

Cyclic Case Study Dec21

Hezar K

2022-11-29

This is an analysis for Cyclic Case Study for Google Data Analytics Course. This is an analysis for December 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)

## Attaching packages
## ✓ ggplot2 3.4.0      ✓ purrr  0.3.0
## ✓ lubridate 3.1.0    ✓ dplyr  0.9.0
## ✓ tidyr  1.2.1      ✓ stringr 1.4.1
## ✓ readr  2.1.3      ✓ forcats 0.5.2
## ✓ Conflicts:
## ✓ dplyr::filter() masks stats::filter()
## ✓ dplyr::lag()   masks stats::lag()

library(lubridate)
```

```
# Loading required package: timechange
## Attaching package: 'lubridate'
##
## The following objects are masked from 'package:base':
##
##   date, intersect, setdiff, union
```

```
library(data.table)

##
## Attaching package: 'data.table'
##
## The following objects are masked from 'package:lubridate':
##
##   hour, isoweek, mday, minute, month, quarter, second, wday, week,
##   yday, year
##
## The following objects are masked from 'package:dplyr':
##
##   between, first, last
##
## The following object is masked from 'package:purrr':
##
##   transpose
```

```
library(ggplot2)
library(anytime)
```

Import data from local drive.

```
Dec21 <- read_csv("C:/Users/thiby/Documents/20212-12ivvy-tripdata.csv")

## Rows: 247549 Columns: 13
## --- Column specification ---
## dbl(1): ride_id, rideable_type, started_at, ended_at, start_station_name, ...
## dbl(4): start_lat, start_lng, end_lat, end_lng
## Use `spec()` to retrieve the full column specification for this data.
## 1 Specify the column types or set `show_col_types` to FALSE to quiet this message.
```

STEP TWO: EXAMINE THE DATA

Examine the dataframe for an overview of the data. Review column names, **colnames()**, dimensions of the dataframe by row and column, **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()**.

```
View(Dec21)

colnames(Dec21)
```

```
## [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"
```

```
nrow(Dec21)
```

```
## [1] 247549
```

```
dim(Dec21)
```

```
## [1] 247549    13
```

```
head(Dec21)
```

```
## # A tibble: 6 x 13
##   ride_id rideable_type started_at ended_at start_station_name end_station_name start_station_id end_station_id start_lat start_lng end_lat end_lng member_casual
##   <dbl> <chr> <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 46F8167220644 electric 12/7/21 12/7/21 Laflin 13387 Morgan TA1397 41.9
## 2 7A7176263088A electric 12/15/ 12/15/ LaSalle KP1785 Claron TA1397 41.9
## 3 4F4245205465 electric 12/7/21 12/7/21 LaSalle KP1785 Claron TA1397 41.9
## 4 278A87896968A classic 12/26/ 12/26/ LaSalle KP1785 Claron TA1397 41.9
## 5 4F4245205465 electric 12/7/21 12/7/21 LaSalle KP1785 Claron TA1397 41.9
## 6 93E8D79490E3A classic 12/7/21 12/7/21 LaSalle KP1785 Claron TA1397 41.9
## # with 4 more variables: start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
## # member_casual <chr>, and abbreviated variable names `rideable_type`,
## # `started_at`, `ended_at`, `start_station_name`, `start_station_id`,
## # `end_station_name`, `end_station_id`, `start_lat`
```

```
tail(Dec21)
```

```
## # A tibble: 6 x 13
##   ride_id rideable_type started_at ended_at start_station_name end_station_name start_station_id end_station_id start_lat start_lng end_lat end_lng member_casual
##   <dbl> <chr> <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 93B8A97091688A electric 12/24/ 12/24/ Canal 13341 <NA> <NA> 41.9
## 2 8A7431F30533A electric 12/27/ 12/27/ Canal 13341 <NA> <NA> 41.9
## 3 4F4245205465 electric 12/7/21 12/7/21 Canal 13341 <NA> <NA> 41.9
## 4 698B96E8F544A electric 12/7/21 12/7/21 Canal 13341 Denro TA1395 41.9
## 5 5141694AC098A electric 12/7/21 12/7/21 Canal 13341 <NA> <NA> 41.9
## 6 37AC57E348267 classic 12/13/ 12/13/ Michig TA1398 Denro TA1395 41.9
## # with 4 more variables: start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
## # member_casual <chr>, and abbreviated variable names `rideable_type`,
## # `started_at`, `ended_at`, `start_station_name`, `start_station_id`,
## # `end_station_name`, `end_station_id`, `start_lat`
```

