

Cyclistic Case Study 2021 All Trips

Hezar K

2022-11-29

This analysis is for Cyclistic Case Study for Google Data Analytics Course. This is an analysis for the year of 2021.

STEP ONE: INSTALL REQUIRED PACKAGES AND IMPORT DATA

Install the required packages. **Tidyverse** package to import and wrangling the data and **ggplot2** package for visualization of the data. **Lubridate** package for date parsing and **anytime** package for the datetime conversion.

- `install.packages("tidyverse")`
- `install.packages("ggplot2")`
- `install.packages("lubridate")`
- `install.packages("anytime")`

```
library(tidyverse)
library(lubridate)
library(data.table)
library(ggplot2)
library(anytime)
```

Import data from local drive.

```
Jan21 <- read_csv("202101-divvy-tripdata.csv")
Feb21 <- read_csv("202102-divvy-tripdata.csv")
Mar21 <- read_csv("202103-divvy-tripdata.csv")
Apr21 <- read_csv("202104-divvy-tripdata.csv")
May21 <- read_csv("202105-divvy-tripdata.csv")
Jun21 <- read_csv("202106-divvy-tripdata.csv")
Jul21 <- read_csv("202107-divvy-tripdata.csv")
Aug21 <- read_csv("202108-divvy-tripdata.csv")
Sep21 <- read_csv("202109-divvy-tripdata.csv")
Oct21 <- read_csv("202110-divvy-tripdata.csv")
Nov21 <- read_csv("202111-divvy-tripdata.csv")
Dec21 <- read_csv("202112-divvy-tripdata.csv")
```

STEP TWO: EXAMINE THE DATA

Examine the dataframe for an overview of the data. Review column names, **colnames()**. Then, we need to combine all data one dataframe. Then we examine dataframes to find dimensions, **dim()**, the first, **head()**, and the last, **tail()**, six rows in the dataframe, the summary, **summary()**, statistics on the columns of the dataframe, and review the data type structure of columns, **str()**. (To reduce cutler I have removed colnames output from Feb21-Dec21, because all tables have the same column names.

```
colnames(Jan21)
##  [1] "ride_id"          "rideable_type"    "started_at"
##  [4] "ended_at"         "start_station_name" "start_station_id"
##  [7] "end_station_name" "end_station_id"   "start_lat"
## [10] "start_lng"        "end_lat"          "end_lng"
## [13] "member_casual"

colnames(Feb21)
colnames(Mar21)
colnames(Apr21)
colnames(May21)
colnames(Jun21)
colnames(Jul21)
colnames(Aug21)
colnames(Sep21)
colnames(Oct21)
colnames(Nov21)
colnames(Dec21)
```

