# Case study data cleaning

*Nancy Chauhan*

*October 31, 2021*

## R Markdown

Library

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.6.3
```

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

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.1      v dplyr   1.0.6
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.1
```

```
## Warning: package 'tibble' was built under R version 3.6.3
```

```
## Warning: package 'tidyr' was built under R version 3.6.3
```

```
## Warning: package 'readr' was built under R version 3.6.3
```

```
## Warning: package 'purrr' was built under R version 3.6.3
```

```
## Warning: package 'dplyr' was built under R version 3.6.3
```

```
## Warning: package 'forcats' was built under R version 3.6.3
```

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

```
library(lubridate)
```

```
## Warning: package 'lubridate' was built under R version 3.6.3
```

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

```
## The following objects are masked from 'package:base':
##
##     date, intersect, setdiff, union
```

```
library(janitor)
```

```
## Warning: package 'janitor' was built under R version 3.6.3
```

```
##
## Attaching package: 'janitor'
```

```
## The following objects are masked from 'package:stats':
##
##     chisq.test, fisher.test
```

```
library(dplyr)
library(skimr)
```

```
## Warning: package 'skimr' was built under R version 3.6.3
```

```
library(tinytex)
```

## Collect data

```
ride_oct20 <- read_csv("cycle/oct20.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()
## )
```

```
ride_nov20 <- read_csv("cycle/nov20.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()
## )
```

```
ride_dec20 <- read_csv("cycle/dec20.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()
## )
```

```
ride_jan21 <- read_csv("cycle/jan21.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()
## )
```

```r
ride_feb21 <- read_csv("cycle/feb21.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()
## )
```

```r
ride_mar21 <- read_csv("cycle/mar21.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()
## )
```

```r
ride_apr21 <- read_csv("cycle/apr21.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()
## )
```

```r
ride_may21 <- read_csv("cycle/may21.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()
## )
```

```r
ride_jun21 <- read_csv("cycle/jun21.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()
## )
```

```r
ride_jul21 <- read_csv("cycle/jul21.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()
## )
```

```r
ride_aug21 <- read_csv("cycle/aug21.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()
## )
```

```r
ride_sep21 <- read_csv("cycle/sep21.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()
## )
```

## Column detail of all file

Column name

```
colnames(ride_oct20)
```

```
##  [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(ride_nov20)
```

```
##  [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(ride_dec20)
```

```
##  [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(ride_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(ride_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(ride_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"
```

```
colnames(ride_apr21)
```

```
##  [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(ride_may21)
```

```
##  [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(ride_jun21)
```

```
##  [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(ride_jul21)
```

```
##  [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(ride_aug21)
```

```
##  [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(ride_sep21)
```

```
##  [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"
```

##Compare columns column datatype

```
compare_df_cols(ride_oct20,ride_nov20,ride_dec20,ride_jan21,ride_feb21,ride_mar21,ride_apr21,ride_may21
```

```
##              column_name      ride_oct20       ride_nov20       ride_dec20
## 1               end_lat         numeric          numeric          numeric
## 2               end_lng         numeric          numeric          numeric
## 3        end_station_id         numeric          numeric        character
## 4      end_station_name       character        character        character
## 5               ended_at POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt
## 6         member_casual       character        character        character
## 7               ride_id       character        character        character
## 8         rideable_type       character        character        character
## 9             start_lat         numeric          numeric          numeric
## 10            start_lng         numeric          numeric          numeric
## 11      start_station_id         numeric          numeric        character
## 12 start_station_name       character        character        character
## 13           started_at POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt
##            ride_jan21      ride_feb21      ride_mar21      ride_apr21
## 1             numeric          numeric          numeric          numeric
## 2             numeric          numeric          numeric          numeric
## 3           character        character        character        character
## 4           character        character        character        character
## 5  POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt
## 6           character        character        character        character
## 7           character        character        character        character
## 8           character        character        character        character
## 9             numeric          numeric          numeric          numeric
## 10            numeric          numeric          numeric          numeric
## 11          character        character        character        character
## 12          character        character        character        character
## 13 POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt
##            ride_may21      ride_jun21      ride_jul21      ride_aug21
## 1             numeric          numeric          numeric          numeric
## 2             numeric          numeric          numeric          numeric
## 3           character        character        character        character
## 4           character        character        character        character
## 5  POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt
## 6           character        character        character        character
## 7           character        character        character        character
## 8           character        character        character        character
## 9             numeric          numeric          numeric          numeric
## 10            numeric          numeric          numeric          numeric
## 11          character        character        character        character
```

