# Excel and SQL Analysis - Notes

Welcome to my comprehensive notes on Excel and SQL analysis techniques. This collection of notes serves as a documentation of my data analysis journey, showcasing my proficiency in Excel and SQL.

## About the Notes

In this document, I've meticulously documented the process and insights derived from my data analysis journey using Excel and SQL. With a focus on the Cylistic Case Study provided under the Google Data Analytics Professional certificate, these notes highlight my approach to:

- **Excel Analysis**: Learn how I harnessed the power of Excel to clean, transform, and visualize the data. From basic descriptive statistics to advanced functions, these notes cover it all.
- **SQL Techniques**: Dive into SQL commands used for data querying, transformations, and analysis. Witness how I employed SQL to unveil hidden patterns and trends within the dataset.

## Highlights

- Detailed step-by-step instructions for replicating the analyses conducted in the project.
- Insights into data cleaning, data transformation, and visualization techniques using Excel.
- SQL commands and strategies for extracting valuable information from the dataset.

## How to Use

Whether you're a beginner or an experienced data enthusiast, these notes provide valuable insights and practical knowledge to elevate your analytical skills. Feel free to explore the notes to enhance your understanding of Excel and SQL analysis.

Your feedback is invaluable as we continue to explore the realm of data analysis and extract insights that drive informed decisions together!
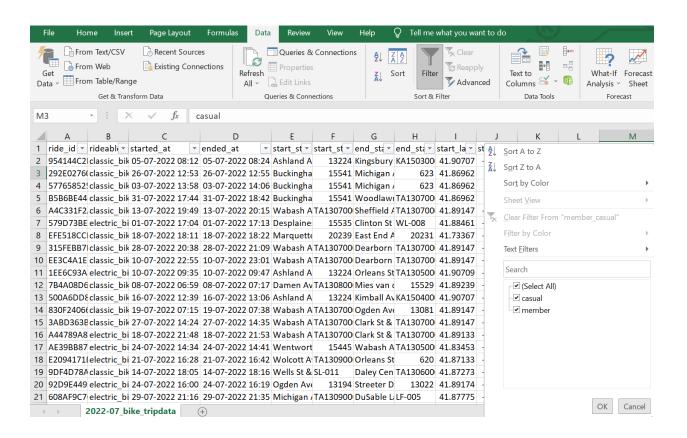
Abhay Dhupar

(Author)

# Excel

First we'll save the csv files to workbook type so that we don't have data loss.

Now analyzing consolidated data for year is difficult on excel as there are too many rows (more than 5 lakhs). So we'll analyze month wise.

## Working on month data

First we'll use the filter option in the data tab to check if there's any wrong value in any field.
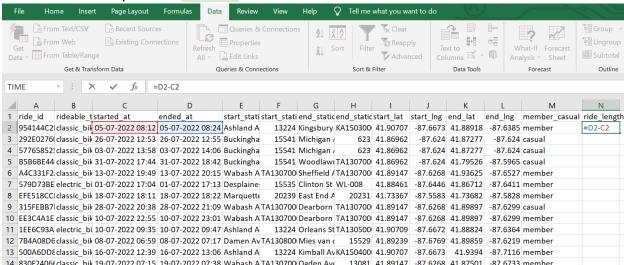


This way, we can check for any missing values and if there are any, we can remove it.
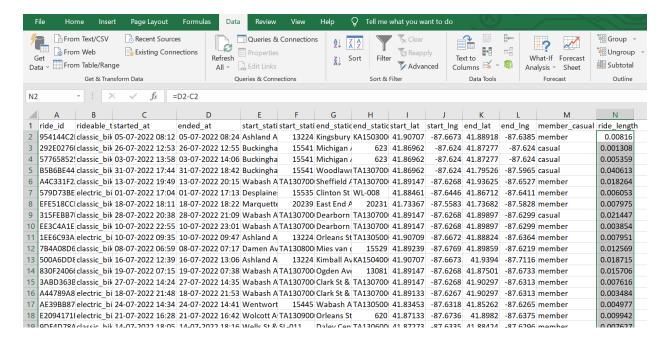
# Add more columns for descriptive analysis

