

GoogleCapstoneTT_Part1_CreateRawDataframe

Tomasz Tomaszewski

1/21/2022

R Markdown

MY GOOGLE CAPSTONE CODE BEGINS HERE

In this project I will follow guidelines and recommendations from Coursera Google Data Analytics Professional Certificate.

All information is available when browsing course number 8, Google Data Analytics Capstone: Complete a Case Study. I decided to choose scenario number one, and dataset from 2019 as i believe this will be more representative and unaffected by covid pandemy etc. The workflow and description will be available in final deliverable.

As recommended in capstone guidelines, the following script was used for further work: Link to capstone script (2022 Feb) (<https://docs.google.com/document/d/1TTj5KNKf4BWvEORGm10oNbpwTRk1hamsWJGj6qRWpul/edit>)

Getting working directory.

```
getwd()
```

```
## [1] "E:/Tomasz/CapstoneGoogle"
```

CSV files with datasets were loaded into directory above with help of RStudio.

Next step will be loading required packages.

```
library(tidyverse)
```

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

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.6      v dplyr   1.0.7
## v tidyr   1.1.4      v stringr 1.4.0
## v readr   2.1.1      v forcats 0.5.1
```

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

```
library(ggplot2)
library(lubridate)
```

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

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

Collecting and loading data:

```
data_q1_2019 <- read_csv("Divvy_Trips_2019_Q1.csv")
```

```
## Rows: 365069 Columns: 12
```

```
## -- Column specification -----  
## Delimiter: ","  
## chr  (4): from_station_name, to_station_name, usertype, gender  
## dbl  (5): trip_id, bikeid, from_station_id, to_station_id, birthyear  
## dtm  (2): start_time, end_time
```

```
##  
## 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.
```

```
data_q2_2019 <- read_csv("Divvy_Trips_2019_Q2.csv")
```

```
## Rows: 1108163 Columns: 12
```

```
## -- Column specification -----  
## Delimiter: ","  
## chr  (4): 03 - Rental Start Station Name, 02 - Rental End Station Name, User...  
## dbl  (5): 01 - Rental Details Rental ID, 01 - Rental Details Bike ID, 03 - R...  
## dtm  (2): 01 - Rental Details Local Start Time, 01 - Rental Details Local En...
```

```
##  
## 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.
```

```
data_q3_2019 <- read_csv("Divvy_Trips_2019_Q3.csv")
```

```
## Rows: 1640718 Columns: 12
```

```
## -- Column specification -----  
## Delimiter: ","  
## chr  (4): from_station_name, to_station_name, usertype, gender  
## dbl  (5): trip_id, bikeid, from_station_id, to_station_id, birthyear  
## dtm  (2): start_time, end_time
```

```
##  
## 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.
```

```
data_q4_2019 <- read_csv("Divvy_Trips_2019_Q4.csv")
```

```
## Rows: 704054 Columns: 12
```

```
## -- Column specification -----  
## Delimiter: ","  
## chr  (4): from_station_name, to_station_name, usertype, gender  
## dbl  (5): trip_id, bikeid, from_station_id, to_station_id, birthyear  
## dtm  (2): start_time, end_time
```

```
##  
## 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.
```

Comparing column names to identify if they match (which is required for merging all four datasets into one).

```
colnames(data_q1_2019) == colnames(data_q1_2019)
```

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

```
colnames(data_q1_2019) == colnames(data_q2_2019)
```

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

```
colnames(data_q1_2019) == colnames(data_q3_2019)
```

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

```
colnames(data_q1_2019) == colnames(data_q4_2019)
```

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

Warning! Column names for quarters 1, 3, 4 are the same however **for Q2 names are different.**

Retrieving column names:

```
colnames(data_q1_2019)
```

```
## [1] "trip_id"      "start_time"   "end_time"
## [4] "bikeid"      "tripduration" "from_station_id"
## [7] "from_station_name" "to_station_id" "to_station_name"
## [10] "usertype"    "gender"       "birthyear"
```