9

```
## 12         character        character        character        character
## 13 POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt POSIXct, POSIXt
##          ride_sep21
## 1          numeric
## 2          numeric
## 3        character
## 4        character
## 5  POSIXct, POSIXt
## 6        character
## 7        character
## 8        character
## 9          numeric
## 10         numeric
## 11       character
## 12       character
## 13 POSIXct, POSIXt
```

## Examine structure of the data

```
str(ride_oct20)
```

```
## spec_tbl_df [388,653 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:388653] "ACB6B40CF5B9044C" "DF450C72FD109C01" "B6396B54A15AC0DF" "44A4A
## $ rideable_type     : chr [1:388653] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at        : POSIXct[1:388653], format: "2020-10-31 19:39:43" "2020-10-31 23:50:08" ...
## $ ended_at          : POSIXct[1:388653], format: "2020-10-31 19:57:12" "2020-11-01 00:04:16" ...
## $ start_station_name: chr [1:388653] "Lakeview Ave & Fullerton Pkwy" "Southport Ave & Waveland Ave"
## $ start_station_id  : num [1:388653] 313 227 102 165 190 359 313 125 NA 174 ...
## $ end_station_name  : chr [1:388653] "Rush St & Hubbard St" "Kedzie Ave & Milwaukee Ave" "University
## $ end_station_id    : num [1:388653] 125 260 423 256 185 53 125 313 199 635 ...
## $ start_lat         : num [1:388653] 41.9 41.9 41.8 42 41.9 ...
## $ start_lng         : num [1:388653] -87.6 -87.7 -87.6 -87.7 -87.7 ...
## $ end_lat           : num [1:388653] 41.9 41.9 41.8 42 41.9 ...
## $ end_lng           : num [1:388653] -87.6 -87.7 -87.6 -87.7 -87.7 ...
## $ member_casual     : chr [1:388653] "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_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()
##  .. )
```

```
str(ride_nov20)
```

```
## spec_tbl_df [259,716 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:259716] "BD0A6FF6FFF9B921" "96A7A7A4BDE4F82D" "C61526D06582BDC5" "E533
## $ rideable_type     : chr [1:259716] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at        : POSIXct[1:259716], format: "2020-11-01 13:36:00" "2020-11-01 10:03:26" ...
## $ ended_at          : POSIXct[1:259716], format: "2020-11-01 13:45:40" "2020-11-01 10:14:45" ...
## $ start_station_name: chr [1:259716] "Dearborn St & Erie St" "Franklin St & Illinois St" "Lake Shore
## $ start_station_id  : num [1:259716] 110 672 76 659 2 72 76 NA 58 394 ...
## $ end_station_name  : chr [1:259716] "St. Clair St & Erie St" "Noble St & Milwaukee Ave" "Federal St
## $ end_station_id    : num [1:259716] 211 29 41 185 2 76 72 NA 288 273 ...
## $ start_lat         : num [1:259716] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng         : num [1:259716] -87.6 -87.6 -87.6 -87.7 -87.6 ...
## $ end_lat           : num [1:259716] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng           : num [1:259716] -87.6 -87.7 -87.6 -87.7 -87.6 ...
## $ member_casual     : chr [1:259716] "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_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()
##   .. )
```

```
str(ride_dec20)
```

```
## spec_tbl_df [131,573 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:131573] "70B6A9A437D4C30D" "158A465D4E74C54A" "5262016E0F1F2F9A" "BE119
## $ rideable_type     : chr [1:131573] "classic_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at        : POSIXct[1:131573], format: "2020-12-27 12:44:29" "2020-12-18 17:37:15" ...
## $ ended_at          : POSIXct[1:131573], format: "2020-12-27 12:55:06" "2020-12-18 17:44:19" ...
## $ start_station_name: chr [1:131573] "Aberdeen St & Jackson Blvd" NA NA NA ...
## $ start_station_id  : chr [1:131573] "13157" NA NA NA ...
## $ end_station_name  : chr [1:131573] "Desplaines St & Kinzie St" NA NA NA ...
## $ end_station_id    : chr [1:131573] "TA1306000003" NA NA NA ...
## $ start_lat         : num [1:131573] 41.9 41.9 41.9 41.9 41.8 ...
## $ start_lng         : num [1:131573] -87.7 -87.7 -87.7 -87.7 -87.6 ...
## $ end_lat           : num [1:131573] 41.9 41.9 41.9 41.9 41.8 ...
## $ end_lng           : num [1:131573] -87.6 -87.7 -87.7 -87.7 -87.6 ...
## $ member_casual     : chr [1:131573] "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()
##    .. )
```