Since all column names are the same. We can combine the data into one dataframe.

```
all_trips <- bind_rows(Jan21, Feb21, Mar21, Apr21, May21, Jun21, Jul21,
Aug21, Sep21, Oct21, Nov21, Dec21)
```

```
View(all_trips)
```

```
nrow(all_trips)
## [1] 5595063
```

```
dim(all_trips)
## [1] 5595063      13
```

```
head(all_trips)
```

```
## # A tibble: 6 × 13
```

```
##   ride_id      ridea...1 start...2 ended...3 start...4 start...5 end_s...6 end_s...7  
start...8
```

```
##   <chr>          <chr>   <chr>   <chr>   <chr>   <chr>   <chr>   <chr>  
<dbl>
```

```
## 1 E19E6F1B8D4C4... electr... 1/23/2... 1/23/2... Califo... 17660   <NA>   <NA>  
41.9
```

```
## 2 DC88F20C2C55F... electr... 1/27/2... 1/27/2... Califo... 17660   <NA>   <NA>  
41.9
```

```
## 3 EC45C94683FE3... electr... 1/21/2... 1/21/2... Califo... 17660   <NA>   <NA>  
41.9
```

```
## 4 4FA453A75AE37... electr... 1/7/20... 1/7/20... Califo... 17660   <NA>   <NA>  
41.9
```

```
## 5 BE5E8EB4E7263... electr... 1/23/2... 1/23/2... Califo... 17660   <NA>   <NA>  
41.9
```

```
## 6 5D8969F88C773... electr... 1/9/20... 1/9/20... Califo... 17660   <NA>   <NA>  
41.9
```

```
## # ... with 4 more variables: start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
```

```
## #   member_casual <chr>, and abbreviated variable names 1rideable_type,
```

```
## #   2started_at, 3ended_at, 4start_station_name, 5start_station_id,
```

```
## #   6end_station_name, 7end_station_id, 8start_lat
```

```
tail(all_trips)
```

```
## # A tibble: 6 × 13
```

```
##   ride_id      ridea...1 start...2 ended...3 start...4 start...5 end_s...6 end_s...7  
start...8
```

```
##   <chr>          <chr>   <chr>   <chr>   <chr>   <chr>   <chr>   <chr>  
<dbl>
```

```
## 1 92BBAB97D1683... electr... 12/24/... 12/24/... Canal ... 13341   <NA>   <NA>  
41.9
```

```
## 2 847431F3D5353... electr... 12/12/... 12/12/... Canal ... 13341   <NA>   <NA>  
41.9
```

```
## 3 CF407BBC3B9FA... electr... 12/6/2... 12/6/2... Canal ... 13341   Kingsb... KA1503...  
41.9
```

```
## 4 60BB69EBF5440... electr... 12/2/2... 12/2/2... Canal ... 13341   Dearbo... TA1305...  
41.9
```

```
## 5 C414F654A2863... electr... 12/13/... 12/13/... Lawnda... 362     <NA>   <NA>  
41.9
```

```
## 6 37AC57E34B2E7... classi... 12/13/... 12/13/... Michig... TA1309... Dearbo... TA1305...
41.9
## # ... with 4 more variables: start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
## #   member_casual <chr>, and abbreviated variable names 1rideable_type,
## #   2started_at, 3ended_at, 4start_station_name, 5start_station_id,
## #   6end_station_name, 7end_station_id, 8start_lat
```

```
summary(all_trips)
##      ride_id      rideable_type      started_at      ended_at
## Length:5595063 Length:5595063 Length:5595063 Length:5595063
## Class :character Class :character Class :character Class :character
## Mode  :character Mode  :character Mode  :character Mode  :character
##
## start_station_name start_station_id end_station_name end_station_id
## Length:5595063 Length:5595063 Length:5595063 Length:5595063
## Class :character Class :character Class :character Class :character
## Mode  :character Mode  :character Mode  :character Mode  :character
##
##      start_lat      start_lng      end_lat      end_lng
## Min.   :41.64 Min.   : -87.84 Min.   :41.39 Min.   : -88.97
## 1st Qu.:41.88 1st Qu.: -87.66 1st Qu.:41.88 1st Qu.: -87.66
## Median :41.90 Median : -87.64 Median :41.90 Median : -87.64
## Mean   :41.90 Mean   : -87.65 Mean   :41.90 Mean   : -87.65
## 3rd Qu.:41.93 3rd Qu.: -87.63 3rd Qu.:41.93 3rd Qu.: -87.63
## Max.   :42.07 Max.   : -87.52 Max.   :42.17 Max.   : -87.49
##
##              NA's      :4771 NA's      :4771
## member_casual
## Length:5595063
## Class :character
## Mode  :character
##
```

```
str(all_trips)
## spc_tbl_ [5,595,063 × 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
```

```

## $ ride_id          : chr [1:5595063] "E19E6F1B8D4C42ED"
"DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA453A75AE377DB" ...

## $ rideable_type    : chr [1:5595063] "electric_bike" "electric_bike"
"electric_bike" "electric_bike" ...

## $ started_at       : chr [1:5595063] "1/23/2021 16:14" "1/27/2021 18:43"
"1/21/2021 22:35" "1/7/2021 13:31" ...

## $ ended_at         : chr [1:5595063] "1/23/2021 16:24" "1/27/2021 18:47"
"1/21/2021 22:37" "1/7/2021 13:42" ...

## $ start_station_name: chr [1:5595063] "California Ave & Cortez St"
"California Ave & Cortez St" "California Ave & Cortez St" "California Ave &
Cortez St" ...

## $ start_station_id  : chr [1:5595063] "17660" "17660" "17660" "17660" ...

## $ end_station_name  : chr [1:5595063] NA NA NA NA ...

## $ end_station_id    : chr [1:5595063] NA NA NA NA ...

## $ start_lat         : num [1:5595063] 41.9 41.9 41.9 41.9 41.9 ...

## $ start_lng         : num [1:5595063] -87.7 -87.7 -87.7 -87.7 -87.7 ...

## $ end_lat          : num [1:5595063] 41.9 41.9 41.9 41.9 41.9 ...

## $ end_lng          : num [1:5595063] -87.7 -87.7 -87.7 -87.7 -87.7 ...

## $ member_casual     : chr [1:5595063] "member" "member" "member" "member"
...

## - attr(*, "spec")=
## .. cols(
## ..   ride_id = col_character(),
## ..   rideable_type = col_character(),
## ..   started_at = col_character(),
## ..   ended_at = col_character(),
## ..   start_station_name = col_character(),
## ..   start_station_id = col_character(),
## ..   end_station_name = col_character(),
## ..   end_station_id = col_character(),
## ..   start_lat = col_double(),
## ..   start_lng = col_double(),
## ..   end_lat = col_double(),
## ..   end_lng = col_double(),
## ..   member_casual = col_character()
## .. )
## - attr(*, "problems")=<externalptr>