```
colnames(data_q2_2019)
```

```
## [1] "01 - Rental Details Rental ID"
## [2] "01 - Rental Details Local Start Time"
## [3] "01 - Rental Details Local End Time"
## [4] "01 - Rental Details Bike ID"
## [5] "01 - Rental Details Duration In Seconds Uncapped"
## [6] "03 - Rental Start Station ID"
## [7] "03 - Rental Start Station Name"
## [8] "02 - Rental End Station ID"
## [9] "02 - Rental End Station Name"
## [10] "User Type"
## [11] "Member Gender"
## [12] "05 - Member Details Member Birthday Year"
```

Note different column names for quarter 2. Requires name change.

```
colnames(data_q3_2019)
```

```
## [1] "trip_id"      "start_time"   "end_time"
## [4] "bikeid"      "tripduration" "from_station_id"
## [7] "from_station_name" "to_station_id" "to_station_name"
## [10] "usertype"    "gender"       "birthyear"
```

```
colnames(data_q4_2019)
```

```
## [1] "trip_id"      "start_time"    "end_time"
## [4] "bikeid"       "tripduration"  "from_station_id"
## [7] "from_station_name" "to_station_id" "to_station_name"
## [10] "usertype"      "gender"        "birthyear"
```

It is confirmed that Q1 Q3 Q4 column names match, but for Q2 column names need to be changed before datasets will be merged. For consistency and better readability and convention column names will be changed according to course R script guidelines.

I will also store datasets in new variable (appended name with ****_v2****).

```
(q1_v2 <- rename(data_q1_2019
  ,ride_id = trip_id
  ,rideable_type = bikeid
  ,started_at = start_time
  ,ended_at = end_time
  ,start_station_name = from_station_name
  ,start_station_id = from_station_id
  ,end_station_name = to_station_name
  ,end_station_id = to_station_id
  ,member_casual = usertype))
```

```
## # A tibble: 365,069 x 12
##   ride_id started_at ended_at rideable_type tripduration
##   <dbl> <dtm>      <dtm>      <dbl>      <dbl>
## 1 21742443 2019-01-01 00:04:37 2019-01-01 00:11:07      2167      390
## 2 21742444 2019-01-01 00:08:13 2019-01-01 00:15:34      4386      441
## 3 21742445 2019-01-01 00:13:23 2019-01-01 00:27:12      1524      829
## 4 21742446 2019-01-01 00:13:45 2019-01-01 00:43:28       252     1783
## 5 21742447 2019-01-01 00:14:52 2019-01-01 00:20:56     1170      364
## 6 21742448 2019-01-01 00:15:33 2019-01-01 00:19:09     2437      216
## 7 21742449 2019-01-01 00:16:06 2019-01-01 00:19:03     2708      177
## 8 21742450 2019-01-01 00:18:41 2019-01-01 00:20:21     2796      100
## 9 21742451 2019-01-01 00:18:43 2019-01-01 00:47:30     6205     1727
## 10 21742452 2019-01-01 00:19:18 2019-01-01 00:24:54     3939      336
## # ... with 365,059 more rows, and 7 more variables: start_station_id <dbl>,
## #   start_station_name <chr>, end_station_id <dbl>, end_station_name <chr>,
## #   member_casual <chr>, gender <chr>, birthyear <dbl>
```

```
(q2_v2 <- rename(data_q2_2019
  ,ride_id = "01 - Rental Details Rental ID"
  ,rideable_type = "01 - Rental Details Bike ID"
  ,started_at = "01 - Rental Details Local Start Time"
  ,ended_at = "01 - Rental Details Local End Time"
  ,start_station_name = "03 - Rental Start Station Name"
  ,start_station_id = "03 - Rental Start Station ID"
  ,end_station_name = "02 - Rental End Station Name"
  ,end_station_id = "02 - Rental End Station ID"
  ,member_casual = "User Type"))
```

