

Cyclistic Case Study Q1_2021

Hezar K

2022-11-29

This is an analysis for Cyclistic Case Study for Google Data Analytics Course. This is an analysis for 2021's first quarter.

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")
```

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()**.

```
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)
```

```
## [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(Mar21)
```

```
## [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"
```

Since all column names are the same. We can combine the data for each month into quarters.

```
q1_2021 <- bind_rows(Jan21, Feb21, Mar21)
```

View(q1_2021)

```
nrow(q1_2021)
```

```
## [1] 374952
```

```
dim(q1_2021)
```

```
## [1] 374952      13
```

```
head(q1_2021)
```

```
## # A tibble: 6 × 13
##   ride_id      ridea...1 started_at      ended_at      start...2 start...3
##   <chr>      <chr>    <dtm>      <dtm>      <chr>    <chr>
## 1 E19E6F1B8D4C4... electr... 2021-01-23 16:14:19 2021-01-23 16:24:44 Califo... 17660
## 2 DC88F20C2C55F... electr... 2021-01-27 18:43:08 2021-01-27 18:47:12 Califo... 17660
## 3 EC45C94683FE3... electr... 2021-01-21 22:35:54 2021-01-21 22:37:14 Califo... 17660
## 4 4FA453A75AE37... electr... 2021-01-07 13:31:13 2021-01-07 13:42:55 Califo... 17660
## 5 BE5E8EB4E7263... electr... 2021-01-23 02:24:02 2021-01-23 02:24:45 Califo... 17660
## 6 5D8969F88C773... electr... 2021-01-09 14:24:07 2021-01-09 15:17:54 Califo... 17660
## # ... with 7 more variables: end_station_name <chr>, end_station_id <chr>,
## #   start_lat <dbl>, start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
## #   member_casual <chr>, and abbreviated variable names 1rideable_type,
## #   2start_station_name, 3start_station_id
```

```
tail(q1_2021)
```

```
## # A tibble: 6 × 13
##   ride_id      ridea...1 started_at      ended_at      start...2 start...3
##   <chr>      <chr>    <dtm>      <dtm>      <chr>    <chr>
## 1 081549DEA616C... electr... 2021-03-14 01:59:38 2021-03-14 03:13:09 Larrab... TA1309...
## 2 9397BDD14798A... docked... 2021-03-20 14:58:56 2021-03-20 17:22:47 Michig... 13042
## 3 BBBEB8D51AAD4... classi... 2021-03-02 11:35:10 2021-03-02 11:43:37 Kingsb... KA1503...
## 4 637FF754DA0BD... classi... 2021-03-09 11:07:36 2021-03-09 11:49:11 Michig... 13042
## 5 F8F43A0B978A7... classi... 2021-03-01 18:11:57 2021-03-01 18:18:37 Kingsb... KA1503...
## 6 3AE64EA5BF43C... electr... 2021-03-26 17:58:14 2021-03-26 18:06:43 <NA>    <NA>
## # ... with 7 more variables: end_station_name <chr>, end_station_id <chr>,
## #   start_lat <dbl>, start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
## #   member_casual <chr>, and abbreviated variable names 1rideable_type,
## #   2start_station_name, 3start_station_id
```

