

Divvy study

Miguel

18/01/2022

Divvy_Exercise_Full_Year_Analysis

This analysis is based on the Divvy case study “‘Sophisticated, Clear, and Polished’: Divvy and Data Visualization” written by Kevin Hartman (found here: <https://artsience.blog/home/divvy-dataviz-case-study>). The purpose of this script is to consolidate downloaded Divvy data into a single dataframe and then conduct simple analysis to help answer the key question:

Install required packages

tidyverse for data import and wrangling

lubridate for date functions

ggplot for visualization

```
library(tidyverse) #helps wrangle data
```

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

```
## v ggplot2 3.3.5      v purrr  0.3.4
## v tibble  3.1.5      v dplyr  1.0.7
## v tidyr   1.1.4      v stringr 1.4.0
## v readr   2.0.2      v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(lubridate) #helps wrangle date attributes
```

```
##
```

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

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

```
##
```

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

```
library(ggplot2) #helps visualize data
getwd() #displays your working directory
```

```
## [1] "C:/Users/maap_/Documents/trip data/TRIP_2021"
```

```
setwd("/Users/maap_/Documents/trip data/TRIP_2021") #sets your working directory to simplify calls
```

```
#COLLECT DATA Trips 2021 Upload Divvy datasets (csv files) here
```

```
trip_01 <- read_csv("202101-divvy-tripdata.csv")
```

```
## Rows: 96834 Columns: 13
```

```

## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_02 <- read_csv("202102-divvy-tripdata.csv")

## Rows: 49622 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_03 <- read_csv("202103-divvy-tripdata.csv")

## Rows: 228496 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_04 <- read_csv("202104-divvy-tripdata.csv")

## Rows: 337230 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_05 <- read_csv("202105-divvy-tripdata.csv")

## Rows: 531633 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at

```

```

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_06 <- read_csv("202106-divvy-tripdata.csv")

## Rows: 729595 Columns: 13
## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_07 <- read_csv("202107-divvy-tripdata.csv")

## Rows: 822410 Columns: 13
## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_08 <- read_csv("202108-divvy-tripdata.csv")

## Rows: 804352 Columns: 13
## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
trip_09 <- read_csv("202109-divvy-tripdata.csv")

## Rows: 756147 Columns: 13
## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl (4): start_lat, start_lng, end_lat, end_lng
## dtm (2): started_at, ended_at
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

```

```

trip_10 <- read_csv("202101-divvy-tripdata.csv")

## Rows: 96834 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr  (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl  (4): start_lat, start_lng, end_lat, end_lng
## dtm  (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

trip_11 <- read_csv("202101-divvy-tripdata.csv")

## Rows: 96834 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr  (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl  (4): start_lat, start_lng, end_lat, end_lng
## dtm  (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

trip_12 <- read_csv("202101-divvy-tripdata.csv")

## Rows: 96834 Columns: 13

## -- Column specification -----
## Delimiter: ","
## chr  (7): ride_id, rideable_type, start_station_name, start_station_id, end...
## dbl  (4): start_lat, start_lng, end_lat, end_lng
## dtm  (2): started_at, ended_at

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

```

WRANGLE DATA AND COMBINE INTO A SINGLE FILE

#Compare

column names each of the files While the names don't have to be in the same order, they DO need to match perfectly before we can use a command to join them into one file.

```

colnames(trip_01)

## [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(trip_02)

```

```
## [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(trip_03)
```

```
## [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(trip_04)
```

```
## [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"
```