```
summary(Dec21)
```

```
##   ride_id      rideable_type      started_at      ended_at
## Length:247549      Length:247549      Length:247549      Length:247549
## 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:247549      Length:247549      Length:247549      Length:247549
## 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.   -141.64   Min.   -187.64   Min.   -141.48   Min.   -87.85
## 1st Qu. -141.88   1st Qu. -187.67   1st Qu. -141.88   1st Qu. -87.87
## Median -141.90   Median -187.64   Median -141.90   Median -87.64
## Mean    -141.90   Mean    -187.65   Mean    -141.90   Mean    -87.65
## 3rd Qu. -141.93   3rd Qu. -187.63   3rd Qu. -141.93   3rd Qu. -87.63
## Max.    -142.07   Max.    -187.52   Max.    -142.07   Max.    -87.52
##
##   member_casual
## Length:247549
## Class:character
## Mode:character
##
##
##   member_casual
## Length:247549
## Class:character
## Mode:character
```

```
str(Dec21)
```

```
## spec_tbl_([247,549 x 13]) (S3: spec_tbl_df/tbl_df/tbl/data.frame)
##   $ ride_id      : chr [1:247549] "46F816722064432P" "7A71762630882P" "4CF42452054F59C5" "3278A878F6983
##   $ rideable_type : chr [1:247549] "electric.bike" "electric.bike" "electric.bike" "classic.bike" ...
##   $ started_at   : chr [1:247549] "12/7/2021 15:18" "12/21/2021 3:43" "12/15/2021 23:30" "12/26/2021 16:1
##   $ ended_at     : chr [1:247549] "12/7/2021 15:19" "12/15/2021 4:18" "12/15/2021 23:33" "12/26/2021 16:3
##   $ start_station_name: chr [1:247549] "Laflin St & Culbertson St" "LaSalle Dr & Huron St" "Halsted St & North B
##   $ start_station_id : chr [1:247549] "13387" "KP1785081826" "KA1504000117" "KA1504000117" ...
##   $ end_station_name : chr [1:247549] "Morgan St & Polk St" "Clarendon Ave & Leland Ave" "Broadway & Barry A
##   $ end_station_id   : chr [1:247549] "TA1387980118" "TA1387980118" "13137" "KP1785091026" ...
##   $ start_lat       : num [1:247549] 41.9 41.9 41.9 41.9 41.9 ...
##   $ start_lng       : num [1:247549] -87.7 -87.6 -87.6 -87.6 -87.7 ...
##   $ end_lat         : num [1:247549] 41.9 42.41 41.9 41.9 41.9 ...
##   $ end_lng         : num [1:247549] -87.7 -87.7 -87.6 -87.6 -87.6 ...
##   $ member_casual   : chr [1:247549] "member" "casual" "member" "member" ...
##   $ attr(,"spec")=
##   .. col()
##   ..   ride_id = 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 convert from character data type to date data type. **str()** syntax confirms changes.

```
Dec21$started_at <- any_time(Dec21$started_at)
Dec21$ended_at <- any_time(Dec21$ended_at)
str(Dec21)
```

```
## spec_tbl_([247,549 x 13]) (S3: spec_tbl_df/tbl_df/tbl/data.frame)
##   $ ride_id      : chr [1:247549] "46F816722064432P" "7A71762630882P" "4CF42452054F59C5" "3278A878F6983
##   $ rideable_type : chr [1:247549] "electric.bike" "electric.bike" "electric.bike" "classic.bike" ...
##   $ started_at   : POSIXct[1:247549], format: "2021-12-07 15:08:00" "2021-12-11 03:43:00" ...
##   $ ended_at     : POSIXct[1:247549], format: "2021-12-07 15:18:00" "2021-12-11 04:10:00" ...
##   $ start_station_name: chr [1:247549] "Laflin St & Culbertson St" "LaSalle Dr & Huron St" "Halsted St & North B
##   $ start_station_id : chr [1:247549] "13387" "KP1785081826" "KA1504000117" "KA1504000117" ...
##   $ end_station_name : chr [1:247549] "Morgan St & Polk St" "Clarendon Ave & Leland Ave" "Broadway & Barry Av
##   $ end_station_id   : chr [1:247549] "TA1387980118" "TA1387980118" "13137" "KP1785091026" ...
##   $ start_lat       : num [1:247549] 41.9 41.9 41.9 41.9 41.9 ...
##   $ start_lng       : num [1:247549] -87.7 -87.6 -87.6 -87.6 -87.7 ...
##   $ end_lat         : num [1:247549] 41.9 42.41 41.9 41.9 41.9 ...
##   $ end_lng         : num [1:247549] -87.7 -87.7 -87.6 -87.6 -87.6 ...
##   $ member_casual   : chr [1:247549] "member" "casual" "member" "member" ...
##   $ attr(,"spec")=
##   .. col()
##   ..   ride_id = 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.