str(ride_jan21)

```
## 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" "Cali:
## $ 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()
##   .. )
```

str(ride_feb21)

```
## 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" ...
```

12

```
## $ 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 S
## $ 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()
##   .. )
```

```
str(ride_mar21)
```

```
## spec_tbl_df [228,496 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:228496] "CFA86D4455AA1030" "30D9DC61227D1AF3" "846D87A15682A284" "994DC
## $ 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()
##     .. )
```

str(ride_apr21)

```
## 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" "Loo
## $ 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()
##   .. )
```

str(ride_may21)

```
## spec_tbl_df [531,633 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:531633] "C809ED75D6160B2A" "DD59FDCE0ACACAF3" "0AB83CB88C43EFC2" "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()
##   .. )
```

```
str(ride_jun21)
```

```
## 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()
##   .. )
```

```
str(ride_jul21)
```

```
## spec_tbl_df [822,410 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id            : chr [1:822410] "0A1B623926EF4E16" "B2D5583A5A5E76EE" "6F264597DDBF427A" "379B5
## $ 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" "Wa
## $ 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()
##   .. )
```

```
str(ride_aug21)
```

```
## 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()
##    .. )
```

```
str(ride_sep21)
```

```
## spec_tbl_df [756,147 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:756147] "9DC7B962304CBFD8" "F930E2C6872D6B32" "6EF72137900BB910" "78D1I
## $ 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()
##   .. )
```

##Change numeric to character station id of oct20 and nov20 are in numeric instead of character datatype

```
ride_oct20 <- mutate(ride_oct20,end_station_id= as.character(end_station_id),start_station_id= as.chara
ride_nov20 <- mutate(ride_nov20,end_station_id= as.character(end_station_id),start_station_id= as.chara
```

