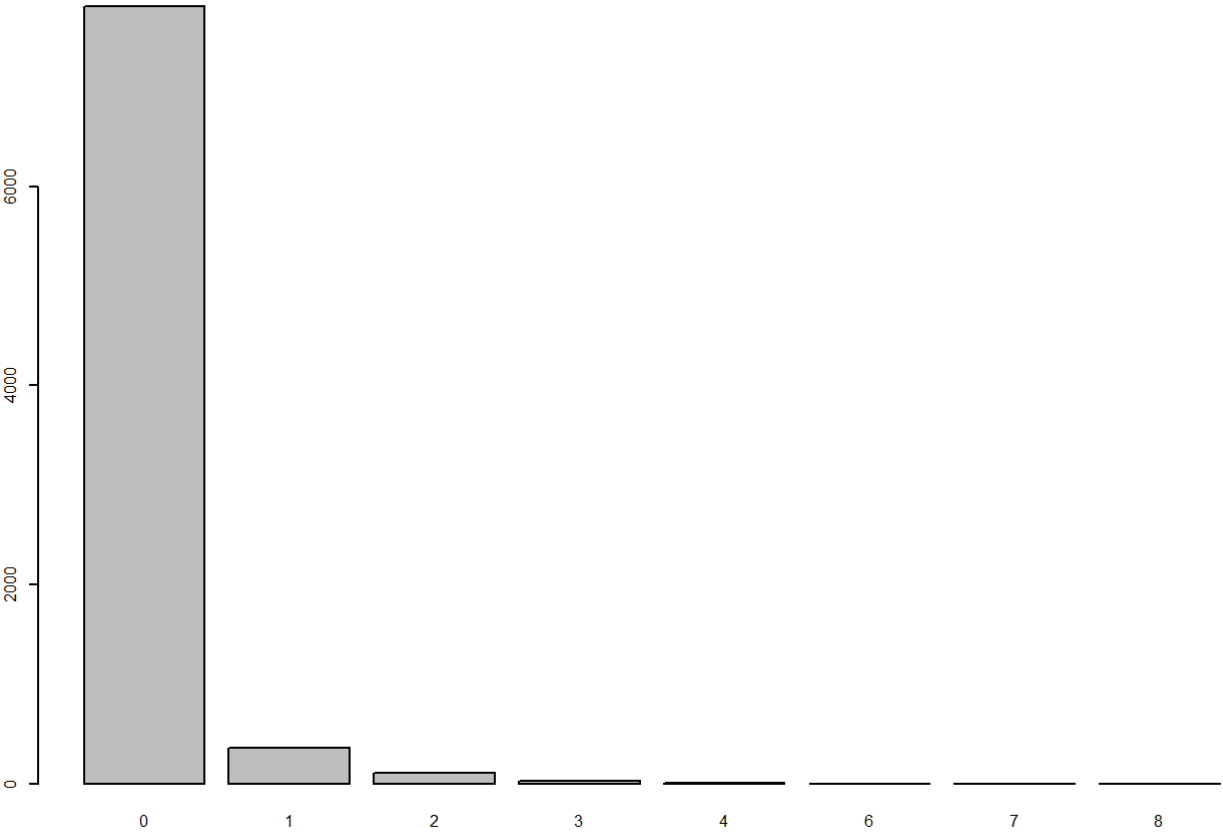
reserved_room_type