Inspect the dataframes and look for incongruencies

```
str(trip_01)
```

```
## spec_tbl_df [96,834 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:96834] "E19E6F1B8D4C42ED" "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA453..."
## $ rideable_type: chr [1:96834] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at   : POSIXct[1:96834], format: "2021-01-23 16:14:19" "2021-01-27 18:43:08" ...
## $ ended_at     : POSIXct[1:96834], format: "2021-01-23 16:24:44" "2021-01-27 18:47:12" ...
## $ start_station_name: chr [1:96834] "California Ave & Cortez St" "California Ave & Cortez St" "California Ave & Cortez St" ...
## $ start_station_id : chr [1:96834] "17660" "17660" "17660" "17660" ...
## $ end_station_name : chr [1:96834] NA NA NA NA ...
## $ end_station_id   : chr [1:96834] NA NA NA NA ...
## $ start_lat       : num [1:96834] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng       : num [1:96834] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ end_lat         : num [1:96834] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng         : num [1:96834] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ member_casual   : chr [1:96834] "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>

str(trip_02)

## spec_tbl_df [49,622 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:49622] "89E7AA6C29227EFF" "0FEFDE2603568365" "E6159D746B2DBB91" "B32D3 ..
## $ rideable_type : chr [1:49622] "classic_bike" "classic_bike" "electric_bike" "classic_bike" ..
## $ started_at    : POSIXct[1:49622], format: "2021-02-12 16:14:56" "2021-02-14 17:52:38" ...
## $ ended_at      : POSIXct[1:49622], format: "2021-02-12 16:21:43" "2021-02-14 18:12:09" ...
## $ start_station_name: chr [1:49622] "Glenwood Ave & Touhy Ave" "Glenwood Ave & Touhy Ave" "Clark St ..
## $ start_station_id : chr [1:49622] "525" "525" "KA1503000012" "637" ...
## $ end_station_name : chr [1:49622] "Sheridan Rd & Columbia Ave" "Bosworth Ave & Howard St" "State ..
## $ end_station_id   : chr [1:49622] "660" "16806" "TA1305000029" "TA1305000034" ...
## $ start_lat       : num [1:49622] 42 42 41.9 41.9 41.8 ...
## $ start_lng       : num [1:49622] -87.7 -87.7 -87.6 -87.7 -87.6 ...
## $ end_lat         : num [1:49622] 42 42 41.9 41.9 41.8 ...
## $ end_lng         : num [1:49622] -87.7 -87.7 -87.6 -87.7 -87.6 ...
## $ member_casual   : chr [1:49622] "member" "casual" "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>

str(trip_03)

## spec_tbl_df [228,496 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:228496] "CFA86D4455AA1030" "30D9DC61227D1AF3" "846D87A15682A284" "994D ..
## $ rideable_type : chr [1:228496] "classic_bike" "classic_bike" "classic_bike" "classic_bike" ..
## $ started_at    : POSIXct[1:228496], format: "2021-03-16 08:32:30" "2021-03-28 01:26:28" ...
## $ ended_at      : POSIXct[1:228496], format: "2021-03-16 08:36:34" "2021-03-28 01:36:55" ...
## $ start_station_name: chr [1:228496] "Humboldt Blvd & Armitage Ave" "Humboldt Blvd & Armitage Ave" ..
## $ start_station_id : chr [1:228496] "15651" "15651" "15443" "TA1308000021" ...
## $ end_station_name : chr [1:228496] "Stave St & Armitage Ave" "Central Park Ave & Bloomingdale Ave ..
## $ end_station_id   : chr [1:228496] "13266" "18017" "TA1308000043" "13323" ...
## $ start_lat       : num [1:228496] 41.9 41.9 41.8 42 42 ...
## $ start_lng       : num [1:228496] -87.7 -87.7 -87.6 -87.7 -87.7 ...
## $ end_lat         : num [1:228496] 41.9 41.9 41.8 42 42.1 ...
## $ end_lng         : num [1:228496] -87.7 -87.7 -87.6 -87.6 -87.7 ...
## $ member_casual   : chr [1:228496] "casual" "casual" "casual" "casual" ...
## - 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>

str(trip_04)

## spec_tbl_df [337,230 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:337230] "6C992BD37A98A63F" "1E0145613A209000" "E498E15508A80BAD" "1887
## $ rideable_type : chr [1:337230] "classic_bike" "docked_bike" "docked_bike" "classic_bike" ...
## $ started_at   : POSIXct[1:337230], format: "2021-04-12 18:25:36" "2021-04-27 17:27:11" ...
## $ ended_at     : POSIXct[1:337230], format: "2021-04-12 18:56:55" "2021-04-27 18:31:29" ...
## $ start_station_name: chr [1:337230] "State St & Pearson St" "Dorchester Ave & 49th St" "Loomis Blv
## $ start_station_id : chr [1:337230] "TA1307000061" "KA1503000069" "20121" "TA1305000034" ...
## $ end_station_name : chr [1:337230] "Southport Ave & Waveland Ave" "Dorchester Ave & 49th St" "Loom
## $ end_station_id   : chr [1:337230] "13235" "KA1503000069" "20121" "13235" ...
## $ start_lat       : num [1:337230] 41.9 41.8 41.7 41.9 41.7 ...
## $ start_lng       : num [1:337230] -87.6 -87.6 -87.7 -87.7 -87.7 ...
## $ end_lat        : num [1:337230] 41.9 41.8 41.7 41.9 41.7 ...
## $ end_lng        : num [1:337230] -87.7 -87.6 -87.7 -87.7 -87.7 ...
## $ member_casual   : chr [1:337230] "member" "casual" "casual" "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>

str(trip_05)

## spec_tbl_df [531,633 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:531633] "C809ED75D6160B2A" "DD59FDCE0ACACAF3" "OAB83CB88C43EFC2" "7881
## $ rideable_type : chr [1:531633] "electric_bike" "electric_bike" "electric_bike" "electric_bike"

