

Cleaning Cylistic Data 2023-04

2023-07-31

Import data

```
data_01 <- read.csv(file="dataset/202304-divvy-tripdata.csv")
```

Check data 01

Check the data type for each meta

```
str(data_01)
```

```
## 'data.frame': 426590 obs. of 13 variables:
## $ ride_id : chr "8FE8F7D9C10E88C7" "34E4ED3ADF1D821B" "5296BF07A2F77CB5" "40759916B76D5D" ...
## $ rideable_type : chr "electric_bike" "electric_bike" "electric_bike" "electric_bike" ...
## $ started_at : chr "2023-04-02 08:37:28" "2023-04-19 11:29:02" "2023-04-19 08:41:22" "2023-04-19 08:43:22" ...
## $ ended_at : chr "2023-04-02 08:41:37" "2023-04-19 11:52:12" "2023-04-19 08:43:22" "2023-04-19 08:43:22" ...
## $ start_station_name: chr "" "" "" "" ...
## $ start_station_id : chr "" "" "" "" ...
## $ end_station_name : chr "" "" "" "" ...
## $ end_station_id : chr "" "" "" "" ...
## $ start_lat : num 41.8 41.9 41.9 41.9 41.9 ...
## $ start_lng : num -87.6 -87.7 -87.7 -87.7 -87.7 ...
## $ end_lat : num 41.8 41.9 41.9 41.9 41.9 ...
## $ end_lng : num -87.6 -87.7 -87.7 -87.7 -87.6 ...
## $ member_casual : chr "member" "member" "member" "member" ...
```

```
summary(data_01)
```

```
## ride_id rideable_type started_at ended_at
## Length:426590 Length:426590 Length:426590 Length:426590
## 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:426590 Length:426590 Length:426590 Length:426590
## 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.      :-88.11
##   1st Qu.:41.88   1st Qu.: -87.66   1st Qu.:41.88   1st Qu.: -87.66
##   Median :41.90   Median : -87.64   Median :41.90   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.    : -87.52   Max.    :42.08   Max.    : -87.53
##                                     NA's    :435    NA's    :435
## member_casual
## Length:426590
## 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[is.na(data_01$ride_id) | data_01$ride_id=="", ])
```

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

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

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

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

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

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

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

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

```
##          n
## 1 63814
```

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

```
##          n
## 1 63814
```

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

```
##          n
## 1 68630
```

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

```
##          n
## 1 68630
```

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

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

Missing value checking result :

ride_id: [0]

rideable_type: [0]

started_at: [0]

ended_at: [0]

start_station_name: [63,814]

start_station_id: [63,814]

end_station_name: [68,630]

end_station_id: [68,630]

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  
## 1 435
```

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

```
##      n  
## 1 435
```

Missing value checking result :

start latitude and longitude : [0]

end latitude and longitude : [435]

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 : 41.90153 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 : -87.64697 -87.52 -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 : 41.902 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 : -87.64723 -87.52823 -88.11NULL
```

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 8FE8F7D9C10E88C7 electric_bike 2023-04-02 08:37:28 2023-04-02 08:41:37
## 2 34E4ED3ADF1D821B electric_bike 2023-04-19 11:29:02 2023-04-19 11:52:12
## 3 5296BF07A2F77CB5 electric_bike 2023-04-19 08:41:22 2023-04-19 08:43:22
## 4 40759916B76D5D52 electric_bike 2023-04-19 13:31:30 2023-04-19 13:35:09
## 5 77A96F460101AC63 electric_bike 2023-04-19 12:05:36 2023-04-19 12:10:26
## 6 8D6A2328E19DC168 electric_bike 2023-04-19 12:17:34 2023-04-19 12:21:38
##   start_station_id end_station_id start_lat start_lng end_lat end_lng
## 1              <NA>          <NA>    41.80   -87.60   41.79   -87.60
## 2              <NA>          <NA>    41.87   -87.65   41.93   -87.68
## 3              <NA>          <NA>    41.93   -87.66   41.93   -87.66
## 4              <NA>          <NA>    41.92   -87.65   41.91   -87.65
## 5              <NA>          <NA>    41.91   -87.65   41.91   -87.63
## 6              <NA>          <NA>    41.91   -87.63   41.92   -87.65
##   member_casual
## 1         member
## 2         member
## 3         member
## 4         member
## 5         member
## 6         member
```

```
str(data_01)
```

```
## 'data.frame':   426155 obs. of  11 variables:
## $ ride_id      : chr  "8FE8F7D9C10E88C7" "34E4ED3ADF1D821B" "5296BF07A2F77CB5" "40759916B76D5D52" ...
## $ rideable_type : chr  "electric_bike" "electric_bike" "electric_bike" "electric_bike" ...
## $ started_at   : chr  "2023-04-02 08:37:28" "2023-04-19 11:29:02" "2023-04-19 08:41:22" "2023-04-19 08:43:22" ...
## $ ended_at     : chr  "2023-04-02 08:41:37" "2023-04-19 11:52:12" "2023-04-19 08:43:22" "2023-04-19 13:35:09" ...
## $ start_station_id: chr  NA NA NA NA ...
## $ end_station_id  : chr  NA NA NA NA ...
## $ start_lat      : num  41.8 41.9 41.9 41.9 41.9 ...
## $ start_lng      : num  -87.6 -87.7 -87.7 -87.7 -87.7 ...
## $ end_lat        : num  41.8 41.9 41.9 41.9 41.9 ...
## $ end_lng        : num  -87.6 -87.7 -87.7 -87.7 -87.6 ...
## $ member_casual  : chr  "member" "member" "member" "member" ...
```

Export clean data into csv

```
# write.csv(data_01, "dataclean/202304-clean.csv", row.names = FALSE)
```