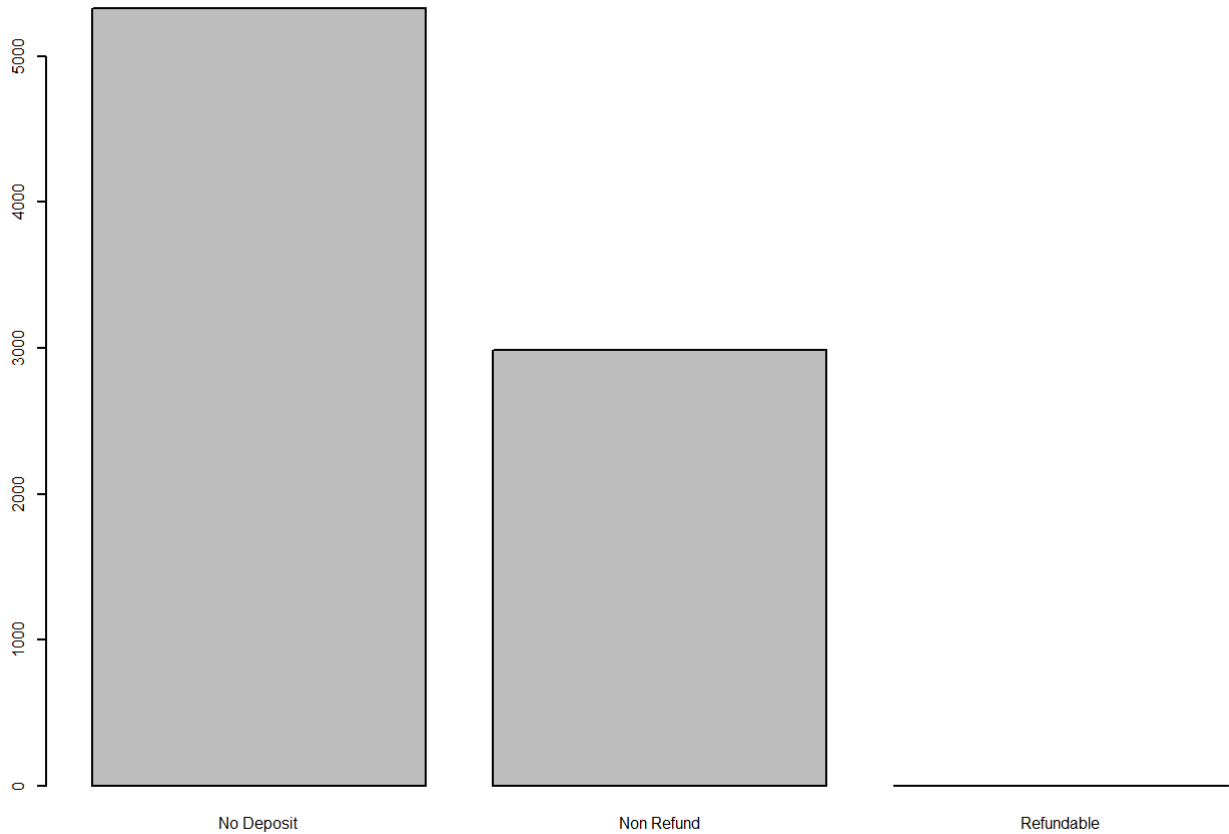
assigned_room_type



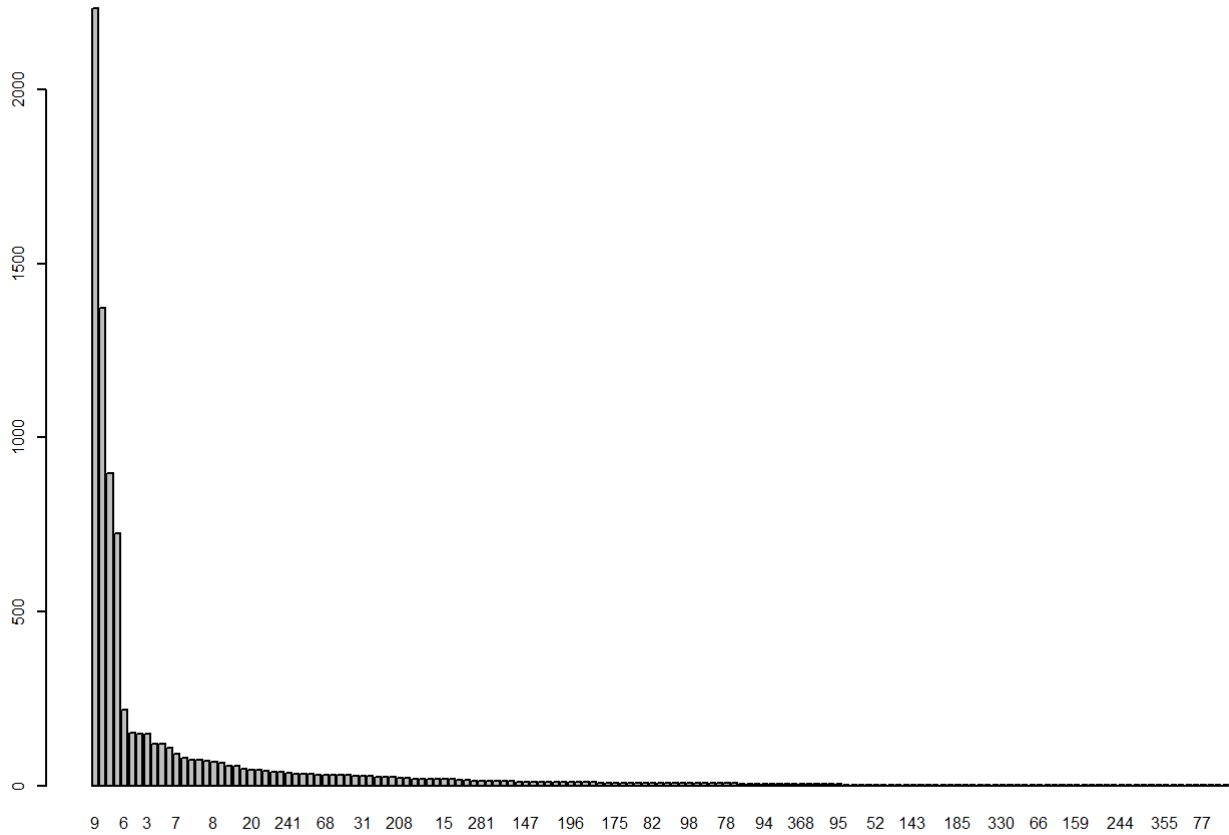
