SUMIT SHAMLAL CHAURE

# Operation Analytics and Investigating Metric Spike
## Trainity Project 2 – Advanced SQL

## Introduction

Operational Analytics is a crucial process that involves analyzing a company's end-to-end operations. This analysis helps identify areas for improvement within the company. As a Data Analyst, you'll work closely with various teams, such as operations, support, and marketing, helping them derive valuable insights from the data they collect.

# Description

Operational Analytics is a crucial process that involves analyzing a company's end-to-end operations. This analysis helps identify areas for improvement within the company. As a Data Analyst, you'll work closely with various teams, such as operations, support, and marketing, helping them derive valuable insights from the data they collect.

One of the key aspects of Operational Analytics is investigating metric spikes. This involves understanding and explaining sudden changes in key metrics, such as a dip in daily user engagement or a drop in sales. As a Data Analyst, you'll need to answer these questions daily, making it crucial to understand how to investigate these metric spikes.

In this project, you'll take on the role of a Lead Data Analyst at a company like Microsoft. You'll be provided with various datasets and tables, and your task will be to derive insights from this data to answer questions posed by different departments within the company. Your goal is to use your advanced SQL skills to analyze the data and provide valuable insights that can help improve the company's operations and understand sudden changes in key metrics.

There are two Case Study in these assignment :

1. Job Data Analysis
2. Investigating Metric Spike

# Requirements -

## 1) Project Description :

The aim of the project is to find the user analytics involved in operations of a job activity pertaining to job reviews done and the time spent on job. The investigation metric spike helps us to identify the sudden changes in key metrics like user engagement activities, drop in sales, productivity etc.

## 2) Approach :

Firstly I created the database and tables from the provided CSV data – for the first case study I manually wrote the query to create and insert data but the case study 2 required to import large data from csv so we made a blank table and then used a special query called 'Inline Load File' to load the data & then change the table to adjust the datatypes where needed. Then using basic query we got the required results for the operations. For the visualization part I made the csv data into excel tables and generated the screenshots of both the query and table for results.For some tables we need to do operations like table alterations too.

The analysis based on query has insights at the bottom of the screengrabs to let the mentor understand the aim of each analysis.

## 3) Tech-Stack Used :

MySQL – For the main query and results part I have relied on MySQL version 8.1.0.

Excel – The tables were generated from the csv data of the query results of sql and imported in excel.

Word – The report is written in word/docx format using MS Word and then exported to pdf.

Drive – To upload all the essential files attached in the report for reference and for pdf upload.

## 4) Insights :

The summary for each query is given with the screenshot but to summarize the overall thing I came to the conclusion that i learnt about the business scenario where multiple tables contains related data which when combined together brings in deeper and more meaningful insights which are necessary to let the company grow and also helps the data analysts to generate insights that can help them predict the future activities and engagements.

## 5) Result :

The project has helped me get the gist of real life examples of database operations from table creation, data insertion management to the complexity involved in connecting various tables to store and get relevant data. The analysis helps to know about the parameter and inputs required to track the activities

. I learnt about joins and select statements more practically with the activity relating to functions which helped me more efficiently do querying operations.

