

Case Study - Cyclistic Bike share (2018-2019)

Mohd Ali Ansari

14/06/2021

Problem Statement

How do annual members and casual riders use Cyclistic bikes differently?

Objective

This documents contains all the cleaning steps taken in order to clean and transform the datasets and preparing it for next step i.e. analysis Cleaning dataset is from the year 2018-2019

About dataset

Cyclistic's historical trip data to analyze and identify trends. Download data from [here](#). (Note: The datasets have a different name because Cyclistic is a fictional company. For the purposes of this case study, the datasets are appropriate and will enable to answer the business questions. The data has been made available by Motivate International Inc. under this license.)

This is public data that can use to explore how different customer types are using Cyclistic bikes. But note that data-privacy issues prohibit from using riders' personally identifiable information. This means that one wouldn't connect pass purchases to credit card numbers to determine if casual riders live in the Cyclistic service area or if they have purchased multiple single passes

```
library(tidyverse)
library(tidyr)
library(dplyr)
library(geosphere)
library(lubridate)
```

Importing the libraries

Loading the dataset The dataset is available in csv format after downloading so we will clean it simultaneously for merging them into one fiscal year Since the companies year starts from April month we will load all the data sets and then after checking for consistency we will merge them to make a complete one year tripdata

```
q2_2018 <- read.csv("Divvy_Trips_2018_Q2.csv")
q3_2018 <- read.csv("Divvy_Trips_2018_Q3.csv")
q4_2018 <- read.csv("Divvy_Trips_2018_Q4.csv")
q1_2019 <- read.csv("Divvy_Trips_2019_Q1.csv")
```

Checking for consistency We have to check for consistency as we have to merge all the datasets into one dataset. So the column names and columns data type should be same for all the datasets

Checking for Column name

```
## [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"

## [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"

## [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"

## [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"
```

There are total 12 columns in every data set and the name of columns are same but not consistent as for 2020 on wards so we will change the column name to 2020-Q1 column names

Renaming columns We will rename the columns into same format as of 2020dataset because it is the latest format

```
q4_2018 <- rename(q4_2018
  ,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)

q3_2018 <- rename(q3_2018
  ,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)

q2_2018 <- rename(q2_2018
  ,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)

q1_2019 <- rename(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)

```

Checking data types of columns

```

## Rows: 1,059,681
## Columns: 12
## $ ride_id          <int> 18000527, 18000528, 18000529, 18000530, 18000531, 1~
## $ started_at       <chr> "2018-04-01 00:04:44", "2018-04-01 00:06:42", "2018~
## $ ended_at         <chr> "2018-04-01 00:13:03", "2018-04-01 00:27:07", "2018~
## $ rideable_type     <int> 3819, 5000, 5165, 3851, 5065, 5962, 4570, 1323, 197~
## $ tripduration     <chr> "499.0", "1,225.0", "960.0", "434.0", "709.0", "659~
## $ start_station_id <int> 22, 157, 106, 241, 228, 244, 128, 130, 130, 121, 12~
## $ start_station_name <chr> "May St & Taylor St", "Lake Shore Dr & Wellington A~
## $ end_station_id   <int> 171, 190, 106, 171, 219, 325, 130, 69, 69, 351, 351~
## $ end_station_name <chr> "May St & Cullerton St", "Southport Ave & Wrightwoo~
## $ member_casual    <chr> "Subscriber", "Subscriber", "Customer", "Subscriber~
## $ gender           <chr> "Male", "Male", "", "Male", "Male", "Male", "Male",~
## $ birthyear        <int> 1994, 1965, NA, 1998, 1983, 1991, 1978, 1991, 1990,~

## Rows: 1,513,570
## Columns: 12
## $ ride_id          <int> 19244622, 19244623, 19244624, 19244625, 19244626, 1~
## $ started_at       <chr> "2018-07-01 00:00:03", "2018-07-01 00:00:13", "2018~
## $ ended_at         <chr> "2018-07-01 23:56:11", "2018-07-01 00:06:39", "2018~
## $ rideable_type     <int> 5429, 93, 2461, 2991, 2851, 5980, 3132, 2281, 3465,~
## $ tripduration     <chr> "86,168.0", "386.0", "1,391.0", "1,386.0", "656.0",~
## $ start_station_id <int> 140, 153, 76, 76, 60, 128, 168, 168, 229, 229, 39, ~
## $ start_station_name <chr> "Dearborn Pkwy & Delaware Pl", "Southport Ave & Wel~
## $ end_station_id   <int> 106, 250, 301, 301, 166, 71, 321, 321, 324, 324, 87~
## $ end_station_name <chr> "State St & Pearson St", "Ashland Ave & Wellington ~
## $ member_casual    <chr> "Customer", "Subscriber", "Subscriber", "Subscriber~
## $ gender           <chr> "", "Male", "Female", "Male", "Male", "Male", "", "~
## $ birthyear        <int> NA, 1986, 1987, 1986, 1961, 1995, NA, NA, NA, NA, N~

## Rows: 642,686
## Columns: 12
## $ ride_id          <int> 20983530, 20983531, 20983532, 20983533, 20983534, 2~
## $ started_at       <chr> "2018-10-01 00:01:17", "2018-10-01 00:03:59", "2018~
## $ ended_at         <chr> "2018-10-01 00:29:35", "2018-10-01 00:10:55", "2018~
## $ rideable_type     <int> 4551, 847, 6188, 6372, 1927, 2392, 308, 1187, 6247,~
## $ tripduration     <chr> "1,698.0", "416.0", "534.0", "778.0", "1,102.0", "2~
## $ start_station_id <int> 85, 13, 59, 328, 93, 229, 148, 374, 268, 125, 110, ~
## $ start_station_name <chr> "Michigan Ave & Oak St", "Wilton Ave & Diversey Pkw~
## $ end_station_id   <int> 166, 144, 197, 419, 159, 318, 11, 130, 289, 175, 28~