booking_changes



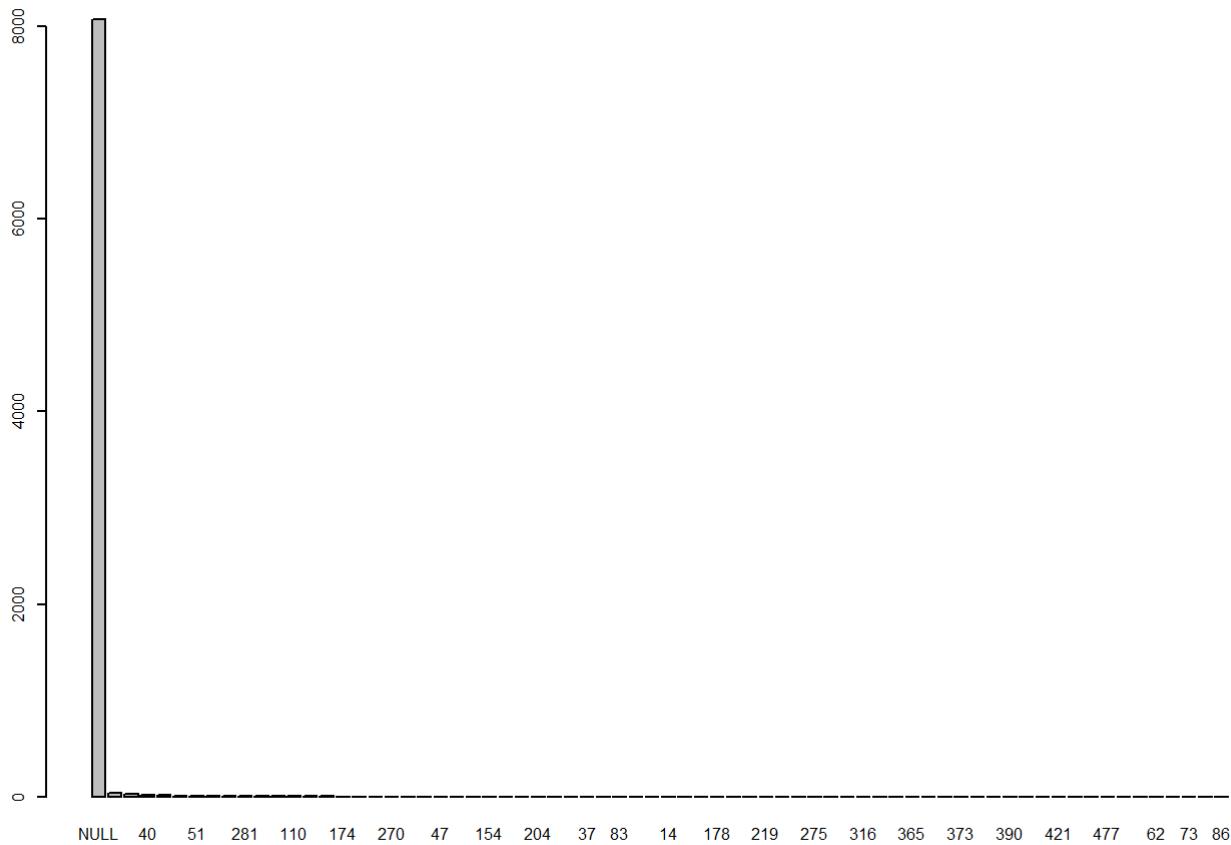
deposit_type