# Main Querying Part & Analysis -

**SQL Task**

_Case Study 1: Job Data Analysis_

    A. Jobs Reviewed Over Time

    B. Throughput Analysis

    C. Language Share Analysis

    D. Duplicate Rows Detection

_Case Study 2: Investigating Metric Spike_

    A. Weekly User Engagement

    B. User Growth Analysis

    C. Weekly Retention Analysis

    D. Weekly Engagement Per Device

    E. Email Engagement Analysis

## A) Case Study 1: Job Data Analysis:

**1.Jobs Reviewed Over Time:** Calculate the number of jobs reviewed per hour for each day in November 2020.

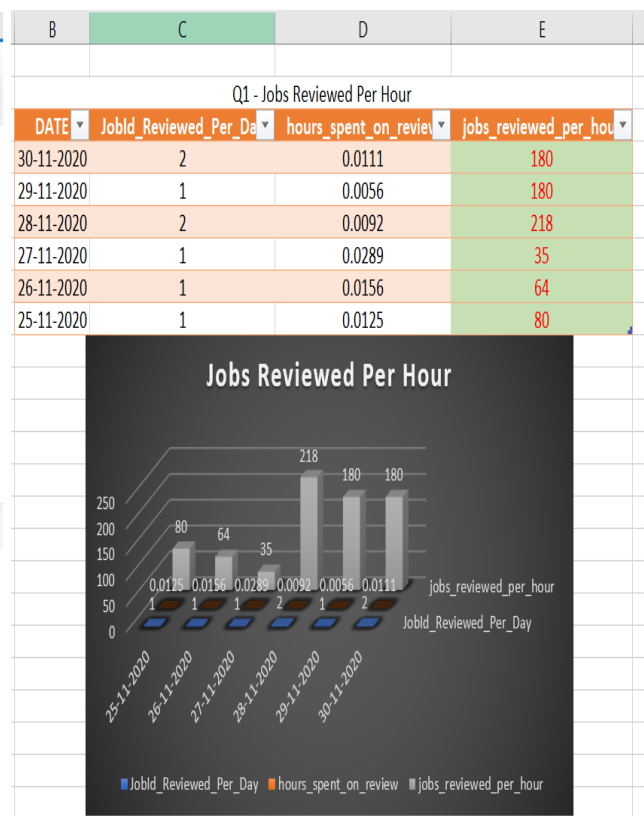*Your Task*: Write an SQL query to calculate the number of jobs reviewed per hour for each day in November 2020.
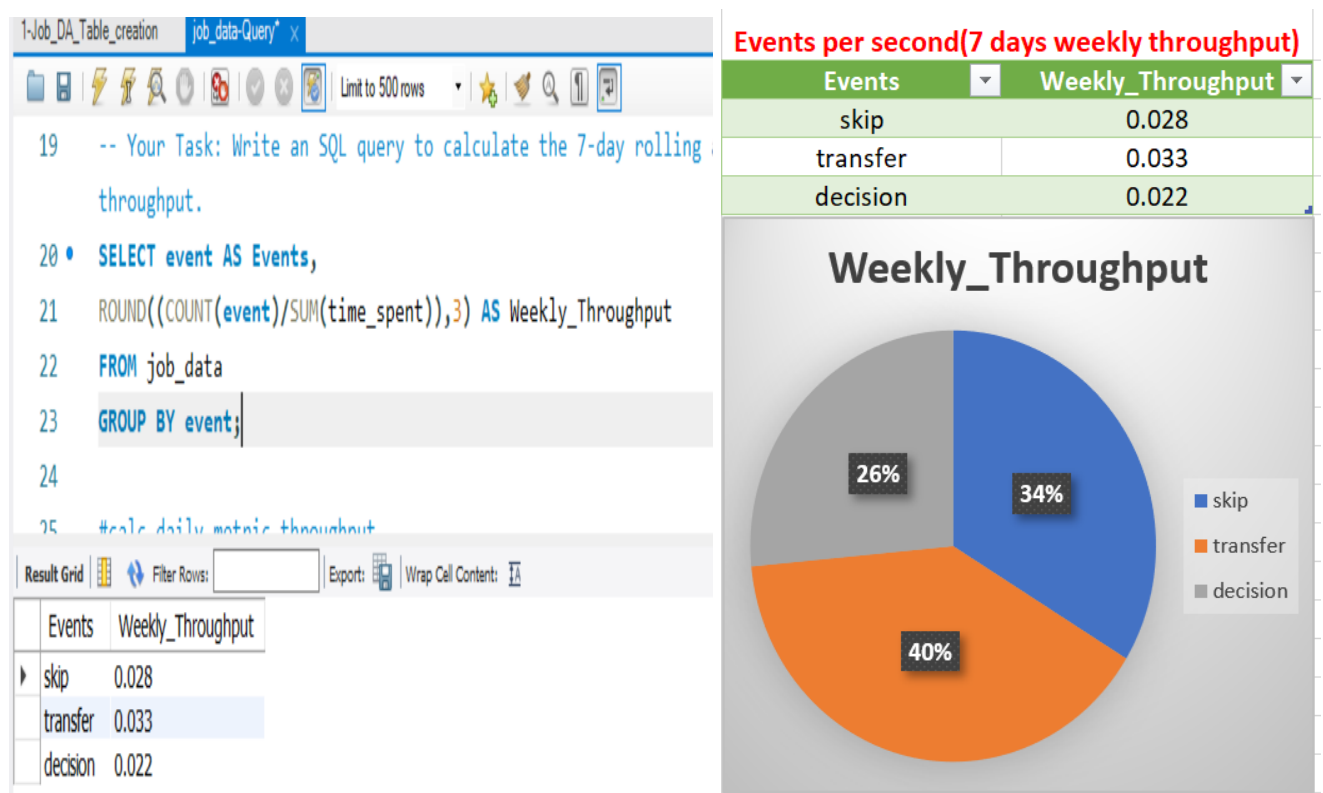
### Query