```

```

## $ started_at      : POSIXct[1:531633], format: "2021-05-30 11:58:15" "2021-05-30 11:29:14" ...
## $ ended_at        : POSIXct[1:531633], format: "2021-05-30 12:10:39" "2021-05-30 12:14:09" ...
## $ start_station_name: chr [1:531633] NA NA NA NA ...
## $ start_station_id  : chr [1:531633] NA NA NA NA ...
## $ end_station_name  : chr [1:531633] NA NA NA NA ...
## $ end_station_id    : chr [1:531633] NA NA NA NA ...
## $ start_lat         : num [1:531633] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng         : num [1:531633] -87.6 -87.6 -87.7 -87.7 -87.7 ...
## $ end_lat           : num [1:531633] 41.9 41.8 41.9 41.9 41.9 ...
## $ end_lng           : num [1:531633] -87.6 -87.6 -87.7 -87.7 -87.7 ...
## $ member_casual     : chr [1:531633] "casual" "casual" "casual" "casual" ...
## - 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>

str(trip_06)

## spec_tbl_df [729,595 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id          : chr [1:729595] "99FEC93BA843FB20" "06048DCFC8520CAF" "9598066F68045DF2" "B03C
## $ rideable_type     : chr [1:729595] "electric_bike" "electric_bike" "electric_bike" "electric_bike
## $ started_at        : POSIXct[1:729595], format: "2021-06-13 14:31:28" "2021-06-04 11:18:02" ...
## $ ended_at          : POSIXct[1:729595], format: "2021-06-13 14:34:11" "2021-06-04 11:24:19" ...
## $ start_station_name: chr [1:729595] NA NA NA NA ...
## $ start_station_id  : chr [1:729595] NA NA NA NA ...
## $ end_station_name  : chr [1:729595] NA NA NA NA ...
## $ end_station_id    : chr [1:729595] NA NA NA NA ...
## $ start_lat         : num [1:729595] 41.8 41.8 41.8 41.8 41.8 ...
## $ start_lng         : num [1:729595] -87.6 -87.6 -87.6 -87.6 -87.6 ...
## $ end_lat           : num [1:729595] 41.8 41.8 41.8 41.8 41.8 ...
## $ end_lng           : num [1:729595] -87.6 -87.6 -87.6 -87.6 -87.6 ...
## $ member_casual     : chr [1:729595] "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>
```

```
str(trip_07)
```

```
## spec_tbl_df [822,410 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:822410] "OA1B623926EF4E16" "B2D5583A5A5E76EE" "6F264597DDBF427A" "379B
## $ rideable_type : chr [1:822410] "docked_bike" "classic_bike" "classic_bike" "classic_bike" ...
## $ started_at   : POSIXct[1:822410], format: "2021-07-02 14:44:36" "2021-07-07 16:57:42" ...
## $ ended_at     : POSIXct[1:822410], format: "2021-07-02 15:19:58" "2021-07-07 17:16:09" ...
## $ start_station_name: chr [1:822410] "Michigan Ave & Washington St" "California Ave & Cortez St" "W
## $ start_station_id : chr [1:822410] "13001" "17660" "SL-012" "17660" ...
## $ end_station_name : chr [1:822410] "Halsted St & North Branch St" "Wood St & Hubbard St" "Rush St
## $ end_station_id   : chr [1:822410] "KA1504000117" "13432" "KA1503000044" "13196" ...
## $ start_lat       : num [1:822410] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng       : num [1:822410] -87.6 -87.7 -87.6 -87.7 -87.7 ...
## $ end_lat         : num [1:822410] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng         : num [1:822410] -87.6 -87.7 -87.6 -87.7 -87.7 ...
## $ member_casual   : chr [1:822410] "casual" "casual" "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>
```

```
str(trip_08)
```

```
## spec_tbl_df [804,352 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:804352] "99103BB87CC6C1BB" "EAFCCCFB0A3FC5A1" "9EF4F46C57AD234D" "5834
## $ rideable_type : chr [1:804352] "electric_bike" "electric_bike" "electric_bike" "electric_bike
## $ started_at   : POSIXct[1:804352], format: "2021-08-10 17:15:49" "2021-08-10 17:23:14" ...
## $ ended_at     : POSIXct[1:804352], format: "2021-08-10 17:22:44" "2021-08-10 17:39:24" ...
## $ start_station_name: chr [1:804352] NA NA NA NA ...
## $ start_station_id : chr [1:804352] NA NA NA NA ...
## $ end_station_name : chr [1:804352] NA NA NA NA ...
## $ end_station_id   : chr [1:804352] NA NA NA NA ...
## $ start_lat       : num [1:804352] 41.8 41.8 42 42 41.8 ...
## $ start_lng       : num [1:804352] -87.7 -87.7 -87.7 -87.7 -87.6 ...
```