```
## # A tibble: 1,108,163 x 12
##   ride_id started_at      ended_at      rideable_type
##   <dbl> <dtm>          <dtm>          <dbl>
## 1 22178529 2019-04-01 00:02:22 2019-04-01 00:09:48      6251
## 2 22178530 2019-04-01 00:03:02 2019-04-01 00:20:30      6226
## 3 22178531 2019-04-01 00:11:07 2019-04-01 00:15:19      5649
## 4 22178532 2019-04-01 00:13:01 2019-04-01 00:18:58      4151
## 5 22178533 2019-04-01 00:19:26 2019-04-01 00:36:13      3270
## 6 22178534 2019-04-01 00:19:39 2019-04-01 00:23:56      3123
## 7 22178535 2019-04-01 00:26:33 2019-04-01 00:35:41      6418
## 8 22178536 2019-04-01 00:29:48 2019-04-01 00:36:11      4513
## 9 22178537 2019-04-01 00:32:07 2019-04-01 01:07:44      3280
## 10 22178538 2019-04-01 00:32:19 2019-04-01 01:07:39      5534
## # ... with 1,108,153 more rows, and 8 more variables:
## #   01 - Rental Details Duration In Seconds Uncapped <dbl>,
## #   start_station_id <dbl>, start_station_name <chr>, end_station_id <dbl>,
## #   end_station_name <chr>, member_casual <chr>, Member Gender <chr>,
## #   05 - Member Details Member Birthday Year <dbl>
```

```
(q3_v2 <- rename(data_q3_2019
  ,ride_id = trip_id
  ,rideable_type = bikeid
  ,started_at = start_time
  ,ended_at = end_time
  ,start_station_name = from_station_name
  ,start_station_id = from_station_id
  ,end_station_name = to_station_name
  ,end_station_id = to_station_id
  ,member_casual = usertype))
```

```
## # A tibble: 1,640,718 x 12
##   ride_id started_at      ended_at      rideable_type tripduration
##   <dbl> <dtm>          <dtm>          <dbl>          <dbl>
## 1 23479388 2019-07-01 00:00:27 2019-07-01 00:20:41      3591      1214
## 2 23479389 2019-07-01 00:01:16 2019-07-01 00:18:44      5353      1048
## 3 23479390 2019-07-01 00:01:48 2019-07-01 00:27:42      6180      1554
## 4 23479391 2019-07-01 00:02:07 2019-07-01 00:27:10      5540      1503
## 5 23479392 2019-07-01 00:02:13 2019-07-01 00:22:26      6014      1213
## 6 23479393 2019-07-01 00:02:21 2019-07-01 00:07:31      4941       310
## 7 23479394 2019-07-01 00:02:24 2019-07-01 00:23:12      3770      1248
## 8 23479395 2019-07-01 00:02:26 2019-07-01 00:28:16      5442      1550
## 9 23479396 2019-07-01 00:02:34 2019-07-01 00:28:57      2957      1583
## 10 23479397 2019-07-01 00:02:45 2019-07-01 00:29:14      6091      1589
## # ... with 1,640,708 more rows, and 7 more variables: start_station_id <dbl>,
## #   start_station_name <chr>, end_station_id <dbl>, end_station_name <chr>,
## #   member_casual <chr>, gender <chr>, birthyear <dbl>
```

```
(q4_v2 <- rename(data_q4_2019
  ,ride_id = trip_id
  ,rideable_type = bikeid
  ,started_at = start_time
  ,ended_at = end_time
  ,start_station_name = from_station_name
  ,start_station_id = from_station_id
  ,end_station_name = to_station_name
  ,end_station_id = to_station_id
  ,member_casual = usertype))
```