```
Dec21$date <- as.Date(Dec21$started_at)
Dec21$month <- format(as.Date(Dec21$date), "%m")
Dec21$day <- format(as.Date(Dec21$date), "%d")
Dec21$year <- format(as.Date(Dec21$date), "%Y")
Dec21$day_of_week <- format(as.Date(Dec21$date), "%A")
Dec21$ride_length <- difftime(Dec21$ended_at, Dec21$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(Dec21$ride_length)
```

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

Recheck ride_length data type.

```
Dec21$ride_length <- as.numeric(as.character(Dec21$ride_length))
is.numeric(Dec21$ride_length)
```

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

STEP THREE: CLEAN DATA

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

```
Dec21 <- na.omit(Dec21)
```

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.

```
Dec21 <- subset(Dec21, nchar(as.character(ride_id)) == 16)
```

Remove rows with the ride_length less than 1 minute.

```
Dec21 <- subset(Dec21, ride_length > "1")
```

STEP FOUR: ANALYZE DATA

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

```
mean(Dec21$ride_length)
```

```
## [1] 860.5358
```

```
median(Dec21$ride_length)
```

```
## [1] 540
```

```
max(Dec21$ride_length)
```

```
## [1] 1824000
```

```
min(Dec21$ride_length)
```

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

Run a statistical summary of the ride_length.

```
summary(Dec21$ride_length)
```

```
##      Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
##    60.0      300.0      540.0      860.5     900.8    1824000.0
```

Compare the members and casual users

```
aggregate(Dec21$ride_length ~ Dec21$member_casual, FUN = mean)
```

```
##   Dec21$member_casual Dec21$ride_length
## 1      casual          1498.1462
## 2      member           646.1371
```

```
aggregate(Dec21$ride_length ~ Dec21$member_casual, FUN = median)
```

```
##   Dec21$member_casual Dec21$ride_length
## 1      casual           720
## 2      member           480
```

```
aggregate(Dec21$ride_length ~ Dec21$member_casual, FUN = max)
```

```
##   Dec21$member_casual Dec21$ride_length
## 1      casual          1824000
## 2      member          73860
```

```
aggregate(Dec21$ride_length ~ Dec21$member_casual, FUN = min)
```

```
##   Dec21$member_casual Dec21$ride_length
## 1      casual           60
## 2      member           60
```

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

```
aggregate(Dec21$ride_length ~ Dec21$member_casual + Dec21$day_of_week, FUN = mean)
```

```
##   Dec21$member_casual Dec21$day_of_week Dec21$ride_length
## 1      casual      Friday          1361.6466
## 2      member      Friday          646.0598
## 3      casual      Monday          1331.2799
## 4      member      Monday          618.0515
## 5      casual      Saturday        1469.7190
## 6      member      Saturday          696.1138
## 7      casual      Sunday          1511.0706
## 8      member      Sunday          691.0534
## 9      casual      Thursday        1518.3089
## 10     member      Thursday          639.2543
## 11     casual      Tuesday          1459.9058
## 12     member      Tuesday          686.0767
## 13     casual      Wednesday        1551.0541
## 14     member      Wednesday          623.0991
```

Sort the days of the week in order.

