

Yearly Bike Sharing Data for Year 2020

Steven Pyae

5/12/2021

Data packages Credited to Divvybikes under this license

Packages Used

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.3      v purrr   0.3.4
## v tibble  3.1.1      v dplyr  1.0.5
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(lubridate)
```

```
##
```

```
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      date, intersect, setdiff, union
```

```
library(ggplot2)
```

```
setwd("D:/R/Data/Google Course Data/CSV Files")
```

- Divvy_Trips_2020_Q1.csv - q1_2020
- 202005-divvy-tripdata.csv - may_2020
- 202006-divvy-tripdata.csv - june_2020
- 202007-divvy-tripdata.csv - july_2020
- 202008-divvy-tripdata.csv - aug_2020
- 202009-divvy-tripdata.csv - sep_2020
- 202010-divvy-tripdata.csv - oct_2020
- 202011-divvy-tripdata.csv - nov_2020
- 202012-divvy-tripdata.csv - dec_2020

Loading Data sets

```
q1_2020 <- read_csv("Divvy_Trips_2020_Q1.csv")
```

```
##
## -- Column specification -----
## 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_double(),
##   end_station_name = col_character(),
##   end_station_id = col_double(),
##   start_lat = col_double(),
##   start_lng = col_double(),
##   end_lat = col_double(),
##   end_lng = col_double(),
##   member_casual = col_character()
## )
```

```
may_2020 <- read_csv("202005-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
##   end_station_name = col_character(),
##   end_station_id = col_double(),
##   start_lat = col_double(),
##   start_lng = col_double(),
##   end_lat = col_double(),
##   end_lng = col_double(),
##   member_casual = col_character()
## )
```

```
june_2020 <- read_csv("202006-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
## end_station_name = col_character(),
## end_station_id = col_double(),
## start_lat = col_double(),
## start_lng = col_double(),
## end_lat = col_double(),
## end_lng = col_double(),
## member_casual = col_character()
## )
```

```
july_2020<- read_csv("202007-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
##   end_station_name = col_character(),
##   end_station_id = col_double(),
##   start_lat = col_double(),
##   start_lng = col_double(),
##   end_lat = col_double(),
##   end_lng = col_double(),
##   member_casual = col_character()
## )
```

```
aug_2020<- read_csv("202008-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
##   end_station_name = col_character(),
##   end_station_id = col_double(),
##   start_lat = col_double(),
##   start_lng = col_double(),
##   end_lat = col_double(),
##   end_lng = col_double(),
##   member_casual = col_character()
## )
```

```
sep_2020<- read_csv("202009-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
##   end_station_name = col_character(),
##   end_station_id = col_double(),
##   start_lat = col_double(),
##   start_lng = col_double(),
##   end_lat = col_double(),
##   end_lng = col_double(),
##   member_casual = col_character()
## )
```

```
oct_2020<- read_csv("202010-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
##   end_station_name = col_character(),
##   end_station_id = col_double(),
##   start_lat = col_double(),
##   start_lng = col_double(),
##   end_lat = col_double(),
##   end_lng = col_double(),
##   member_casual = col_character()
## )
```

```
nov_2020 <- read_csv("202011-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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_double(),
```

```
## end_station_name = col_character(),
## end_station_id = col_double(),
## start_lat = col_double(),
## start_lng = col_double(),
## end_lat = col_double(),
## end_lng = col_double(),
## member_casual = col_character()
## )
```

```
dec_2020 <- read_csv("202012-divvy-tripdata.csv")
```

```
##
## -- Column specification -----
## 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()
## )
```

July_2020's Dates are formatted into POSIXct to bind with other datasets

```
july_2020 <- mutate(july_2020, started_at = as.POSIXct(format(strptime(started_at, "%d/%m/%Y %H:%M"), "%d/%m/%Y %H:%M"), "%d/%m/%Y %H:%M"),
july_2020 <- mutate(july_2020, ended_at = as.POSIXct(format(strptime(ended_at, "%d/%m/%Y %H:%M"), "%d/%m/%Y %H:%M"), "%d/%m/%Y %H:%M"))
```

Binding all the datasets into one

```
all_trips <- bind_rows(q1_2020, may_2020, june_2020, july_2020, aug_2020, sep_2020, oct_2020, nov_2020)
```

Data Cleaning

```
# Remove the lat, lng, lat, lng
all_trips <- all_trips %>%
  select(-c(start_lat, start_lng, end_lat, end_lng))
# check for consistent name conventions
table(all_trips$member_casual)
```

```
##
## casual member
## 1312867 2012467
```

```
#add columns that list the date, month, day and year of each ride
all_trips$date <- as.Date(all_trips$started_at) #The default format is yyyy-mm-dd
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")
```

Adding additional columns: ride_length, day_of_week