```
## # A tibble: 704,054 x 12
##   ride_id started_at ended_at rideable_type tripduration
##   <dbl> <dtm>      <dtm>      <dbl>      <dbl>
## 1 25223640 2019-10-01 00:01:39 2019-10-01 00:17:20      2215      940
## 2 25223641 2019-10-01 00:02:16 2019-10-01 00:06:34      6328      258
## 3 25223642 2019-10-01 00:04:32 2019-10-01 00:18:43      3003      850
## 4 25223643 2019-10-01 00:04:32 2019-10-01 00:43:43      3275     2350
## 5 25223644 2019-10-01 00:04:34 2019-10-01 00:35:42      5294     1867
## 6 25223645 2019-10-01 00:04:38 2019-10-01 00:10:51      1891      373
## 7 25223646 2019-10-01 00:04:52 2019-10-01 00:22:45      1061     1072
## 8 25223647 2019-10-01 00:04:57 2019-10-01 00:29:16      1274     1458
## 9 25223648 2019-10-01 00:05:20 2019-10-01 00:29:18      6011     1437
## 10 25223649 2019-10-01 00:05:20 2019-10-01 02:23:46      2957     8306
## # ... with 704,044 more rows, and 7 more variables: start_station_id <dbl>,
## #   start_station_name <chr>, end_station_id <dbl>, end_station_name <chr>,
## #   member_casual <chr>, gender <chr>, birthyear <dbl>
```

Inspecting data frames for type of data in each column:

```
str(q1_v2)
```

```
## spec_tbl_df [365,069 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id      : num [1:365069] 21742443 21742444 21742445 21742446 21742447 ...
## $ started_at   : POSIXct[1:365069], format: "2019-01-01 00:04:37" "2019-01-01 00:08:13" ...
## $ ended_at     : POSIXct[1:365069], format: "2019-01-01 00:11:07" "2019-01-01 00:15:34" ...
## $ rideable_type : num [1:365069] 2167 4386 1524 252 1170 ...
## $ tripduration : num [1:365069] 390 441 829 1783 364 ...
## $ start_station_id : num [1:365069] 199 44 15 123 173 98 98 211 150 268 ...
## $ start_station_name: chr [1:365069] "Wabash Ave & Grand Ave" "State St & Randolph St" "Racine Ave & 18th S
t" "California Ave & Milwaukee Ave" ...
## $ end_station_id   : num [1:365069] 84 624 644 176 35 49 49 142 148 141 ...
## $ end_station_name : chr [1:365069] "Milwaukee Ave & Grand Ave" "Dearborn St & Van Buren St (*)" "Western
Ave & Fillmore St (*)" "Clark St & Elm St" ...
## $ member_casual    : chr [1:365069] "Subscriber" "Subscriber" "Subscriber" "Subscriber" ...
## $ gender           : chr [1:365069] "Male" "Female" "Female" "Male" ...
## $ birthyear        : num [1:365069] 1989 1990 1994 1993 1994 ...
## - attr(*, "spec")=
## .. cols(
## ..   trip_id = col_double(),
## ..   start_time = col_datetime(format = ""),
## ..   end_time = col_datetime(format = ""),
## ..   bikeid = col_double(),
## ..   tripduration = col_number(),
## ..   from_station_id = col_double(),
## ..   from_station_name = col_character(),
## ..   to_station_id = col_double(),
## ..   to_station_name = col_character(),
## ..   usertype = col_character(),
## ..   gender = col_character(),
## ..   birthyear = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
str(q2_v2)
```