```
summary(q1_2021)
```

```
##   ride_id      rideable_type      started_at
##   Length:374952   Length:374952   Min.    :2021-01-01 00:02:05.0
##   Class :character Class :character   1st Qu.:2021-01-29 19:27:50.0
##   Mode  :character Mode  :character   Median :2021-03-08 16:10:06.5
##                                     Mean  :2021-02-26 11:12:08.9
##                                     3rd Qu.:2021-03-21 13:36:25.0
##                                     Max.   :2021-03-31 23:59:08.0
##
##   ended_at      start_station_name start_station_id
##   Min.    :2021-01-01 00:08:39.00   Length:374952   Length:374952
##   1st Qu.:2021-01-29 19:40:06.50   Class :character Class :character
##   Median :2021-03-08 16:32:02.50   Mode  :character Mode  :character
##   Mean  :2021-02-26 11:33:15.50
##   3rd Qu.:2021-03-21 14:06:50.25
##   Max.   :2021-04-06 11:00:11.00
##
##   end_station_name end_station_id      start_lat      start_lng
##   Length:374952   Length:374952   Min.    :41.64   Min.    :-87.78
##   Class :character Class :character   1st Qu.:41.88   1st Qu.: -87.66
##   Mode  :character Mode  :character   Median :41.90   Median : -87.64
##                                     Mean  :41.90   Mean  : -87.65
##                                     3rd Qu.:41.93   3rd Qu.: -87.63
##                                     Max.   :42.07   Max.   : -87.53
##
##   end_lat      end_lng      member_casual
##   Min.    :41.54   Min.    :-88.07   Length:374952
##   1st Qu.:41.88   1st Qu.: -87.66   Class :character
##   Median :41.90   Median : -87.64   Mode  :character
##   Mean  :41.90   Mean  : -87.65
##   3rd Qu.:41.93   3rd Qu.: -87.63
##   Max.   :42.08   Max.   : -87.51
##   NA's   :484     NA's   :484
```

```
str(q1_2021)
```

```
## spc_tbl_ [374,952 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:374952] "E19E6F1B8D4C42ED" "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA453A75AE377
DB" ...
## $ rideable_type : chr [1:374952] "electric_bike" "electric_bike" "electric_bike" "electric_bike" ...
## $ started_at   : POSIXct[1:374952], format: "2021-01-23 16:14:19" "2021-01-27 18:43:08" ...
## $ ended_at     : POSIXct[1:374952], format: "2021-01-23 16:24:44" "2021-01-27 18:47:12" ...
## $ start_station_name: chr [1:374952] "California Ave & Cortez St" "California Ave & Cortez St" "California Av
e & Cortez St" "California Ave & Cortez St" ...
## $ start_station_id : chr [1:374952] "17660" "17660" "17660" "17660" ...
## $ end_station_name : chr [1:374952] NA NA NA NA ...
## $ end_station_id   : chr [1:374952] NA NA NA NA ...
## $ start_lat        : num [1:374952] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng        : num [1:374952] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ end_lat          : num [1:374952] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng          : num [1:374952] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ member_casual    : chr [1:374952] "member" "member" "member" "member" ...
## - attr(*, "spec")=
## .. cols(
## ..   ride_id = col_character(),
## ..   rideable_type = col_character(),
## ..   started_at = col_datetime(format = ""),
## ..   ended_at = col_datetime(format = ""),
## ..   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.

```
q1_2021$date <- as.Date(q1_2021$started_at)
q1_2021$month <- format(as.Date(q1_2021$date), "%m")
q1_2021$month <- month.name[as.numeric(q1_2021$month)]
q1_2021$day <- format(as.Date(q1_2021$date), "%d")
q1_2021$year <- format(as.Date(q1_2021$date), "%Y")
q1_2021$day_of_week <- format(as.Date(q1_2021$date), "%A")
q1_2021$ride_length <- difftime(q1_2021$ended_at, q1_2021$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(q1_2021$ride_length)
```

```
## [1] FALSE
```

Recheck *ride_length* data type.

```
q1_2021$ride_length <- as.numeric(as.character(q1_2021$ride_length))
is.numeric(q1_2021$ride_length)
```

```
## [1] TRUE
```

STEP THREE: CLEAN DATA

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

```
q1_2021 <- na.omit(q1_2021)
```

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.

```
q1_2021 <- subset(q1_2021, nchar(as.character(ride_id)) == 16)
```

Remove rows with the *ride_length* less than 60 seconds or 1 minute.

```
q1_2021 <- subset (q1_2021, ride_length > 59)
```

STEP FOUR: ANALYZE DATA

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

```
mean(q1_2021$ride_length)
```

```
## [1] 1245.537
```

```
median(q1_2021$ride_length)
```

```
## [1] 689
```

```
max(q1_2021$ride_length)
```

```
## [1] 1900899
```

```
min(q1_2021$ride_length)
```

```
## [1] 60
```

Run a statistical summary of the *ride_length*.

```
summary(q1_2021$ride_length)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       60     398     689    1246    1260 1900899
```