Here we will add two more columns **ride_length** and **day_of_week** to make our descriptive analysis more easy.
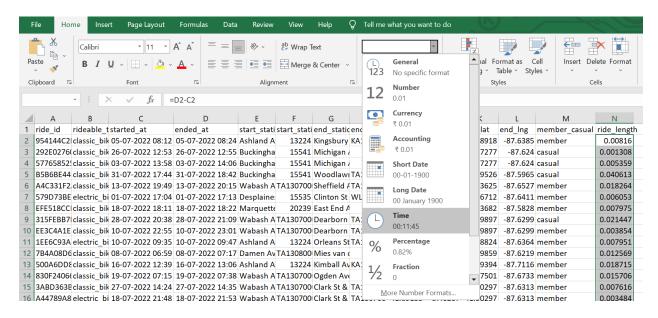
To calculate **ride_length,** we'll find the difference between columns **ended_at** and **started_at.** So we'll use the expression '=D2-C2' to find the difference.



After than we can press enter and use the fill handle to apply for all rows. Or we can also select the rows column to fill and press CTRL+D to make the fill handle on column operation.

Now the data type we have for ride length is in decimal. So let's convert it to time. For this first we'll select the data and then go to Home>Number>Dropdown to select time.



This will update the data type of all values to be shown in time format (i.e. hh:mm:ss)

Now let's add the **day_of_week** column to our data.

We'll use the WEEKDAY() to do so. There are two argument we'll give. FIrst will be the date from which we'll find our weekday. In this case, we'll use Column C, which is **started_at**
And then we'll give the second argument as **1,** to indicate that the series as (1-Sunday, 2-Monday .. and so on)



After that we'll press enter and use the fill handle or select all and click CTRL + D to apply to all rows.

## Doing some descriptive analysis

Now we can find **mean, median** of the **ride length** and **mode** of **day_of_week** to get some simple insights about our data.

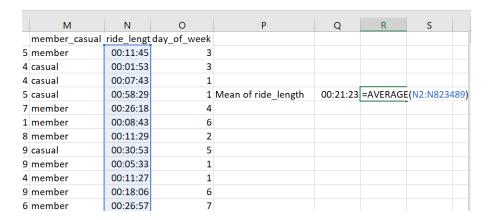| | M | N | O | P | Q | R | S |
|---|---|---|---|---|---|---|---|
| | member_casual | ride_lengt | day_of_week | | | | |
| 5 | member | 00:11:45 | 3 | | | | |
| 4 | casual | 00:01:53 | 3 | | | | |
| 4 | casual | 00:07:43 | 1 | | | | |
| 5 | casual | 00:58:29 | 1 | Mean of ride_length | 00:21:23 | =AVERAGE(N2:N823489) | |
| 7 | member | 00:26:18 | 4 | | | | |
| 1 | member | 00:08:43 | 6 | | | | |
| 8 | member | 00:11:29 | 2 | | | | |
| 9 | casual | 00:30:53 | 5 | | | | |
| 9 | member | 00:05:33 | 1 | | | | |
| 4 | member | 00:11:27 | 1 | | | | |
| 9 | member | 00:18:06 | 6 | | | | |
| 6 | member | 00:26:57 | 7 | | | | |

This way we'll find the average or mean.

