

Data Exploratory Analysis

Simone Zanetti

19/6/2018

The following section provides an exploratory analysis of the data available for this project. At first, a generic analysis of the company will be performed to observe the situation and address eventual trends. After that, aggregate data will be created in order to obtain the dependent variable of interest which is the ratio between deliveries and number of worked hours. Each dependent variable is referred to a driver, that is the statistical unit of this analysis. In conclusion, an analysis of the relationship between the y - dependent variable - and each independent variable will be performed in order to analyse the situation.

The following codes need a series of packages to be installed in order to perform the analysis:

```
library(tidyr)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##   filter, lag
```

```
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

```
library(chron)
library(ggmap)
```

```
## Loading required package: ggplot2
```

```
library(ggplot2)
library(ggplot2)
library(ggthemes)
```

```
## Warning: package 'ggthemes' was built under R version 3.4.4
```

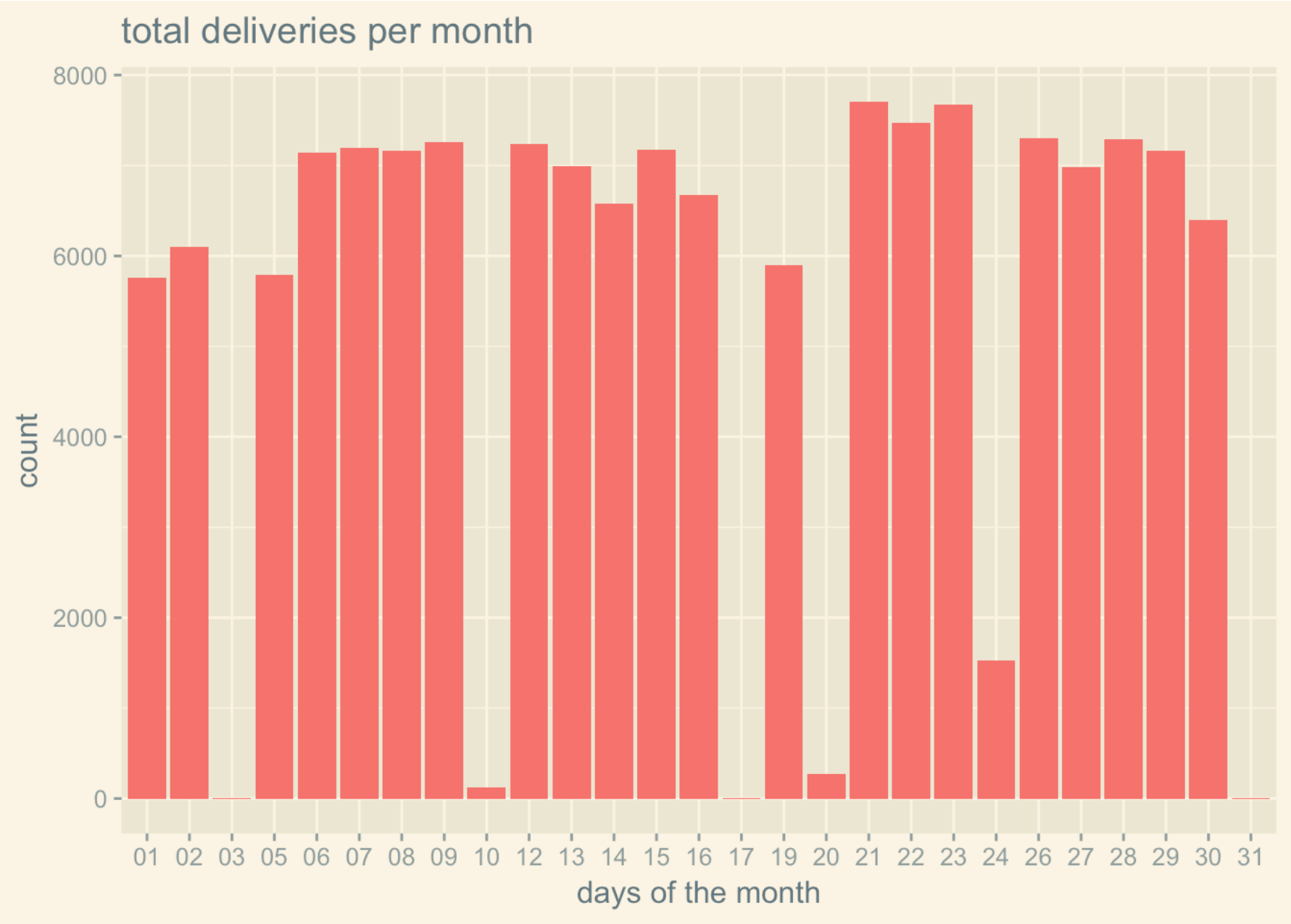
```
load("Data_exploratory.RData")
```

GENERIC TREND OF THE COMPANY

DAYS OF THE MONTH

The data frame and the plot obtained by the codes below can be summarised as following:

##	driver_code	day_deliv	tot_deliveries
##	Length:1882	Length:1882	Min. : 1.00
##	Class :character	Class :character	1st Qu.: 72.00
##	Mode :character	Mode :character	Median : 82.00
##			Mean : 78.05
##			3rd Qu.: 90.00
##			Max. :143.00



The bar chart shows two aspect of the analysis that need to be fixed: 1) For future analysis based on the day of the week there are one Tuesday and One Friday more than any other weekday. 2) The Tuesday 20 corresponds to a strike that took place in Brescia. As a consequence, deliveries that day have been limited and apparently redistributed in the following days causing an increase on the number of deliveries the following days.

This issue can cause problems in the moment where an analysis of the day of the week will be performed and as a consequence solutions will be provided in the following section.

```
stronger_day <- deliveries_day_driver %>% group_by(day_deliv) %>% summarise(tot_de
liveries = sum(tot_deliveries))
stronger_day <- stronger_day %>% arrange(desc(tot_deliveries))

head(stronger_day)
```

```
## # A tibble: 6 x 2
##   day_deliv tot_deliveries
##   <chr>      <dbl>
## 1 21          7704.
## 2 23          7676.
## 3 22          7469.
## 4 26          7303.
## 5 28          7292.
## 6 09          7258.
```

```
tail(stronger_day)
```

```
## # A tibble: 6 x 2
##   day_deliv tot_deliveries
##   <chr>      <dbl>
## 1 24          1529.
## 2 20           277.
## 3 10           118.
## 4 17             9.
## 5 31             3.
## 6 03             2.
```

```
summary(stronger_day)
```

```
##   day_deliv      tot_deliveries
## Length:27      Min.   : 2
## Class :character 1st Qu.:5776
## Mode  :character Median :6984
##                  Mean    :5441
##                  3rd Qu.:7216
##                  Max.    :7704
```

The chart shows the 21, 23 and 22 to be the days with more deliveries and this is apparently coherent with the fact that the 20 March a strike paralyzed the normal activity of the company. The Mean of deliveries per day is 5441 with a Median of 6984 suggesting a distribution skewed to the left.

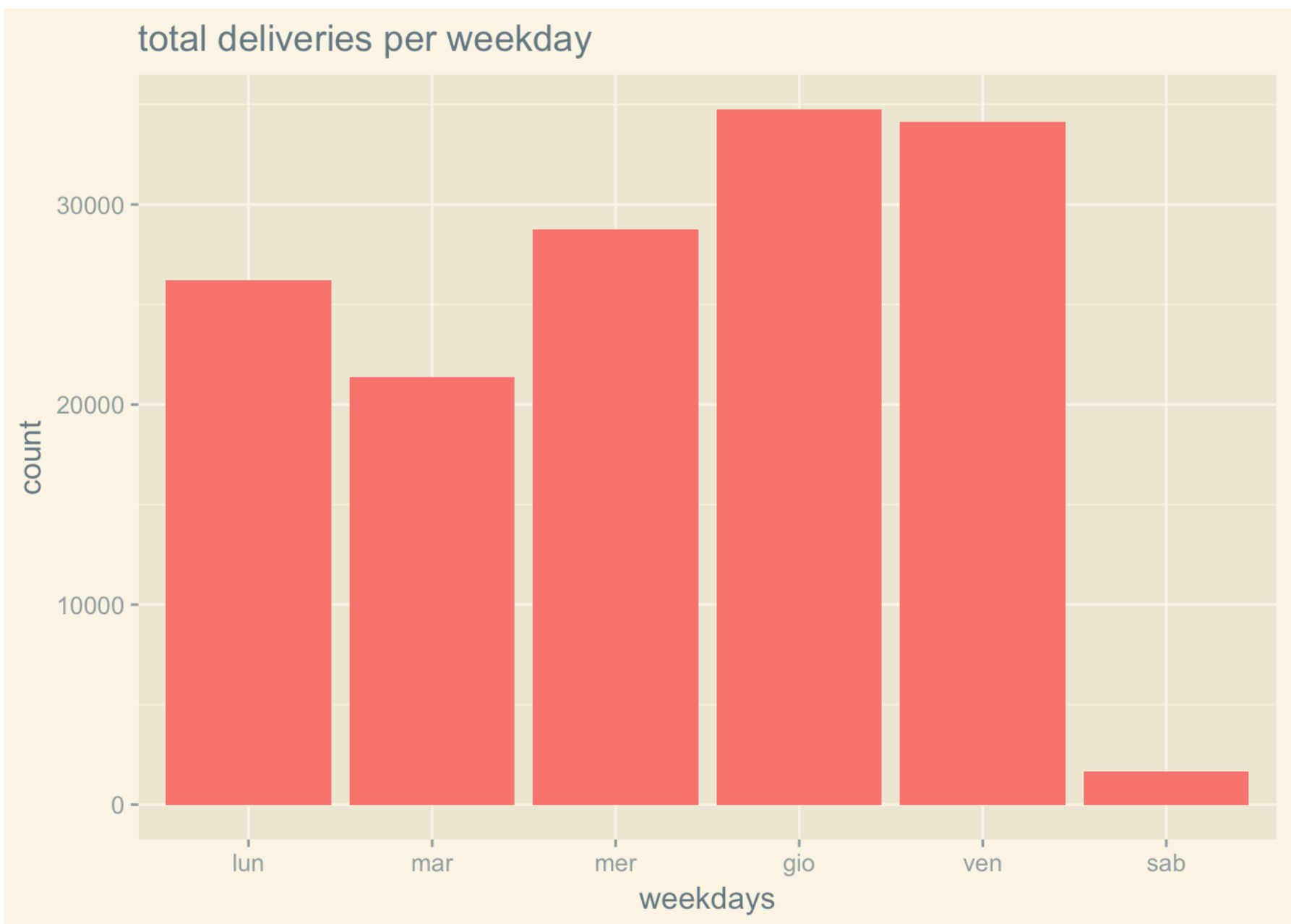
ANALYSIS OF WEEKDAYS

```
data_exploratory$weekday_deliv <- factor(x= data_exploratory$weekday_deliv, levels = c("lun", "mar", "mer", "gio", "ven", "sab"))

stronger_weekday_driver <- data_exploratory %>% group_by(driver_code, weekday_deliv) %>% summarise(tot_delivery = sum(delivery))

stronger_weekday <- data_exploratory %>% group_by(weekday_deliv) %>% summarise(tot_delivery = sum(delivery))

weekday_plot <- ggplot(data = data_exploratory, aes(x = weekday_deliv, fill = "indianred2")) +
  theme_solarized_2() + labs(title = "total deliveries per weekday", x = "weekdays") + guides(fill = FALSE)
weekday_plot + geom_bar()
```

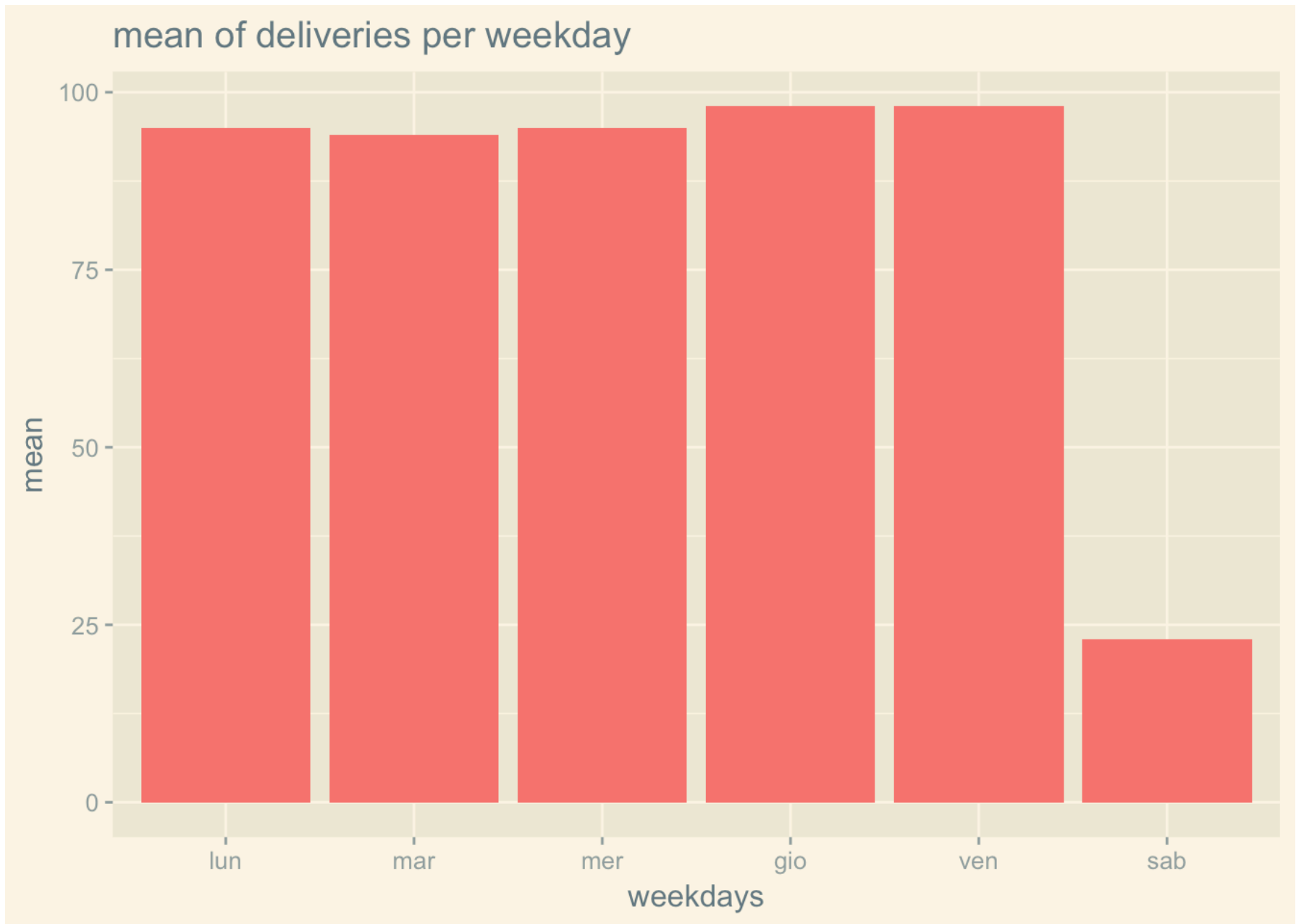


As aforementioned described, these data can be misleading due to the reasons previously described. To overcome this, two possibilities:

1) Use the Mean

