# Exploratory Data Analysis in Microsoft Excel

- I created new column called "ride length" and set up values in "time" HH: MM: SS 37:30:55

- I created new column called "day of week" Where Sunday is 1 and Saturday is 7

- I deleted columns which I will not use for the analysis, making tables easy to read and understand.

- Using the function "find and replace", I located all blanks in spreadsheet and deleted them.

- Calculated mean(average) and max values in column "ride length".

- Calculated mode in column "day of week", to discover the most frequently occurring value that appears in this column.

- Using Custom Filter function, I found outliers and deleted all the rows that are showing ride length <01.00 min and 24.00 h>

- Made pivot tables for each of twelve files to get initial insights of how these two types of riders use bikes differently.

- To be specific, for each month, I calculated average ride length, average ride length by day, and number of rides per day, and displayed results through charts to make it easier for stakeholders to understand.

## Cleaning Process

Removing duplicates:



- Made sure that all the tables are consistent (column names etc.)

- Removed all the extra spaces using TRIM function. (Where applicable)

- Made sure that all dates are in the same format.

- Used FILTER function to be sure that there is not any error or unusual values.



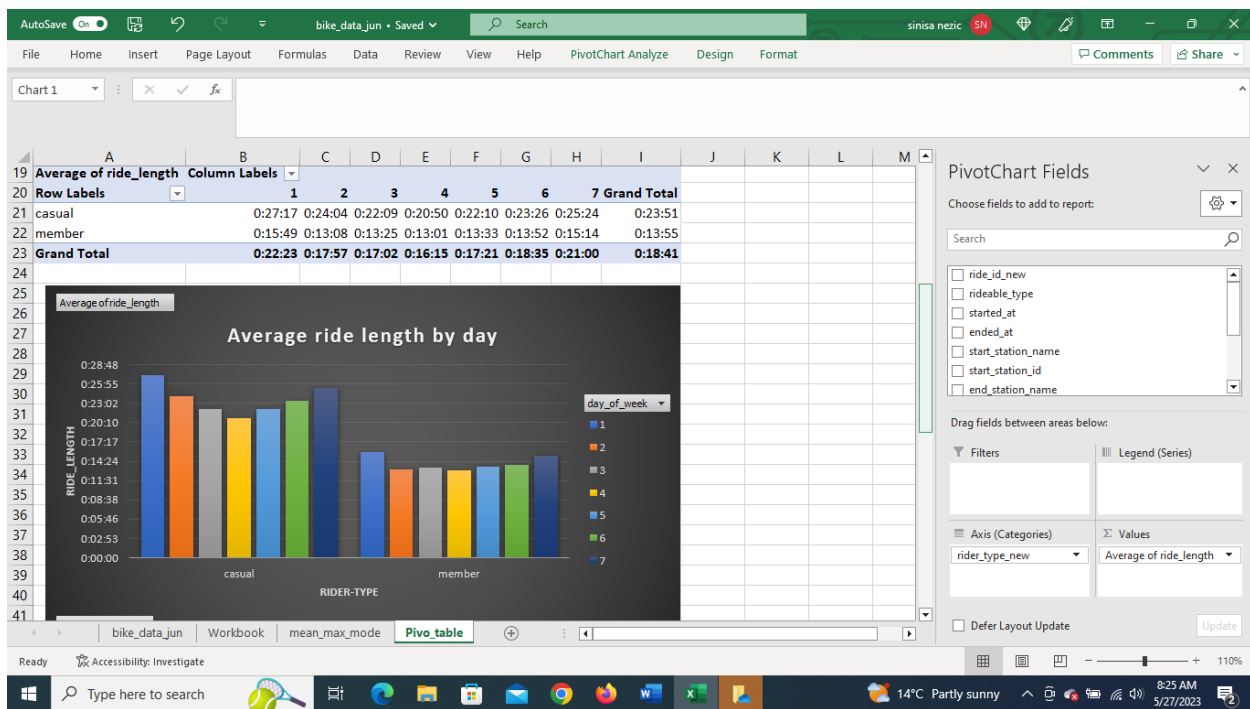Average ride length:

## Average ride length by day (Where #1 is Sunday, #7 is Saturday)



| Average of ride_length | Column Labels | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Row Labels | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Grand Total |
| casual | 0:27:17 | 0:24:04 | 0:22:09 | 0:20:50 | 0:22:10 | 0:23:26 | 0:25:24 | 0:23:51 |
| member | 0:15:49 | 0:13:08 | 0:13:25 | 0:13:01 | 0:13:33 | 0:13:52 | 0:15:14 | 0:13:55 |
| Grand Total | 0:22:23 | 0:17:57 | 0:17:02 | 0:16:15 | 0:17:21 | 0:18:35 | 0:21:00 | 0:18:41 |

## Number of rides per day (Where #1 is Sunday, #7 is Saturday)



| Count of ride_id_new | Column Labels | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Row Labels | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Grand Total |
| casual | 64443 | 36199 | 37923 | 47369 | 56704 | 54677 | 63733 | 361048 |
| member | 47973 | 45845 | 53843 | 67372 | 72002 | 56172 | 48588 | 391795 |
| Grand Total | 112416 | 82044 | 91766 | 114741 | 128706 | 110849 | 112321 | 752843 |