**Insights :** From the above query we have come to know that if we take the hours spent every hour on a job review then do the sum of time spent on each job for the day we get the result of our question. The bar chart shows us the distribution of tasks per day and the number of task done per hour according to the time taken for each review (*The highlighted column is our desired output*).

**2.Throughput Analysis:** Calculate the 7-day rolling average of throughput (number of events per second) .

*Your Task*: Write an SQL query to calculate the 7-day rolling average of throughput. Additionally, explain whether you prefer using the daily metric or the 7-day rolling average for throughput, and why.

***Query***



***Img -*** *Weekly Throughput shows a smooth distribution of events done for the week and distribution is more readable. Weekly metric helps us to understand which events holds account for the most work/reviews.*
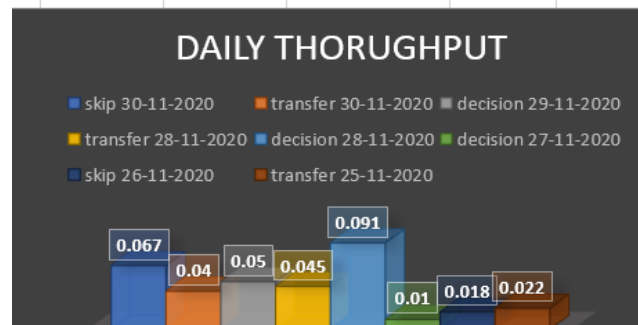
```
25    #Daily Metric Throughput(Daywise & event grouping)
26 •  select event,ds as date, round((count(event)/sum(
      time_spent)),3) as daily_metric
27    from job_data group by event,date;
28
```

| event | date | daily_metric |
|---|---|---|
| ▶ skip | 2020-11-30 | 0.067 |
| transfer | 2020-11-30 | 0.040 |
| decision | 2020-11-29 | 0.050 |
| transfer | 2020-11-28 | 0.045 |
| decision | 2020-11-28 | 0.091 |
| decision | 2020-11-27 | 0.010 |
| skip | 2020-11-26 | 0.018 |
| transfer | 2020-11-25 | 0.022 |

Events per second(Daily throughput)

| event | date | daily_metric |
|---|---|---|
| skip | ######### | 0.067 |
| transfer | ######### | 0.04 |
| decision | ######### | 0.05 |
| transfer | ######### | 0.045 |
| decision | ######### | 0.091 |
| decision | ######### | 0.01 |
| skip | ######### | 0.018 |
| transfer | ######### | 0.022 |

**DAILY THORUGHPUT**

■ skip 30-11-2020  ■ transfer 30-11-2020  ■ decision 29-11-2020
■ transfer 28-11-2020  ■ decision 28-11-2020  ■ decision 27-11-2020
■ skip 26-11-2020  ■ transfer 25-11-2020

*Img - The daily throughput based on the events has distributed the work over daily basis which is not that readable for the insights if we look at the event wise distributions too.*

**Insights :** From the above charts and values we can infer the weekly throughput of the various events based on a weekly metric over events and daily metric over the events and date distribution.