agent



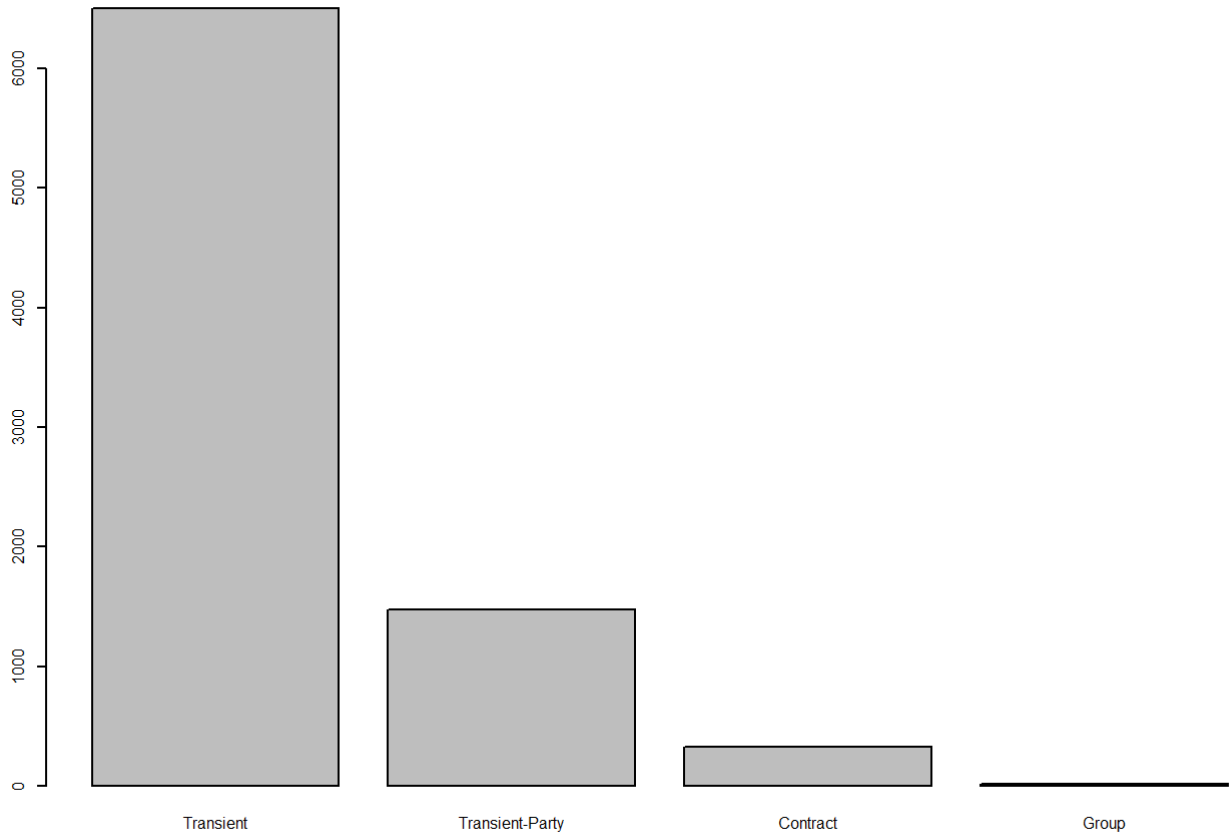
company



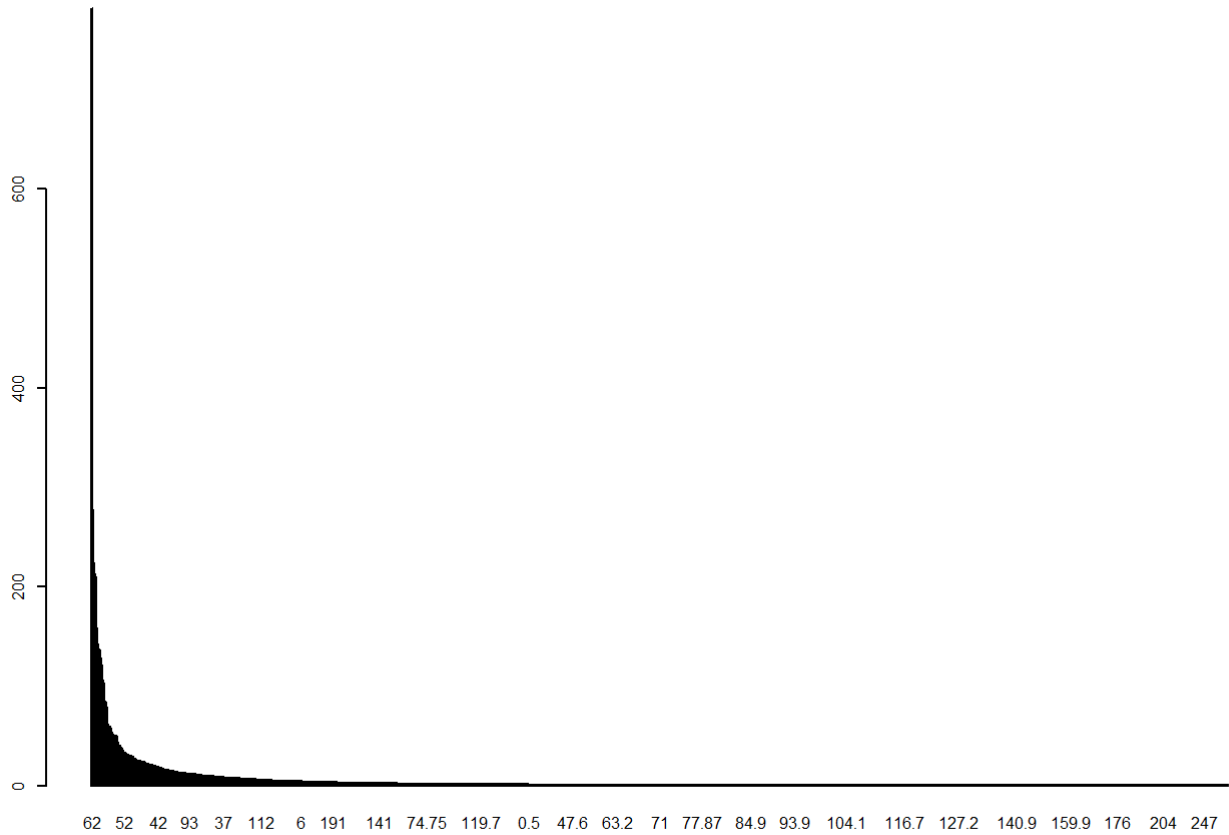