```

```
## $ end_station_name <chr> "Ashland Ave & Wrightwood Ave", "Larrabee St & Webs~
## $ member_casual <chr> "Subscriber", "Subscriber", "Subscriber", "Subscrib~
## $ gender <chr> "Male", "Female", "Male", "Female", "Female", "Male~
## $ birthyear <int> 1992, 1982, 1986, 1960, 1993, 1992, 1997, 1992, 198~

## Rows: 365,069
## Columns: 12
## $ ride_id <int> 21742443, 21742444, 21742445, 21742446, 21742447, 2~
## $ started_at <chr> "2019-01-01 00:04:37", "2019-01-01 00:08:13", "2019~
## $ ended_at <chr> "2019-01-01 00:11:07", "2019-01-01 00:15:34", "2019~
## $ rideable_type <int> 2167, 4386, 1524, 252, 1170, 2437, 2708, 2796, 6205~
## $ tripduration <chr> "390.0", "441.0", "829.0", "1,783.0", "364.0", "216~
## $ start_station_id <int> 199, 44, 15, 123, 173, 98, 98, 211, 150, 268, 299, ~
## $ start_station_name <chr> "Wabash Ave & Grand Ave", "State St & Randolph St",~
## $ end_station_id <int> 84, 624, 644, 176, 35, 49, 49, 142, 148, 141, 295, ~
## $ end_station_name <chr> "Milwaukee Ave & Grand Ave", "Dearborn St & Van Bur~
## $ member_casual <chr> "Subscriber", "Subscriber", "Subscriber", "Subscrib~
## $ gender <chr> "Male", "Female", "Female", "Male", "Male", "Female~
## $ birthyear <int> 1989, 1990, 1994, 1993, 1994, 1983, 1984, 1990, 199~
```

ride_id is in integer form but it should be in character form for consistency

```
q4_2018 <- mutate(q4_2018, ride_id = as.character(ride_id)
,rideable_type = as.character(rideable_type))
q3_2018 <- mutate(q3_2018, ride_id = as.character(ride_id)
,rideable_type = as.character(rideable_type))
q2_2018 <- mutate(q2_2018, ride_id = as.character(ride_id)
,rideable_type = as.character(rideable_type))
q1_2019 <- mutate(q1_2019, ride_id = as.character(ride_id)
,rideable_type = as.character(rideable_type))
```