```
weekday_mean <- stronger_weekday_driver %>% group_by(weekday_deliv) %>% summarise(
  mean = mean(n()))

weekday_mean_plot <- ggplot(data = weekday_mean, aes(x = weekday_deliv, y = mean,
  fill = "indianred2")) +
  theme_solarized_2()
weekday_mean_plot + geom_col()+ labs(title = "mean of deliveries per weekday", x =
"weekdays") + guides(fill = FALSE)
```



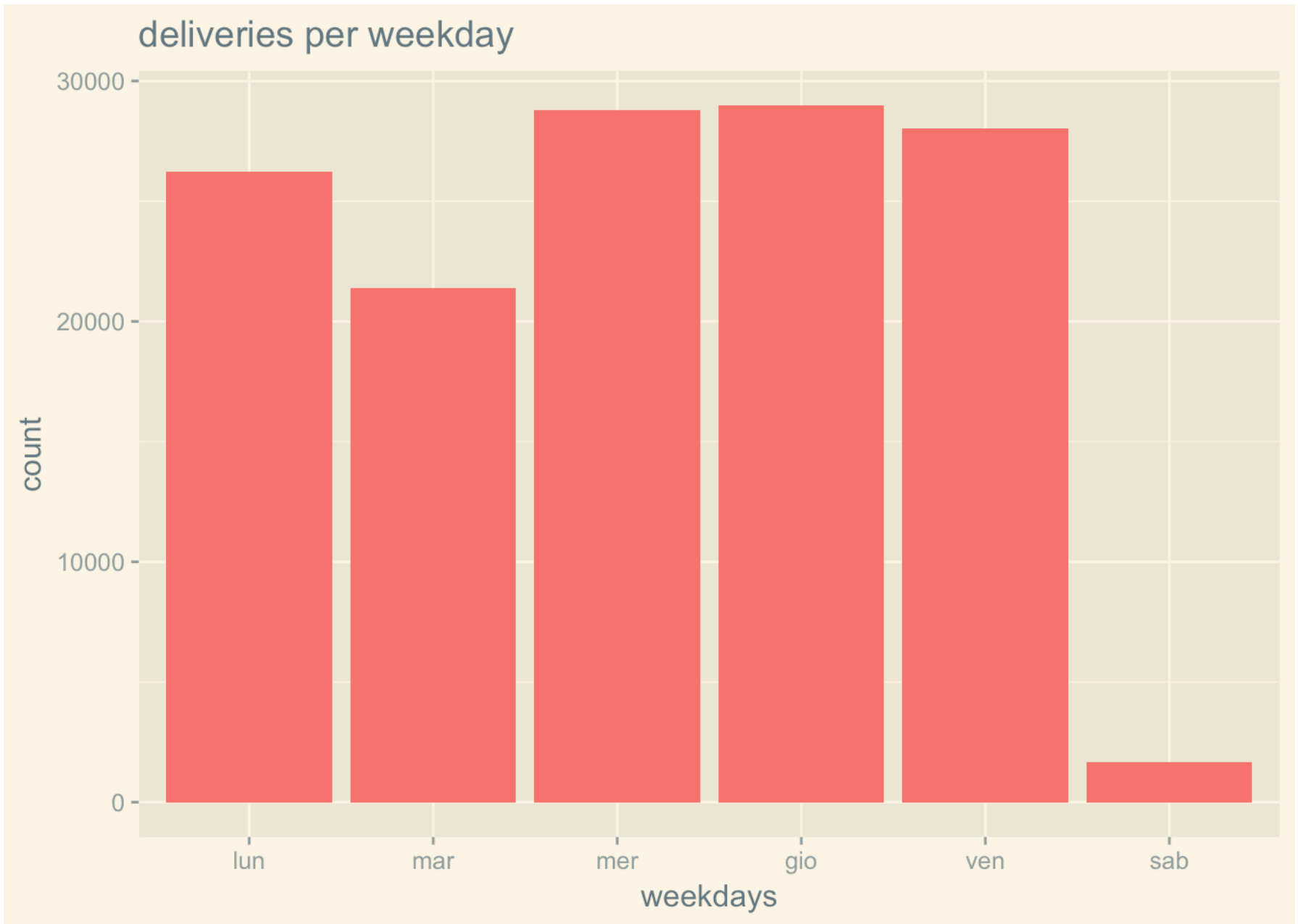
The new graph shows a more leveled situation. The Saturday results the day with less deliveries, as expectable by the fact that the company works only the morning.

2) WORK WITH COMPLETE WEEKS ONLY

I erase the first two days of the month in order to obtain just complete weeks, from Monday to Saturday.

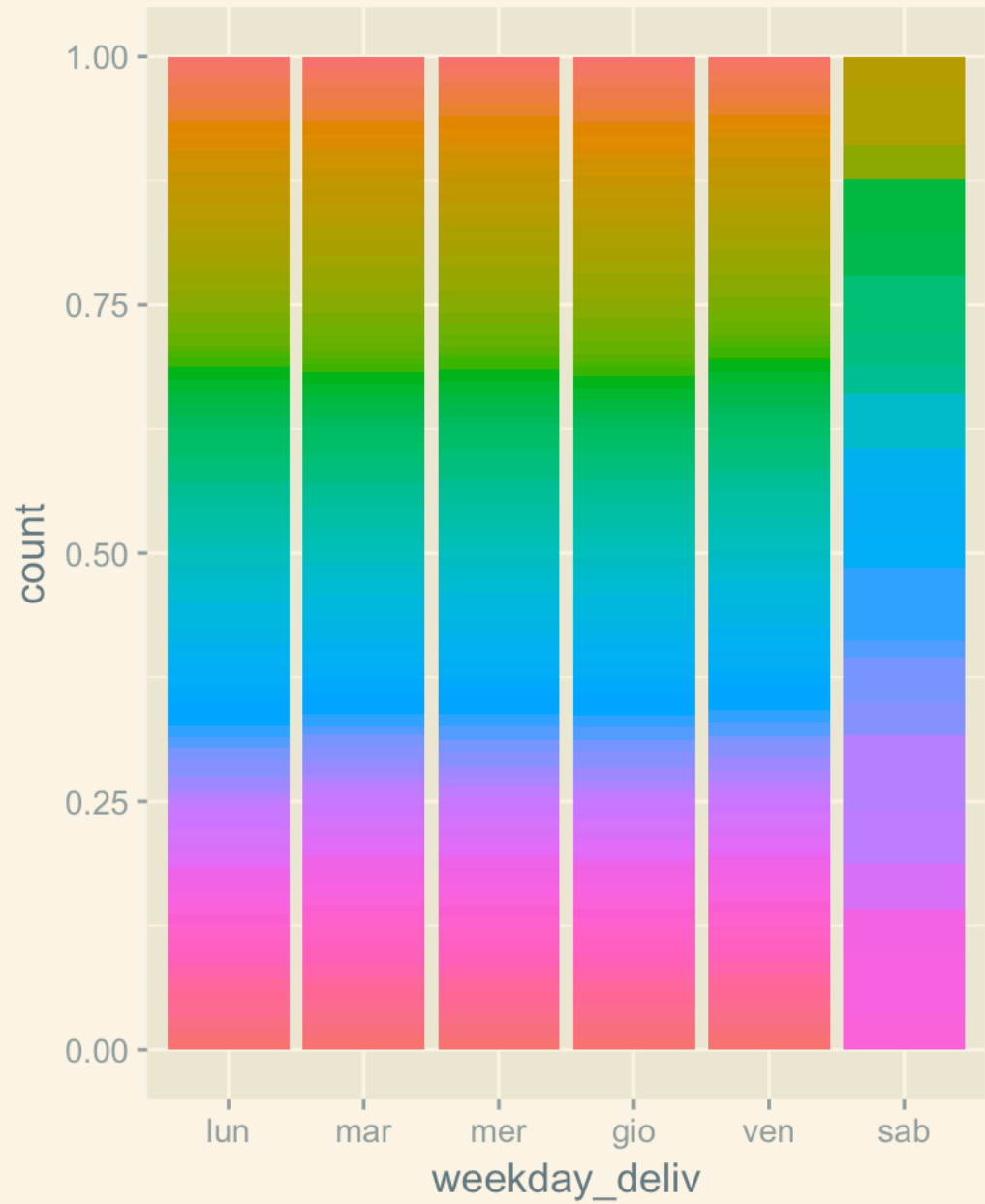
```
data_strongerweekdays <- data_exploratory %>% filter(data_exploratory$day_deliv !=
"01",
                                                    data_exploratory$day_deliv !=
"02",
                                                    data_exploratory$day_deliv !=
"03")

weekday2_plot <- ggplot(data = data_strongerweekdays, aes(x = weekday_deliv, fill
= "indianred2")) +
  theme_solarized_2()
weekday2_plot + geom_bar()+ labs(title = "deliveries per weekday", x = "weekdays")
+ guides(fill = FALSE)
```



The following codes aimed to deepen the analysis to observe if weekdays can be significantly influenced by other variables such as the driver, the day of delivery and the range of hour.

```
weekday2_plot + aes(fill = driver_code) + geom_bar(position = "fill")
```

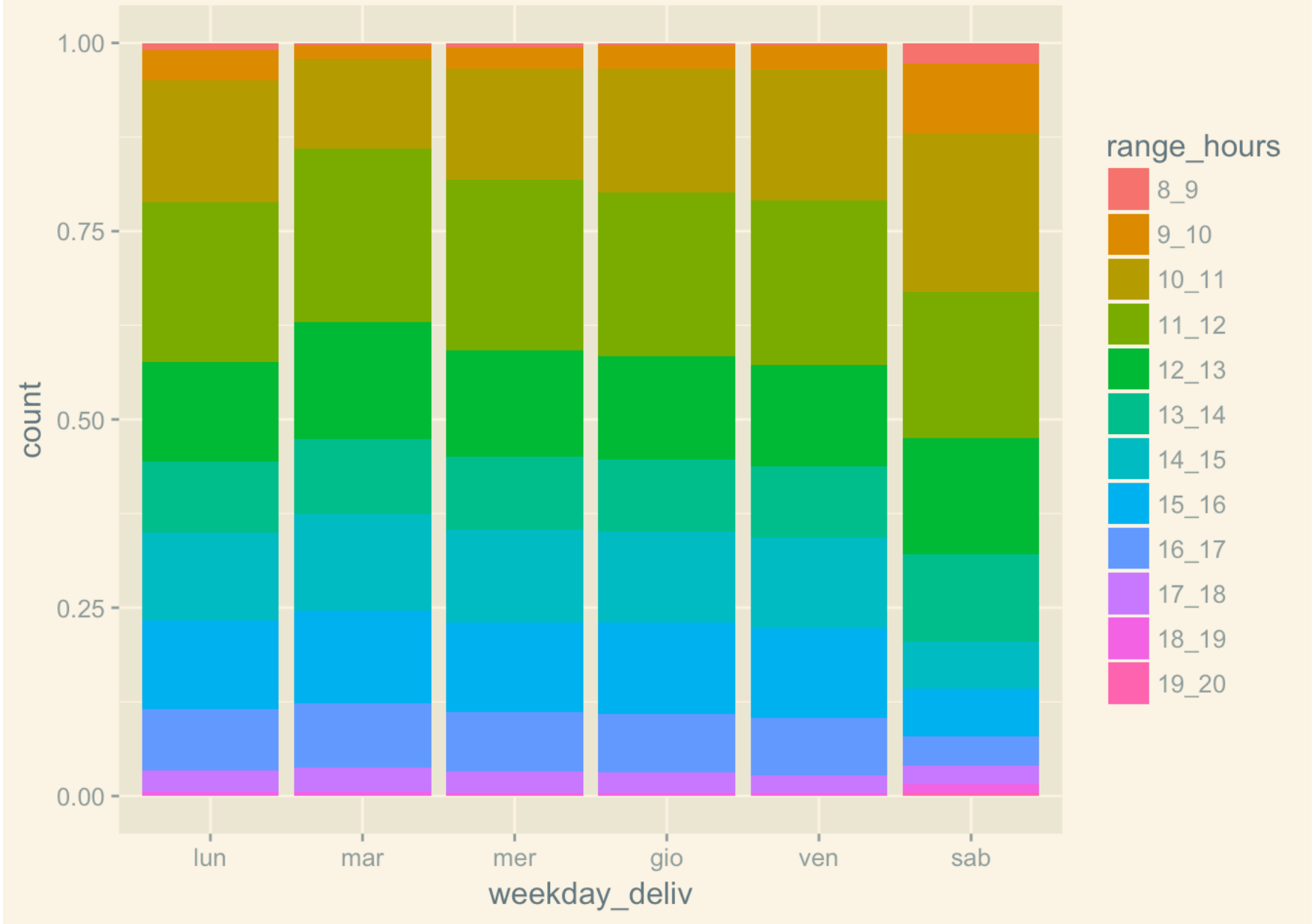


113	306	425	608	678
114	307	426	611	679
115	308	427	612	760
116	309	428	613	761
117	312	429	618	762
119	321	430	621	770
120	322	431	629	771
121	374	432	650	773
122	380	433	651	774
123	381	435	653	775
125	384	436	656	776
126	387	437	657	803
127	388	438	658	804
128	401	529	661	850
129	414	530	662	857
131	415	531	663	858
202	420	532	664	884
266	421	533	666	888
291	422	598	669	972
303	424	599	677	994

```
weekday2_plot + aes(fill = day_deliv) + geom_bar(position = "fill")
```



```
weekday2_plot + aes(fill = range_hours) + geom_bar(position = "fill")
```

ANALYSIS OF THE TIME OF DELIVERY

Creating of Range Hours from the variable pickup_time

```

data_exploratory = data_exploratory %>% mutate(range_hours = sub(pattern = "^08.*"
, replacement = "8_9",x = pickup_time),
range_hours = sub(pattern = "^09.*"
, replacement = "9_10",x = range_hours),
range_hours = sub(pattern = "^10.*"
, replacement = "10_11",x = range_hours),
range_hours = sub(pattern = "^11.*"
, replacement = "11_12",x = range_hours),
range_hours = sub(pattern = "^12.*"
, replacement = "12_13",x = range_hours),
range_hours = sub(pattern = "^13.*"
, replacement = "13_14",x = range_hours),
range_hours = sub(pattern = "^14.*"
, replacement = "14_15",x = range_hours),
range_hours = sub(pattern = "^15.*"
, replacement = "15_16",x = range_hours),
range_hours = sub(pattern = "^16.*"
, replacement = "16_17",x = range_hours),
range_hours = sub(pattern = "^17.*"
, replacement = "17_18",x = range_hours),
range_hours = sub(pattern = "^18.*"
, replacement = "18_19",x = range_hours),
range_hours = sub(pattern = "^19.*"
, replacement = "19_20",x = range_hours),
range_hours = sub(pattern = "^20.*"
, replacement = "20_21",x = range_hours))

```

```
unique(data_exploratory$range_hours)
```

```

## [1] "15_16" "16_17" "11_12" "14_15" "13_14" "12_13" "17_18" "10_11"
## [9] "9_10"  "18_19" "8_9"   "19_20"

```

```
data_exploratory = data_exploratory %>% filter(range_hours != "20_21" & range_hours != "06:51:00")
```

```
data_exploratory$range_hours <- factor(x = data_exploratory$range_hours, levels = c("8_9","9_10","10_11","11_12","12_13","13_14","14_15","15_16","16_17","17_18","18_19", "19_20"))
```