To find median of ride_length

| M | N | O | P | Q | R | S |
|---|---|---|---|---|---|---|
| member_casual | ride_lengt | day_of_week | | | | |
| member | 00:11:45 | 3 | | | | |
| casual | 00:01:53 | 3 | | | | |
| casual | 00:07:43 | 1 | | | | |
| casual | 00:58:29 | 1 | Mean of ride_length | 00:21:23 | | |
| member | 00:26:18 | 4 | Median of ride_length | 00:11:42 | =MEDIAN(N2:N823490) | |
| member | 00:08:43 | 6 | | | | |
| member | 00:11:29 | 2 | | | | |
| casual | 00:30:53 | 5 | | | | |
| member | 00:05:33 | 1 | | | | |
| member | 00:11:27 | 1 | | | | |
| member | 00:18:06 | 6 | | | | |
| member | 00:26:57 | 7 | | | | |

To find mode of day_of_week

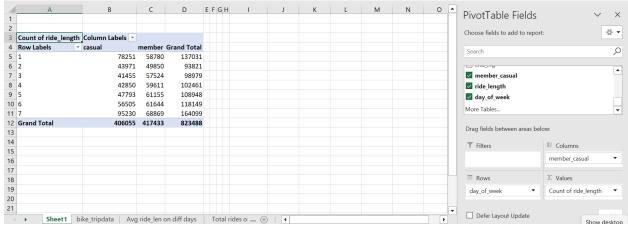| M | N | O | P | Q | R | S |
|---|---|---|---|---|---|---|
| member_casual | ride_lengt | day_of_week | | | | |
| member | 00:11:45 | 3 | | | | |
| casual | 00:01:53 | 3 | | | | |
| casual | 00:07:43 | 1 | | | | |
| casual | 00:58:29 | 1 | Mean of ride_length | 00:21:23 | | |
| member | 00:26:18 | 4 | Median of ride_length | 00:11:42 | 00:11:42 | |
| member | 00:08:43 | 6 | Mode of day_of_week | | 7 | =MODE(O2:O823489) |
| member | 00:11:29 | 2 | | | | |
| casual | 00:30:53 | 5 | | | | |
| member | 00:05:33 | 1 | | | | |
| member | 00:11:27 | 1 | | | | |
| member | 00:18:06 | 6 | | | | |
| member | 00:26:57 | 7 | | | | |

This is the frequency of **day_of_week** is 7.

# Creating Pivot table and inserting visuals

Select all the rows and insert a pivot table on a blank sheet.
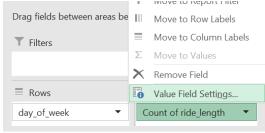


Now let's calculate the avg ride length on each weekday by each customer type.
For this, let's put day_of_week to column, and member_casual to rows. And ride length to values in the **pivot table fields pane.**
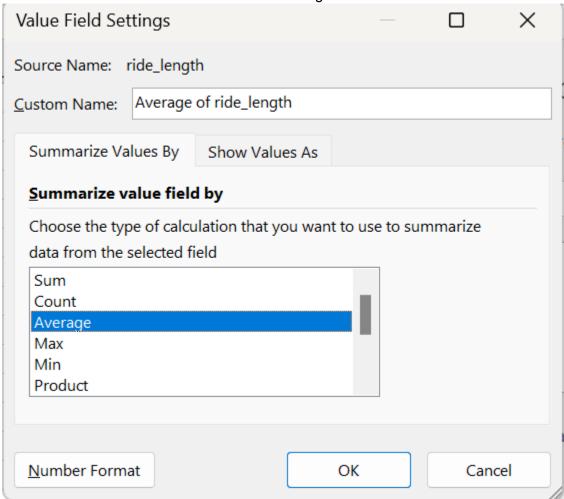


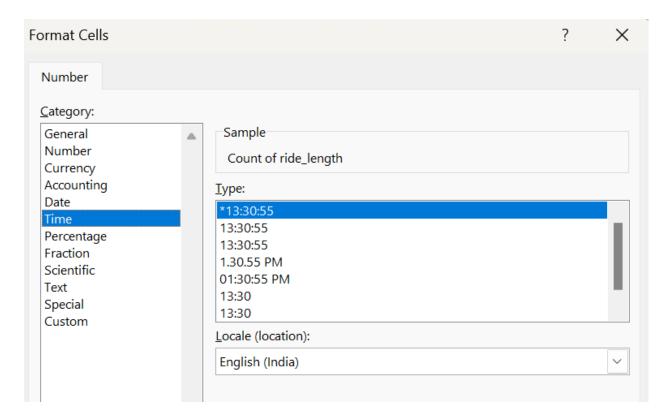This will generate a summary table in the sheets.