```

Columns *started_at* and *ended_at* need to be converted from character data type to date data type. **Str()** syntax confirms changes.

```
all_trips$started_at <- mdy_hm(all_trips$started_at)
all_trips$ended_at <- mdy_hm(all_trips$ended_at)
str(all_trips)

## spc_tbl_ [5,595,063 × 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
##   $ ride_id           : chr [1:5595063] "E19E6F1B8D4C42ED"
##   "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA453A75AE377DB" ...
##   $ rideable_type      : chr [1:5595063] "electric_bike" "electric_bike"
##   "electric_bike" "electric_bike" ...
##   $ started_at         : POSIXct[1:5595063], format: "2021-01-23 16:14:00"
##   "2021-01-27 18:43:00" ...
##   $ ended_at           : POSIXct[1:5595063], format: "2021-01-23 16:24:00"
##   "2021-01-27 18:47:00" ...
##   $ start_station_name: chr [1:5595063] "California Ave & Cortez St"
##   "California Ave & Cortez St" "California Ave & Cortez St" "California Ave &
##   Cortez St" ...
##   $ start_station_id   : chr [1:5595063] "17660" "17660" "17660" "17660" ...
##   $ end_station_name   : chr [1:5595063] NA NA NA NA ...
##   $ end_station_id     : chr [1:5595063] NA NA NA NA ...
##   $ start_lat          : num [1:5595063] 41.9 41.9 41.9 41.9 41.9 ...
##   $ start_lng          : num [1:5595063] -87.7 -87.7 -87.7 -87.7 -87.7 ...
##   $ end_lat            : num [1:5595063] 41.9 41.9 41.9 41.9 41.9 ...
##   $ end_lng            : num [1:5595063] -87.7 -87.7 -87.7 -87.7 -87.7 ...
##   $ member_casual      : chr [1:5595063] "member" "member" "member" "member"
##   ...
##   - attr(*, "spec")=
##     .. cols(
##       ..   ride_id = col_character(),
##       ..   rideable_type = col_character(),
##       ..   started_at = col_character(),
##       ..   ended_at = col_character(),
##       ..   start_station_name = col_character(),
##       ..   start_station_id = col_character(),
##       ..   end_station_name = col_character(),
##       ..   end_station_id = col_character(),
##       ..   start_lat = col_double(),
##       ..   start_lng = col_double(),
```

```
## .. end_lat = col_double(),
## .. end_lng = col_double(),
## .. member_casual = col_character()
## .. )
## - attr(*, "problems")=<externalptr>
```

Create new columns as for *date*, *month*, *day*, *year*, *day_of_week*, and *ride_length* in seconds.

```
all_trips$date <- as.Date(all_trips$started_at)
all_trips$month <- format(as.Date(all_trips$date), "%m")
all_trips$day <- format(as.Date(all_trips$date), "%d")
all_trips$year <- format(as.Date(all_trips$date), "%Y")
all_trips$day_of_week <- format(as.Date(all_trips$date), "%A")
all_trips$ride_length <- difftime(all_trips$ended_at, all_trips$started_at)
```

Convert *ride_length* column to numeric in order to run calculations on the data. First, check to see if the data type is numeric, and then convert if needed.

```
is.numeric(all_trips$ride_length)
## [1] FALSE
```

Recheck *ride_length* data type.

```
all_trips$ride_length <- as.numeric(as.character(all_trips$ride_length))
is.numeric(all_trips$ride_length)
## [1] TRUE
```

STEP THREE: CLEAN DATA

na.omit() will remove all NA from the dataframe.

```
all_trips <- na.omit(all_trips)
```

Remove rows with the *ride_id* column character length is not 16. This will remove all the scientific ride ids that we noticed while examining the data.

```
all_trips <- subset(all_trips, nchar(as.character(ride_id)) == 16)
```

Remove rows with the *ride_length* less than 1 minute.

```
all_trips <- subset (all_trips, ride_length > "1")
```

STEP FOUR: ANALYZE DATA

Analyze the dataframe by find the **mean**, **median**, **max** (maximum), and **min** (minimum) of *ride_length*.

```
mean(all_trips$ride_length)
## [1] 1318.707
```

```
median(all_trips$ride_length)
## [1] 720
```

```
max(all_trips$ride_length)
## [1] 3356640
```

```
min(all_trips$ride_length)
## [1] 60
```

Run a statistical summary of the *ride_length*.

```
summary(all_trips$ride_length)
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       60    420     720    1319    1320 3356640
```

Compare the members and casual users

```
aggregate(all_trips$ride_length ~ all_trips$member_casual, FUN = mean)
##   all_trips$member_casual all_trips$ride_length
## 1                    casual          1961.3398
## 2                    member           798.4492
aggregate(all_trips$ride_length ~ all_trips$member_casual, FUN = median)
##   all_trips$member_casual all_trips$ride_length
## 1                    casual             1020
## 2                    member              600
aggregate(all_trips$ride_length ~ all_trips$member_casual, FUN = max)
##   all_trips$member_casual all_trips$ride_length
## 1                    casual          3356640
## 2                    member          89700
aggregate(all_trips$ride_length ~ all_trips$member_casual, FUN = min)
##   all_trips$member_casual all_trips$ride_length
## 1                    casual              60
```