```

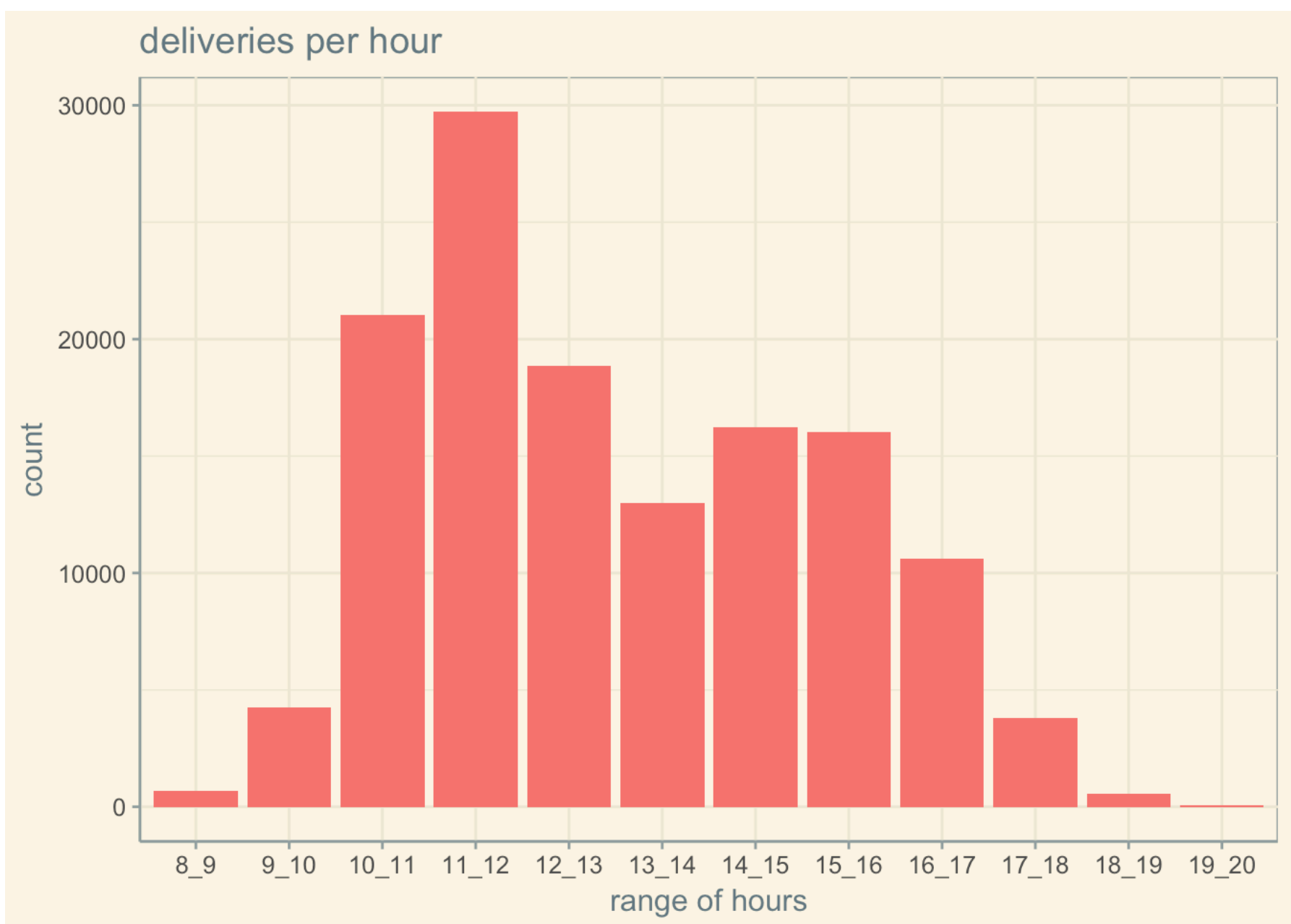
stronger_hour <- data_exploratory %>% group_by(driver_code,range_hours) %>% summar
ise(tot_deliveries = sum(delivery))

data_strongerhours <- data_exploratory %>% filter(data_exploratory$day_deliv != "0
1",
                                                    data_exploratory$day_deliv != "0
2",
                                                    data_exploratory$day_deliv != "2
0")

strongerhour_plot <- ggplot(data = data_strongerhours, aes(x = range_hours, fill=
"indianared2")) + theme_solarized()+
  labs(title = "deliveries per hour", x = "range of hours")

strongerhour_plot + geom_bar() + guides(fill = FALSE)

```



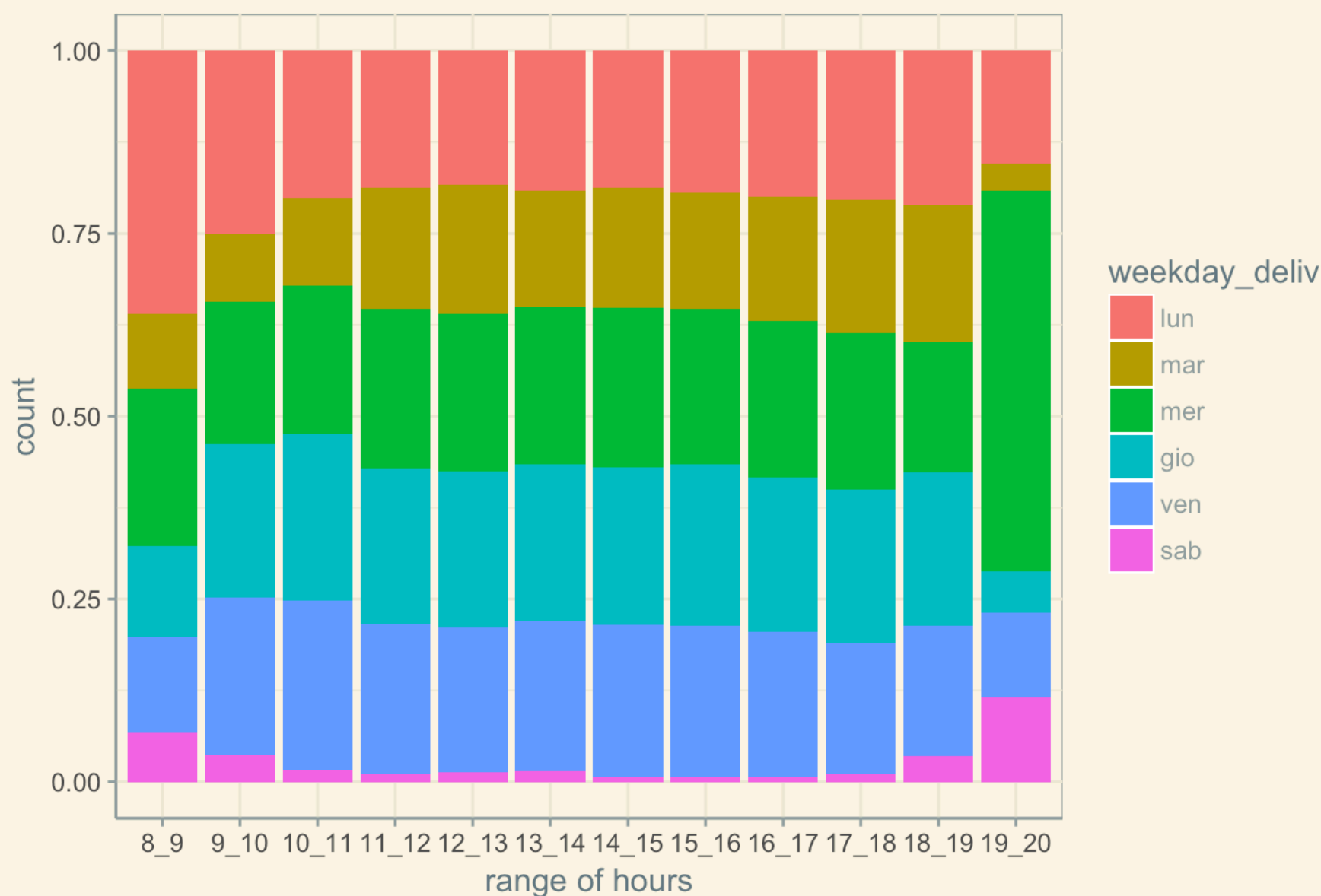
The plot shows a peak of deliveries on the range between 11 am and 12 pm. After that range, the curve regularly decreases until the end of the day with the exception of the range between 13 and 14, where the number of deliveries decreases before to increase again on the following range. From this point of view, the morning between 10 and 13 the majority of the deliveries are done.

```

strongerhour_plot + aes(fill = weekday_deliv) + geom_bar(position = "fill")

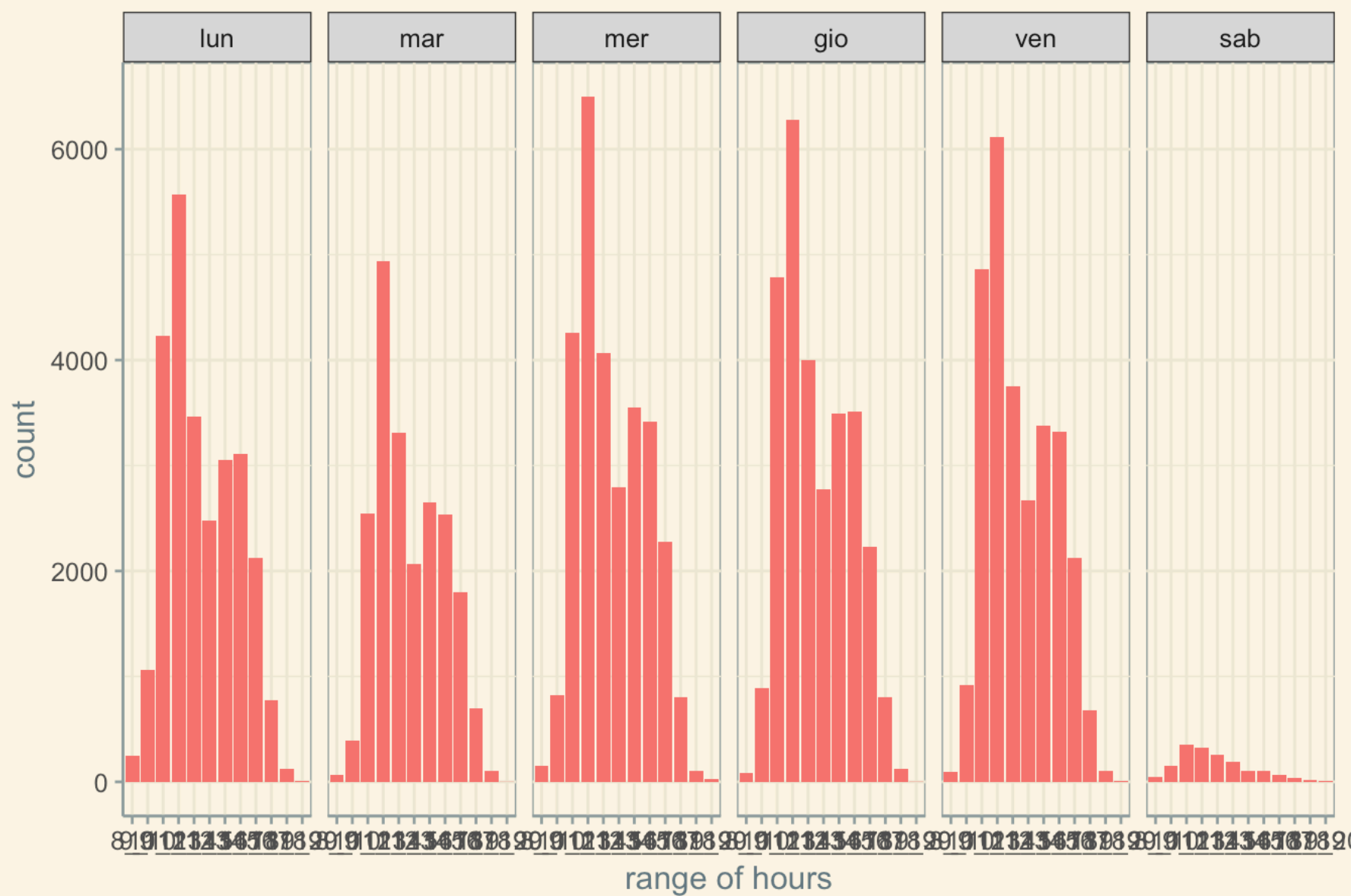
```

deliveries per hour

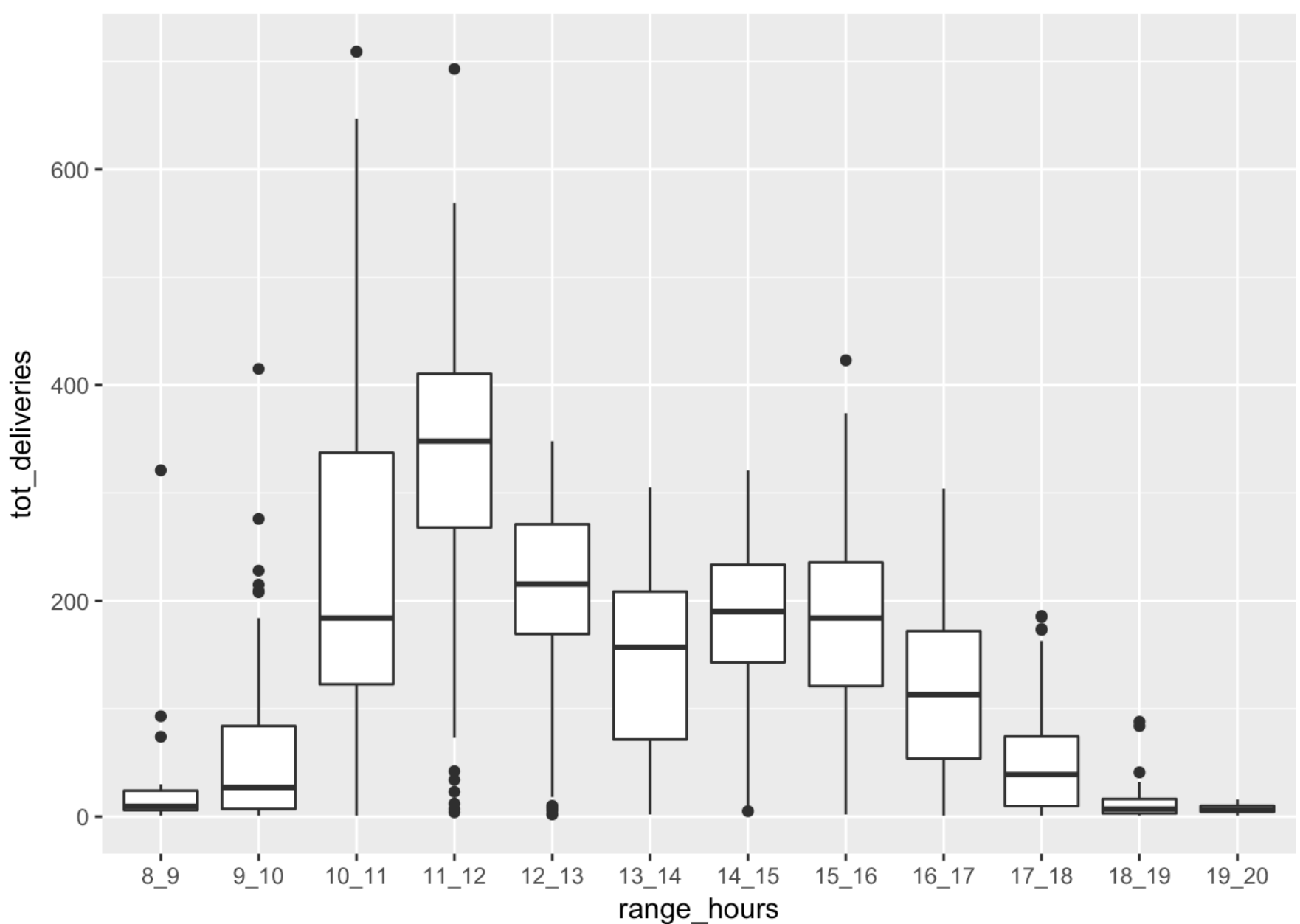