Since we want the average ride length, we'd go the the value field settings of the Count of ride length in the values section.
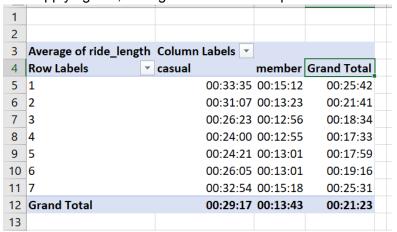
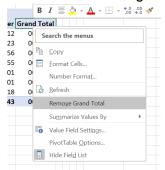Now we'll set the Summarize value field to average.



Also we'll go to the number format section in the lower left and set it to time. This will change the type of our result to the following format 'hh/mm/ss'

## Format Cells

**Number**

Category:
- General
- Number
- Currency
- Accounting
- Date
- **Time**
- Percentage
- Fraction
- Scientific
- Text
- Special
- Custom

Sample

Count of ride_length

Type:

*13:30:55
13:30:55
13:30:55
1.30.55 PM
01:30:55 PM
13:30
13:30

Locale (location):

English (India)

After applying this, we'll get our desired output.

| | Average of ride_length | Column Labels | |
|---|---|---|---|
| | Row Labels | casual | member | Grand Total |
| 1 | | 00:33:35 | 00:15:12 | 00:25:42 |
| 2 | | 00:31:07 | 00:13:23 | 00:21:41 |
| 3 | | 00:26:23 | 00:12:56 | 00:18:34 |
| 4 | | 00:24:00 | 00:12:55 | 00:17:33 |
| 5 | | 00:24:21 | 00:13:01 | 00:17:59 |
| 6 | | 00:26:05 | 00:13:01 | 00:19:16 |
| 7 | | 00:32:54 | 00:15:18 | 00:25:31 |
| Grand Total | | 00:29:17 | 00:13:43 | 00:21:23 |

We can remove the grand totals by right clicking on the name and removing it.

Search the menus
- Copy
- Format Cells...
- Number Format...
- Refresh
- Remove Grand Total
- Summarize Values By ▸
- Value Field Settings...
- PivotTable Options...
- Hide Field List

Now we can visualize it, select the table then follow the steps to choose the chart to visualize. You can choose your preferred chart. In this example, I'm using a side by side bar chart.



The results would be like this,

Let's now count the total rides by each member on each weekday.

Regenerate the pivot table insertion steps again,and reach upto this mark.



Now let's put the day_of_week and member_casual in the same places that we did before.
This time rather than putting ride_length to values, we'd use **ride_id.**



This ride_id is already set to count. And we'd get desired result.

| Count of ride_id | Column Labels | |
|---|---|---|
| Row Labels | casual | member |
| 1 | 78251 | 58780 |
| 2 | 43971 | 49850 |
| 3 | 41455 | 57524 |
| 4 | 42850 | 59611 |
| 5 | 47793 | 61155 |
| 6 | 56505 | 61644 |
| 7 | 95230 | 68869 |

Now let's put the stacked bar chart to analyze this visually.
To get a stack bar chart, you can perform following operations

Select the table, and follow the steps.



We'd get this



We can analyze with other months data like this, and see the trends for each month.

# Working on Year data

We can easily merge the data by following these steps:



Now we can put all our files in a folder to give the location of the folder.

And now we can select the folder where all our files resides.



If the data is clean and does not contain any error or change in number of columns, we can directly do **Combine and load,** But if we think that the data needs some cleanliness and change of data types, we can select **Combine and transform.** This will open the power query editor and we can perform the operations.