```
Dec21$day_of_week <- ordered(Dec21$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(Dec21$ride_length ~ Dec21$member_casual + Dec21$day_of_week, FUN = mean)
```

```
head(x)
```

```
##   Dec21$member_casual Dec21$day_of_week Dec21$ride_length
## 1      casual      Sunday          1361.6466
## 2      member      Sunday          691.0534
## 3      casual      Monday          1331.2799
## 4      member      Monday          618.0515
## 5      casual      Tuesday          1459.9058
## 6      member      Tuesday          686.0767
```

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

```
y <- Dec21 %>%
  mutate(weekday = weekday(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 x 4
##   member_casual weekday number_of_rides average_duration
##   <chr>          <int>      <int>
## 1 casual      Sunday          3          5867
## 2 casual      Monday         2          4899
## 3 casual      Tuesday         3          5865
## 4 casual      Wednesday        6          6622
## 5 casual      Thursday         5          8992
## 6 casual      Friday         6          8320
```

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

```
table(Dec21$member_casual)
```

```
##
## casual member
## 44844 139086
```

```
table(Dec21$rideable_type)
```

```
##
## classic_bike docked_bike electric_bike
## 99519      4851      76560
```

STEP FIVE: VISUALIZATION

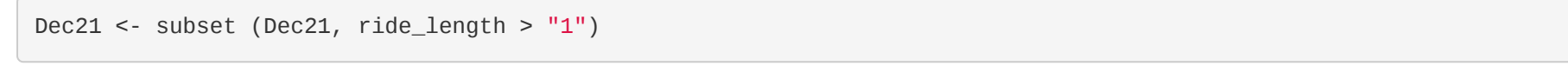
Display full digits instead of scientific number.

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

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

```
Dec21 %>%
  mutate(day_of_week = "Sunday") %>%
  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_bar(position = "dodge") +
  labs(x = "Day Of The Week",
       y = "Number of Rides",
       title = "Days of the Week")
```

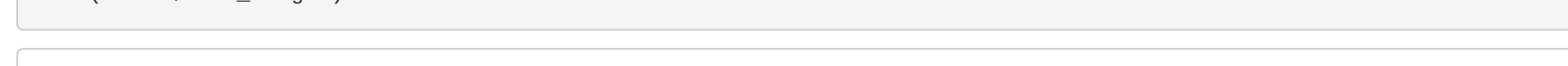
Days of the Week



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

```
Dec21 %>%
  mutate(day_of_week = "Sunday") %>%
  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_bar(position = "dodge") +
  labs(x = "Day Of The 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 weekdays trends vs weekend trends.

```
Dec <- as.data.frame(table(Dec21$day_of_week, Dec21$member_casual))
```

Rename columns

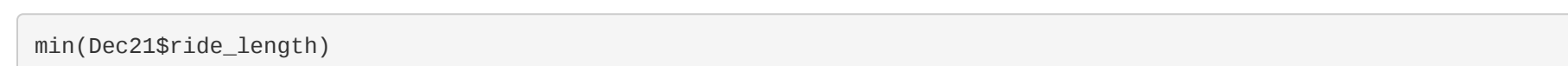
```
Dec <- rename(Dec, day_of_week = Var1, member_casual = Var2)
head(Dec)
```

```
##   day_of_week member_casual Freq
## 1 Sunday      casual      5907
## 2 Monday      casual      4899
## 3 Tuesday      casual      3865
## 4 Wednesday      casual      6622
## 5 Thursday      casual      8992
## 6 Friday      casual      8320
```

Weekday trends (Monday through Friday)

```
Dec %>%
  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")
```

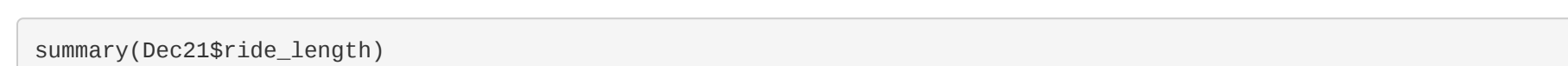
Weekdays Trends



Weekend trends (Sunday and Saturday)

```
Dec %>%
  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")
```

Weekends Trends



Create dataframe for member and casual riders vs ride type

```
rt <- as.data.frame(table(Dec21$rideable_type, Dec21$member_casual))
```

Rename columns

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

```
##   rideable_type member_casual Freq
## 1 classic_bike      casual      19600
## 2 docked_bike      casual      4851
## 3 electric_bike      casual      39393
## 4 classic_bike      member      79919
## 5 docked_bike      member       9
## 6 electric_bike      member      9507
```

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")
```

Riders and Ride Types



STEP SIX: EXPORT ANALYZED DATA

Save the analyzed data as a new file. **write_csv(Dec21, Dec21.csv)**