## 2	member	60
------	--------	----

Aggregate the average ride length by each day of the week for members and users.

```
aggregate(all_trips$ride_length ~ all_trips$member_casual +
all_trips$day_of_week, FUN = mean)
```

##	all_trips\$member_casual	all_trips\$day_of_week	all_trips\$ride_length
## 1	casual	Friday	1865.4044
## 2	member	Friday	774.7142
## 3	casual	Monday	1969.3185
## 4	member	Monday	770.6198
## 5	casual	Saturday	2103.7133
## 6	member	Saturday	898.2403
## 7	casual	Sunday	2268.3352
## 8	member	Sunday	921.2654
## 9	casual	Thursday	1690.0091
## 10	member	Thursday	747.7804
## 11	casual	Tuesday	1737.8340
## 12	member	Tuesday	749.7625
## 13	casual	Wednesday	1705.6794
## 14	member	Wednesday	754.1646

Sort the days of the week in order.

```
all_trips$day_of_week <- ordered(all_trips$day_of_week, levels=c("Sunday",
"Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday"))
```

Assign the aggregate the average ride length by each day of the week for members and users to x.

```
x <- aggregate(all_trips$ride_length ~ all_trips$member_casual +
all_trips$day_of_week, FUN = mean)
```

```
head(x)
```

##	all_trips\$member_casual	all_trips\$day_of_week	all_trips\$ride_length
## 1	casual	Sunday	2268.3352
## 2	member	Sunday	921.2654
## 3	casual	Monday	1969.3185
## 4	member	Monday	770.6198
## 5	casual	Tuesday	1737.8340
## 6	member	Tuesday	749.7625

Find the average ride length of member riders and casual riders per day and assign it to y.

```
y <- all_trips %>%
  mutate(weekday = wday(started_at)) %>%
  group_by(member_casual, weekday) %>%
  summarise(number_of_rides = n(),
            average_duration = mean(ride_length), .groups = 'drop') %>%
  arrange(member_casual, weekday)
```

```
head(y)
```

```
## # A tibble: 6 × 4
##   member_casual weekday number_of_rides average_duration
##   <chr>          <int>          <int>          <dbl>
## 1 casual         1          401470          2268.
## 2 casual         2          227603          1969.
## 3 casual         3          213707          1738.
## 4 casual         4          216912          1706.
## 5 casual         5          222919          1690.
## 6 casual         6          288411          1865.
```

Analyze the dataframe to find the frequency of member riders, casual riders, classic bikes, docked bikes, and electric bikes.

```
table(all_trips$member_casual)
```

```
##
##   casual  member
## 2036760 2515844
```

```
table(all_trips$rideable_type)
```

```
##
##   classic_bike  docked_bike electric_bike
##      3216339      310815      1025450
```

STEP FIVE: VISUALIZATION

Display full digits instead of scientific number.