Combining the data may take time. And analyzing it will also make a lot of time, as excel is suitable for less data and processing and analysis with big data is not possible with excel.



**Now that the data is so big (about 5 million rows), it can't be proceed in a Excel file.**
**So now we'll switch to BigQuery**

# BigQuery

First I uploaded my files on the dataset. Greater than 100 mb is not available with the upload option.
So I rather use the drive option. I uploaded it on drive and use the link to give table.

## Create table

**Select Drive URI ***
https://drive.google.com/open?id=1aWBh8D5xjfd8-kv4YFS6dtlBtpVoVXBR

**File format**
CSV

## Destination

**Project ***
cyclist-394306

**Data set ***
previous_12_month_data

**Table ***
2023-06

Unicode letters, marks, numbers, connectors, dashes or spaces allowed.

**Table type**
External table

## Schema

☑ Auto-detect

**CREATE TABLE**    CANCEL

## Working on Month data

In this below screenshot, I created a temp_table to store the following results from the query I'll use again and again to analyze from.
Here I added two columns, one for **ride_length** and the other for **day_of_week.**
Then I counted the total people grouped my their member_casual.

```sql
 1  WITH temp_table AS (
 2    SELECT
 3      ride_id,
 4      rideable_type,
 5      member_casual,
 6      TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
 7      EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
 8    FROM `cyclist-394306.previous_12_month_data.2022-07`
 9  )
10
11  SELECT
12    member_casual,
13    COUNT(*) AS total
14  FROM temp_table
15  GROUP BY
16    member_casual
```

## Query results

| | JOB INFORMATION | RESULTS | JSON | EXECUTI |
|---|---|---|---|---|

| Row | member_casual ▼ | total ▼ |
|---|---|---|
| 1 | member | 417433 |
| 2 | casual | 406055 |

**Total rideable_type each member_casual prefer**

Now let's calculate the total number of rideable_type each member_casual prefer.

```
1   WITH temp_table AS (
2     SELECT
3       ride_id,
4       rideable_type,
5       member_casual,
6       TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
7       EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
8     FROM `cyclist-394306.previous_12_month_data.2022-07`
9   )
10
11  SELECT
12    member_casual,
13    rideable_type,
14    COUNT(rideable_type) AS Total
15  FROM
16    temp_table
17  GROUP BY
18    member_casual,
19    rideable_type
20  ORDER BY
21    member_casual,
22    rideable_type
23
```

## Query results

⬇ SAV

| JOB INFORMATION | RESULTS | JSON | EXECUTION DETAILS | EXECUT |
|---|---|---|---|---|

| Row | member_casual ▼ | rideable_type ▼ | Total ▼ |
|---|---|---|---|
| 1 | casual | classic_bike | 156095 |
| 2 | casual | docked_bike | 31055 |
| 3 | casual | electric_bike | 218905 |
| 4 | member | classic_bike | 217078 |
| 5 | member | electric_bike | 200355 |

**Total rides by each member_casual on weekends.**

Now let's calculate the total rides count by each member_casual on weekends.

```sql
1   WITH temp_table AS (
2     SELECT
3       ride_id,
4       rideable_type,
5       member_casual,
6       TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
7       EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
8     FROM `cyclist-394306.previous_12_month_data.2022-07`
9   )
10
11  SELECT
12    member_casual,
13    SUM(CASE WHEN day_of_week = 1 OR day_of_week = 7 THEN 1 ELSE 0 END) AS weekend_rides_count
14  FROM
15    temp_table
16  GROUP BY
17    member_casual;
18
```

## Query results

| | JOB INFORMATION | RESULTS | JSON | EXECU |
|---|---|---|---|---|

| Row | member_casual ▼ | weekend_rides_coun |
|---|---|---|
| 1 | member | 127649 |
| 2 | casual | 173481 |