Compare the members and casual users

```
aggregate(q1_2021$ride_length ~ q1_2021$member_casual, FUN = mean)
```

```
##      q1_2021$member_casual q1_2021$ride_length
## 1                      casual          2264.1109
## 2                      member           811.2135
```

```
aggregate(q1_2021$ride_length ~ q1_2021$member_casual, FUN = median)
```

```
##      q1_2021$member_casual q1_2021$ride_length
## 1                      casual             1084
## 2                      member              586
```

```
aggregate(q1_2021$ride_length ~ q1_2021$member_casual, FUN = max)
```

```
##      q1_2021$member_casual q1_2021$ride_length
## 1                      casual          1900899
## 2                      member           88461
```

```
aggregate(q1_2021$ride_length ~ q1_2021$member_casual, FUN = min)
```

```
##      q1_2021$member_casual q1_2021$ride_length
## 1                      casual              60
## 2                      member              60
```

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

```
aggregate(q1_2021$ride_length ~ q1_2021$member_casual + q1_2021$day_of_week, FUN = mean)
```

```
##      q1_2021$member_casual q1_2021$day_of_week q1_2021$ride_length
## 1          casual          Friday          1942.7303
## 2          member          Friday           754.7240
## 3          casual          Monday          2510.5686
## 4          member          Monday           801.6531
## 5          casual          Saturday         2613.8336
## 6          member          Saturday           910.1541
## 7          casual          Sunday          2390.5358
## 8          member          Sunday           922.3624
## 9          casual          Thursday         1618.0957
## 10         member          Thursday           727.4440
## 11         casual          Tuesday          2147.4895
## 12         member          Tuesday           793.4528
## 13         casual          Wednesday        1731.7417
## 14         member          Wednesday         772.0138
```

Sort the days of the week in order.

```
q1_2021$day_of_week <- ordered(q1_2021$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(q1_2021$ride_length ~ q1_2021$member_casual + q1_2021$day_of_week, FUN = mean)

head(x)
```

```
##      q1_2021$member_casual q1_2021$day_of_week q1_2021$ride_length
## 1          casual          Sunday          2390.5358
## 2          member          Sunday           922.3624
## 3          casual          Monday          2510.5686
## 4          member          Monday           801.6531
## 5          casual          Tuesday          2147.4895
## 6          member          Tuesday           793.4528
```

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

```
y <- q1_2021 %>%
  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            19343           2391.
## 2 casual          2            12764           2511.
## 3 casual          3            11523           2147.
## 4 casual          4            10201           1732.
## 5 casual          5             7492           1618.
## 6 casual          6            10216           1943.
```

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

```
table(q1_2021$member_casual)
```

```
##
## casual member
## 98150 230181
```

```
table(q1_2021$rideable_type)
```

```
##
## classic_bike  docked_bike electric_bike
##      245207      18915      64209
```

```
table(q1_2021$day_of_week)
```

```
##
## Sunday Monday Tuesday Wednesday Thursday Friday Saturday
## 46703 46163 46557 45911 37325 42463 63209
```

```
table(q1_2021$month)
```

```
##
## February January March
## 42301 82622 203408
```