```
## spec_tbl_df [1,108,163 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id : num [1:1108163] 22178529 22178530 22178531 22178532 221
78533 ...
## $ started_at : POSIXct[1:1108163], format: "2019-04-01 00:02:22" "2019
-04-01 00:03:02" ...
## $ ended_at : POSIXct[1:1108163], format: "2019-04-01 00:09:48" "2019
-04-01 00:20:30" ...
## $ rideable_type : num [1:1108163] 6251 6226 5649 4151 3270 ...
## $ 01 - Rental Details Duration In Seconds Uncapped: num [1:1108163] 446 1048 252 357 1007 ...
## $ start_station_id : num [1:1108163] 81 317 283 26 202 420 503 260 211 211
...
## $ start_station_name : chr [1:1108163] "Daley Center Plaza" "Wood St & Taylor
St" "LaSalle St & Jackson Blvd" "McClurg Ct & Illinois St" ...
## $ end_station_id : num [1:1108163] 56 59 174 133 129 426 500 499 211 211
...
## $ end_station_name : chr [1:1108163] "Desplaines St & Kinzie St" "Wabash Ave
& Roosevelt Rd" "Canal St & Madison St" "Kingsbury St & Kinzie St" ...
## $ member_casual : chr [1:1108163] "Subscriber" "Subscriber" "Subscriber"
"Subscriber" ...
## $ Member Gender : chr [1:1108163] "Male" "Female" "Male" "Male" ...
## $ 05 - Member Details Member Birthday Year : num [1:1108163] 1975 1984 1990 1993 1992 ...
## - attr(*, "spec")=
## .. cols(
## .. `01 - Rental Details Rental ID` = col_double(),
## .. `01 - Rental Details Local Start Time` = col_datetime(format = ""),
## .. `01 - Rental Details Local End Time` = col_datetime(format = ""),
## .. `01 - Rental Details Bike ID` = col_double(),
## .. `01 - Rental Details Duration In Seconds Uncapped` = col_number(),
## .. `03 - Rental Start Station ID` = col_double(),
## .. `03 - Rental Start Station Name` = col_character(),
## .. `02 - Rental End Station ID` = col_double(),
## .. `02 - Rental End Station Name` = col_character(),
## .. `User Type` = col_character(),
## .. `Member Gender` = col_character(),
## .. `05 - Member Details Member Birthday Year` = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
str(q3_v2)
```

```
## spec_tbl_df [1,640,718 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id          : num [1:1640718] 23479388 23479389 23479390 23479391 23479392 ...
## $ started_at       : POSIXct[1:1640718], format: "2019-07-01 00:00:27" "2019-07-01 00:01:16" ...
## $ ended_at         : POSIXct[1:1640718], format: "2019-07-01 00:20:41" "2019-07-01 00:18:44" ...
## $ rideable_type     : num [1:1640718] 3591 5353 6180 5540 6014 ...
## $ tripduration     : num [1:1640718] 1214 1048 1554 1503 1213 ...
## $ start_station_id  : num [1:1640718] 117 381 313 313 168 300 168 313 43 43 ...
## $ start_station_name: chr [1:1640718] "Wilton Ave & Belmont Ave" "Western Ave & Monroe St" "Lakeview Ave & Fullerton Pkwy" "Lakeview Ave & Fullerton Pkwy" ...
## $ end_station_id    : num [1:1640718] 497 203 144 144 62 232 62 144 195 195 ...
## $ end_station_name  : chr [1:1640718] "Kimball Ave & Belmont Ave" "Western Ave & 21st St" "Larrabee St & Webster Ave" "Larrabee St & Webster Ave" ...
## $ member_casual     : chr [1:1640718] "Subscriber" "Customer" "Customer" "Customer" ...
## $ gender            : chr [1:1640718] "Male" NA NA NA ...
## $ birthyear         : num [1:1640718] 1992 NA NA NA NA ...
## - attr(*, "spec")=
## .. cols(
## ..   trip_id = col_double(),
## ..   start_time = col_datetime(format = ""),
## ..   end_time = col_datetime(format = ""),
## ..   bikeid = col_double(),
## ..   tripduration = col_number(),
## ..   from_station_id = col_double(),
## ..   from_station_name = col_character(),
## ..   to_station_id = col_double(),
## ..   to_station_name = col_character(),
## ..   usertype = col_character(),
## ..   gender = col_character(),
## ..   birthyear = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
str(q4_v2)
```