```
all_trips$ride_length <- difftime(all_trips$ended_at, all_trips$started_at)
# Convert "ride_length" from Factor to numeric so we can run calculations on the data
is.factor(all_trips$ride_length)
```

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

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

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

Data contains negative ride_length and non-applicable numbers, Filtering

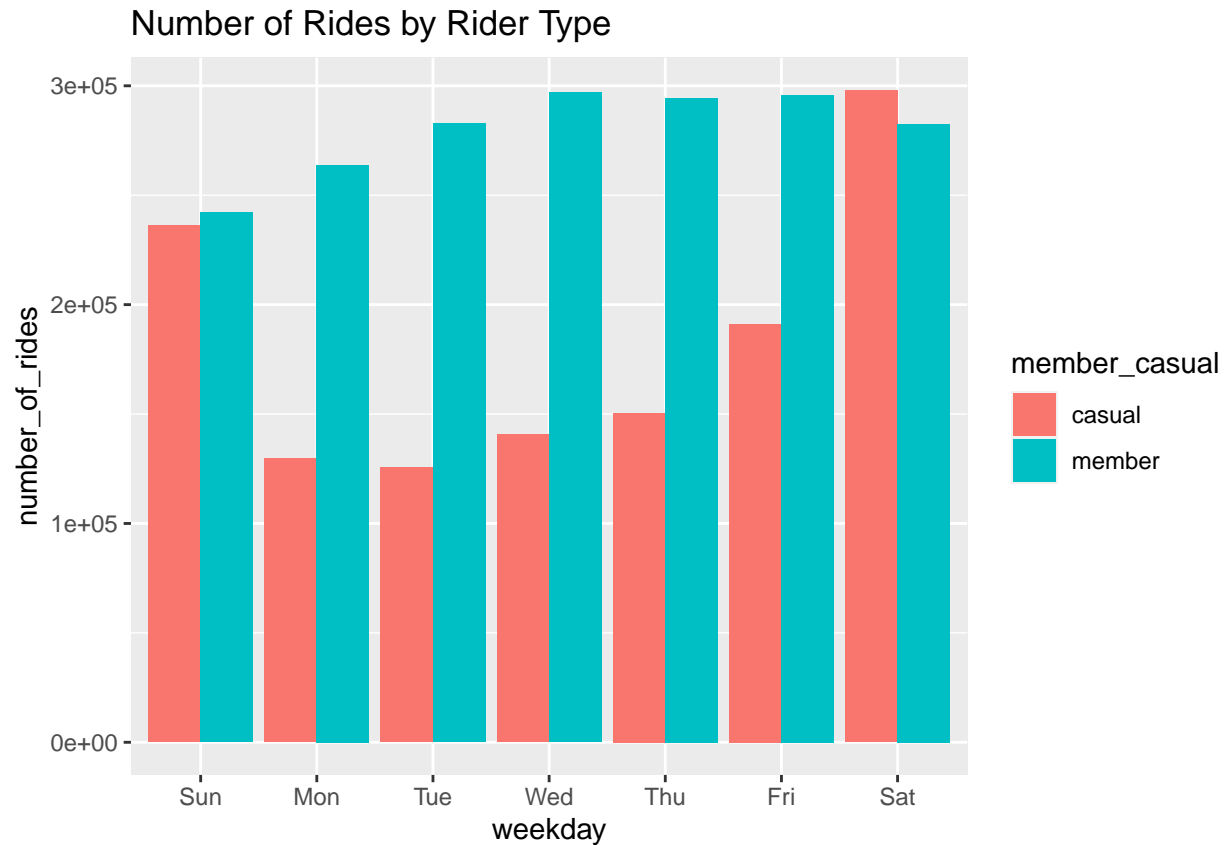
```
all_trips_v2 <- all_trips[!(all_trips$start_station_name == "HQ QR" | all_trips$ride_length < 0),]
# Remove non existent ride lengths
all_trips_v2 <- all_trips_v2[!is.na(all_trips_v2$ride_length),]
```

Plots

Analysing the results

```
all_trips_v2 %>%
  mutate(weekday = wday(started_at, label = TRUE)) %>%
  group_by(member_casual, weekday) %>%
  summarise(number_of_rides = n()
            , average_duration = mean(ride_length)) %>%
  arrange(member_casual, weekday) %>%
  ggplot(aes(x = weekday, y = number_of_rides, fill = member_casual)) +
  geom_col(position = "dodge") + labs(title = "Number of Rides by Rider Type")
```

```
## 'summarise()' has grouped output by 'member_casual'. You can override using the '.groups' argument.
```



```
all_trips_v2 %>%
  mutate(weekday = wday(started_at, label = TRUE)) %>%
  group_by(member_casual, weekday) %>%
  summarise(number_of_rides = n()
            ,average_duration = mean(ride_length)) %>%
  arrange(member_casual, weekday) %>%
  ggplot(aes(x = weekday, y = average_duration, fill = member_casual)) +
  geom_col(position = "dodge") + labs(title = "Average Duration by Rider Type")
```

'summarise()' has grouped output by 'member_casual'. You can override using the '.groups' argument.