```
strongerhour_plot + geom_bar() + facet_grid(.~data_strongerhours$weekday_deliv)+
guides(fill = FALSE)
```

deliveries per hour



```
ggplot(data = stronger_hour, aes(x = range_hours, y = tot_deliveries))+ geom_boxplot()
```



The trend seems to suggest a different attitude of deliveries during the day based on the type of day of the week. In particular, Monday is the day where the majority of deliveries between 8 and 9 am are performed.

ANALYSIS OF THE AREA OF DELIVERY

```
stronger_area <- data_exploratory %>% group_by(postal_code) %>% summarise(cnt = n(
)) %>%
  arrange(desc(cnt))

summary(stronger_area)
```

```
## postal_code      cnt
## Length:76      Min.   :    1.0
## Class :character 1st Qu.: 312.8
## Mode  :character Median : 1024.5
##                  Mean    : 1932.8
##                  3rd Qu.: 2036.5
##                  Max.    :31846.0
```

```
head(stronger_area)
```

```
## # A tibble: 6 x 2
##   postal_code    cnt
##   <chr>        <int>
## 1 25121/25136 31846
## 2 25020        9111
## 3 25080        9103
## 4 25010        6531
## 5 25030        5775
## 6 25040        4819
```

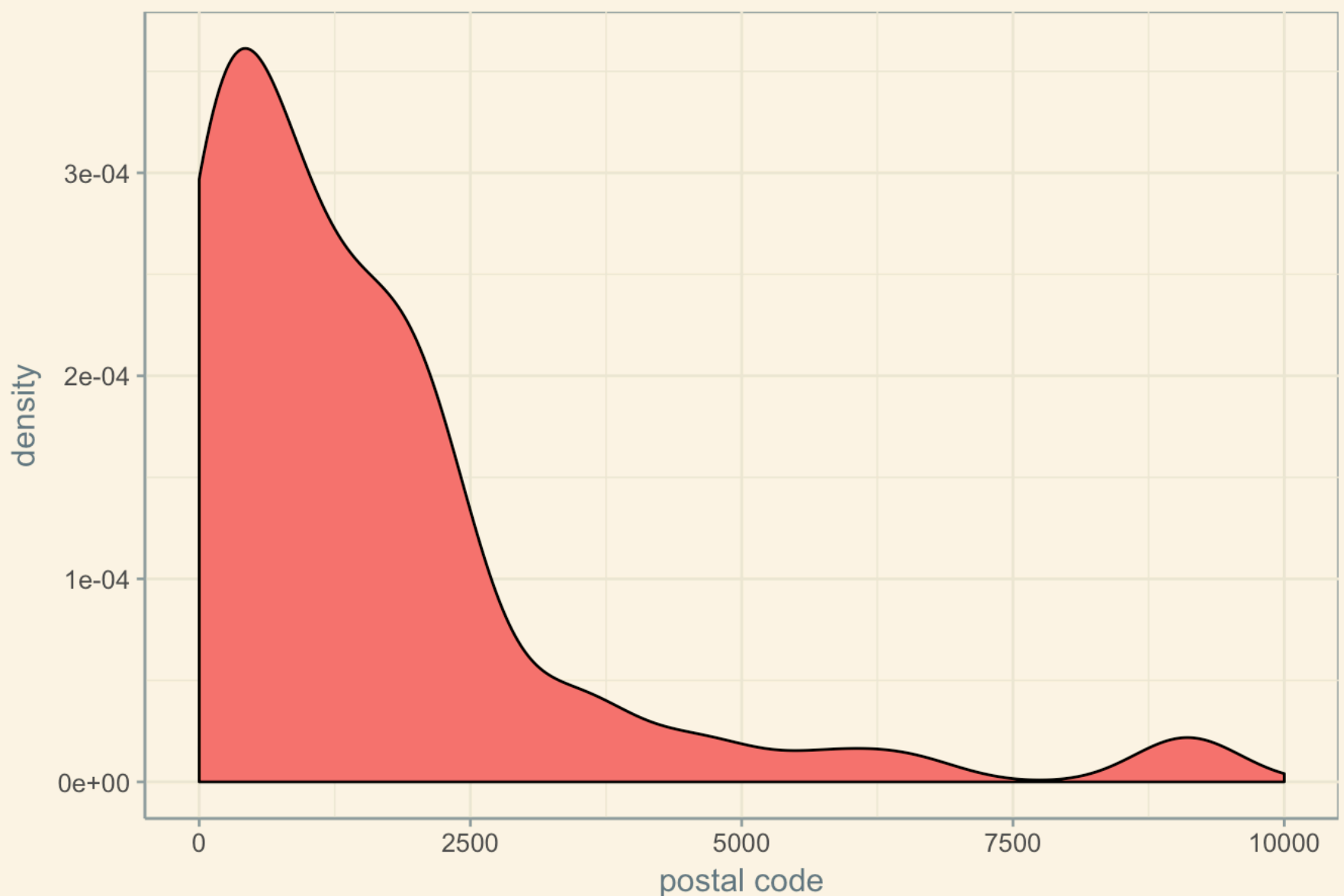
```
tail(stronger_area)
```

```
## # A tibble: 6 x 2
##   postal_code    cnt
##   <chr>        <int>
## 1 25129         6
## 2 25127         4
## 3 unknown       4
## 4 25126         3
## 5 25128         3
## 6 25046         1
```

```
strongerarea_plot <- ggplot(data = stronger_area, aes(x = cnt, fill= "indianared2"
)) + theme_solarized()+
  labs(title = "deliveries per area", x = "postal code")
strongerarea_plot + geom_density()+guides(fill = F)+ scale_x_continuous(limits = c
(0,10000))
```

```
## Warning: Removed 1 rows containing non-finite values (stat_density).
```

deliveries per area



As shown below, we can observe a right skewed distribution whose Mean of deliveries per Area is 1932.8. However, due to the fact that for the center of Brescia one postal code because was not possible to obtain, it has been identified with the observation “25121/25136”. This represents an important outlier of this distribution. For this, the number of monthly deliveries is 31846, against the second postal code with maximum number of deliveries which is 9111. This postal code is (an area south of Brescia which includes different important district of the city). Third postal code is 25080 which is an area close to the Lake Garda, where commercial activities are really frequent. In this case, due to the aforementioned outlier the median of 1024.5 is more indicated to describe the tendency of this trend.

SPLIT AREA

The aim of this paragraph is to split postal codes into 3 different areas based on the distance between them and the center of Brescia. As imaginable, one of the factor that can really influence the performance of the drivers is the presence of traffic, and the distance of the points of delivery between each others. From this point of view, it will be expected to have more traffic and poins of delivery closer to each other in the center of Brescia and the opposite further from the city center.

— — —

- At first, I will have to obtain the coordinates of each postal code.
- Secondly, I will obtain the distances between points and the center of Brescia.

ANALYSIS OF INDEPENDENT VARIABLE

In the following section I will aggregate data in order to obtain the independent variable of this analysis, which is the relationship between deliveries and hours worked for each driver, that represent my statistical unit. By the end of this paragraph a data frame named `aggregate_data_last` will be obtained in order to conduct a machine learning analysis later on.

OBTAIN THE HOURS WORKED FOR EACH DRIVER AND RATIO DELIVERIES PER DAY WORKED

```
data_exploratory <- na.omit(data_exploratory)
data_exploratory <- data_exploratory %>% filter(driver_code != "208" & driver_code
!= "234"& driver_code != "260" & driver_code != "336" & driver_code != "404" & dri
ver_code != "534" & driver_code != "535" & driver_code != "623"  & driver_code !=
"132")
aggregate_data <- data_exploratory %>% group_by(driver_code, day_deliv, pickup_tim
e, delivery) %>% summarise()

# aggregate_data <- aggregate_data %>% filter(driver_code != "208" & driver_code !=
"234"& driver_code != "260" & driver_code != "336" & driver_code != "404" & driv
er_code != "534" & driver_code != "535" & driver_code != "623"  & driver_code != "
132")

aggregate_data$pickup_time <- as.character(aggregate_data$pickup_time)
aggregate_data$pickup_time = as.POSIXlt(aggregate_data$pickup_time, format = "%H:%
M:%S")

class(aggregate_data$pickup_time)
```

```
## [1] "POSIXlt" "POSIXt"
```

```

time_delivery <- rep(0,nrow(aggregate_data))

# Turning time into the difference of minutes between every point of delivery for
every driver

for (i in 1:(nrow(aggregate_data)-1)) {
  if(aggregate_data$driver_code[i] == aggregate_data$driver_code[i+1])
    {time_delivery[i] = difftime(aggregate_data$pickup_time[i], aggregate_data$picku
p_time[i+1], units = "mins" ) }
}

time_delivery=ifelse(time_delivery>0,0,time_delivery)
aggregate_data=aggregate_data[,-3]
aggregate_data=cbind(aggregate_data,time_delivery)
aggregate_data$time_delivery <- abs(aggregate_data$time_delivery)

# sum the worked minutes per driver(minutes worked) and create a new variable base
d on hours (hours worked per driver)

aggregate_data2 <- aggregate_data %>% group_by(driver_code) %>% summarise(tot_deli
veries = sum(delivery), tot_minutes = sum(time_delivery))
hours_delivery <- aggregate_data2$tot_minutes/60
aggregate_data2 <- cbind(aggregate_data2, hours_delivery)

y = aggregate_data2$tot_deliveries / aggregate_data2$hours_delivery
aggregate_data2 <- cbind(aggregate_data2, y)

# aggregate_data2 <- aggregate_data2 %>% filter(driver_code != "208" & driver_code
!= "234"& driver_code != "260" & driver_code != "336" & driver_code != "404" & dri
ver_code != "534" & driver_code != "535" & driver_code != "623" & driver_code !=
"132")

# Now that I have the total of deliveries per driver and the Independent variable
I can plot them to analyse the trend.

summary(aggregate_data2$y)

```

```

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.4167  8.7378  9.8758  9.4411 10.7959 13.8572

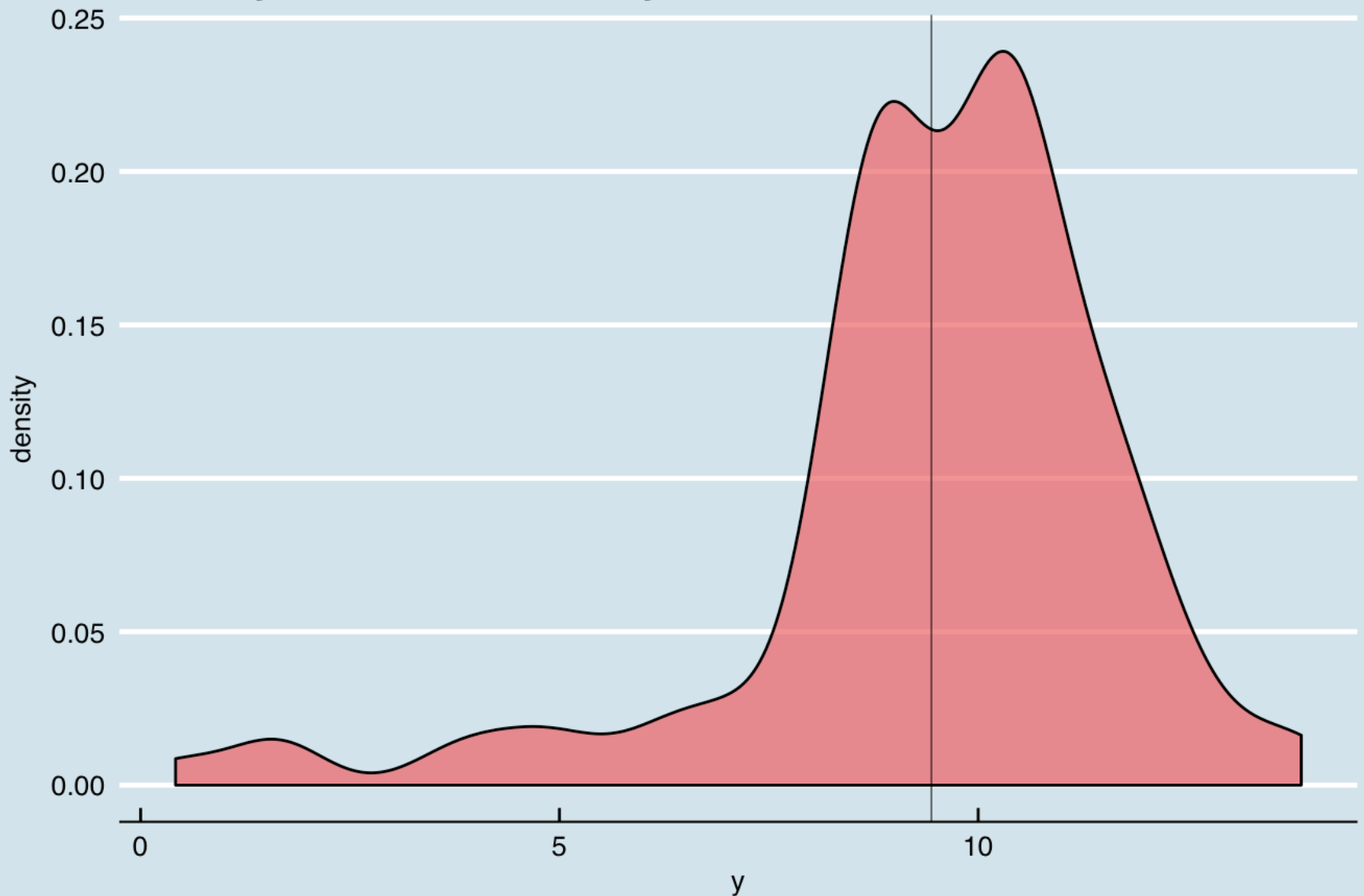
```

```

y_plot <- ggplot(data = aggregate_data2, aes(x = y))
y_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist() + labs
(title = "delivery hours ratio density")+ geom_vline(xintercept = 9.4412, size = 0
.2)

```

delivery hours ratio density

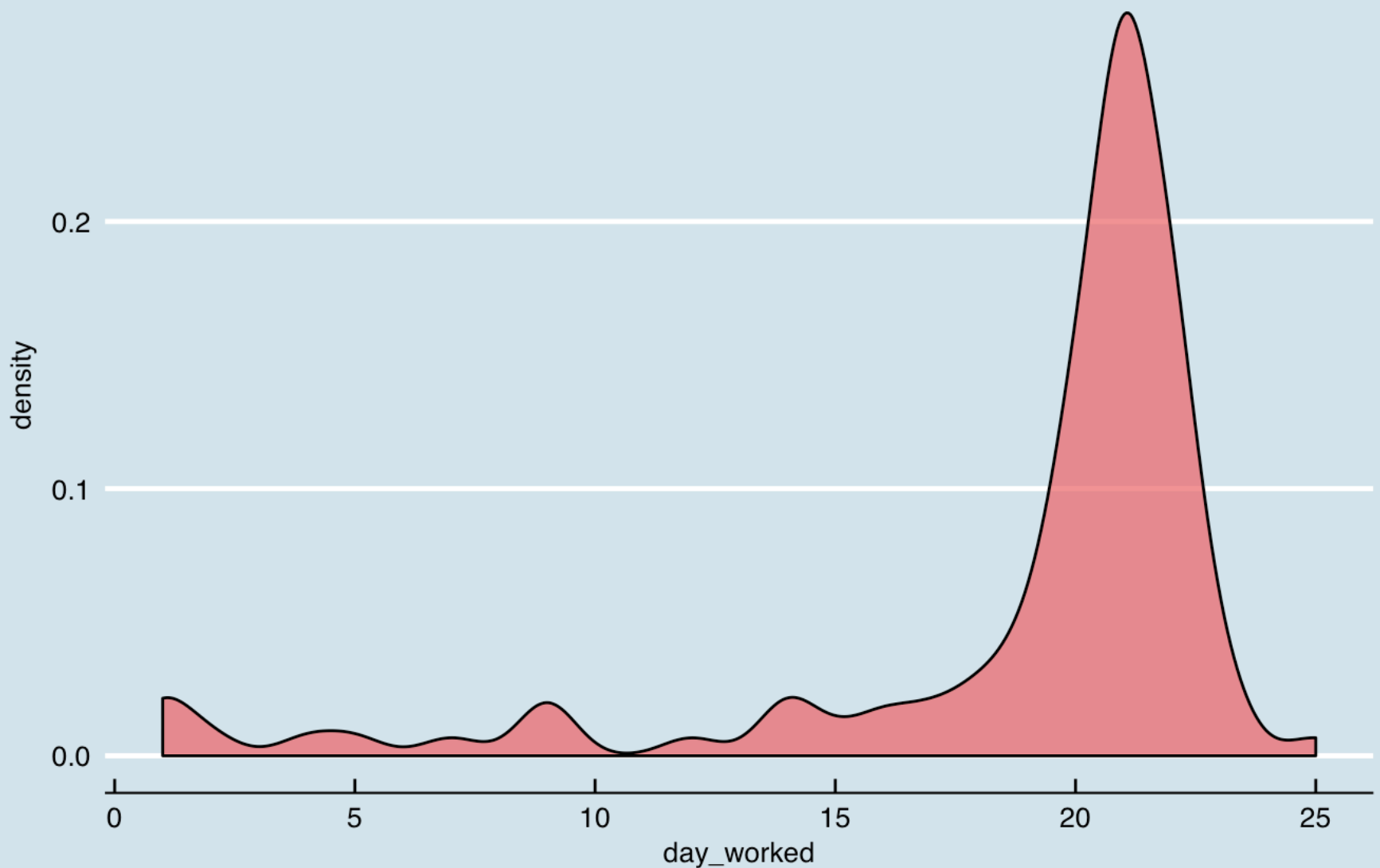


The dependent variable y is represented by a left tailed distribution, with a mean of 9.4412 deliveries per hours. 68% of the drivers owns a delivery hours ratio between ...