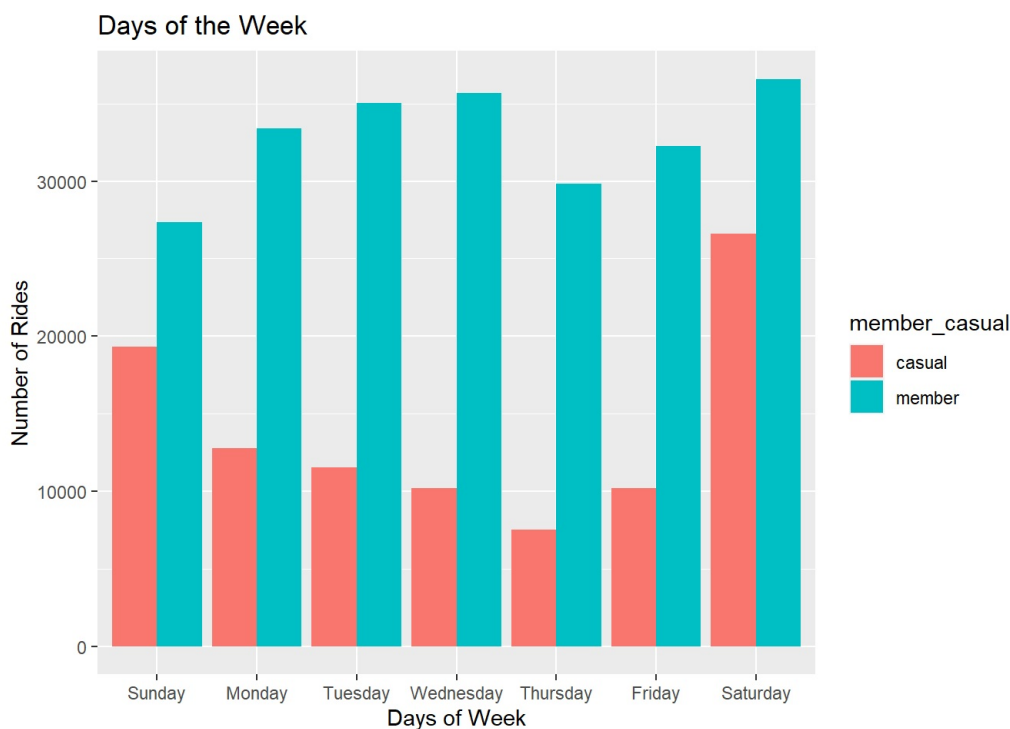
STEP FIVE: VISUALIZATION

Display full digits instead of scientific number.

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

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

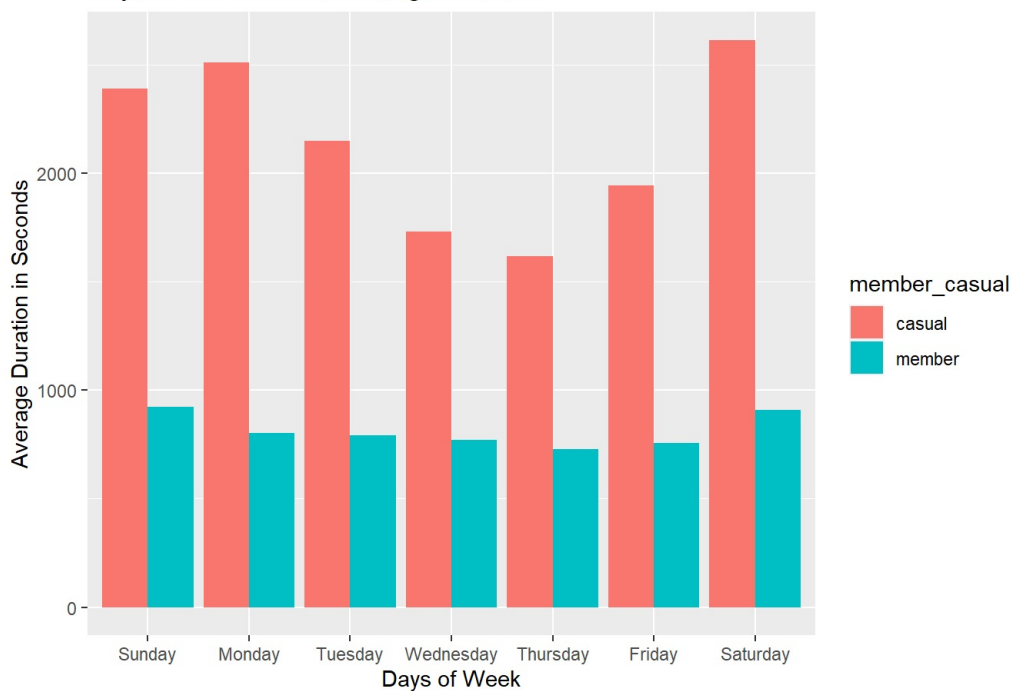
```
q1_2021 %>%
  mutate(day_of_week) %>%
  group_by(member_casual, day_of_week) %>%
  summarise(number_of_rides = n(), average_duration = mean(ride_length), .groups = 'drop') %>%
  arrange(member_casual, day_of_week) %>%
  ggplot(aes(x = day_of_week, y = number_of_rides, fill = member_casual)) +
  geom_col(position = "dodge")+
  labs(x = "Days of Week",
       y = "Number of Rides",
       title= "Days of the Week")
```