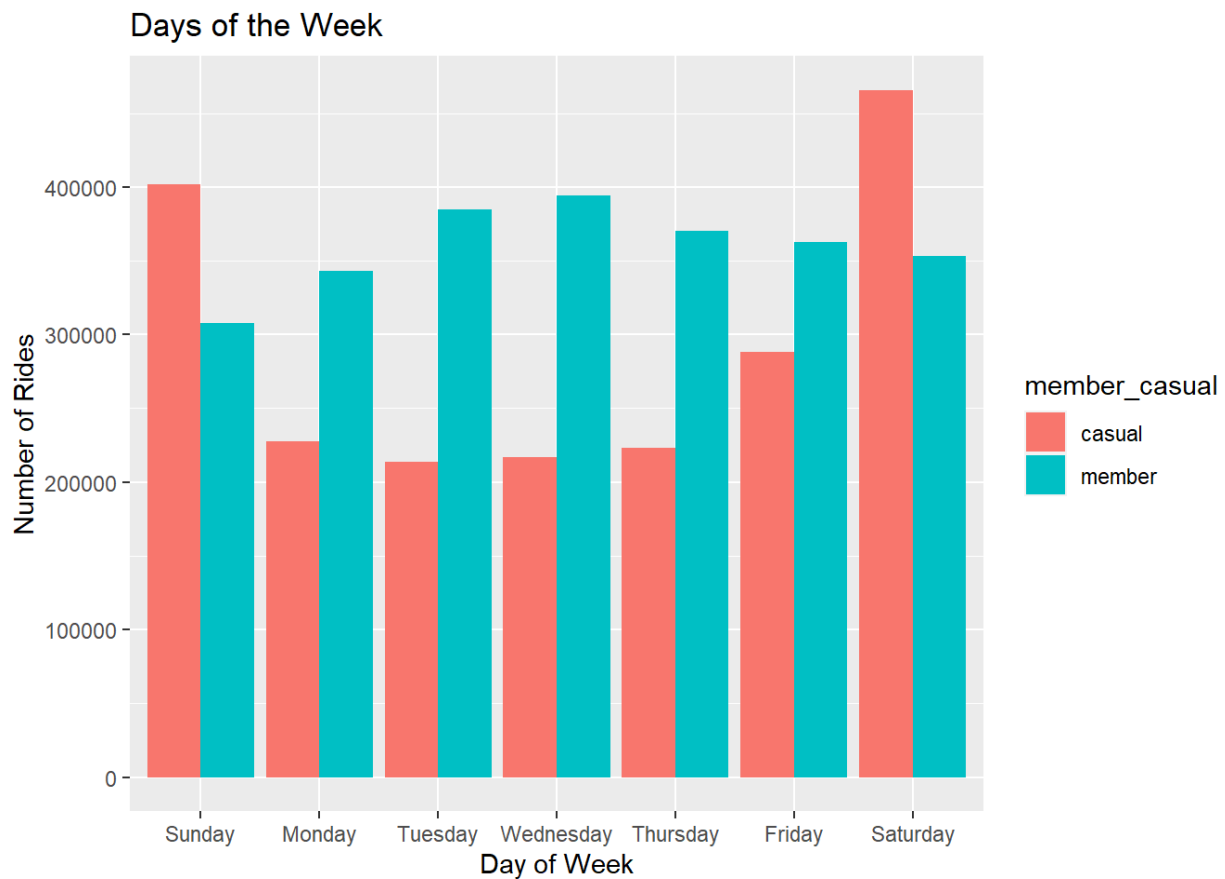
```
options(scipen=999)
```

Plot the number of rides by user type during the week.

```

all_trips %>%
  mutate(day_of_week) %>%
  group_by(member_casual, day_of_week) %>%
  summarise(number_of_rides = n(), average_duration = mean(ride_length),
    .groups = 'drop') %>%
  arrange(member_casual, day_of_week) %>%
  ggplot(aes(x = day_of_week, y = number_of_rides, fill = member_casual)) +
  geom_col(position = "dodge") +
  labs(x = "Day of Week",
    y = "Number of Rides",
    title = "Days of the Week")

```



Plot the duration of the ride by user type during the week.

```

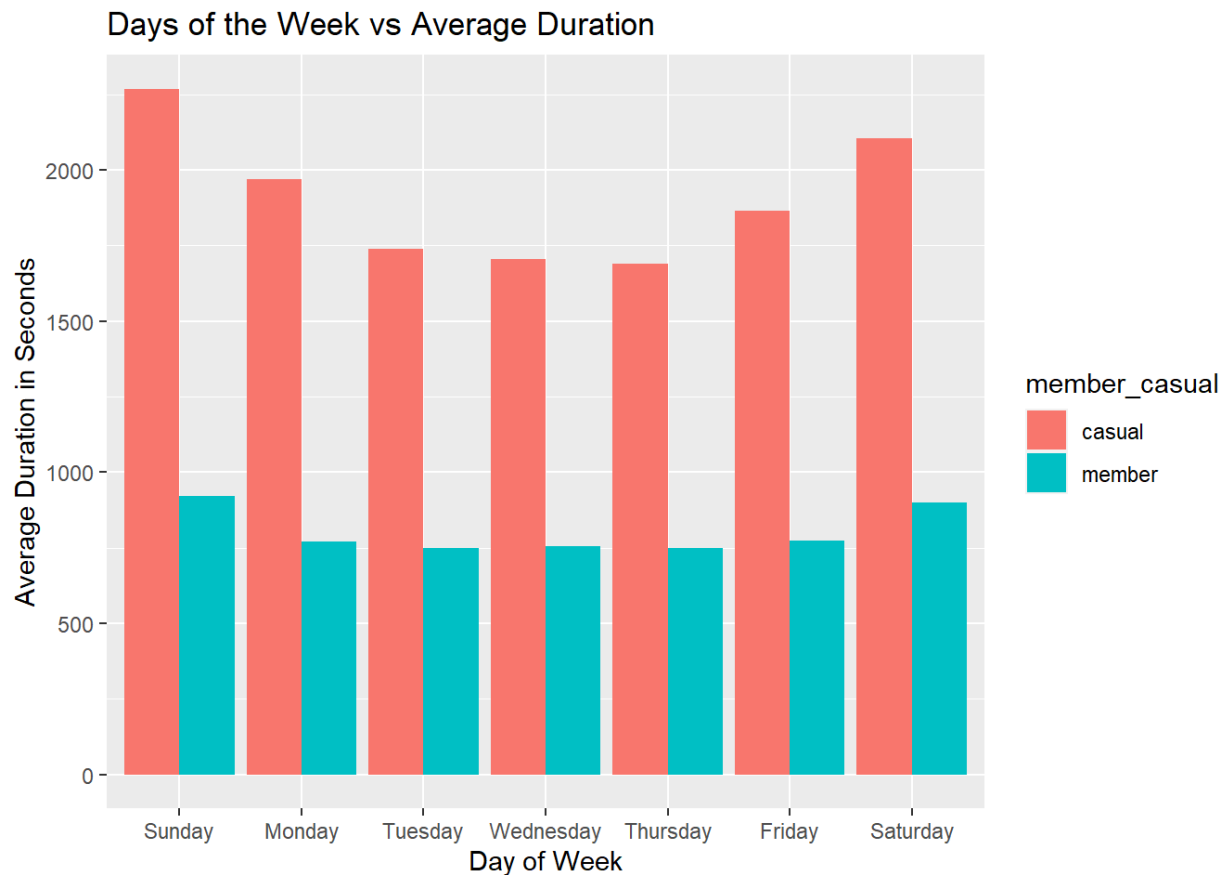
all_trips %>%
  mutate(day_of_week) %>%
  group_by(member_casual, day_of_week) %>%