```
aggregate_data$day_deliv <- as.numeric(aggregate_data$day_deliv)
aggregate_data3 <- data_exploratory %>% group_by(driver_code) %>% summarise(day_worked = n_distinct(day_deliv))
# aggregate_data3 <- aggregate_data3 %>% filter(driver_code != "208" & driver_code != "234" & driver_code != "260" & driver_code != "336" & driver_code != "404" & driver_code != "534" & driver_code != "535" & driver_code != "623" & driver_code != "132")
aggregate_data2 <- cbind(aggregate_data2, aggregate_data3$day_worked)

dayworked_plot <- ggplot(data = aggregate_data3, aes(x = day_worked))
dayworked_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist() + labs(title = "day worked density")
```

day worked density



The curve ...

TOTAL PACK LOADED PER DRIVER

```
load("aggregate_data4.Rdata") # Load a dataset previously cleaned during the wrangling phase. Dataset provided by the client later on
```

```
# clean the data for the needs of this analysis
```

```
aggregate_data4 <- aggregate_data4 %>% filter(driver_code != "208" & driver_code != "234" & driver_code != "260" & driver_code != "336" & driver_code != "404" & driver_code != "534" & driver_code != "535" & driver_code != "623" & driver_code != "132",
```

```
driver_code != "1000", driver_code != "101", driver_code != "402")
```

```
aggregate_data4 <- aggregate_data4 %>% mutate(driver_code = gsub(pattern = "8421", replacement = "421", x = driver_code),
```

```
driver_code = gsub(pattern = "8618", replacement = "618", x = driver_code),
```

```
driver_code = gsub(pattern = "8678", replacement = "678", x = driver_code),
```

```
driver_code = gsub(pattern = "8679", replacement = "679", x = driver_code),
```

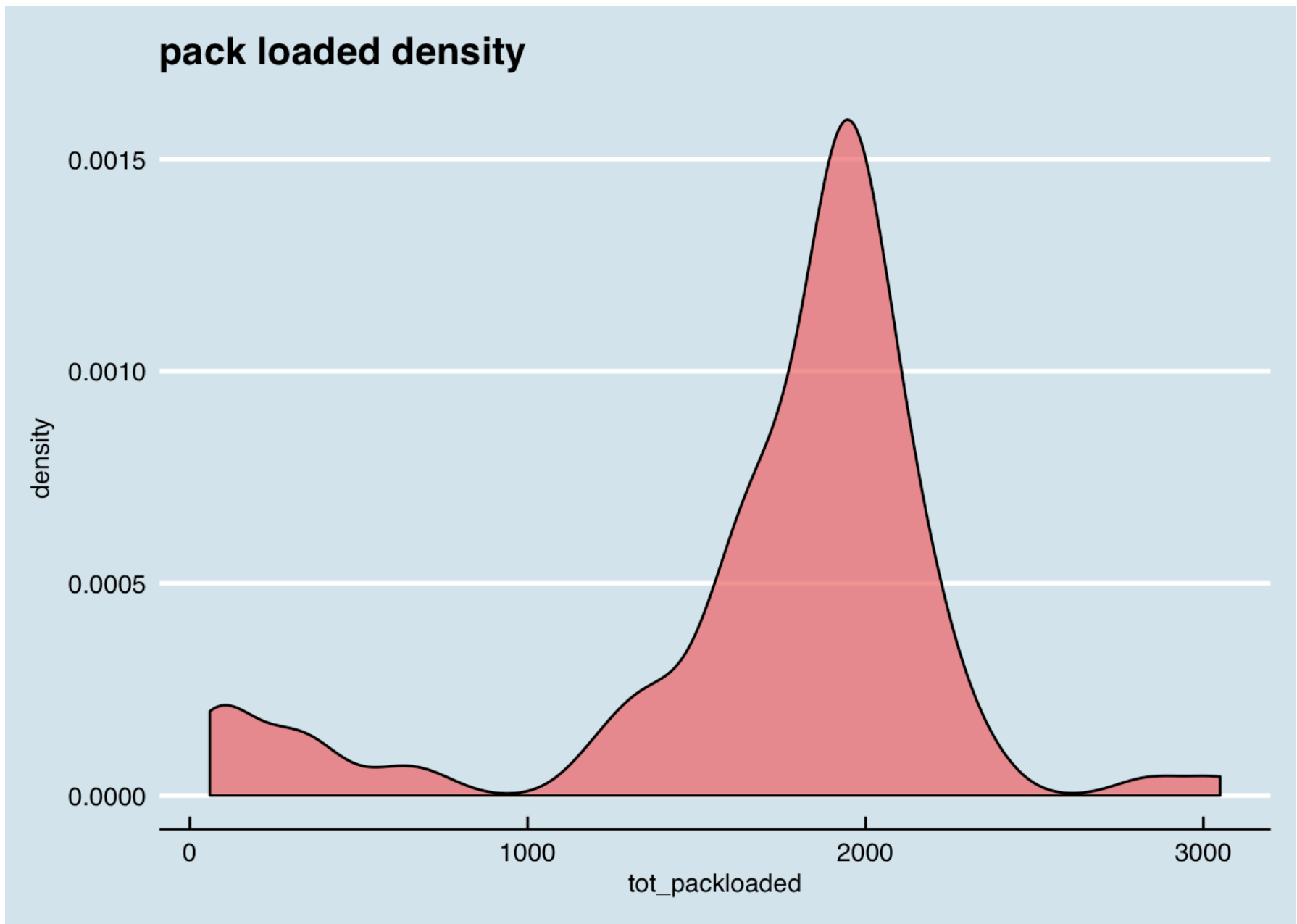
```
driver_code = gsub(pattern = "8531", replacement = "531", x = driver_code))
```

```
aggregate_data5 <- aggregate_data4 %>% group_by(driver_code ) %>% summarise(tot_packloaded = sum(pack_loaded))
setdiff(aggregate_data5$driver_code, aggregate_data2$driver_code)
```

```
## [1] "805" "851"
```

```
aggregate_data5 <- aggregate_data5 %>% filter(tot_packloaded > 0, driver_code != "851", driver_code != "805")
```

```
packloaded_plot <- ggplot(data = aggregate_data5, aes(x = tot_packloaded))
packloaded_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist() + labs(title = "pack loaded density")
```



The curve...

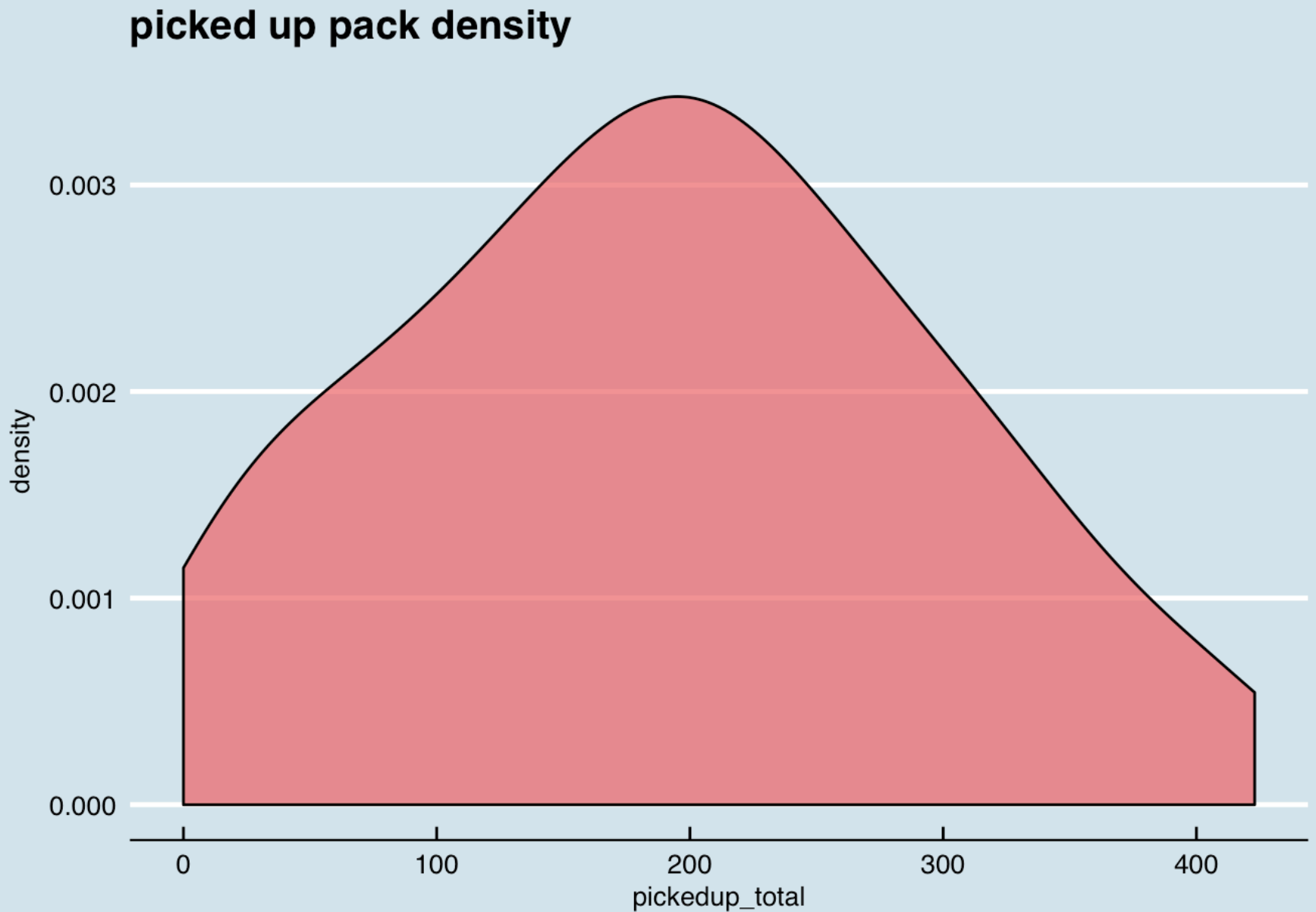
TOTAL PICKED UP PACKS PER DRIVER

```

aggregate_data6 <- aggregate_data4 %>% group_by(driver_code) %>% summarise(pickedu
p_total = sum(pickup_services))
aggregate_data6 <- aggregate_data6 %>% filter( driver_code != "851", driver_code !
= "805")

pickedup_plot <- ggplot(data = aggregate_data6, aes(x = pickedup_total ))
pickedup_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist()
+ labs(title = "picked up pack density")

```



TOTAL ARRIVED PACKS PER DRIVER

```

aggregate_data7 <- aggregate_data4 %>% group_by(driver_code) %>% summarise(packarr
ived_total = sum(pack_arrived))
setdiff(aggregate_data7$driver_code, aggregate_data2$driver_code)

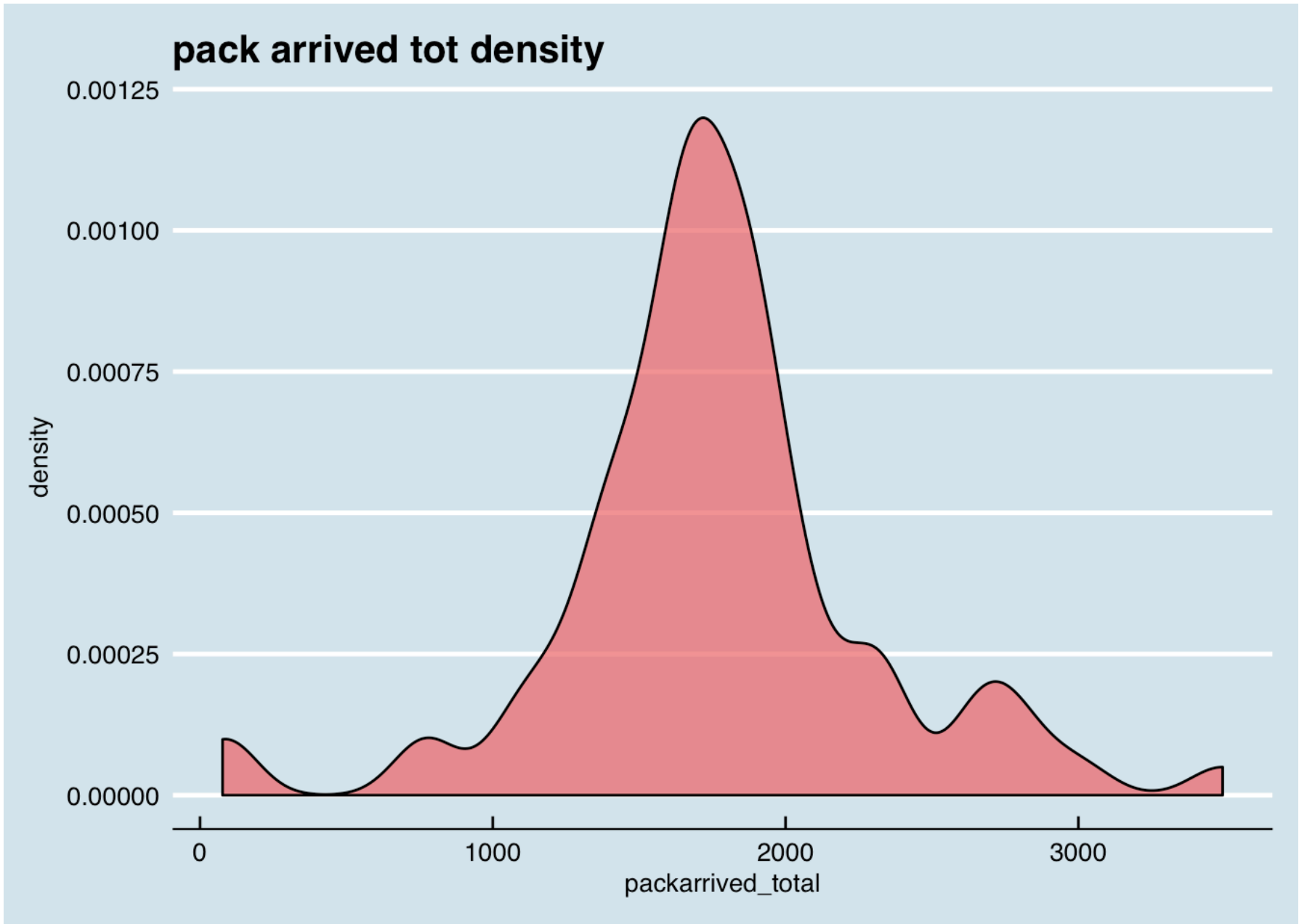
```

```
## [1] "805" "851"
```

```
aggregate_data7 <- aggregate_data7 %>% filter( driver_code != "851", driver_code != "805")
```