*Throughput is the number of events happening per second.* For calculating the 7-day rolling average, we used the number of events per second for each day & divided it by the week to get the average needed.From the insights I have inferred that it is better to take a 7 day rolling average over the daily metric as the daily data shows fluctuations and does not give estimated data of events average but the weekly throughput helps us to understand the events handled on seconds basis for the week much more efficiently.
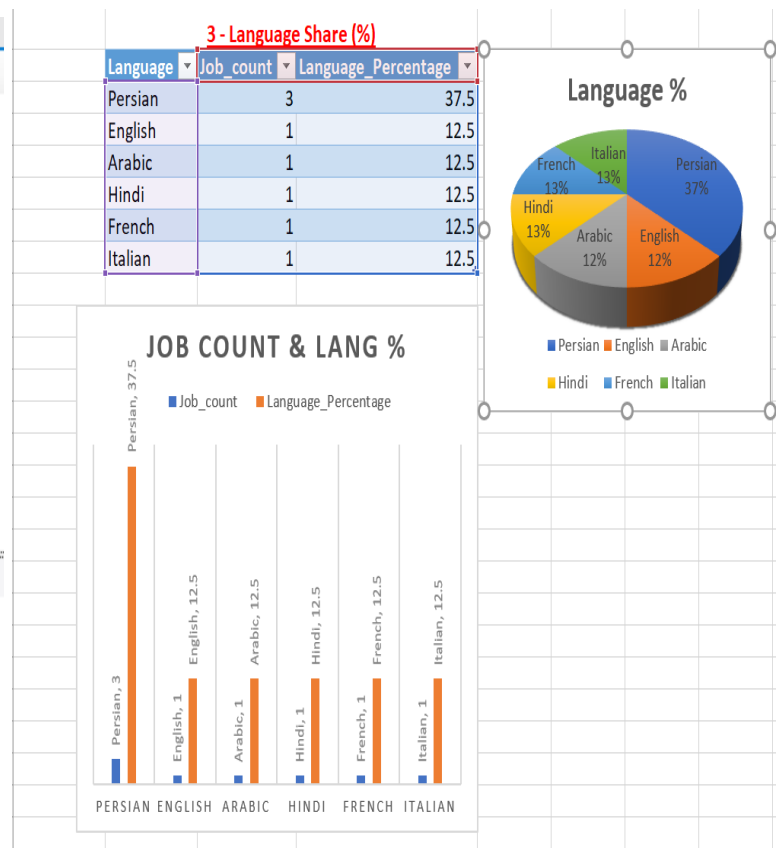
**3. Language Share Analysis** : Calculate the percentage share of each language in the last 30 days .

*Your Task*: Write an SQL query to calculate the percentage share of each language over the last 30 days .
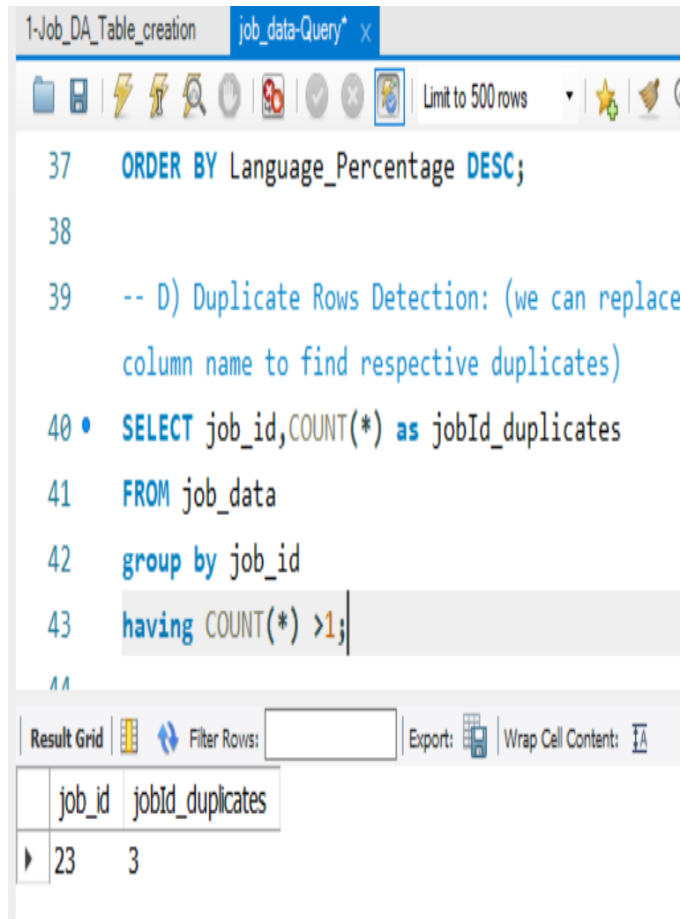
***Query***



***Img - Language Percentage Shares***

*Insight :* *The above insight shows us that* **Persian** *language has the highest percentage share of all the 8 languages in our job id followed by all the other language having 1 counts from total count.*

**4.Duplicate Rows Detection** : Identify duplicate rows in the data.

*Your Task*: Write an SQL query to display duplicate rows from the job_data table.