## Total No. of rides on each weekday

Now let us calculate the number of rides on each weekday and find out what weekday has the highest number of rides.

```sql
1   WITH temp_table AS (
2     SELECT
3       ride_id,
4       rideable_type,
5       member_casual,
6       TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
7       EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
8     FROM `cyclist-394306.previous_12_month_data.2022-07`
9   )
10
11  SELECT
12    day_of_week,
13    COUNT(*)  AS frequency_of_weekday
14  FROM
15    temp_table
16  GROUP BY
17    day_of_week
18  ORDER BY
19    frequency_of_weekday DESC
20
```

## Query results

| | JOB INFORMATION | **RESULTS** | JSO |
|---|---|---|---|

| Row | day_of_week ▼ | frequency_of_weekd |
|---|---|---|
| 1 | 7 | 164099 |
| 2 | 1 | 137031 |
| 3 | 6 | 118149 |
| 4 | 5 | 108948 |
| 5 | 4 | 102461 |
| 6 | 3 | 98979 |
| 7 | 2 | 93821 |

Now we can conclude that weekday 7 or Saturday has the highest number of rides.

# Working on Year data

### Combining each month data.
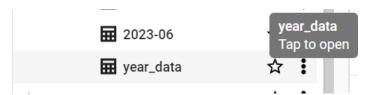
Now we will work on yearly data.
But we only have data in separate tables. So let's create a new table for year data.

```
 1    CREATE TABLE `cyclist-394306.previous_12_month_data.year_data` AS
 2    WITH temp_table AS (
 3      SELECT * FROM `cyclist-394306.previous_12_month_data.2022-07`
 4      UNION ALL
 5      SELECT * FROM `cyclist-394306.previous_12_month_data.2022-08`
 6      UNION ALL
 7      SELECT * FROM `cyclist-394306.previous_12_month_data.2022-09`
 8      UNION ALL
 9      SELECT * FROM `cyclist-394306.previous_12_month_data.2022-10`
10      UNION ALL
11      SELECT * FROM `cyclist-394306.previous_12_month_data.2022-11`
12      UNION ALL
13      SELECT * FROM `cyclist-394306.previous_12_month_data.2022-12`
14      UNION ALL
15      SELECT * FROM `cyclist-394306.previous_12_month_data.2023-01`
16      UNION ALL
17      SELECT * FROM `cyclist-394306.previous_12_month_data.2023-02`
18      UNION ALL
19      SELECT * FROM `cyclist-394306.previous_12_month_data.2023-03`
20      UNION ALL
21      SELECT * FROM `cyclist-394306.previous_12_month_data.2023-04`
22      UNION ALL
23      SELECT * FROM `cyclist-394306.previous_12_month_data.2023-05`
24      UNION ALL
25      SELECT * FROM `cyclist-394306.previous_12_month_data.2023-06`
26    )
27
28    SELECT * FROM temp_table
```

The UNION ALL will combine all rows from the different tables.
And the SELECT * FROM temp_table with flash the data then it will be stored in the table.

Then table will be created at the specified location



This will make query faster. Now we can use this data for yearly data analysis.

## Total rides each month

Let's analyze total rides done in each month.

```
1   SELECT
2     EXTRACT(YEAR FROM started_at) AS year,
3     EXTRACT(MONTH FROM started_at) AS month,
4     COUNT(*) AS total_rides
5   FROM
6     `cyclist-394306.previous_12_month_data.year_data`
7
8   GROUP BY
9     year, month
10  ORDER BY
11    year, month
```

## Query results

| Row | year | month | total_rides |
|-----|------|-------|-------------|
| 1 | 2022 | 7 | 823488 |
| 2 | 2022 | 8 | 785932 |
| 3 | 2022 | 9 | 701339 |
| 4 | 2022 | 10 | 558685 |
| 5 | 2022 | 11 | 337735 |
| 6 | 2022 | 12 | 181806 |
| 7 | 2023 | 1 | 190301 |
| 8 | 2023 | 2 | 190445 |
| 9 | 2023 | 3 | 258678 |
| 10 | 2023 | 4 | 426590 |
| 11 | 2023 | 5 | 604827 |
| 12 | 2023 | 6 | 719618 |