##Recheck the structure

```
str(ride_oct20)
```

```
## spec_tbl_df [388,653 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id            : chr [1:388653] "ACB6B40CF5B9044C" "DF450C72FD109C01" "B6396B54A15AC0DF" "44A4/
## $ rideable_type      : chr [1:388653] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at         : POSIXct[1:388653], format: "2020-10-31 19:39:43" "2020-10-31 23:50:08" ...
## $ ended_at           : POSIXct[1:388653], format: "2020-10-31 19:57:12" "2020-11-01 00:04:16" ...
## $ start_station_name : chr [1:388653] "Lakeview Ave & Fullerton Pkwy" "Southport Ave & Waveland Ave"
## $ start_station_id   : chr [1:388653] "313" "227" "102" "165" ...
## $ end_station_name   : chr [1:388653] "Rush St & Hubbard St" "Kedzie Ave & Milwaukee Ave" "University
## $ end_station_id     : chr [1:388653] "125" "260" "423" "256" ...
## $ start_lat          : num [1:388653] 41.9 41.9 41.8 42 41.9 ...
## $ start_lng          : num [1:388653] -87.6 -87.7 -87.6 -87.7 -87.7 ...
## $ end_lat            : num [1:388653] 41.9 41.9 41.8 42 41.9 ...
## $ end_lng            : num [1:388653] -87.6 -87.7 -87.6 -87.7 -87.7 ...
## $ member_casual      : chr [1:388653] "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_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()
##   .. )
```

```
str(ride_nov20)
```

```
## spec_tbl_df [259,716 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id            : chr [1:259716] "BD0A6FF6FFF9B921" "96A7A7A4BDE4F82D" "C61526D06582BDC5" "E533I
## $ rideable_type      : chr [1:259716] "electric_bike" "electric_bike" "electric_bike" "electric_bike"
## $ started_at         : POSIXct[1:259716], format: "2020-11-01 13:36:00" "2020-11-01 10:03:26" ...
## $ ended_at           : POSIXct[1:259716], format: "2020-11-01 13:45:40" "2020-11-01 10:14:45" ...
## $ start_station_name : chr [1:259716] "Dearborn St & Erie St" "Franklin St & Illinois St" "Lake Shore
## $ start_station_id   : chr [1:259716] "110" "672" "76" "659" ...
## $ end_station_name   : chr [1:259716] "St. Clair St & Erie St" "Noble St & Milwaukee Ave" "Federal S
## $ end_station_id     : chr [1:259716] "211" "29" "41" "185" ...
## $ start_lat          : num [1:259716] 41.9 41.9 41.9 41.9 41.9 ...
## $ start_lng          : num [1:259716] -87.6 -87.6 -87.6 -87.7 -87.6 ...
## $ end_lat            : num [1:259716] 41.9 41.9 41.9 41.9 41.9 ...
## $ end_lng            : num [1:259716] -87.6 -87.7 -87.6 -87.7 -87.6 ...
## $ member_casual      : chr [1:259716] "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_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()
##    ..  )
```

## Merge all the tables

```
total_ride <- bind_rows(ride_oct20,ride_nov20,ride_dec20,ride_jan21,ride_feb21,ride_mar21,ride_apr21,ri
```

## Delete column remove dirty data

```
total_ride= total_ride%>%
  select(-c(start_lat,start_lng,end_lat,end_lng))
```

## Rename columns

```
total_ride <- total_ride %>%
  rename(ride_type = rideable_type,
         start_time = started_at,
         end_time = ended_at,
         user_type = member_casual)
```

## Change to character

```
total_ride <- mutate(total_ride,start_time=as.character(start_time),end_time= as.character(end_time))
```

## New column for day,month, year and day of the week

```
total_ride$date <- as.Date(total_ride$start_time)
total_ride$day <- format(as.Date(total_ride$date),format="%d")
total_ride$month <- format(as.Date(total_ride$date),format="%m")
total_ride$year <- format(as.Date(total_ride$date),format="%Y")
total_ride$day_of_week <- format(as.Date(total_ride$date),"%A")
```

## Total time taken by each trip

```
total_ride$ride_length <-difftime(total_ride$end_time,total_ride$start_time)
```

## skim

```
skim(total_ride)
```

Table 1: Data summary

| Name | total_ride |
|------|-----------|
| Number of rows | 5136261 |
| Number of columns | 15 |
| | |
| Column type frequency: | |
| character | 13 |
| Date | 1 |
| difftime | 1 |
| | |
| Group variables | None |

**Variable type: character**

| skim_variable | n_missing | complete_rate | min | max | empty | n_unique | whitespace |
|---------------|-----------|---------------|-----|-----|-------|----------|------------|
| ride_id | 0 | 1.00 | 16 | 16 | 0 | 5136052 | 0 |
| ride_type | 0 | 1.00 | 11 | 13 | 0 | 3 | 0 |
| start_time | 0 | 1.00 | 19 | 19 | 0 | 4301706 | 0 |
| end_time | 0 | 1.00 | 19 | 19 | 0 | 4291553 | 0 |
| start_station_name | 523467 | 0.90 | 3 | 53 | 0 | 784 | 0 |
| start_station_id | 523781 | 0.90 | 1 | 36 | 0 | 1299 | 0 |
| end_station_name | 567268 | 0.89 | 10 | 53 | 0 | 781 | 0 |
| end_station_id | 567501 | 0.89 | 1 | 36 | 0 | 1299 | 0 |
| user_type | 0 | 1.00 | 6 | 6 | 0 | 2 | 0 |
| day | 0 | 1.00 | 2 | 2 | 0 | 31 | 0 |
| month | 0 | 1.00 | 2 | 2 | 0 | 12 | 0 |
| year | 0 | 1.00 | 4 | 4 | 0 | 2 | 0 |
| day_of_week | 0 | 1.00 | 6 | 9 | 0 | 7 | 0 |