```
## spec_tbl_df [704,054 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id          : num [1:704054] 25223640 25223641 25223642 25223643 25223644 ...
## $ started_at       : POSIXct[1:704054], format: "2019-10-01 00:01:39" "2019-10-01 00:02:16" ...
## $ ended_at         : POSIXct[1:704054], format: "2019-10-01 00:17:20" "2019-10-01 00:06:34" ...
## $ rideable_type     : num [1:704054] 2215 6328 3003 3275 5294 ...
## $ tripduration     : num [1:704054] 940 258 850 2350 1867 ...
## $ start_station_id : num [1:704054] 20 19 84 313 210 156 84 156 156 336 ...
## $ start_station_name: chr [1:704054] "Sheffield Ave & Kingsbury St" "Throop (Loomis) St & Taylor St" "Milwaukee Ave & Grand Ave" "Lakeview Ave & Fullerton Pkwy" ...
## $ end_station_id   : num [1:704054] 309 241 199 290 382 226 142 463 463 336 ...
## $ end_station_name : chr [1:704054] "Leavitt St & Armitage Ave" "Morgan St & Polk St" "Wabash Ave & Grand Ave" "Kedzie Ave & Palmer Ct" ...
## $ member_casual    : chr [1:704054] "Subscriber" "Subscriber" "Subscriber" "Subscriber" ...
## $ gender           : chr [1:704054] "Male" "Male" "Female" "Male" ...
## $ birthyear        : num [1:704054] 1987 1998 1991 1990 1987 ...
## - attr(*, "spec")=
## .. cols(
## ..   trip_id = col_double(),
## ..   start_time = col_datetime(format = ""),
## ..   end_time = col_datetime(format = ""),
## ..   bikeid = col_double(),
## ..   tripduration = col_number(),
## ..   from_station_id = col_double(),
## ..   from_station_name = col_character(),
## ..   to_station_id = col_double(),
## ..   to_station_name = col_character(),
## ..   usertype = col_character(),
## ..   gender = col_character(),
## ..   birthyear = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