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

```
q1_2021 %>%
  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 = "Days of Week",
       y = "Average Duration in Seconds",
       title= "Days of the Week vs Average Duration")
```

Days of the Week vs Average Duration



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

```
mc<- as.data.frame(table(q1_2021$day_of_week,q1_2021$member_casual))
```

Rename columns

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

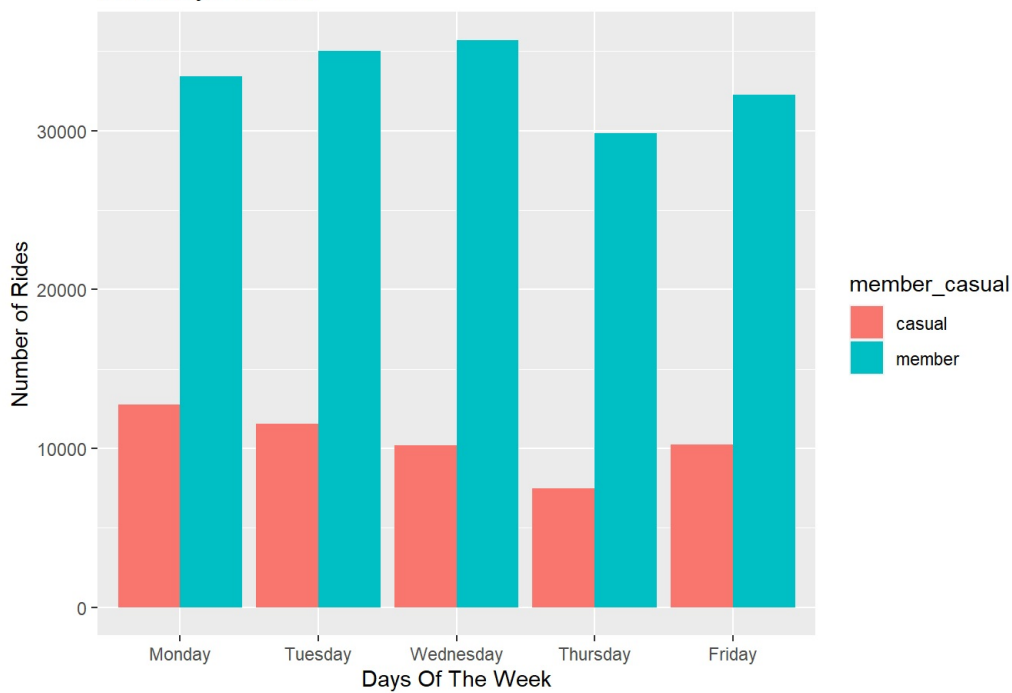
```
head(mc)
```

```
##  day_of_week member_casual  Freq
## 1    Sunday          casual 19343
## 2    Monday          casual 12764
## 3    Tuesday          casual 11523
## 4   Wednesday          casual 10201
## 5   Thursday          casual  7492
## 6    Friday          casual 10216
```