```

```

  summarise(number_of_rides = n(), average_duration = mean(ride_length),
.groups = 'drop') %>%
  arrange(member_casual, day_of_week) %>%
  ggplot(aes(x = day_of_week, y = average_duration, fill = member_casual)) +
  geom_col(position = "dodge") +
  labs(x = "Day of Week",
       y= "Average Duration in Seconds",
       title= "Days of the Week vs Average Duration")

```



Create new dataframe for plots for weekday trends vs weekend trends.

```
mc<- as.data.frame(table(all_trips$day_of_week,all_trips$member_casual))
```

Rename columns

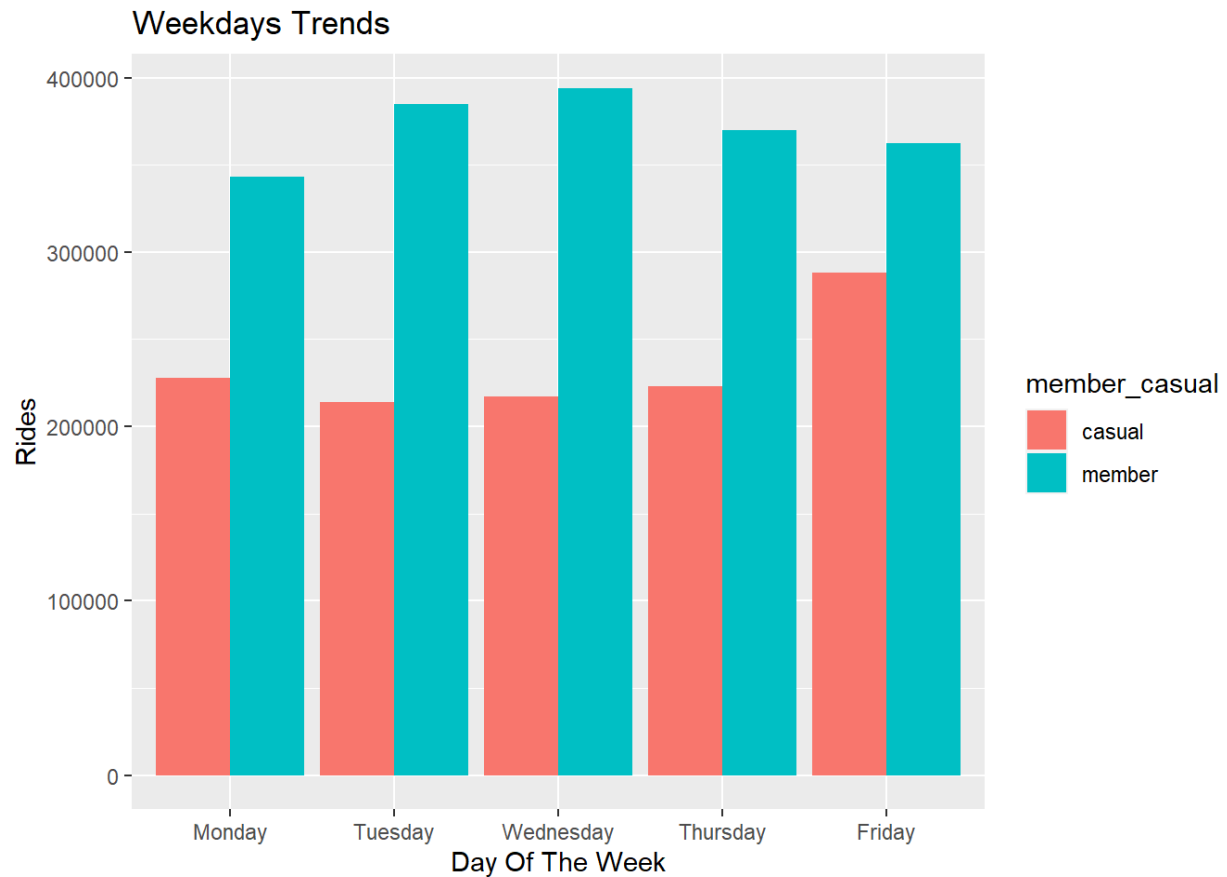
```
mc<-rename(mc, day_of_week = Var1, member_casual = Var2)
head(mc)
```

```
##   day_of_week member_casual   Freq
## 1      Sunday          casual 401470
```

```
## 2      Monday      casual 227603
## 3      Tuesday     casual 213707
## 4      Wednesday   casual 216912
## 5      Thursday    casual 222919
## 6      Friday      casual 288411
```