***Query***



*Img -* *1st image considers job_id column while 2nd image is considering every column*

*Insights :* *If we want to find the duplicate rows we can do a count and use the column name on which we want to find the duplicates in above query if we replace the job_id with say event or with actor_id it will query for any duplicates present in that column as we have not been given any specific primary key or column we can check the different duplicates as needed.In 2nd query we tried checking for whole*

*columns if they are repeated so we don't get any results as any specific row does not contain a completely duplicate row.*

---

### B) Case Study 2: Investigating Metric Spike:

**1.Weekly User Engagement:** Measure the activeness of users on a weekly basis.

*Your Task*: Write an SQL query to calculate the weekly user engagement.

*Query*



*Insights :* *From the above query we get the user engagement that is unique user visit according to the week numbers.This demograph helps us to gather information about the events density and user uniqueness at the events.*

**2.User Growth Analysis:** Analyze the growth of users over time for a product.

*Your Task*: Write an SQL query to calculate the user growth for the product.

***Query***

```
15      -- 2.User Growth Analysis:
16      -- Your Task:  Write an SQL query to calculate the user growth for the product.
17 •    select week_num, year_num,
18      sum(active_users) over (order by week_num, year_num
19      rows between unbounded preceding and current row) as cumulative_sum
20      from (
21      select extract(week from activated_at) as week_num,
22      extract(year from activated_at) as year_num,
23      count(distinct user_id) as active_users from users
24      where state= "active"
25      group by year_num, week_num
26      order by year_num, week_num) as alias;
27
28
```

*Insights :* *From the above query we get to find out about the number of users on the platform and the growing numbers according to week and year wise when we see at the active state of users column.*

*This data would help us determine strategy for future growth planning and ad or other revenue generation activities.*

**3.Weekly Retention Analysis:** Analyze the retention of users on a weekly basis after signing up for a product.

*Your Task*: Write an SQL query to calculate the weekly retention of users based on their sign-up cohort..

*Query*

```
29    -- 3.Weekly Retention Analysis:

30    -- Your Task:  Write an SQL query to calculate the weekly retentio

31 •  select

32    extract(week from occurred_at) as Weeks,

33    count(distinct user_id) as No_of_RetainedUsers from events

34    where event_type="signup_flow" and event_name="complete_signup"

35    group by weeks order by Weeks;

36
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: 𝐈𝐀

| Weeks | No_of_Retainedusers |
|---|---|
| 17 | 72 |
| 18 | 163 |
| 19 | 185 |
| 20 | 176 |

Result 7 ✕

**3- Retained Users**

| Weeks | No_of_Retainedusers |
|---|---|
| 17 | 72 |
| 18 | 163 |
| 19 | 185 |
| 20 | 176 |
| 21 | 183 |
| 22 | 196 |
| 23 | 196 |
| 24 | 229 |
| 25 | 207 |
| 26 | 201 |
| 27 | 222 |
| 28 | 215 |
| 29 | 221 |
| 30 | 238 |
| 31 | 193 |
| 32 | 245 |
| 33 | 261 |
| 34 | 259 |
| 35 | 18 |

Retained User Week Wise

*Insight ;* *From above query we have found the number of users who have signed up into our cohort program by completing the signup flow and complete the signup process.*

**4.Weekly Engagement Per Device:**  Measure the activeness of users on a weekly basis per device.

*Your Task*:  Write an SQL query to calculate the weekly engagement per device.

***Query***