```
## $ end_lat          : num [1:804352] 41.8 41.8 42 42 41.8 ...
## $ end_lng          : num [1:804352] -87.7 -87.6 -87.7 -87.7 -87.6 ...
## $ member_casual    : chr [1:804352] "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>
```

```
str(trip_09)
```

```
## spec_tbl_df [756,147 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id          : chr [1:756147] "9DC7B962304CBFD8" "F930E2C6872D6B32" "6EF72137900BB910" "78D11
## $ rideable_type     : chr [1:756147] "electric_bike" "electric_bike" "electric_bike" "electric_bike
## $ started_at       : POSIXct[1:756147], format: "2021-09-28 16:07:10" "2021-09-28 14:24:51" ...
## $ ended_at         : POSIXct[1:756147], format: "2021-09-28 16:09:54" "2021-09-28 14:40:05" ...
## $ start_station_name: chr [1:756147] NA NA NA NA ...
## $ start_station_id  : chr [1:756147] NA NA NA NA ...
## $ end_station_name  : chr [1:756147] NA NA NA NA ...
## $ end_station_id    : chr [1:756147] NA NA NA NA ...
## $ start_lat         : num [1:756147] 41.9 41.9 41.8 41.8 41.9 ...
## $ start_lng         : num [1:756147] -87.7 -87.6 -87.7 -87.7 -87.7 ...
## $ end_lat          : num [1:756147] 41.9 42 41.8 41.8 41.9 ...
## $ end_lng          : num [1:756147] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ member_casual     : chr [1:756147] "casual" "casual" "casual" "casual" ...
## - 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>
```

```
str(trip_11)
```

```
## spec_tbl_df [96,834 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:96834] "E19E6F1B8D4C42ED" "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA45..."
## $ rideable_type : chr [1:96834] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at   : POSIXct[1:96834], format: "2021-01-23 16:14:19" "2021-01-27 18:43:08" ...
## $ ended_at     : POSIXct[1:96834], format: "2021-01-23 16:24:44" "2021-01-27 18:47:12" ...
## $ start_station_name: chr [1:96834] "California Ave & Cortez St" "California Ave & Cortez St" "California Ave & Cortez St" ...
## $ start_station_id : chr [1:96834] "17660" "17660" "17660" "17660" ...
## $ end_station_name : chr [1:96834] NA NA NA NA ...
## $ end_station_id   : chr [1:96834] NA NA NA NA ...
## $ start_lat       : num [1:96834] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng       : num [1:96834] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ end_lat         : num [1:96834] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng         : num [1:96834] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ member_casual   : chr [1:96834] "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>
```

```
str(trip_12)
```

```
## spec_tbl_df [96,834 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:96834] "E19E6F1B8D4C42ED" "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA45..."
## $ rideable_type : chr [1:96834] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at   : POSIXct[1:96834], format: "2021-01-23 16:14:19" "2021-01-27 18:43:08" ...
## $ ended_at     : POSIXct[1:96834], format: "2021-01-23 16:24:44" "2021-01-27 18:47:12" ...
## $ start_station_name: chr [1:96834] "California Ave & Cortez St" "California Ave & Cortez St" "California Ave & Cortez St" ...
## $ start_station_id : chr [1:96834] "17660" "17660" "17660" "17660" ...
## $ end_station_name : chr [1:96834] NA NA NA NA ...
## $ end_station_id   : chr [1:96834] NA NA NA NA ...
## $ start_lat       : num [1:96834] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng       : num [1:96834] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ end_lat         : num [1:96834] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng         : num [1:96834] -87.7 -87.7 -87.7 -87.7 -87.7 ...
## $ member_casual   : chr [1:96834] "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>
```

