

Cleaning Cylistic Data 2023-01

2023-07-31

Import data

```
data_01 <- read.csv(file="dataset/202301-divvy-tripdata.csv")
```

Check data 01

Check the data type for each meta

```
str(data_01)
```

```
## 'data.frame': 190301 obs. of 13 variables:
## $ ride_id : chr "F96D5A74A3E41399" "13CB7EB698CEDB88" "BD88A2E670661CE5" "C90792D034FED9" ...
## $ rideable_type : chr "electric_bike" "classic_bike" "electric_bike" "classic_bike" ...
## $ started_at : chr "2023-01-21 20:05:42" "2023-01-10 15:37:36" "2023-01-02 07:51:57" "2023-01-02 08:05:11" ...
## $ ended_at : chr "2023-01-21 20:16:33" "2023-01-10 15:46:05" "2023-01-02 08:05:11" "2023-01-02 08:05:11" ...
## $ start_station_name: chr "Lincoln Ave & Fullerton Ave" "Kimbark Ave & 53rd St" "Western Ave & Lunt Ave" "Western Ave & Lunt Ave" ...
## $ start_station_id : chr "TA1309000058" "TA1309000037" "RP-005" "TA1309000037" ...
## $ end_station_name : chr "Hampden Ct & Diversey Ave" "Greenwood Ave & 47th St" "Valli Produce - E" "Valli Produce - E" ...
## $ end_station_id : chr "202480.0" "TA1308000002" "599" "TA1308000002" ...
## $ start_lat : num 41.9 41.8 42 41.8 41.8 ...
## $ start_lng : num -87.6 -87.6 -87.7 -87.6 -87.6 ...
## $ end_lat : num 41.9 41.8 42 41.8 41.8 ...
## $ end_lng : num -87.6 -87.6 -87.7 -87.6 -87.6 ...
## $ member_casual : chr "member" "member" "casual" "member" ...
```

```
summary(data_01)
```

```
## ride_id rideable_type started_at ended_at
## Length:190301 Length:190301 Length:190301 Length:190301
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
##
##
## start_station_name start_station_id end_station_name end_station_id
## Length:190301 Length:190301 Length:190301 Length:190301
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
```

```
##
##
##
##   start_lat   start_lng   end_lat   end_lng
##   Min.   :41.65   Min.   :-87.83   Min.   :41.65   Min.   :-87.84
##   1st Qu.:41.88   1st Qu.: -87.66   1st Qu.:41.88   1st Qu.: -87.66
##   Median :41.89   Median : -87.64   Median :41.89   Median : -87.64
##   Mean   :41.90   Mean   : -87.65   Mean   :41.90   Mean   : -87.65
##   3rd Qu.:41.93   3rd Qu.: -87.63   3rd Qu.:41.93   3rd Qu.: -87.63
##   Max.   :42.07   Max.   : -87.53   Max.   :42.08   Max.   : -87.53
##                                     NA's   :127   NA's   :127
## member_casual
## Length:190301
## Class :character
## Mode  :character
##
##
##
##
```

From meta check we know that data type of column “started_at” and “end_at” should be datetime

Check duplicate data 01

```
print(data_01[duplicated(data_01), ])
```

```
## [1] ride_id      rideable_type started_at    ended_at
## [5] start_station_name start_station_id end_station_name end_station_id
## [9] start_lat     start_lng     end_lat      end_lng
## [13] member_casual
## <0 rows> (or 0-length row.names)
```

Duplicate data checking result : no data duplicate in data_01

Remove duplicate data

Remove Duplicate data result : No data to remove

Check missing value data in character data type

```
count(data_01[data_01$ride_id=="", ])
```

```
##   n
## 1 0
```

```
count(data_01[data_01$rideable_type=="", ])
```

```
##      n
## 1 0
```

```
count(data_01[data_01$started_at=="", ])
```

```
##      n
## 1 0
```

```
count(data_01[data_01$ended_at=="", ])
```

```
##      n
## 1 0
```

```
count(data_01[data_01$start_station_name=="", ])
```

```
##      n
## 1 26721
```

```
count(data_01[data_01$start_station_id=="", ])
```

```
##      n
## 1 26721
```

```
count(data_01[data_01$end_station_name=="", ])
```

```
##      n
## 1 27840
```

```
count(data_01[data_01$end_station_id=="", ])
```

```
##      n
## 1 27840
```

```
count(data_01[data_01$member_casual=="", ])
```

```
##      n
## 1 0
```

Missing value checking result :

ride_id: [0] rideable_type: [0] started_at: [0] ended_at: [0] start_station_name: [26,721] start_station_id: [26,721] end_station_name: [27,840] end_station_id: [27,840] member_casual: [0]

Fill Missing value with NA

Missing value (empty data) in start_station_name, start_station_id, end_station_name, end_station_id will be filling with NA

```
data_01 <- replace(data_01, data_01 == "", NA)
```

Fill missing value result : empty data was replace with NA

Check missing value data

```
count(data_01[is.na(data_01$start_lat) | data_01$start_lat=="", ])
```

```
##      n  
## 1 0
```

```
count(data_01[is.na(data_01$start_lng) | data_01$start_lng=="", ])
```

```
##      n  
## 1 0
```

```
count(data_01[is.na(data_01$end_lat) | data_01$end_lat=="", ])
```

```
##      n  
## 1 127
```

```
count(data_01[is.na(data_01$end_lng) | data_01$end_lng=="", ])
```

```
##      n  
## 1 127
```