Converting ride_id and rideable_type from numerical value to character value

```
q1_v2 <- mutate(q1_v2, ride_id = as.character(ride_id),rideable_type = as.character(rideable_type))

str(q1_v2)
```

```
## spec_tbl_df [365,069 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id           : chr [1:365069] "21742443" "21742444" "21742445" "21742446" ...
## $ started_at       : POSIXct[1:365069], format: "2019-01-01 00:04:37" "2019-01-01 00:08:13" ...
## $ ended_at         : POSIXct[1:365069], format: "2019-01-01 00:11:07" "2019-01-01 00:15:34" ...
## $ rideable_type     : chr [1:365069] "2167" "4386" "1524" "252" ...
## $ tripduration     : num [1:365069] 390 441 829 1783 364 ...
## $ start_station_id : num [1:365069] 199 44 15 123 173 98 98 211 150 268 ...
## $ start_station_name: chr [1:365069] "Wabash Ave & Grand Ave" "State St & Randolph St" "Racine Ave & 18th S
t" "California Ave & Milwaukee Ave" ...
## $ end_station_id   : num [1:365069] 84 624 644 176 35 49 49 142 148 141 ...
## $ end_station_name : chr [1:365069] "Milwaukee Ave & Grand Ave" "Dearborn St & Van Buren St (*)" "Western
Ave & Fillmore St (*)" "Clark St & Elm St" ...
## $ member_casual    : chr [1:365069] "Subscriber" "Subscriber" "Subscriber" "Subscriber" ...
## $ gender           : chr [1:365069] "Male" "Female" "Female" "Male" ...
## $ birthyear        : num [1:365069] 1989 1990 1994 1993 1994 ...
## - attr(*, "spec")=
## .. cols(
## ..   trip_id = col_double(),
## ..   start_time = col_datetime(format = ""),
## ..   end_time = col_datetime(format = ""),
## ..   bikeid = col_double(),
## ..   tripduration = col_number(),
## ..   from_station_id = col_double(),
## ..   from_station_name = col_character(),
## ..   to_station_id = col_double(),
## ..   to_station_name = col_character(),
## ..   usertype = col_character(),
## ..   gender = col_character(),
## ..   birthyear = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
q2_v2 <- mutate(q2_v2, ride_id = as.character(ride_id),rideable_type = as.character(rideable_type))
```

```
str(q2_v2)
```

```
## spec_tbl_df [1,108,163 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id : chr [1:1108163] "22178529" "22178530" "22178531" "22178
532" ...
## $ started_at : POSIXct[1:1108163], format: "2019-04-01 00:02:22" "2019
-04-01 00:03:02" ...
## $ ended_at : POSIXct[1:1108163], format: "2019-04-01 00:09:48" "2019
-04-01 00:20:30" ...
## $ rideable_type : chr [1:1108163] "6251" "6226" "5649" "4151" ...
## $ 01 - Rental Details Duration In Seconds Uncapped: num [1:1108163] 446 1048 252 357 1007 ...
## $ start_station_id : num [1:1108163] 81 317 283 26 202 420 503 260 211 211
...
## $ start_station_name : chr [1:1108163] "Daley Center Plaza" "Wood St & Taylor
St" "LaSalle St & Jackson Blvd" "McClurg Ct & Illinois St" ...
## $ end_station_id : num [1:1108163] 56 59 174 133 129 426 500 499 211 211
...
## $ end_station_name : chr [1:1108163] "Desplaines St & Kinzie St" "Wabash Ave
& Roosevelt Rd" "Canal St & Madison St" "Kingsbury St & Kinzie St" ...
## $ member_casual : chr [1:1108163] "Subscriber" "Subscriber" "Subscriber"
"Subscriber" ...
## $ Member Gender : chr [1:1108163] "Male" "Female" "Male" "Male" ...
## $ 05 - Member Details Member Birthday Year : num [1:1108163] 1975 1984 1990 1993 1992 ...
## - attr(*, "spec")=
## .. cols(
## .. `01 - Rental Details Rental ID` = col_double(),
## .. `01 - Rental Details Local Start Time` = col_datetime(format = ""),
## .. `01 - Rental Details Local End Time` = col_datetime(format = ""),
## .. `01 - Rental Details Bike ID` = col_double(),
## .. `01 - Rental Details Duration In Seconds Uncapped` = col_number(),
## .. `03 - Rental Start Station ID` = col_double(),
## .. `03 - Rental Start Station Name` = col_character(),
## .. `02 - Rental End Station ID` = col_double(),
## .. `02 - Rental End Station Name` = col_character(),
## .. `User Type` = col_character(),
## .. `Member Gender` = col_character(),
## .. `05 - Member Details Member Birthday Year` = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
q3_v2 <- mutate(q3_v2, ride_id = as.character(ride_id),rideable_type = as.character(rideable_type))