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= "Days Of The Week",
       y = "Number of Rides")
```

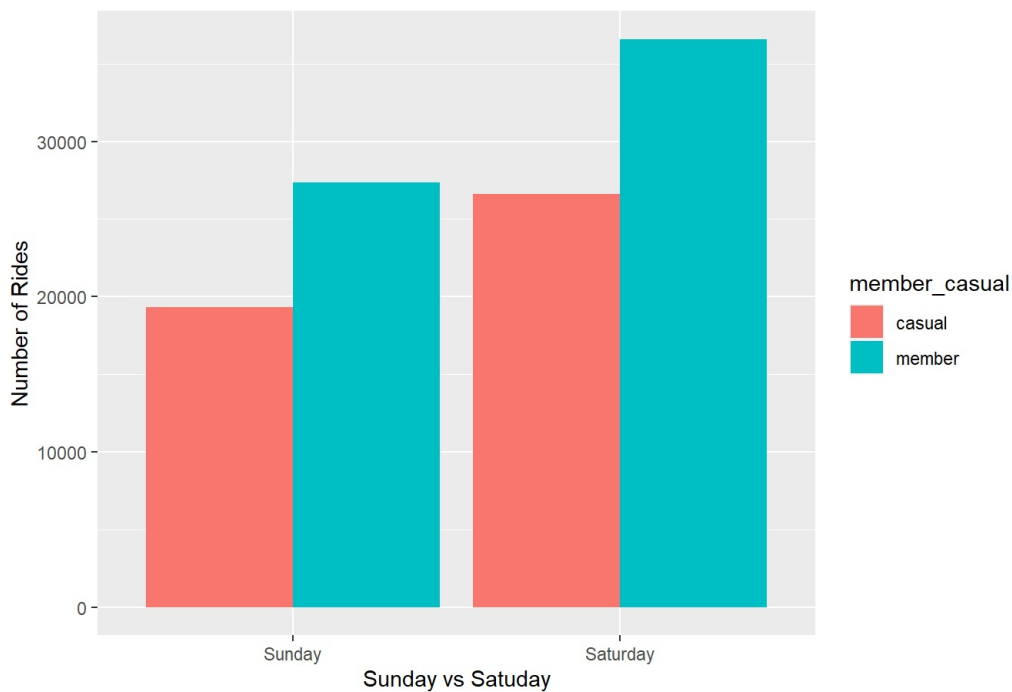
Weekdays Trends



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 = "Sunday vs Saturday",
       y = "Number of Rides")
```

Weekends Trends



Create dataframe for member and casual riders vs ride type

```
rt<- as.data.frame(table(q1_2021$rideable_type,q1_2021$member_casual))
```

Rename columns.