```
packarrived_plot <-ggplot(data = aggregate_data7, aes(x = packarrived_total ))  
packarrived_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist()  
+ labs(title = "pack arrived tot density")
```

```
## Warning: Removed 26 rows containing non-finite values (stat_density).
```



TOTAL WEIGHT OF PACKS DELIVERED BY DRIVER

```
data_exploratory <- data_exploratory %>% mutate(weight_pack = gsub(pattern = ",", replacement = ".", x = weight_pack))
data_exploratory$weight_pack <- as.double(data_exploratory$weight_pack)
aggregate_data8 <- data_exploratory %>% group_by(driver_code, day_deliv) %>% summarise(sum(weight_pack))
aggregate_data9 <- data_exploratory %>% group_by(driver_code) %>% summarise(tot_weight_pack = sum(weight_pack))

summary(aggregate_data9)
```

```
## driver_code      tot_weight_pack
## Length:100      Min.   : 135.7
## Class :character 1st Qu.: 7128.3
## Mode  :character Median : 9589.8
##                  Mean    : 9662.5
##                  3rd Qu.:11392.9
##                  Max.    :35663.4
```

```
weightpack_plot <-ggplot(data = aggregate_data9, aes(x = tot_weight_pack ))
weightpack_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist
() + labs(title = "total weight pack density")
```

total weight pack density



TOTAL PACKS NOT DELIVERED PER DRIVER

```
aggregate_data10 <- aggregate_data4 %>% group_by(driver_code) %>% summarise(packnotdelivered_total = sum(not_delivered))
setdiff(aggregate_data10$driver_code, aggregate_data2$driver_code)
```

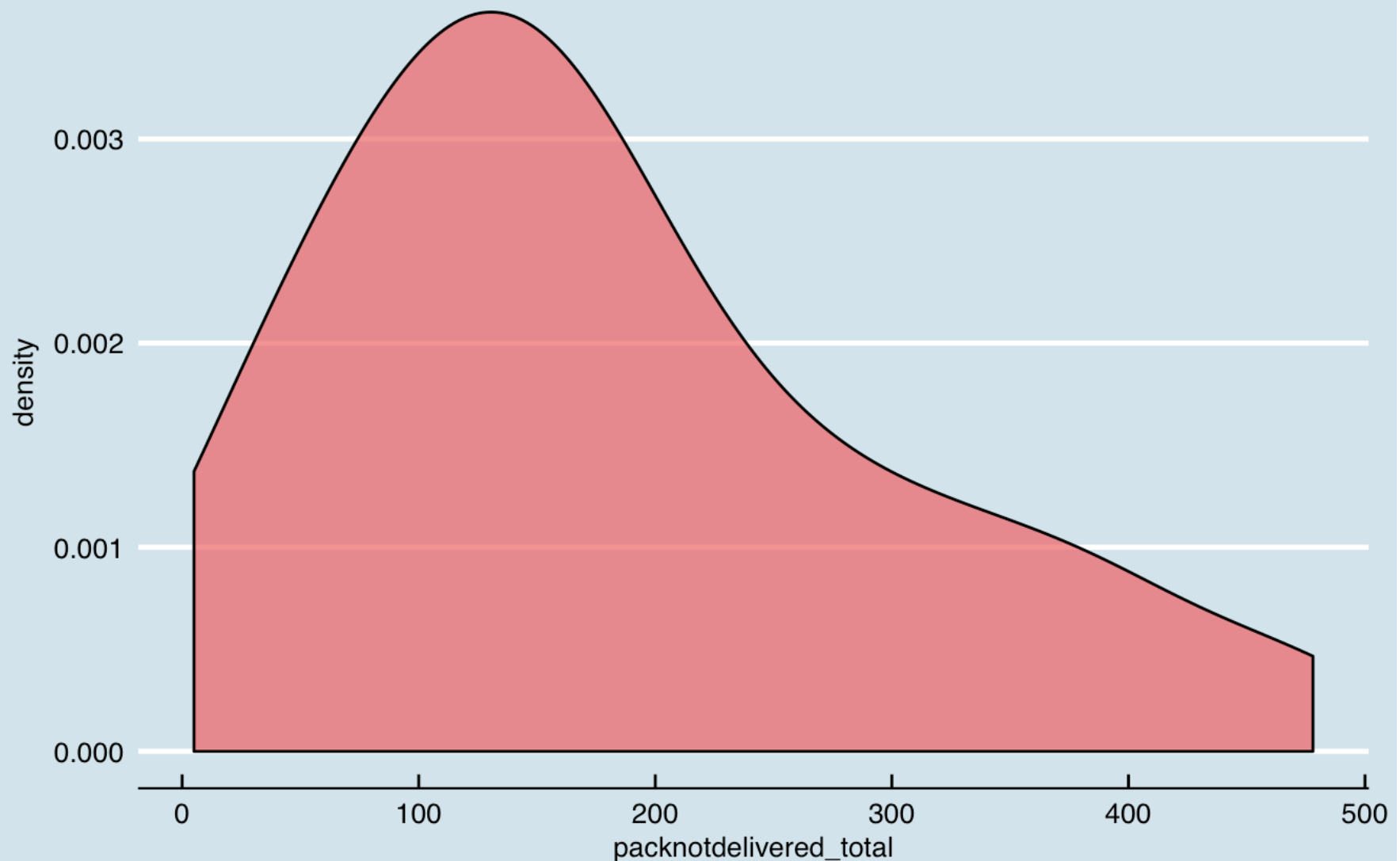
```
## [1] "805" "851"
```

```
aggregate_data10 <- aggregate_data10 %>% filter( driver_code != "851", driver_code != "805")
```

```
nondelivered_plot <-ggplot(data = aggregate_data10, aes(x = packnotdelivered_total ))
nondelivered_plot + geom_density(fill = "indianred2", alpha = 0.7) + theme_economist() + labs(title = "total non delivered packs density")
```

```
## Warning: Removed 28 rows containing non-finite values (stat_density).
```

total non delivered packs density



OBTAINING AN AGGREGATE DATA FRAME WITH ALL THE NECESSARY VALUES

```
aggregate_data_last <- left_join(aggregate_data2, aggregate_data3)
```

```
## Joining, by = "driver_code"
```

```
aggregate_data_last <- left_join(aggregate_data_last, aggregate_data3)
```

```
## Joining, by = c("driver_code", "day_worked")
```

```
aggregate_data_last <- left_join(aggregate_data_last, aggregate_data5)
```

```
## Joining, by = "driver_code"
```

```
aggregate_data_last <- left_join(aggregate_data_last, aggregate_data6)
```

```
## Joining, by = "driver_code"
```

```
aggregate_data_last <- left_join(aggregate_data_last, aggregate_data7)
```

```
## Joining, by = "driver_code"
```

```
aggregate_data_last <- left_join(aggregate_data_last, aggregate_data9)
```

```
## Joining, by = "driver_code"
```

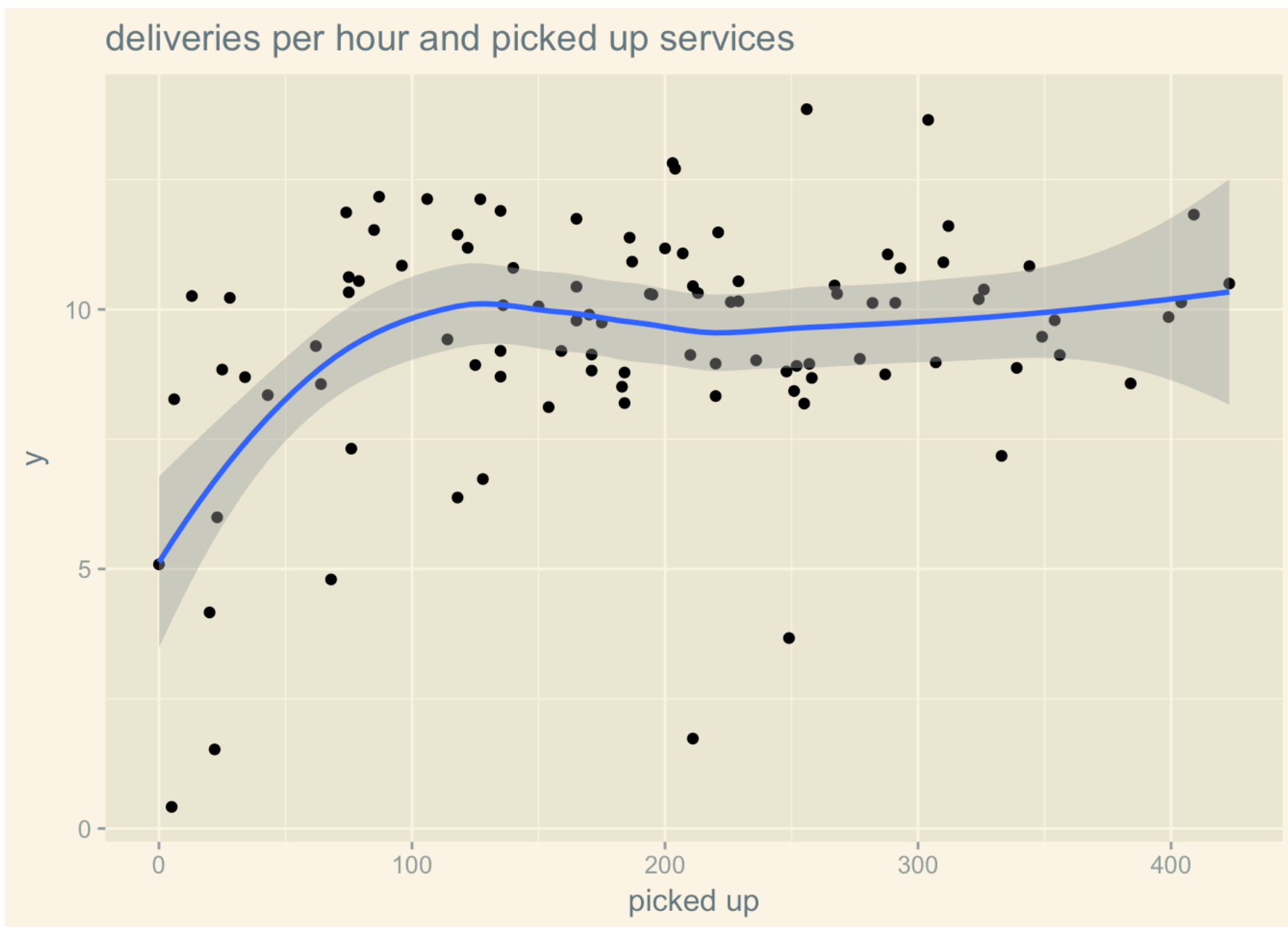
```
aggregate_data_last$`aggregate_data3$day_worked` <- NULL
```

ANALYSIS OF THE DEPENDENT VARIABLE COMPARED TO THE INDEPENDENT VARIABLES

Y AND PICKED UP PACKS

```
y_pickedup_plot<- ggplot(data = aggregate_data_last, aes(x = pickedup_total,y = y))  
  + theme_solarized_2()+  
    labs(title = "deliveries per hour and picked up services", x = "picked up" )  
  
#  
  
y_pickedup_plot + geom_point()+ geom_smooth()
```

```
## `geom_smooth()` using method = 'loess'
```



Y AND DAY WORKED

```
y_dayworked_plot<- ggplot(data = aggregate_data_last, aes(x = day_worked,y = y))+  
  theme_solarized_2()+  
  labs(title = "deliveries per hour and day worked", x = "day worked" )  
  
y_dayworked_plot + geom_point() + geom_smooth() + scale_x_continuous(limits = c(15  
,25))
```

```
## `geom_smooth()` using method = 'loess'
```

```
## Warning: Removed 14 rows containing non-finite values (stat_smooth).
```

```
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric =  
## parametric, : pseudoinverse used at 21
```

```
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric =  
## parametric, : neighborhood radius 1
```

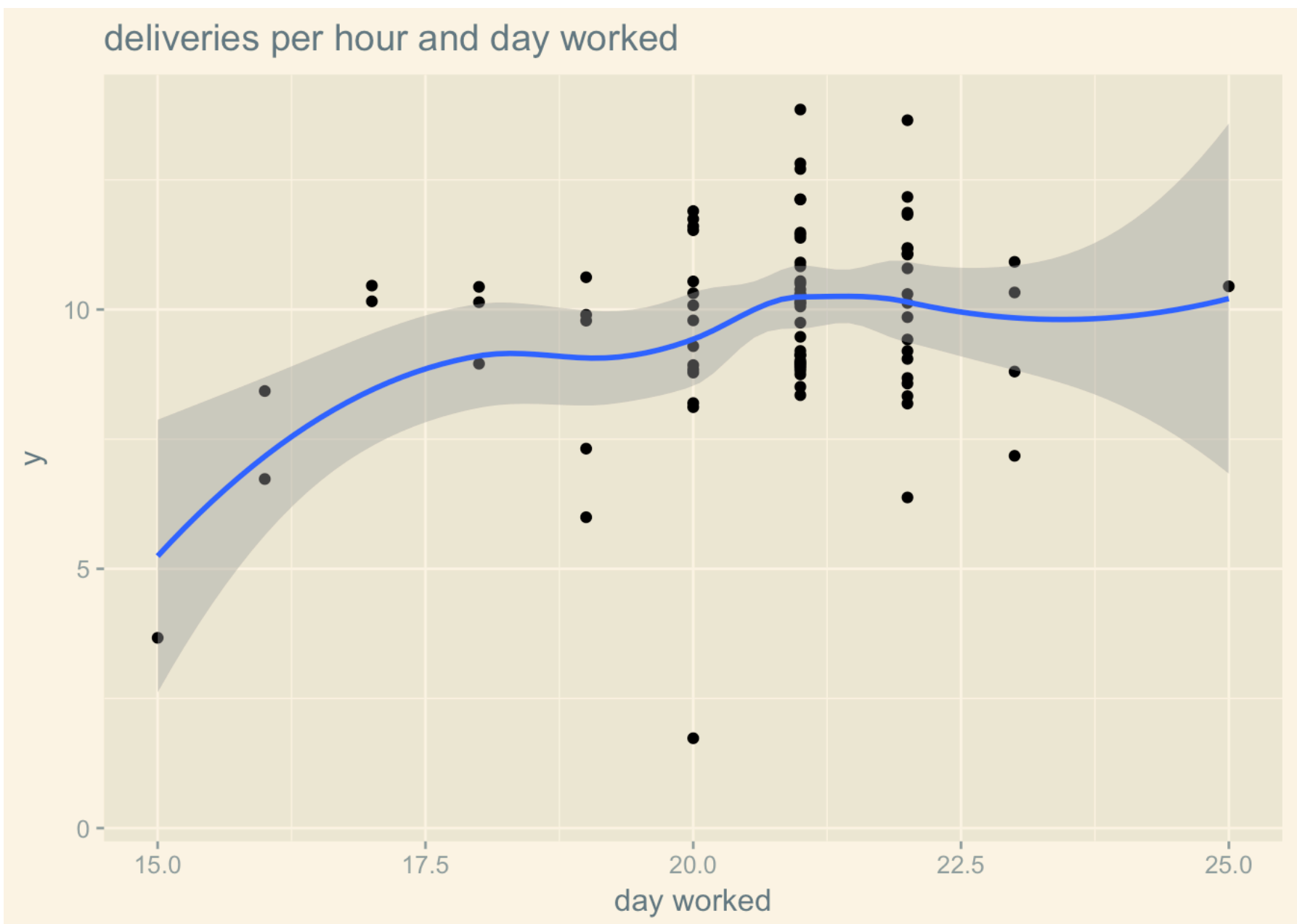
```
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric =  
## parametric, : reciprocal condition number 0
```