Let's order the above results in descending to know which month has the highest number of rides.

```
1   SELECT
2     EXTRACT(YEAR FROM started_at) AS year,
3     EXTRACT(MONTH FROM started_at) AS month,
4     COUNT(*) AS total_rides
5   FROM
6     `cyclist-394306.previous_12_month_data.year_data`
7
8   GROUP BY
9     year, month
10  ORDER BY
11    total_rides DESC
```

## Query results

| JOB INFORMATION | RESULTS | JSON | EXECUTION D |
|---|---|---|---|

| Row | year ▼ | month ▼ | total_rides ▼ |
|---|---|---|---|
| 1 | 2022 | 7 | 823488 |
| 2 | 2022 | 8 | 785932 |
| 3 | 2023 | 6 | 719618 |
| 4 | 2022 | 9 | 701339 |
| 5 | 2023 | 5 | 604827 |
| 6 | 2022 | 10 | 558685 |
| 7 | 2023 | 4 | 426590 |
| 8 | 2022 | 11 | 337735 |
| 9 | 2023 | 3 | 258678 |
| 10 | 2023 | 2 | 190445 |
| 11 | 2023 | 1 | 190301 |
| 12 | 2022 | 12 | 181806 |

This data concludes that July, 2022 has highest number of rides.
While December, 2022 accounts for lowest number of rides.

**Total no. of each member_casual in an year data**

Now let's count the total members or casual customer in an year

```
1   SELECT
2     member_casual,
3     COUNT(*) AS total
4   FROM
5     `cyclist-394306.previous_12_month_data.year_data`
6
7   GROUP BY
8     member_casual
9
```

## Query results

| | JOB INFORMATION | RESULTS | JSON | EXE |
|---|---|---|---|---|

| Row | member_casual ▼ | total ▼ |
|---|---|---|
| 1 | member | 3535192 |
| 2 | casual | 2244252 |

This data shows that there are more members riding bikes than the casual members.

**Adding two more columns to perform more descriptive analysis**

w

Now let's create two more columns **ride_length** and **day_of_week** to do more descriptive analysis

```
1  SELECT
2    ride_id,
3    rideable_type,
4    member_casual,
5    TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
6    EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
7  FROM
8    `cyclist-394306.previous_12_month_data.year_data`
9
```

## Query results

⬇ SAVE RESULTS ▾          📈 EXPLORE D

| JOB INFORMATION | RESULTS | JSON | EXECUTION DETAILS | EXECUTION GRAPH |

| Row | ride_id ▾ | rideable_type ▾ | member_casual ▾ | ride_length ▾ |
|---|---|---|---|---|
| 1 | FBE9EDA91114E989 | electric_bike | member | 00:03:14 |
| 2 | BDA67BDB30E46204 | electric_bike | casual | 00:15:19 |
| 3 | 3EF41562883C9EE4 | electric_bike | member | 00:14:42 |
| 4 | A701E2DFF37F2657 | electric_bike | casual | 00:04:49 |
| 5 | B6A875A9BE72EA71 | electric_bike | member | 00:31:01 |
| 6 | EBA20773A363823A | electric_bike | casual | 00:29:10 |
| 7 | 31B177DB41A1F426 | electric_bike | casual | 00:11:33 |
| 8 | 009184CFF2D9B1AC | electric_bike | casual | 00:09:25 |
| 9 | 752831C94187383D | electric_bike | member | 00:09:37 |
| 10 | 5C1221590EABEC13 | electric_bike | casual | 00:08:29 |
| 11 | 65A08F0C6F458485 | electric_bike | member | 00:05:17 |
| 12 | 699F3FED065A4B9A | electric_bike | member | 00:00:59 |
| 13 | D482BDAD2B8076E7 | electric_bike | casual | 00:13:15 |

Now we can make this addition of two more variable as temp_table so we don't need to write it again and again for our analysis.