Stack individual months data frames into one annual big data frame

```
Trips_2021 <- bind_rows(trip_01,trip_02,trip_03,trip_04,trip_05,trip_06,trip_07,trip_08,trip_09,trip_10)
```

Remove lat, long.

```
Trips_2021 <- Trips_2021 %>% select(-c(start_lat, start_lng, end_lat, end_lng))
```

CLEAN UP AND

ADD DATA TO PREPARE FOR ANALYSIS

```
colnames(Trips_2021) #List of column names
```

```
## [1] "ride_id"          "rideable_type"    "started_at"
## [4] "ended_at"         "start_station_name" "start_station_id"
## [7] "end_station_name" "end_station_id"   "member_casual"
```

```
nrow(Trips_2021) #How many rows are in data frame?
```

```
## [1] 4646821
```

```
dim(Trips_2021) #Dimensions of the data frame?
```

```
## [1] 4646821      9
```

```
head(Trips_2021) #See the first 6 rows of data frame. Also tail(Trips_2021)
```

```
## # A tibble: 6 x 9
##   ride_id rideable_type started_at      ended_at      start_station_n-
##   <chr>   <chr>         <dtm>         <dtm>         <chr>
## 1 E19E6F~ electric_bike 2021-01-23 16:14:19 2021-01-23 16:24:44 California Ave ~
## 2 DC88F2~ electric_bike 2021-01-27 18:43:08 2021-01-27 18:47:12 California Ave ~
## 3 EC45C9~ electric_bike 2021-01-21 22:35:54 2021-01-21 22:37:14 California Ave ~
## 4 4FA453~ electric_bike 2021-01-07 13:31:13 2021-01-07 13:42:55 California Ave ~
## 5 BE5E8E~ electric_bike 2021-01-23 02:24:02 2021-01-23 02:24:45 California Ave ~
## 6 5D8969~ electric_bike 2021-01-09 14:24:07 2021-01-09 15:17:54 California Ave ~
## # ... with 4 more variables: start_station_id <chr>, end_station_name <chr>,
## #   end_station_id <chr>, member_casual <chr>
```

```
str(Trips_2021) #See list of columns and data types (numeric, character, etc)
```

```
## tibble [4,646,821 x 9] (S3: tbl_df/tbl/data.frame)
## $ ride_id      : chr [1:4646821] "E19E6F1B8D4C42ED" "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA..."
## $ rideable_type : chr [1:4646821] "electric_bike" "electric_bike" "electric_bike" "electric_bike" ...
## $ started_at    : POSIXct[1:4646821], format: "2021-01-23 16:14:19" "2021-01-27 18:43:08" ...
## $ ended_at      : POSIXct[1:4646821], format: "2021-01-23 16:24:44" "2021-01-27 18:47:12" ...
## $ start_station_name: chr [1:4646821] "California Ave & Cortez St" "California Ave & Cortez St" "California Ave & Cortez St" ...
## $ start_station_id  : chr [1:4646821] "17660" "17660" "17660" "17660" ...
## $ end_station_name  : chr [1:4646821] NA NA NA NA ...
## $ end_station_id    : chr [1:4646821] NA NA NA NA ...
## $ member_casual    : chr [1:4646821] "member" "member" "member" "member" ...
```

```
summary(Trips_2021) #Statistical summary of data. Mainly for numerics
```

```
##      ride_id      rideable_type      started_at
## Length:4646821 Length:4646821 Min.      :2021-01-01 00:02:05
## Class :character Class :character 1st Qu.:2021-05-12 20:10:39
## Mode  :character Mode  :character Median :2021-07-03 11:52:13
##                                     Mean  :2021-06-20 21:19:03
##                                     3rd Qu.:2021-08-15 22:31:12
##                                     Max.   :2021-09-30 23:59:48
##      ended_at      start_station_name start_station_id
## Min.      :2021-01-01 00:08:39 Length:4646821 Length:4646821
## 1st Qu.:2021-05-12 20:28:33 Class :character Class :character
## Median :2021-07-03 12:20:21 Mode  :character Mode  :character
## Mean    :2021-06-20 21:41:54
## 3rd Qu.:2021-08-15 22:53:50
## Max.    :2021-10-01 22:55:35
##      end_station_name end_station_id      member_casual
## Length:4646821 Length:4646821 Length:4646821
## Class :character Class :character Class :character
## Mode  :character Mode  :character Mode  :character
##
##
##
```

There are a few problems we will need to fix:

We will want to add some additional columns of data – such as day, month, year – that provide additional opportunities to aggregate the data.

```
Trips_2021$date <- as.Date(Trips_2021$started_at) #The default format is yyyy-mm-dd
Trips_2021$month <- format(as.Date(Trips_2021$date), "%m")
Trips_2021$day <- format(as.Date(Trips_2021$date), "%d")
Trips_2021$year <- format(as.Date(Trips_2021$date), "%Y")
Trips_2021$day_of_week <- format(as.Date(Trips_2021$date), "%A")
```

Add a “ride_length” calculation to all_trips (in seconds)

<https://stat.ethz.ch/R-manual/R-devel/library/base/html/difftime.html>

```
Trips_2021$ride_length <- difftime(Trips_2021$ended_at, Trips_2021$started_at)
```