```
## Warning in predLoess(object$y, object$x, newx = if  
## (is.null(newdata)) object$x else if (is.data.frame(newdata))  
## as.matrix(model.frame(delete.response(terms(object))), : pseudoinverse used  
## at 21
```

```
## Warning in predLoess(object$y, object$x, newx = if  
## (is.null(newdata)) object$x else if (is.data.frame(newdata))  
## as.matrix(model.frame(delete.response(terms(object))), : neighborhood radius  
## 1
```

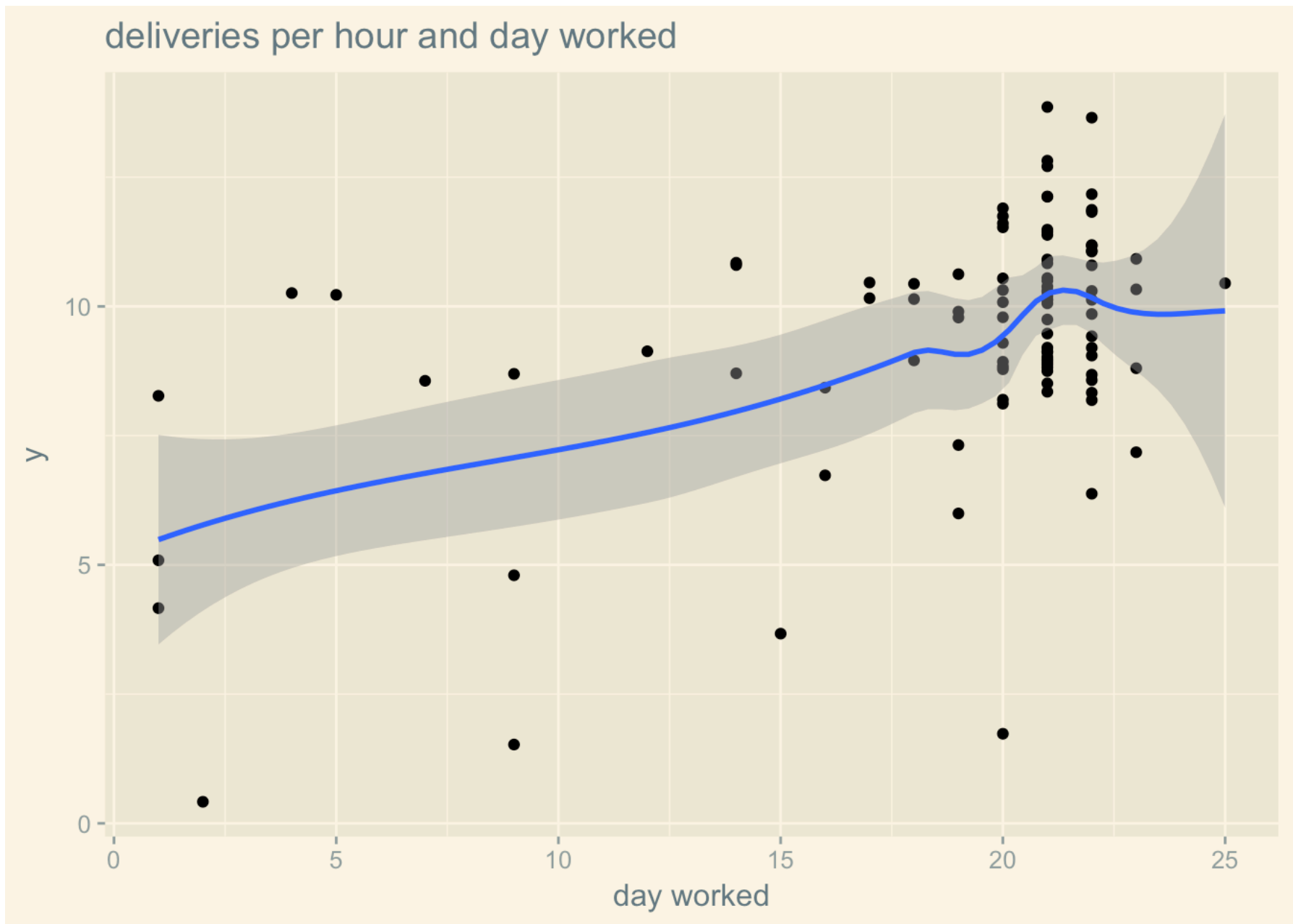
```
## Warning in predLoess(object$y, object$x, newx = if  
## (is.null(newdata)) object$x else if (is.data.frame(newdata))  
## as.matrix(model.frame(delete.response(terms(object))), : reciprocal  
## condition number 0
```

```
## Warning: Removed 14 rows containing missing values (geom_point).
```



```
y_dayworked_plot + geom_point() + geom_smooth()
```

```
## `geom_smooth()` using method = 'loess'
```



Y AND TOT PACK LOADED

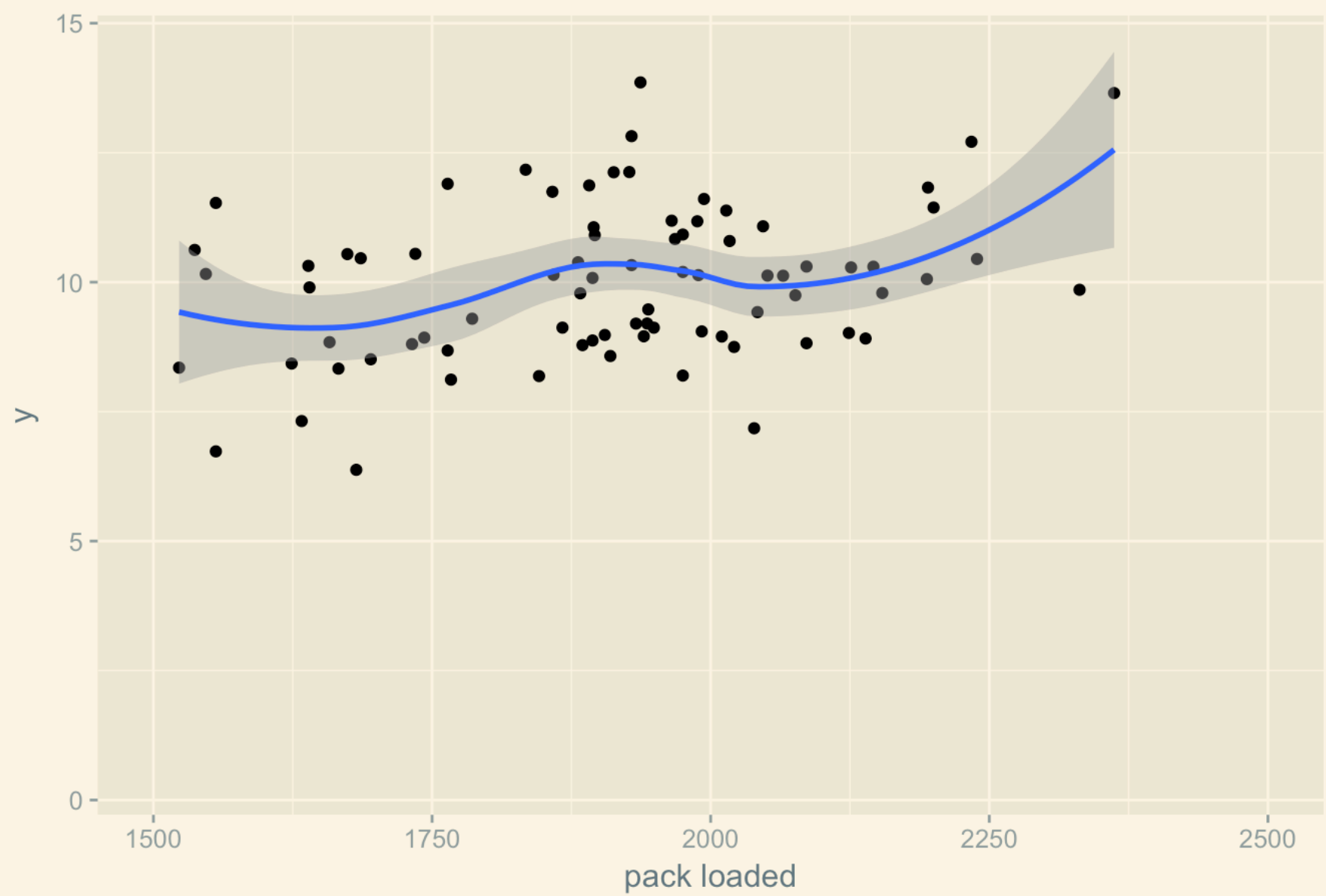
```
y_packloaded_plot<- ggplot(data = aggregate_data_last, aes(x = tot_packloaded,y = y))+ theme_solarized_2()+  
  labs(title = "deliveries per hour and total pack loaded", x = "pack loaded" )  
  
y_packloaded_plot + geom_point() + geom_smooth()+ scale_x_continuous(limits = c(1500, 2500))
```

```
## `geom_smooth()` using method = 'loess'
```

```
## Warning: Removed 20 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 20 rows containing missing values (geom_point).
```

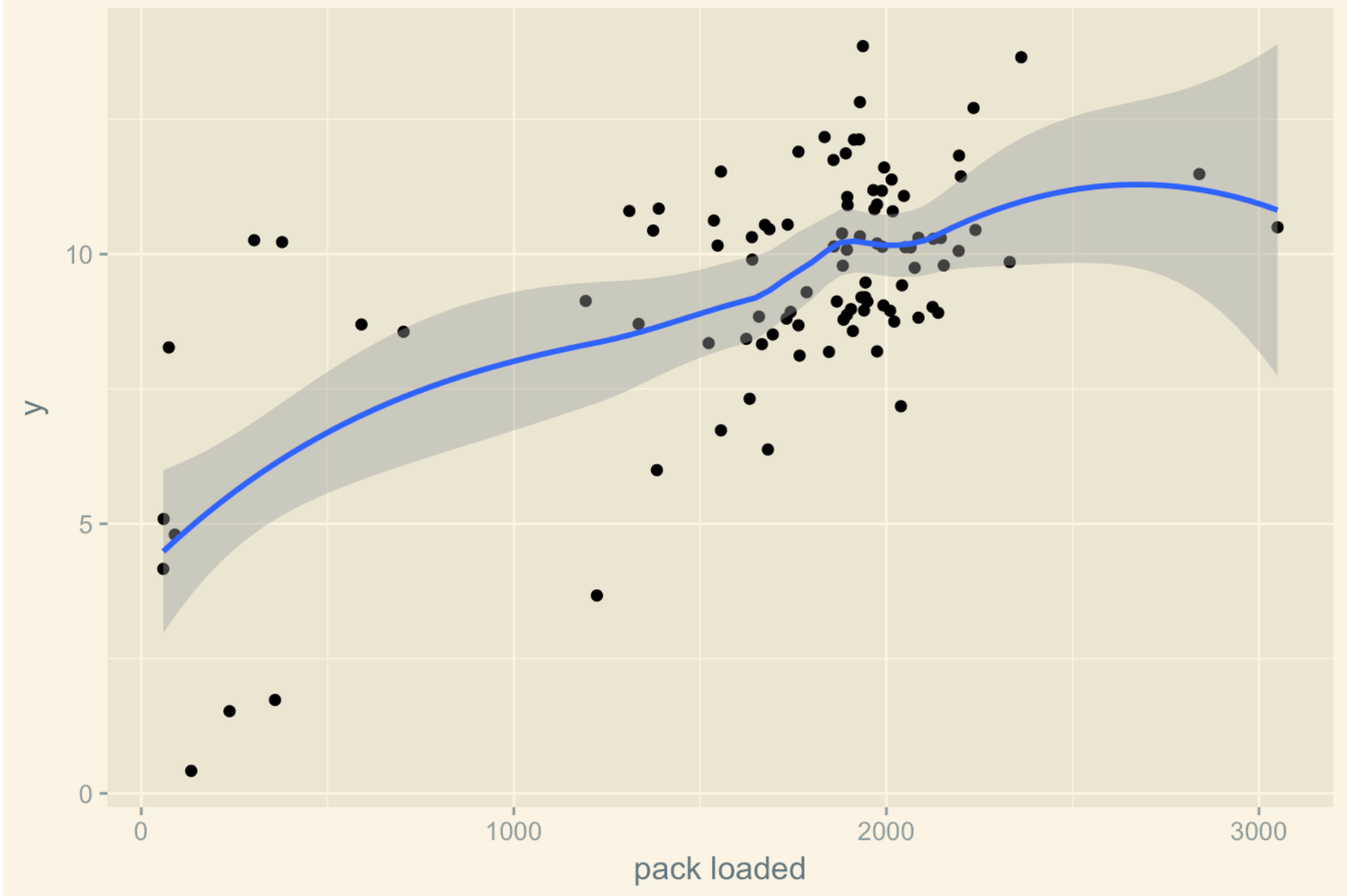
deliveries per hour and total pack loaded



```
y_packloaded_plot + geom_point() + geom_smooth()
```

```
## `geom_smooth()` using method = 'loess'
```

deliveries per hour and total pack loaded



Y AND PACK ARRIVED TOTAL

```
y_packarrived_plot<- ggplot(data = aggregate_data_last, aes(x = packarrived_total,
y = y))+ theme_solarized_2()+
  labs(title = "deliveries per hour and total pack arrived", x = "pack arrived" )