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

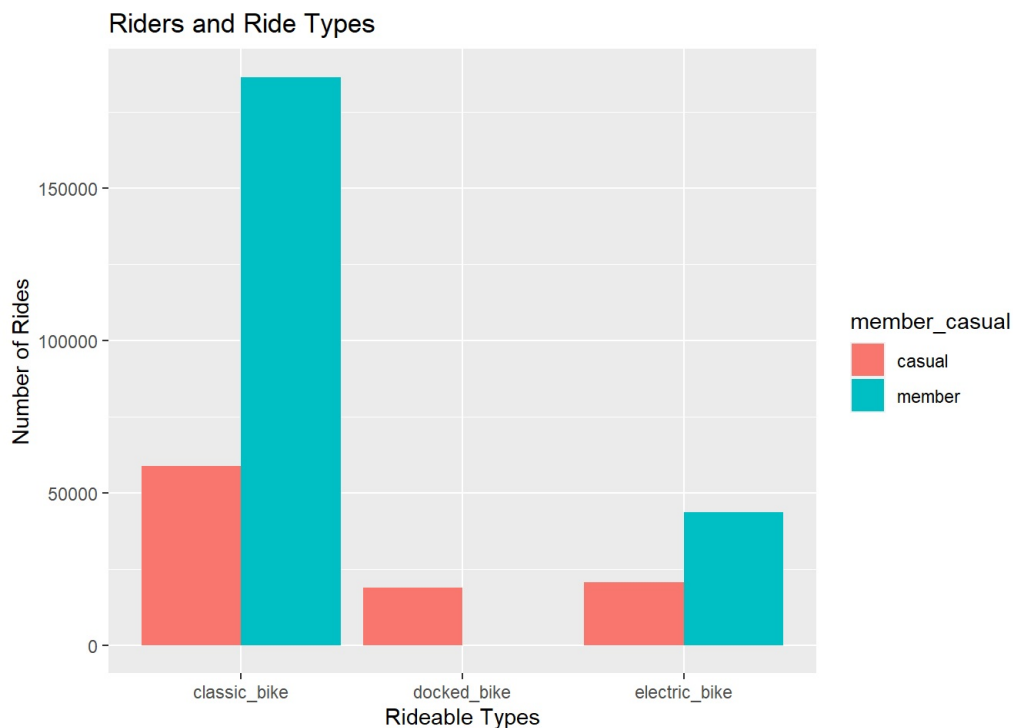
head(rt)
```



```
## rideable_type member_casual Freq
## 1 classic_bike casual 58671
## 2 docked_bike casual 18914
## 3 electric_bike casual 20565
## 4 classic_bike member 186536
## 5 docked_bike member 1
## 6 electric_bike member 43644
```

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 = "Rideable Types",
       y = "Number of Rides")
```



Create vector of month names for Q1 2021

```
q1_months <- c("January", "February", "March")
```

Subset month.name to include only Q1 2021 months

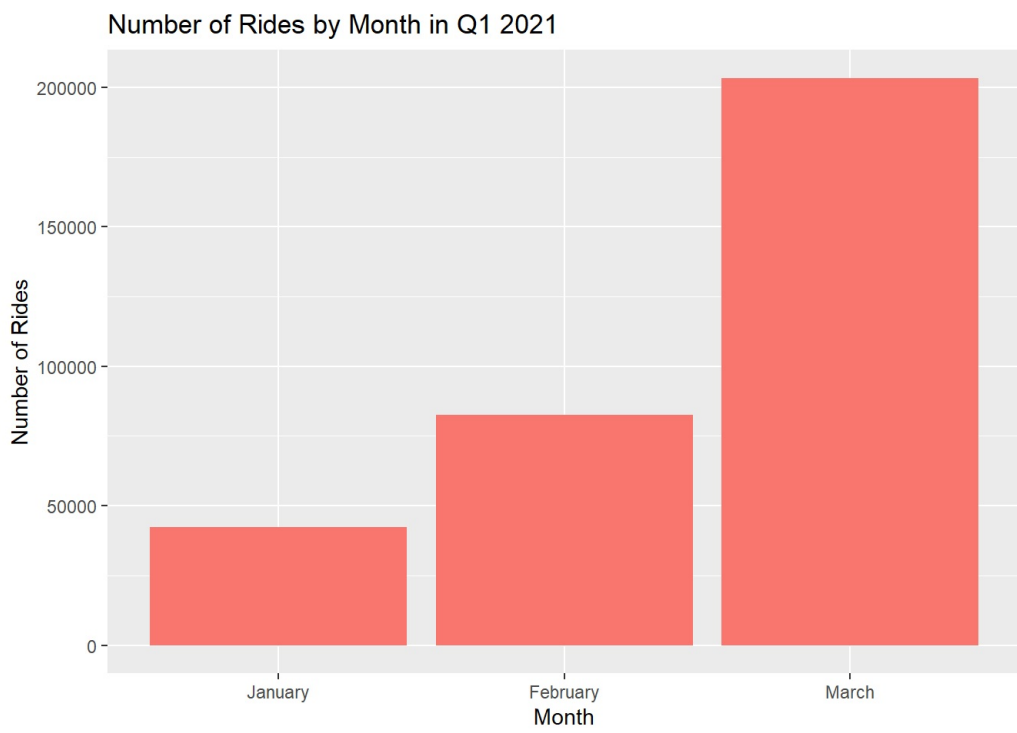
```
q1_month_names <- month.name[match(q1_months, month.name)]
```

Create trips_by_month dataframe with only Q1 2021 months

```
trips_by_month <- data.frame(month = q1_month_names, count = table(q1_2021$month))
```

Set the levels of the month variable in the trips_by_month dataframe

```
trips_by_month$month <- factor(trips_by_month$month, levels = c("January", "February", "March"))
ggplot(trips_by_month, aes(x = month, y = count.Freq)) +
  geom_bar(stat = "identity", fill = "#F8766D") +
  labs(x = "Month", y = "Number of Rides", title = "Number of Rides by Month in Q1 2021")
```



STEP SIX: EXPORT ANALYZED DATA

Save the analyzed data as a new file. `fwrite(q1_2021, "q1_2021.csv")`