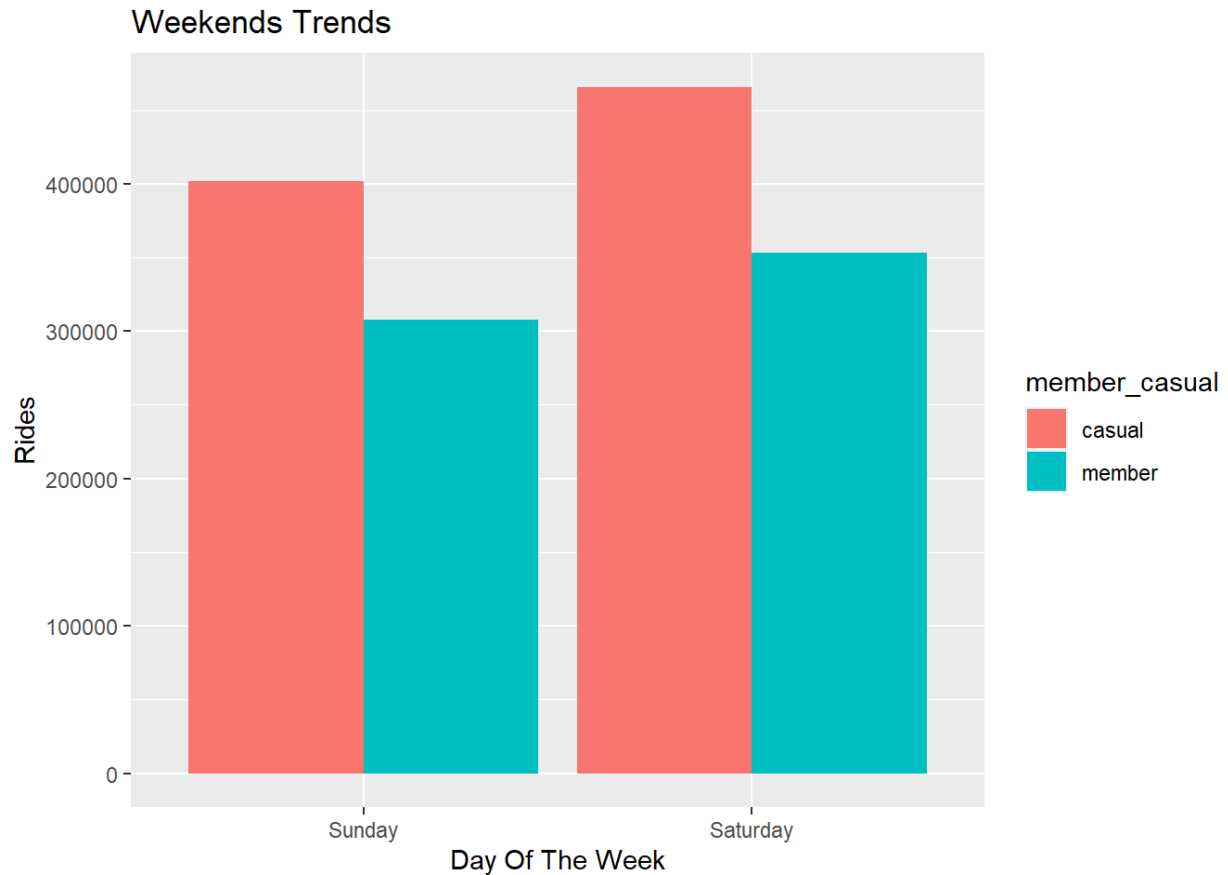
Weekday trends (Monday through Friday).

```
mc %>%
  filter(day_of_week == "Monday" |
         day_of_week == "Tuesday" |
         day_of_week == "Wednesday" |
         day_of_week == "Thursday" |
         day_of_week == "Friday") %>%
  ggplot(aes(x = day_of_week, y = Freq, fill = member_casual))+
  geom_bar(stat = "identity" , position = "dodge") +
  labs(title = "Weekdays Trends",
       x= "Day Of The Week",
       y = "Rides")
```