**Variable type: Date**

| skim_variable | n_missing | complete_rate | min | max | median | n_unique |
|---------------|-----------|---------------|-----|-----|--------|----------|
| date | 0 | 1 | 2020-10-01 | 2021-09-30 | 2021-06-21 | 365 |

**Variable type: difftime**

| skim_variable | n_missing | complete_rate | min | max | median | n_unique |
|---------------|-----------|---------------|-----|-----|--------|----------|
| ride_length | 0 | 1 | -1742998 secs | 3356649 secs | 756 secs | 25804 |

##Day Intial half of the month are quite popular for ride.

```
table(format(total_ride$day))
```

```
## 
##     01     02     03     04     05     06     07     08     09     10
## 157436 166459 177708 178677 188693 183630 181690 171598 175258 166749
##     11     12     13     14     15     16     17     18     19     20
```

```
## 163912 164812 181438 174777 153042 165595 176828 170011 169663 163007
##     21     22     23     24     25     26     27     28     29     30
## 172565 183893 159259 155732 151698 150280 170642 144140 151709 161765
##     31
## 103595
```

## Number of weekday

It shows weekend has more rider than the weekday (approx:1.4times). Saturday=1.4Monday

```
table(format(total_ride$day_of_week))
```

```
##
## Friday    Monday   Saturday  Sunday    Thursday  Tuesday   Wednesday
##    745341    642310    924170    790346    695259    658756    680079
```

##Bike type used by user Classic bike are 4 times the docked bike and 1.6 times the electric bike

```
table(format(total_ride$ride_type))
```

```
##
## classic_bike  docked_bike   electric_bike
##      2750831       677980        1707450
```

## Number of each user type

member are 1.2 times the casual(non-memeber)

```
table(format(total_ride$user_type))
```

```
##
##   casual   member
## 2358287 2777974
```

##Month Month of July has maxinum number of rider (96834),and February has minimun number(96864)

```
table(format(total_ride$month))
```

```
##
##     01     02     03     04     05     06     07     08     09     10
##  96834  49622 228496 337230 531633 729595 822410 804352 756147 388653
##     11     12
## 259716 131573
```

##summary

```
summary(total_ride)
```

```
##     ride_id            ride_type          start_time
## Length:5136261    Length:5136261     Length:5136261
## Class :character   Class :character   Class :character
## Mode  :character   Mode  :character   Mode  :character
##
##
##
##     end_time        start_station_name start_station_id
## Length:5136261    Length:5136261     Length:5136261
## Class :character   Class :character   Class :character
## Mode  :character   Mode  :character   Mode  :character
##
##
##
## end_station_name    end_station_id      user_type
## Length:5136261    Length:5136261     Length:5136261
## Class :character   Class :character   Class :character
## Mode  :character   Mode  :character   Mode  :character
##
##
##
##      date                day                month
## Min.   :2020-10-01   Length:5136261     Length:5136261
## 1st Qu.:2021-04-11   Class :character   Class :character
## Median :2021-06-21   Mode  :character   Mode  :character
## Mean   :2021-05-25
## 3rd Qu.:2021-08-11
## Max.   :2021-09-30
##     year           day_of_week        ride_length
## Length:5136261    Length:5136261     Length:5136261
## Class :character   Class :character   Class :difftime
## Mode  :character   Mode  :character   Mode  :numeric
##
##
##
```