Inspect the structure of the columns

```
str(Trips_2021)

## tibble [4,646,821 x 15] (S3: tbl_df/tbl/data.frame)
##  $ ride_id          : chr [1:4646821] "E19E6F1B8D4C42ED" "DC88F20C2C55F27F" "EC45C94683FE3F27" "4FA
##  $ rideable_type     : chr [1:4646821] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
##  $ started_at        : POSIXct[1:4646821], format: "2021-01-23 16:14:19" "2021-01-27 18:43:08" ...
##  $ ended_at          : POSIXct[1:4646821], format: "2021-01-23 16:24:44" "2021-01-27 18:47:12" ...
##  $ start_station_name: chr [1:4646821] "California Ave & Cortez St" "California Ave & Cortez St" "Ca
##  $ start_station_id  : chr [1:4646821] "17660" "17660" "17660" "17660" ...
##  $ end_station_name  : chr [1:4646821] NA NA NA NA ...
##  $ end_station_id    : chr [1:4646821] NA NA NA NA ...
##  $ member_casual     : chr [1:4646821] "member" "member" "member" "member" ...
##  $ date              : Date[1:4646821], format: "2021-01-23" "2021-01-27" ...
##  $ month             : chr [1:4646821] "01" "01" "01" "01" ...
##  $ day              : chr [1:4646821] "23" "27" "21" "07" ...
##  $ year              : chr [1:4646821] "2021" "2021" "2021" "2021" ...
##  $ day_of_week       : chr [1:4646821] "Saturday" "Wednesday" "Thursday" "Thursday" ...
##  $ ride_length       : 'difftime' num [1:4646821] 625 244 80 702 ...
##  ..- attr(*, "units")= chr "secs"
```

Convert “ride_length” from Factor to numeric so we can run calculations on the data

```
is.factor(Trips_2021$ride_length)

## [1] FALSE

Trips_2021$ride_length <- as.numeric(as.character(Trips_2021$ride_length))
is.numeric(Trips_2021$ride_length)

## [1] TRUE
```

Remove “bad” data

The dataframe includes a few hundred entries when bikes were taken out of docks and checked for quality by Divvy or ride_length was negative We will create a new version of the dataframe (v2) since data is being removed <https://www.datasciencemadesimple.com/delete-or-drop-rows-in-r-with-conditions-2/>

```
trips_v2 <- na.omit(Trips_2021[!(Trips_2021$ride_length < 0),])
```

DESCRIPTIVE ANALYSIS

Ride_length (all figures in seconds)

```
mean(trips_v2$ride_length) #straight average (total ride length / rides)

## [1] 1361.688

median(trips_v2$ride_length) #midpoint number in the ascending array of ride lengths

## [1] 761
```

```
max(trips_v2$ride_length) #longest ride
```

```
## [1] 3356649
```

```
min(trips_v2$ride_length) #shortest ride
```

```
## [1] 0
```

condense the four lines above to one line using summary() on the specific attribute

```
summary(trips_v2$ride_length)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         0      433     761    1362    1375 3356649
```

Compare members and casual users

```
aggregate(trips_v2$ride_length ~ trips_v2$member_casual, FUN = mean)
```

```
##      trips_v2$member_casual trips_v2$ride_length
## 1                          casual          2015.9429
## 2                          member           815.0345
```

```
aggregate(trips_v2$ride_length ~ trips_v2$member_casual, FUN = median)
```

```
##      trips_v2$member_casual trips_v2$ride_length
## 1                          casual           1031
## 2                          member            603
```

```
aggregate(trips_v2$ride_length ~ trips_v2$member_casual, FUN = max)
```

```
##      trips_v2$member_casual trips_v2$ride_length
## 1                          casual          3356649
## 2                          member           89738
```

```
aggregate(trips_v2$ride_length ~ trips_v2$member_casual, FUN = min)
```

```
##      trips_v2$member_casual trips_v2$ride_length
## 1                          casual                0
## 2                          member                0
```

average ride time by each day for members vs casual users

```
aggregate(trips_v2$ride_length ~ trips_v2$member_casual + trips_v2$day_of_week, FUN = mean)
```

```
##      trips_v2$member_casual trips_v2$day_of_week trips_v2$ride_length
## 1                          casual      Friday          1922.1796
## 2                          member      Friday           788.7082
## 3                          casual     Monday          2012.0333
## 4                          member     Monday           787.8024
## 5                          casual     Saturday          2166.9403
## 6                          member     Saturday           909.6612
## 7                          casual      Sunday          2317.4768
## 8                          member      Sunday           934.8033
## 9                          casual    Thursday          1731.7716
## 10                         member    Thursday           763.1619
## 11                         casual     Tuesday          1801.6003
## 12                         member     Tuesday           769.2630
## 13                         casual    Wednesday          1770.7528
## 14                         member    Wednesday           772.9548
```