Missing value checking result :

start latitude and longitude : [0] end latitude and longitude : [127]

Remove Missing value with NA

Missing value in end_lat, end_lng will be delete by remove the row

```
# remove missing value data in this other data if there are also missing values  
# data_01 <- data_01[!is.na(data_01$rideable_type), ]  
# data_01 <- data_01[!is.na(data_01$started_at), ]  
# data_01 <- data_01[!is.na(data_01$ended_at), ]  
# data_01 <- data_01[!is.na(data_01$member_casual), ]
```

```
data_01 <- data_01[!is.na(data_01$end_lat), ]  
data_01 <- data_01[!is.na(data_01$end_lng), ]
```

```
count(data_01[is.na(data_01$end_lat) | data_01$end_lat=="", ])
```

```
##      n  
## 1 0
```

```
count(data_01[is.na(data_01$end_lng) | data_01$end_lng=="", ])
```

```
##      n  
## 1 0
```

Remove missing value result : Row with missing value data was removed

Check outliers in coordinate data

```
print(cat("start_lat : mean max min : ",  
        mean(data_01$start_lat),  
        max(data_01$start_lat),  
        min(data_01$start_lat)))
```

```
## start_lat : mean max min :  41.8971 42.07 41.6485NULL
```

```
print(cat("start_lng : mean max min : ",  
        mean(data_01$start_lng), max(data_01$start_lng), min(data_01$start_lng)))
```

```
## start_lng : mean max min : -87.64717 -87.52823 -87.83NULL
```

```
print(cat("end_lat : mean max min : ",  
        mean(data_01$end_lat), max(data_01$end_lat), min(data_01$end_lat)))
```

```
## end_lat : mean max min :  41.89716 42.08 41.6485NULL
```

```
print(cat("end_lng : mean max min : ",  
        mean(data_01$end_lng), max(data_01$end_lng), min(data_01$end_lng)))
```

```
## end_lng : mean max min : -87.64726 -87.52823 -87.84NULL
```

Outliers checking result : no outliers in coordinate data, max and min value for each data doesnt far from average value

Remove useless column data

According to the bussines task, start_station_name and end_station_name will be remove

```
data_01 <- data_01[, -which(names(data_01) == "start_station_name")]  
data_01 <- data_01[, -which(names(data_01) == "end_station_name")]  
  
head(data_01)
```

```
##          ride_id rideable_type      started_at      ended_at
## 1 F96D5A74A3E41399 electric_bike 2023-01-21 20:05:42 2023-01-21 20:16:33
## 2 13CB7EB698CEDB88 classic_bike 2023-01-10 15:37:36 2023-01-10 15:46:05
## 3 BD88A2E670661CE5 electric_bike 2023-01-02 07:51:57 2023-01-02 08:05:11
## 4 C90792D034FED968 classic_bike 2023-01-22 10:52:58 2023-01-22 11:01:44
## 5 3397017529188E8A classic_bike 2023-01-12 13:58:01 2023-01-12 14:13:20
## 6 58E68156DAE3E311 electric_bike 2023-01-31 07:18:03 2023-01-31 07:21:16
##   start_station_id end_station_id start_lat start_lng end_lat  end_lng
## 1      TA1309000058      202480.0  41.92407 -87.64628 41.93000 -87.64000
## 2      TA1309000037      TA1308000002  41.79957 -87.59475 41.80983 -87.59938
## 3          RP-005          599  42.00857 -87.69048 42.03974 -87.69941
## 4      TA1309000037      TA1308000002  41.79957 -87.59475 41.80983 -87.59938
## 5      TA1309000037      TA1308000002  41.79957 -87.59475 41.80983 -87.59938
## 6      TA1309000019      202480.0  41.92607 -87.63886 41.93000 -87.64000
##   member_casual
## 1      member
## 2      member
## 3      casual
## 4      member
## 5      member
## 6      member
```

```
str(data_01)
```

```
## 'data.frame':   190174 obs. of  11 variables:
## $ ride_id      : chr  "F96D5A74A3E41399" "13CB7EB698CEDB88" "BD88A2E670661CE5" "C90792D034FED968" ...
## $ rideable_type : chr  "electric_bike" "classic_bike" "electric_bike" "classic_bike" ...
## $ started_at   : chr  "2023-01-21 20:05:42" "2023-01-10 15:37:36" "2023-01-02 07:51:57" "2023-01-22 10:52:58" ...
## $ ended_at     : chr  "2023-01-21 20:16:33" "2023-01-10 15:46:05" "2023-01-02 08:05:11" "2023-01-22 11:01:44" ...
## $ start_station_id: chr  "TA1309000058" "TA1309000037" "RP-005" "TA1309000037" ...
## $ end_station_id : chr  "202480.0" "TA1308000002" "599" "TA1308000002" ...
## $ start_lat     : num  41.9 41.8 42 41.8 41.8 ...
## $ start_lng     : num  -87.6 -87.6 -87.7 -87.6 -87.6 ...
## $ end_lat       : num  41.9 41.8 42 41.8 41.8 ...
## $ end_lng       : num  -87.6 -87.6 -87.7 -87.6 -87.6 ...
## $ member_casual : chr  "member" "member" "casual" "member" ...
```

Export clean data into csv

```
# write.csv(data_01, "dataclean/202301-clean.csv", row.names = FALSE)
```