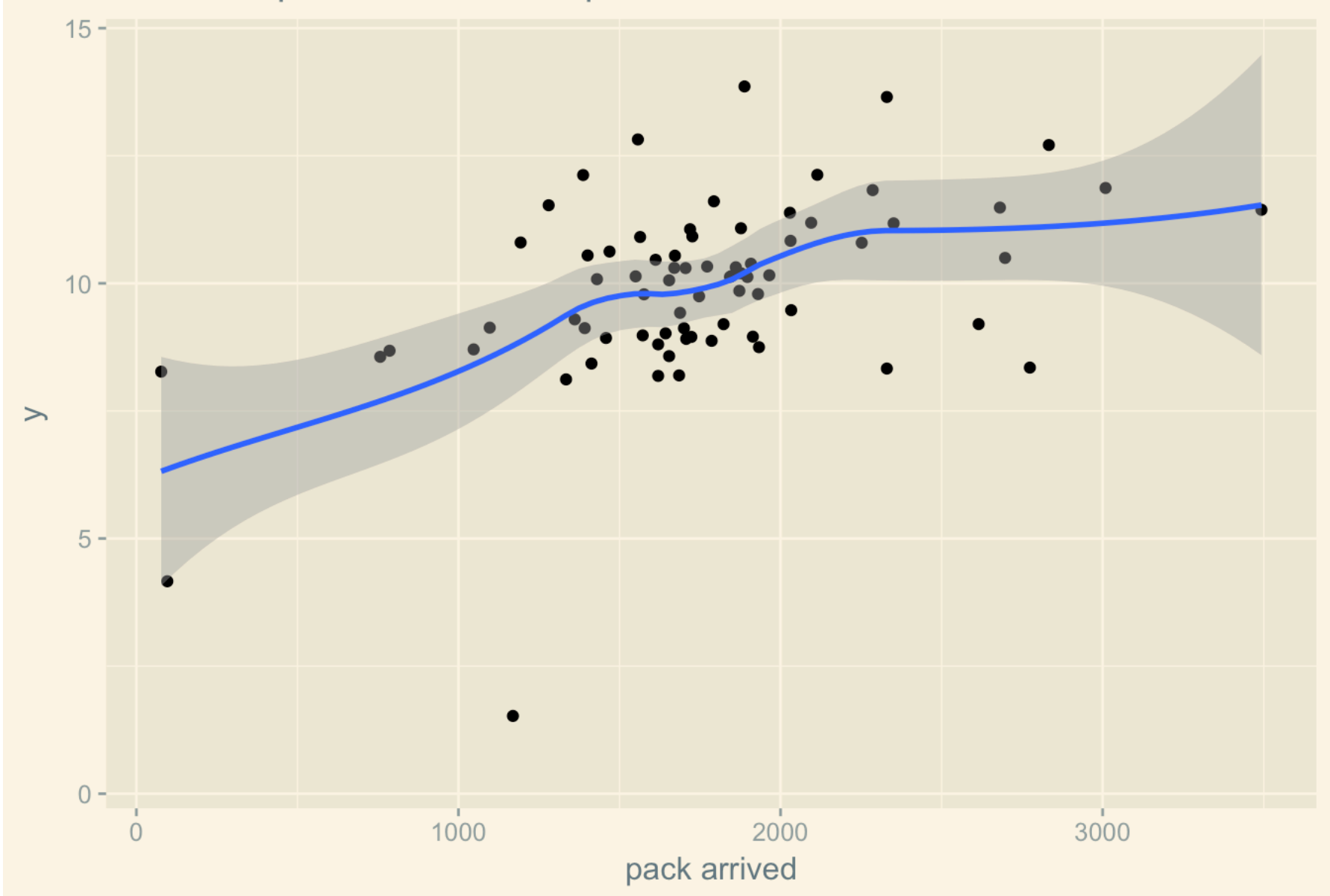
y_packarrived_plot + geom_point() + geom_smooth()
```

```
## `geom_smooth()` using method = 'loess'
```

```
## Warning: Removed 26 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 26 rows containing missing values (geom_point).
```

deliveries per hour and total pack arrived



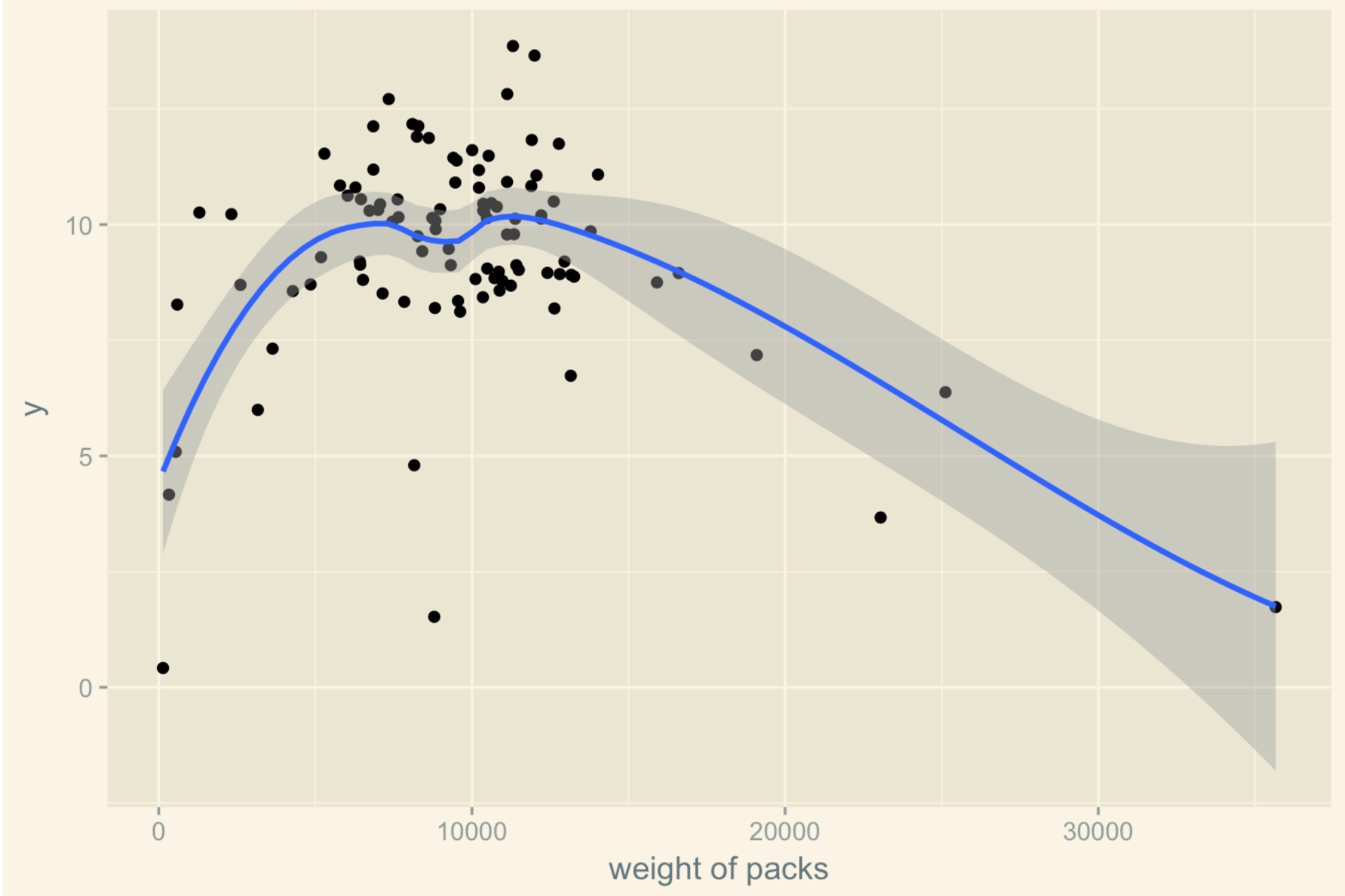
Y AND WEIGHT TOTAL

```
y_weight_plot<- ggplot(data = aggregate_data_last, aes(x = tot_weight_pack,y = y))
+ theme_solarized_2()+
  labs(title = "deliveries per hour and total weight of packs", x = "weight of pa
cks" )
```

```
y_weight_plot + geom_point() + geom_smooth()
```

```
## `geom_smooth()` using method = 'loess'
```


deliveries per hour and total weight of packs



Some driver result in a lower weight due to a lower amount of worked day. As a consequence, I should divide the weight for the number of day worked.

```
y_weight_dayworked_plot<- ggplot(data = aggregate_data_last, aes(x = tot_weight_pack/day_worked,y = y))+ theme_solarized_2()+
  labs(title = "deliveries per hour and total weight of packs", x = "weight of packs" )
```

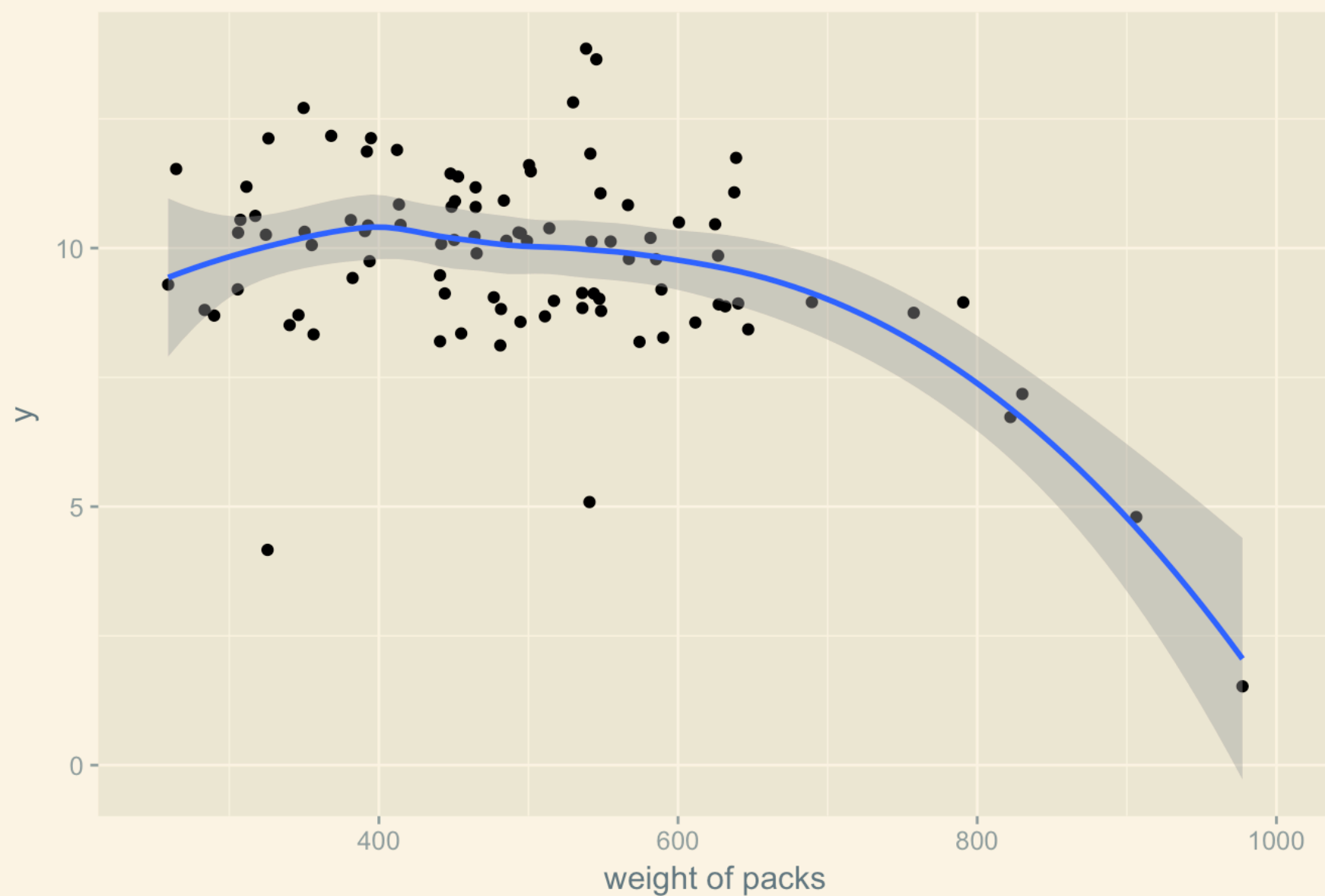
```
y_weight_dayworked_plot + geom_point() + geom_smooth()+ scale_x_continuous(limits = c(250,1000))
```

```
## `geom_smooth()` using method = 'loess'
```

```
## Warning: Removed 6 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 6 rows containing missing values (geom_point).
```

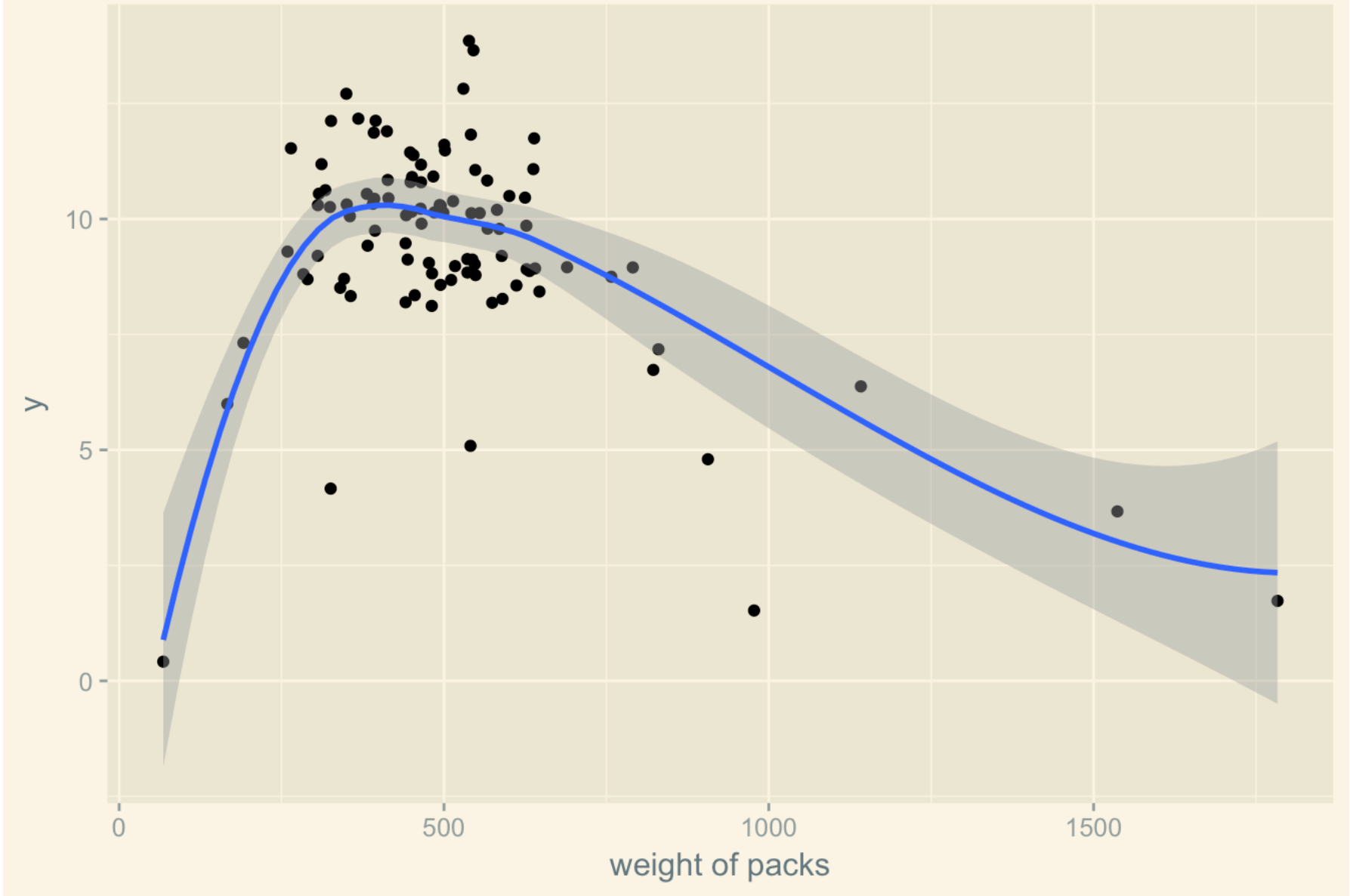
deliveries per hour and total weight of packs



```
y_weight_dayworked_plot + geom_point() + geom_smooth()
```

```
## `geom_smooth()` using method = 'loess'
```

deliveries per hour and total weight of packs



Y AND AREA OF DELIVERY