In order

```
trips_v2$day_of_week <- ordered(trips_v2$day_of_week, levels=c("Sunday", "Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday", "Sunday"))

aggregate(trips_v2$ride_length ~ trips_v2$member_casual + trips_v2$day_of_week, FUN = mean)

##      trips_v2$member_casual trips_v2$day_of_week trips_v2$ride_length
## 1          casual      Sunday          2317.4768
## 2          member      Sunday           934.8033
## 3          casual     Monday          2012.0333
## 4          member     Monday           787.8024
## 5          casual    Tuesday          1801.6003
## 6          member    Tuesday           769.2630
## 7          casual   Wednesday          1770.7528
## 8          member   Wednesday           772.9548
## 9          casual   Thursday          1731.7716
## 10         member   Thursday           763.1619
## 11         casual    Friday          1922.1796
## 12         member    Friday           788.7082
## 13         casual   Saturday          2166.9403
## 14         member   Saturday           909.6612
```

Analyze ridership data by type and weekday

```
trips_v2 %>%
  mutate(weekday = wday(started_at, label = TRUE)) %>% #creates weekday field using wday()
  group_by(member_casual, weekday) %>% #groups by usertype and weekday
  summarise(number_of_rides = n() #calculates the number of rides and average
            ,average_duration = mean(ride_length)) %>% # calculates the average duration
  arrange(member_casual, weekday) # sorts

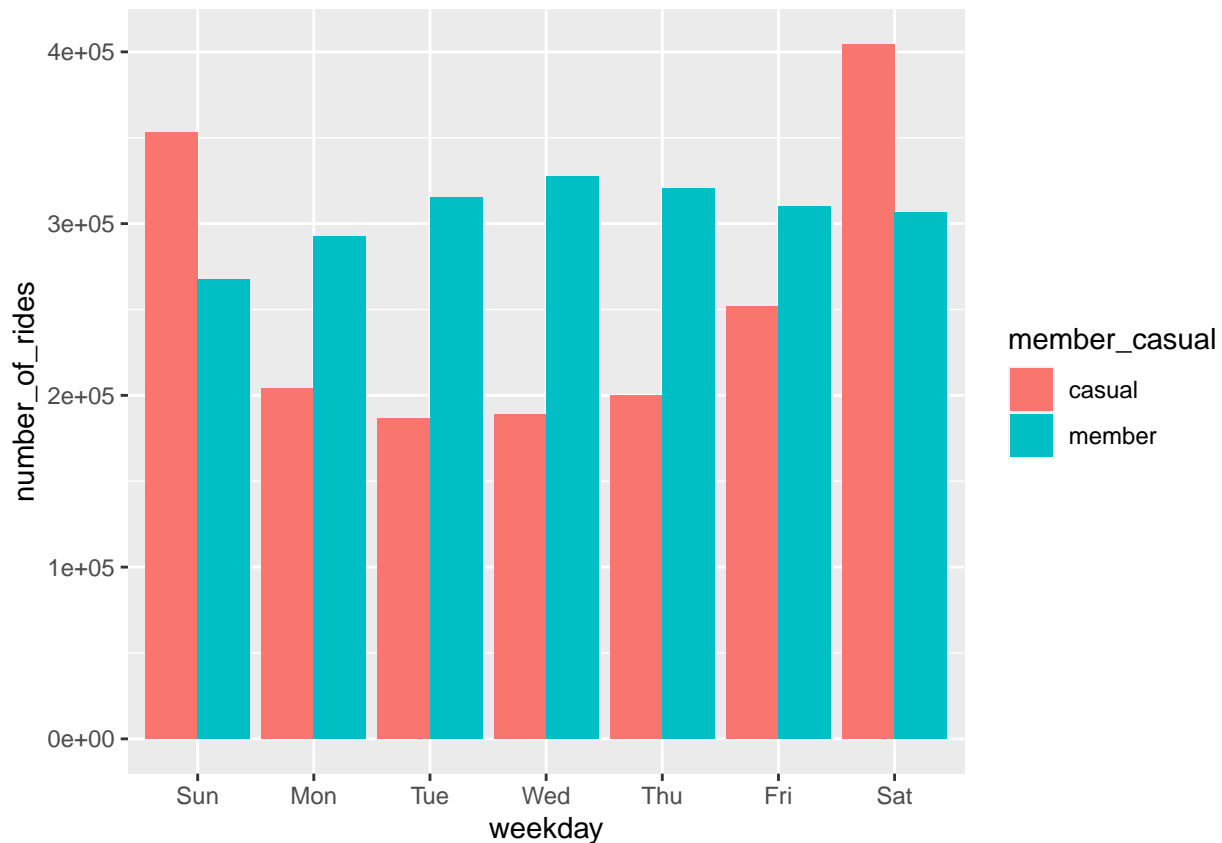
## 'summarise()' has grouped output by 'member_casual'. You can override using the '.groups' argument.

## # A tibble: 14 x 4
## # Groups:   member_casual [2]
##   member_casual weekday number_of_rides average_duration
##   <chr>         <ord>         <int>         <dbl>
## 1 casual      Sun           352824         2317.
## 2 casual      Mon           204044         2012.
## 3 casual      Tue           186326         1802.
## 4 casual      Wed           189154         1771.
## 5 casual      Thu           199758         1732.
## 6 casual      Fri           251615         1922.
## 7 casual      Sat           404533         2167.
## 8 member      Sun           267292           935.
## 9 member      Mon           292731           788.
## 10 member     Tue           315520           769.
## 11 member     Wed           327604           773.
## 12 member     Thu           320792           763.
## 13 member     Fri           309982           789.
## 14 member     Sat           306328           910.
```