```
37    -- 4.Weekly Engagement Per Device:
38    -- Your Task:  Write an SQL query to calculate the
39 •  select * from events;
40 •  select device AS Device,
41    extract(week from occurred_at) as Weeks,
42    count(distinct user_id) as No_Of_Users
43    from events
44    where event_type="engagement"
45    group by Device, Weeks order by Weeks;
```

| Device | Weeks | No_Of_Users |
|--------|-------|-------------|
| ▶ acer aspire desktop | 17 | 9 |
| acer aspire notebook | 17 | 20 |
| amazon fire phone | 17 | 4 |
| asus chromebook | 17 | 21 |

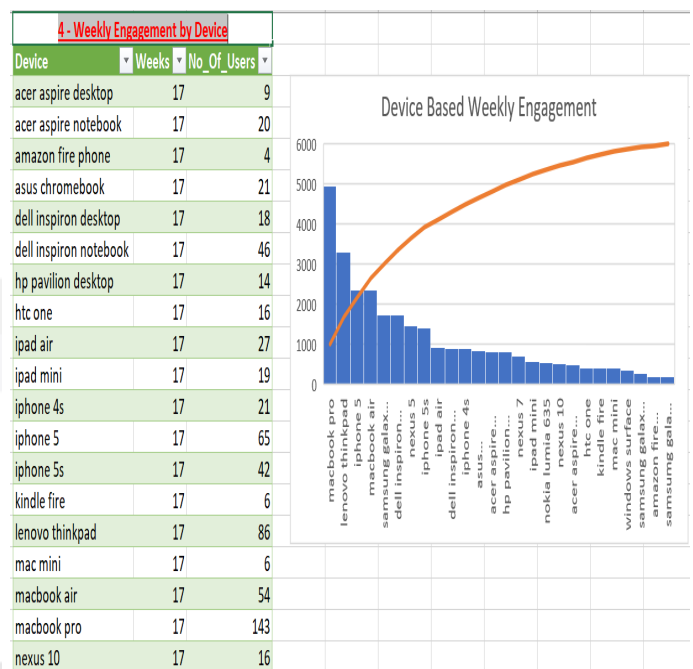| 4 - Weekly Engagement by Device | | |
|--------|-------|-------------|
| Device | Weeks | No_Of_Users |
| acer aspire desktop | 17 | 9 |
| acer aspire notebook | 17 | 20 |
| amazon fire phone | 17 | 4 |
| asus chromebook | 17 | 21 |
| dell inspiron desktop | 17 | 18 |
| dell inspiron notebook | 17 | 46 |
| hp pavilion desktop | 17 | 14 |
| htc one | 17 | 16 |
| ipad air | 17 | 27 |
| ipad mini | 17 | 19 |
| iphone 4s | 17 | 21 |
| iphone 5 | 17 | 65 |
| iphone 5s | 17 | 42 |
| kindle fire | 17 | 6 |
| lenovo thinkpad | 17 | 86 |
| mac mini | 17 | 6 |
| macbook air | 17 | 54 |
| macbook pro | 17 | 143 |
| nexus 10 | 17 | 16 |



Device Based Weekly Engagement

*Insights:* The above insight shows us the weekly demographic of the various devices used to access the events page and the number of user using the device in a certain weeks period.These data helps us to optimize our site/page/app to suite to majority of the devices/OS and thereby increasing our reach.

The data shows that majority of the users prefer to use **macbook pro.**

**5.Email Engagement Analysis:**  Analyze how users are engaging with the email service
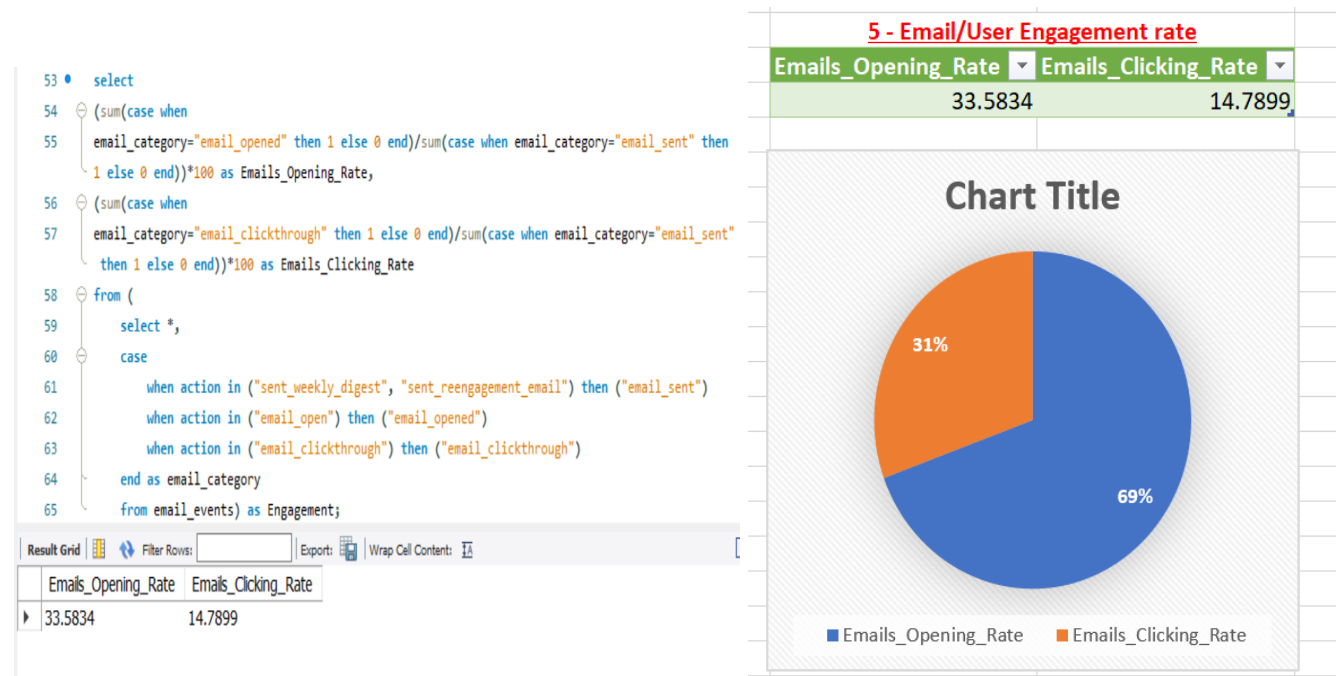
*Your Task*:  Write an SQL query to calculate the email engagement metrics.

**Query**