Weekend trends (Sunday and Saturday).

```
mc %>%
  filter(day_of_week == "Sunday" |
         day_of_week == "Saturday") %>%
  ggplot(aes(x = day_of_week, y = Freq, fill = member_casual)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(title = "Weekends Trends",
       x = "Day Of The Week",
       y = "Rides")
```



Create dataframe for member and casual riders vs ride type

```
rt<- as.data.frame(table(all_trips$rideable_type,all_trips$member_casual))
```

Rename columns.

```
rt<-rename(rt, rideable_type = Var1, member_casual = Var2)
head(rt)
```

```
##  rideable_type member_casual    Freq
## 1  classic_bike        casual 1254201
## 2  docked_bike        casual  310814
## 3 electric_bike        casual  471745
## 4  classic_bike        member 1962138
## 5  docked_bike        member         1
## 6 electric_bike        member  553705
```

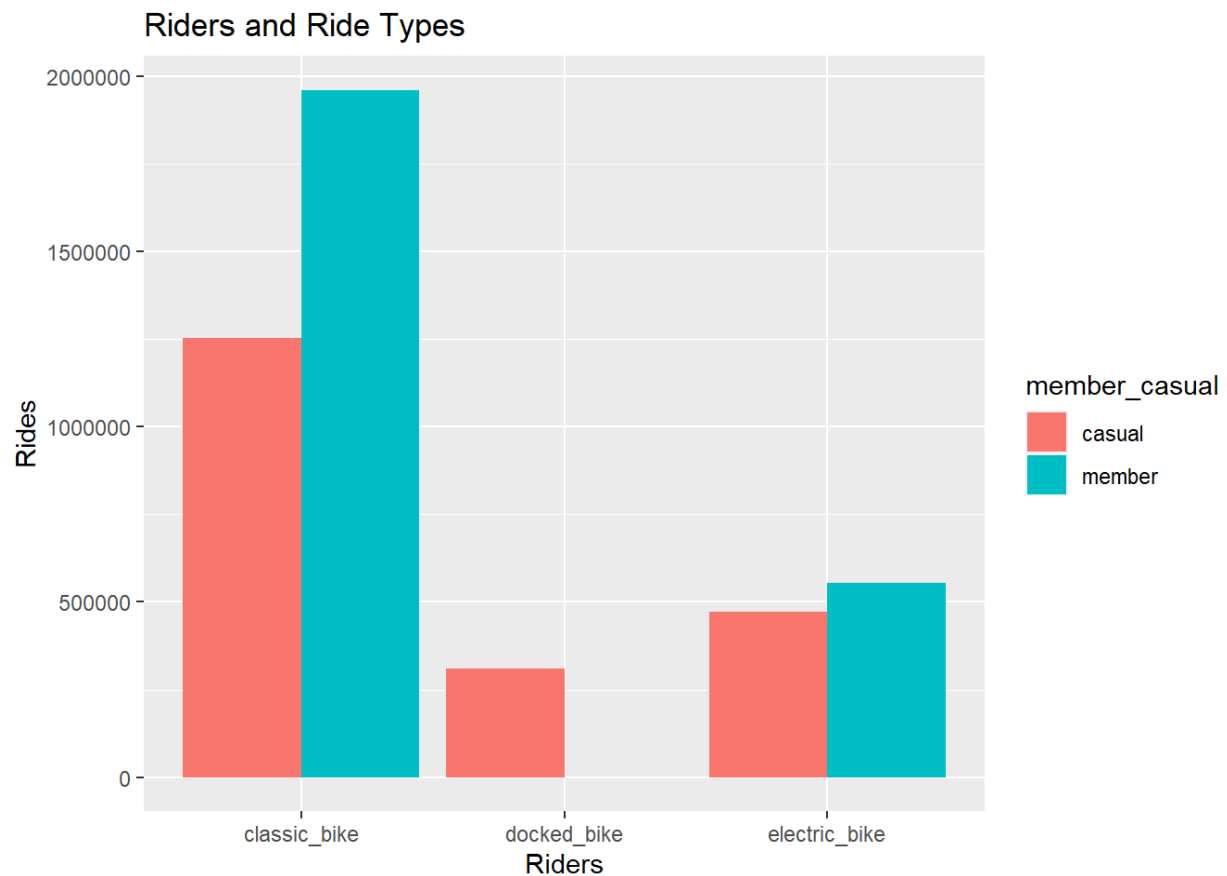
Plot for bike user vs bike type.

```
rt %>%
  filter(member_casual == "member" |
```

```

member_casual == "casual") %>%
ggplot(aes(x = rideable_type, y = Freq, fill = member_casual))+
geom_bar(stat = "identity", position = "dodge") +
labs(title = "Riders and Ride Types",
      x= "Riders",
      y = "Rides")

```



STEP SIX: EXPORT ANALYZED DATA

Save the analyzed data as a new file.

```
fwrite(all_trips, "all_trips.csv")
```