days_in_waiting_list



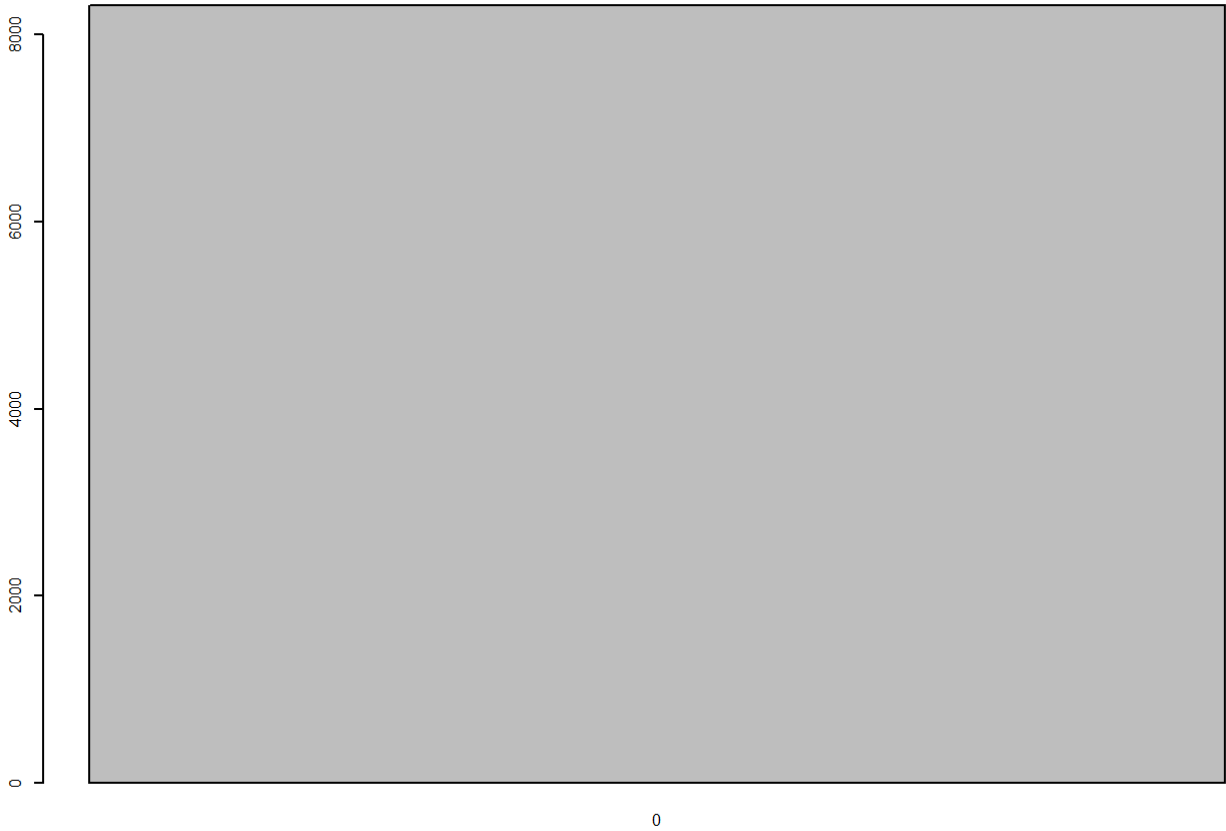
customer_type



adr



required_car_parking_spaces



total_of_special_requests