Changing data type *Checking one dataset if the conversion happens or not*

```
glimpse(q3_2018)
```

```
## Rows: 1,513,570
## Columns: 12
## $ ride_id <chr> "19244622", "19244623", "19244624", "19244625", "19~
## $ started_at <chr> "2018-07-01 00:00:03", "2018-07-01 00:00:13", "2018~
## $ ended_at <chr> "2018-07-01 23:56:11", "2018-07-01 00:06:39", "2018~
## $ rideable_type <chr> "5429", "93", "2461", "2991", "2851", "5980", "3132~
## $ tripduration <chr> "86,168.0", "386.0", "1,391.0", "1,386.0", "656.0",~
## $ start_station_id <int> 140, 153, 76, 76, 60, 128, 168, 168, 229, 229, 39, ~
## $ start_station_name <chr> "Dearborn Pkwy & Delaware Pl", "Southport Ave & Wel~
## $ end_station_id <int> 106, 250, 301, 301, 166, 71, 321, 321, 324, 324, 87~
## $ end_station_name <chr> "State St & Pearson St", "Ashland Ave & Wellington ~
## $ member_casual <chr> "Customer", "Subscriber", "Subscriber", "Subscriber~
## $ gender <chr> "", "Male", "Female", "Male", "Male", "Male", "", "~
## $ birthyear <int> NA, 1986, 1987, 1986, 1961, 1995, NA, NA, NA, NA, N~
```

Merging the datasets Now all the datasets is in the same order so we can now combine them into one single dataset to do the further cleaning and transform

```
all_trips <- bind_rows(q2_2018, q3_2018, q4_2018, q1_2019)
```

Inspect the dataset Since now the dataset of all the month is transform into one dataset we will inspect it and then clean and process it to make ready for analysis

```
#checking column names
colnames(all_trips)
```

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

```
#checking data types
glimpse(all_trips)
```

```
## Rows: 3,581,006
## Columns: 12
## $ ride_id          <chr> "18000527", "18000528", "18000529", "18000530", "18~
## $ started_at       <chr> "2018-04-01 00:04:44", "2018-04-01 00:06:42", "2018~
## $ ended_at         <chr> "2018-04-01 00:13:03", "2018-04-01 00:27:07", "2018~
## $ rideable_type     <chr> "3819", "5000", "5165", "3851", "5065", "5962", "45~
## $ tripduration     <chr> "499.0", "1,225.0", "960.0", "434.0", "709.0", "659~
## $ start_station_id <int> 22, 157, 106, 241, 228, 244, 128, 130, 130, 121, 12~
## $ start_station_name <chr> "May St & Taylor St", "Lake Shore Dr & Wellington A~
## $ end_station_id   <int> 171, 190, 106, 171, 219, 325, 130, 69, 69, 351, 351~
## $ end_station_name <chr> "May St & Cullerton St", "Southport Ave & Wrightwoo~
## $ member_casual    <chr> "Subscriber", "Subscriber", "Customer", "Subscriber~
## $ gender           <chr> "Male", "Male", "", "Male", "Male", "Male", "Male",~
## $ birthyear        <int> 1994, 1965, NA, 1998, 1983, 1991, 1978, 1991, 1990,~
```

Remove unnecessary columns Removing unnecessary columns for consistency

```
all_trips <- all_trips %>%
  select(-c(birthyear, gender, tripduration))
```

Converting data types Convert started_at and ended_at to date and time

```
all_trips$started_at<-ymd_hms(all_trips$started_at)
all_trips$ended_at <- ymd_hms(all_trips$ended_at)
```

Removing inconsistency There are four unique values in member_casual subscriber, member, customer, casual but 2020 onwards these member has been changed into two unique values member, casual

```
all_trips <- all_trips %>%
  mutate(member_casual = recode(member_casual
                                , "Subscriber" = "member"
                                , "Customer" = "casual"))
```

Ride length (new column) ride_length is the distance between started time and ended time

```
all_trips$ride_length <- difftime(all_trips$ended_at,all_trips$started_at,units = "mins")
head(all_trips$ride_length)
```

```
## Time differences in mins
## [1] 8.316667 20.416667 16.000000 7.233333 11.816667 10.983333
```

Also we will convert the ride_length into numeric for further calculations

```
all_trips$ride_length <- round(as.numeric(as.character(all_trips$ride_length)),2)
```