Let's visualize the number of rides by rider type

```
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")
```

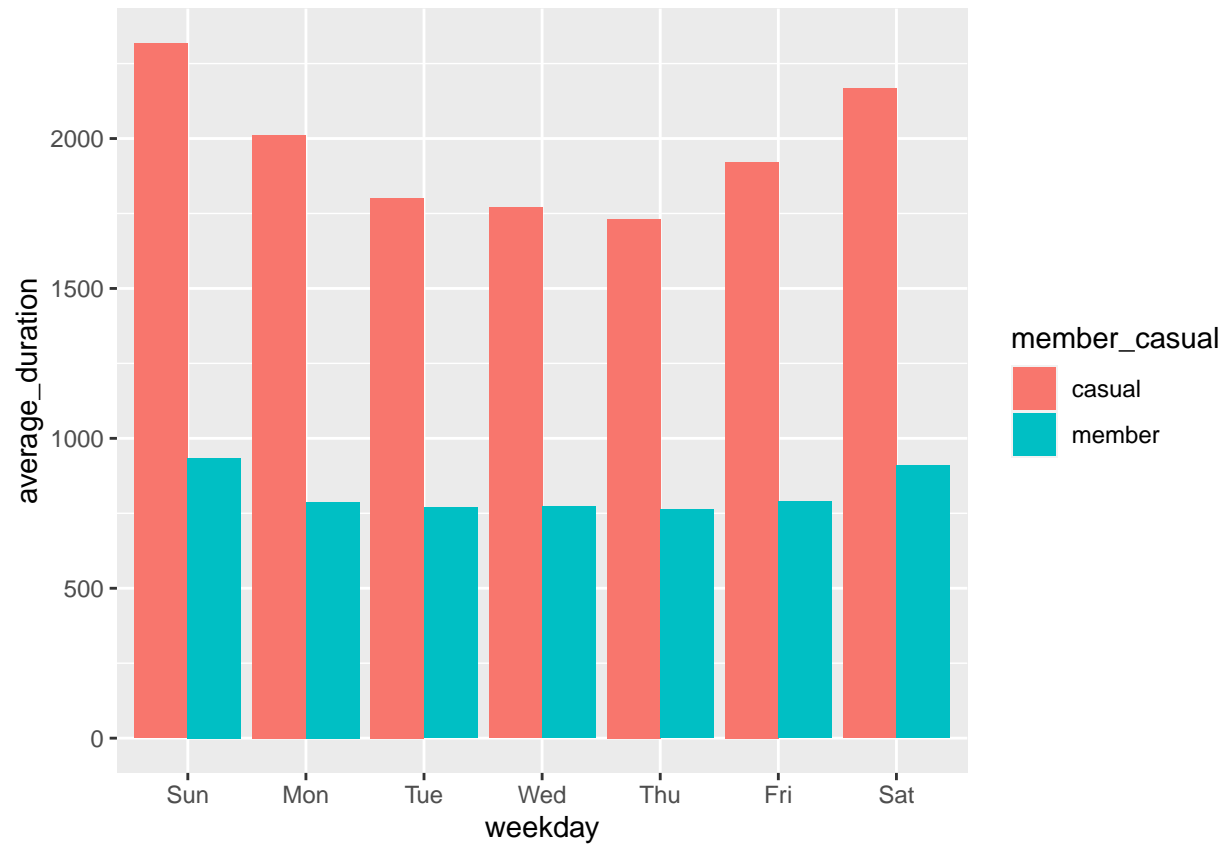
'summarise()' has grouped output by 'member_casual'. You can override using the '.groups' argument.



Let's create a visualization for average duration

```
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")
```

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



EXPORT SUMMARY FILE

<https://datatofish.com/export-dataframe-to-csv-in-r/>

```
write.csv(trips_v2,"C:\Users\maap_\Documents\trip data\TRIP_2021\Year_trips_2021.csv", row.names
= TRUE )
```