str(q3_v2)
```

```
## spec_tbl_df [1,640,718 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id          : chr [1:1640718] "23479388" "23479389" "23479390" "23479391" ...
## $ started_at       : POSIXct[1:1640718], format: "2019-07-01 00:00:27" "2019-07-01 00:01:16" ...
## $ ended_at         : POSIXct[1:1640718], format: "2019-07-01 00:20:41" "2019-07-01 00:18:44" ...
## $ rideable_type     : chr [1:1640718] "3591" "5353" "6180" "5540" ...
## $ tripduration     : num [1:1640718] 1214 1048 1554 1503 1213 ...
## $ start_station_id : num [1:1640718] 117 381 313 313 168 300 168 313 43 43 ...
## $ start_station_name: chr [1:1640718] "Wilton Ave & Belmont Ave" "Western Ave & Monroe St" "Lakeview Ave & Fullerton Pkwy" "Lakeview Ave & Fullerton Pkwy" ...
## $ end_station_id   : num [1:1640718] 497 203 144 144 62 232 62 144 195 195 ...
## $ end_station_name : chr [1:1640718] "Kimball Ave & Belmont Ave" "Western Ave & 21st St" "Larrabee St & Webster Ave" "Larrabee St & Webster Ave" ...
## $ member_casual    : chr [1:1640718] "Subscriber" "Customer" "Customer" "Customer" ...
## $ gender           : chr [1:1640718] "Male" NA NA NA ...
## $ birthyear        : num [1:1640718] 1992 NA NA NA NA ...
## - attr(*, "spec")=
## .. cols(
## ..   trip_id = col_double(),
## ..   start_time = col_datetime(format = ""),
## ..   end_time = col_datetime(format = ""),
## ..   bikeid = col_double(),
## ..   tripduration = col_number(),
## ..   from_station_id = col_double(),
## ..   from_station_name = col_character(),
## ..   to_station_id = col_double(),
## ..   to_station_name = col_character(),
## ..   usertype = col_character(),
## ..   gender = col_character(),
## ..   birthyear = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
q4_v2 <- mutate(q4_v2, ride_id = as.character(ride_id),rideable_type = as.character(rideable_type))

str(q4_v2)
```

```
## spec_tbl_df [704,054 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ride_id          : chr [1:704054] "25223640" "25223641" "25223642" "25223643" ...
## $ started_at       : POSIXct[1:704054], format: "2019-10-01 00:01:39" "2019-10-01 00:02:16" ...
## $ ended_at         : POSIXct[1:704054], format: "2019-10-01 00:17:20" "2019-10-01 00:06:34" ...
## $ rideable_type     : chr [1:704054] "2215" "6328" "3003" "3275" ...
## $ tripduration     : num [1:704054] 940 258 850 2350 1867 ...
## $ start_station_id : num [1:704054] 20 19 84 313 210 156 84 156 156 336 ...
## $ start_station_name: chr [1:704054] "Sheffield Ave & Kingsbury St" "Throop (Loomis) St & Taylor St" "Milwaukee Ave & Grand Ave" "Lakeview Ave & Fullerton Pkwy" ...
## $ end_station_id   : num [1:704054] 309 241 199 290 382 226 142 463 463 336 ...
## $ end_station_name : chr [1:704054] "Leavitt St & Armitage Ave" "Morgan St & Polk St" "Wabash Ave & Grand Ave" "Kedzie Ave & Palmer Ct" ...
## $ member_casual    : chr [1:704054] "Subscriber" "Subscriber" "Subscriber" "Subscriber" ...
## $ gender           : chr [1:704054] "Male" "Male" "Female" "Male" ...
## $ birthyear        : num [1:704054] 1987 1998 1991 1990 1987 ...
## - attr(*, "spec")=
## .. cols(
## ..   trip_id = col_double(),
## ..   start_time = col_datetime(format = ""),
## ..   end_time = col_datetime(format = ""),
## ..   bikeid = col_double(),
## ..   tripduration = col_number(),
## ..   from_station_id = col_double(),
## ..   from_station_name = col_character(),
## ..   to_station_id = col_double(),
## ..   to_station_name = col_character(),
## ..   usertype = col_character(),
## ..   gender = col_character(),
## ..   birthyear = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

Merging quarterly dataframes into one yearly dataframe for 2019.

```
data_2019_base <- bind_rows(q1_v2, q2_v2, q3_v2, q4_v2)
```

export merged dataframe to a csv file

```
#Comment made to avoid repeating of exporting the file. Uncomment if required#
#write.csv(data_2019_base, file = 'E:/Tomasz/CapstoneGoogle/capstone_bike_2019_base.csv')
```

Due to large file and dataset I will continue work in different R markdown file.