Round trip (new column) We will produce a new column named `round_trip` = "Yes" where `start_station_name` is equal to `end_station_name`

```
all_trips <- all_trips %>%
  mutate(round_trip=case_when(
    start_station_name==end_station_name ~ "Yes",
    start_station_name!=end_station_name ~ "No"
  ))
head(all_trips$round_trip)
```

```
## [1] "No" "No" "Yes" "No" "No" "No"
```

Day (new column) calculating the day using the `started_date` column

```
all_trips$day <- day(all_trips$started_at)
head(all_trips$day)
```

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

```
all_trips$day_of_week <- weekdays(all_trips$started_at)
head(all_trips$day_of_week)
```

Day of the week (new column)

```
## [1] "Sunday" "Sunday" "Sunday" "Sunday" "Sunday" "Sunday"
```

Month of the year (new column) calculating month using the `started_date` column

```
all_trips$month <- months.Date(all_trips$started_at)
head(all_trips$month)
```

```
## [1] "April" "April" "April" "April" "April" "April"
```

Year (new column) finally year column for summarizing the data by year

```
all_trips$year <- year(all_trips$started_at)
head(all_trips$year)
```

```
## [1] 2018 2018 2018 2018 2018 2018
```

Deleting/Filtering bad data The `start_station_name` = "DIVVY CASSETTE REPAIR MOBILE STATION" is not relevant because it is the maintenance station for the bike so we have to remove it

Also the negative `ride_length` is not good for analysis as the `ended_time` is less than the started time which is simply a bad data

```
all_trips<- all_trips %>%
  filter(!(all_trips$ride_length<0 | (all_trips$start_station_name=="DIVVY CASSETTE REPAIR MOBILE STATION")))
```

Saving the transform data Finally saved the transform data for analysis

```
write.csv(all_trips,row.names=F,"Bike_sharing_clean/2018-19_tripdatas.csv")
```

Aggregating the file After cleaning, merging and saving all the file its time to aggregate them because the file size is too large to work with them so it is a must to agregate them into most suitable form

We will use `ride_length` for aggregating the data since it is a numerical column and it is most important for our analysis

Loading the cleaned data

```
trip <- read.csv("Bike_sharing_clean/2018-19_tripdatas.csv")
head(trip)
```

```
##      ride_id      started_at      ended_at rideable_type
## 1 18000527 2018-04-01 00:04:44 2018-04-01 00:13:03      3819
## 2 18000528 2018-04-01 00:06:42 2018-04-01 00:27:07      5000
## 3 18000529 2018-04-01 00:07:19 2018-04-01 00:23:19      5165
## 4 18000530 2018-04-01 00:07:33 2018-04-01 00:14:47      3851
## 5 18000531 2018-04-01 00:10:23 2018-04-01 00:22:12      5065
## 6 18000532 2018-04-01 00:11:29 2018-04-01 00:22:28      5962
##      start_station_id      start_station_name end_station_id
## 1              22      May St & Taylor St          171
## 2             157 Lake Shore Dr & Wellington Ave          190
## 3             106      State St & Pearson St          106
## 4             241      Morgan St & Polk St          171
## 5             228      Damen Ave & Melrose Ave          219
## 6             244 Ravenswood Ave & Irving Park Rd          325
##      end_station_name member_casual ride_length round_trip day
## 1      May St & Cullerton St      member          8.32      No    1
## 2 Southport Ave & Wrightwood Ave      member         20.42      No    1
## 3      State St & Pearson St      casual          16.00     Yes    1
## 4      May St & Cullerton St      member          7.23      No    1
## 5      Damen Ave & Cortland St      member         11.82      No    1
## 6 Clark St & Winnemac Ave (Temp)      member         10.98      No    1
##      day_of_week month year
## 1      Sunday April 2018
## 2      Sunday April 2018
## 3      Sunday April 2018
## 4      Sunday April 2018
## 5      Sunday April 2018
## 6      Sunday April 2018
```