**Most preferred day for riding**

Now let's find out which day is most preferred for riding bikes in an year data.

```
1   WITH temp_table AS (
2     SELECT
3       ride_id,
4       rideable_type,
5       member_casual,
6       TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
7       EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
8     FROM
9       `cyclist-394306.previous_12_month_data.year_data`
10  )
11
12  SELECT
13    day_of_week,
14    COUNT(*) AS total_rides
15  FROM
16    temp_table
17  GROUP BY
18    day_of_week
19  ORDER BY
20    total_rides DESC
21
```

## Query results

| | JOB INFORMATION | RESULTS | JSON |
|---|---|---|---|

| Row | day_of_week | total_rides |
|---|---|---|
| 1 | 7 | 924388 |
| 2 | 6 | 865994 |
| 3 | 5 | 863962 |
| 4 | 4 | 848891 |
| 5 | 3 | 806146 |
| 6 | 1 | 739351 |
| 7 | 2 | 730712 |

This data reflects that Saturday is the most preferred day.
Now let's find out total count of different member_casual on each day.

**Find out total member_casual rides each weekday**

```sql
1   WITH temp_table AS (
2     SELECT
3       ride_id,
4       rideable_type,
5       member_casual,
6       TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
7       EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
8     FROM
9       `cyclist-394306.previous_12_month_data.year_data`
10  )
11
12  SELECT
13    day_of_week,
14    member_casual,
15    COUNT(*) AS total
16  FROM
17    temp_table
18  GROUP BY
19    day_of_week, member_casual
20  ORDER BY
21    day_of_week, member_casual
22
```

## Query results

| JOB INFORMATION | RESULTS | JSON | EXECUTION DETAILS |
|---|---|---|---|

| Row | day_of_week ▾ | member_casual ▾ | total ▾ |
|---|---|---|---|
| 1 | 1 | casual | 351412 |
| 2 | 1 | member | 387939 |
| 3 | 2 | casual | 253088 |
| 4 | 2 | member | 477624 |
| 5 | 3 | casual | 257108 |
| 6 | 3 | member | 549038 |
| 7 | 4 | casual | 277263 |
| 8 | 4 | member | 571628 |
| 9 | 5 | casual | 298282 |
| 10 | 5 | member | 565680 |
| 11 | 6 | casual | 347112 |
| 12 | 6 | member | 518882 |
| 13 | 7 | casual | 459987 |
| 14 | 7 | member | 464401 |

**Total Ride length in hours each member has in the year**

Now let's find out how much ride_length in total each member has.

```
1   WITH temp_table AS (
2     SELECT
3       ride_id,
4       started_at,
5       ended_at,
6       rideable_type,
7       member_casual,
8       TIME(TIMESTAMP_SECONDS(TIMESTAMP_DIFF(ended_at, started_at, SECOND))) AS ride_length,
9       EXTRACT(DAYOFWEEK FROM started_at) as day_of_week
10    FROM
11      `cyclist-394306.previous_12_month_data.year_data`
12    WHERE
13      started_at<ended_at
14  )
15
16  SELECT
17    member_casual,
18    ROUND(SUM(Extract(HOUR FROM ride_length)*3600 + Extract(MINUTE FROM ride_length) * 60 + Extract(SECOND FROM
      ride_length))/3600, 2) AS total_hours
19  FROM
20    temp_table
21  GROUP BY
22    member_casual
23
```

## Query results

| | JOB INFORMATION | RESULTS | JSON | EXE |
|---|---|---|---|---|

| Row | member_casual ▼ | total_hours ▼ |
|---|---|---|
| 1 | member | 710718.82 |
| 2 | casual | 787721.39 |

This data shows that even the casual customers are less, they use the bikes the most.

**The following analysis done give a good understanding of the data and some of the patterns in the data. Now we'll switch to Tableau for visualization.**