```
47    -- 5.Email Engagement Analysis:
48    -- Your Task:  Write an SQL query to calculate the email engagement metrics.
49    -- To get the actions which are related to the user activities
50 •  select count(action) as Action_Counter, action
51    from email_events
52    group by action;
```

| Action_Counter | action |
|---|---|
| 57267 | sent_weekly_digest |
| 20459 | email_open |
| 9010 | email_clickthrough |
| 3653 | sent_reengagement_email |

```
53 •  select
54    ⊖ (sum(case when
55      email_category="email_opened" then 1 else 0 end)/sum(case when email_category="email_sent" then
         1 else 0 end))*100 as Emails_Opening_Rate,
56    ⊖ (sum(case when
57      email_category="email_clickthrough" then 1 else 0 end)/sum(case when email_category="email_sent"
         then 1 else 0 end))*100 as Emails_Clicking_Rate
58    ⊖ from (
59        select *,
60    ⊖    case
61          when action in ("sent_weekly_digest", "sent_reengagement_email") then ("email_sent")
62          when action in ("email_open") then ("email_opened")
63          when action in ("email_clickthrough") then ("email_clickthrough")
64        end as email_category
65        from email_events) as Engagement;
```

| Emails_Opening_Rate | Emails_Clicking_Rate |
|---|---|
| 33.5834 | 14.7899 |

**5 - Email/User Engagement rate**

| Emails_Opening_Rate | Emails_Clicking_Rate |
|---|---|
| 33.5834 | 14.7899 |

**Chart Title**

31%

69%

■ Emails_Opening_Rate   ■ Emails_Clicking_Rate

**Q1 -** *Query to check various user activities involved.* **Q2 -** *Shows the actual engagement activity based on various activities we got from 1st querry.*

***Insights*** : *From the above query we can say that from the emails sent to user the email opening rate is approximately 33.6% and the user engagement with the emails is around 15%. This shows that half of the user who opens the mail dont engage into an action or engagement activity and the overall opening rate shows that the marketing team and creative department needs to give focus into bringing the customers on board and making them interact with our platform more.*

## THANK YOU