Checking the data type of data

```
glimpse(trip)

## Rows: 3,580,968
## Columns: 15
## $ ride_id      <int> 18000527, 18000528, 18000529, 18000530, 18000531, 1~
## $ started_at   <chr> "2018-04-01 00:04:44", "2018-04-01 00:06:42", "2018~
## $ ended_at     <chr> "2018-04-01 00:13:03", "2018-04-01 00:27:07", "2018~
## $ rideable_type <int> 3819, 5000, 5165, 3851, 5065, 5962, 4570, 1323, 197~
## $ start_station_id <int> 22, 157, 106, 241, 228, 244, 128, 130, 130, 121, 12~
## $ start_station_name <chr> "May St & Taylor St", "Lake Shore Dr & Wellington A~
## $ end_station_id <int> 171, 190, 106, 171, 219, 325, 130, 69, 69, 351, 351~
## $ end_station_name <chr> "May St & Cullerton St", "Southport Ave & Wrightwoo~
## $ member_casual <chr> "member", "member", "casual", "member", "member", "~
## $ ride_length   <dbl> 8.32, 20.42, 16.00, 7.23, 11.82, 10.98, 3.97, 5.88, ~
## $ round_trip    <chr> "No", "No", "Yes", "No", "No", "No", "No", "No", "N~
```

```
## $ day          <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ day_of_week  <chr> "Sunday", "Sunday", "Sunday", "Sunday", "Sunday", "~
## $ month        <chr> "April", "April", "April", "April", "April", "April~
## $ year         <int> 2018, 2018, 2018, 2018, 2018, 2018, 2018, 2018, 201~
```

Changing the data_type We will convert year data type since we will consider it as the categorical data when we will merge all the dataset when aggregating

```
trip$year <- as.character(trip$year)
```

Checking the Statistics Since we will aggregate the data based on ride_length, its important to check its statistics to decide the aggregate parameter but we have already aggregate the 2016-17 data based on median due to skewed column, we will consider this parameter for all the aggregation for consistency

```
summary(trip$ride_length)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
##      1.02     6.60    11.27    23.56    20.28 224220.33
```

We can easily see that the *ride_length* is right-skewed since there is a BIG difference between Q3 and Max value. So we will use median instead of mean for aggregation which is more ideal in this case of skewed column

Aggregating the data

```
trip_data<-aggregate(trip$ride_length~trip$member_casual+trip$round_trip+trip$day_of_week+trip$month+tr
head(trip_data)
```

```
##   trip$member_casual trip$round_trip trip$day_of_week trip$month trip$year
## 1                casual             No             Friday    April    2018
## 2                member             No             Friday    April    2018
## 3                casual             Yes             Friday    April    2018
## 4                member             Yes             Friday    April    2018
## 5                casual             No              Monday    April    2018
## 6                member             No              Monday    April    2018
##   trip$ride_length
## 1                25.800
## 2                 8.500
## 3                36.250
## 4                10.430
## 5                27.725
## 6                 8.820
```

Saving the aggregate

Finally, last step is to save the data so we can use this data to merge all other aggregates data

We will merge the data with the old data we saved while aggregating

```
trip_old <- read.csv("Bike_sharing_clean/tripdata_aggregate.csv")
head(trip_old)
```

```
##   trip.member_casual trip.round_trip trip.day_of_week trip.month trip.year
## 1                casual             No             Friday    April    2016
## 2                member             No             Friday    April    2016
## 3                casual             Yes             Friday    April    2016
## 4                member             Yes             Friday    April    2016
## 5                casual             No              Monday    April    2016
## 6                member             No              Monday    April    2016
##   trip.ride_length
```



```
## 1          19
## 2           9
## 3          22
## 4          10
## 5          22
## 6          10
```

Transforming the old data to merge perfectly We have to make column name and type consistent before merging

```
trip_old$trip.year <- as.character(trip_old$trip.year)
trip_old <- rename(trip_old,
  "trip$member_casual"=trip.member_casual,
  "trip$round_trip"=trip.round_trip,
  "trip$day_of_week"=trip.day_of_week,
  "trip$month"=trip.month,
  "trip$year"=trip.year,
  "trip$ride_length"=trip.ride_length
)
```

```
trip_merged <- bind_rows(trip_old, trip_data)
write.csv(trip_merged, row.names = F, "Bike_sharing_clean/tripdata_aggregate.csv")
```