# PROJECT – III

## Operation Analytics and Investigating Metric Spike

## PROJECT DESCRIPTION

**Operation Analytics** is the analysis done for the complete end to end operations of a company. With the help of this, the company then finds the areas on which it must improve upon. You work closely with the ops team, support team, marketing team, etc and help them derive **insights** out of the data they collect.

Being one of the most important parts of a company, this kind of analysis is further used to predict the overall growth or decline of a company's fortune. It means better automation, better understanding between cross-functional teams, and more effective workflows.

Investigating **metric spike** is also an important part of operation analytics as being a Data Analyst you must be able to understand or make other teams understand questions like- Why is there a **dip** in daily engagement? Why have sales taken a dip? Etc. Questions like these must be answered daily and for that it's very important to investigate metric spike.

You (I) are (am) required to provide a detailed report for the two operations by mentioning the answers for the related questions that are asked in Case Study-1 and Case Study-2.

The Project first requires us to create a **database** and import it any SQL platform. And then it requires us to give answers to the questions that they have asked.

I would be using MYSQL in the further report.

I know this project will help me clear my doubts regarding **SQL programming** and **databases**. Having strong concepts, initially, about any topic is very important for further strengthening of the concepts and through this project I will learn to create a database and then import in MYSQL platform and then further learn to code in the form of making **queries.**

# APPROACH

I will be following the steps as directed for completing this project . Learning at every step will be my goal. I have watched all the videos pertaining to this project . I have understood the concepts at every step, and I will be applying that knowledge here to finish this project with great respect for the mentors. I have practiced few questions at my level to test my knowledge before jumping into the project.

Hopefully with all the knowledge that I have gained from the videos and my practice would help me in solving the questions and finish the project.

# TECH – STACK USED

**FOLLOWING ARE THE TECH STACKS USED BY ME FOR COMPLETEING THS PROJECT**

• MySQL Workbench 8.0 CE

• Excel

# INSIGHTS: CASE STUDY – 1 (Job Data)

Below is the structure of the table with the definition of each column that we have to work on:

**Table-1: job_data**

**job_id:** unique identifier of jobs

**actor_id:** unique identifier of actor

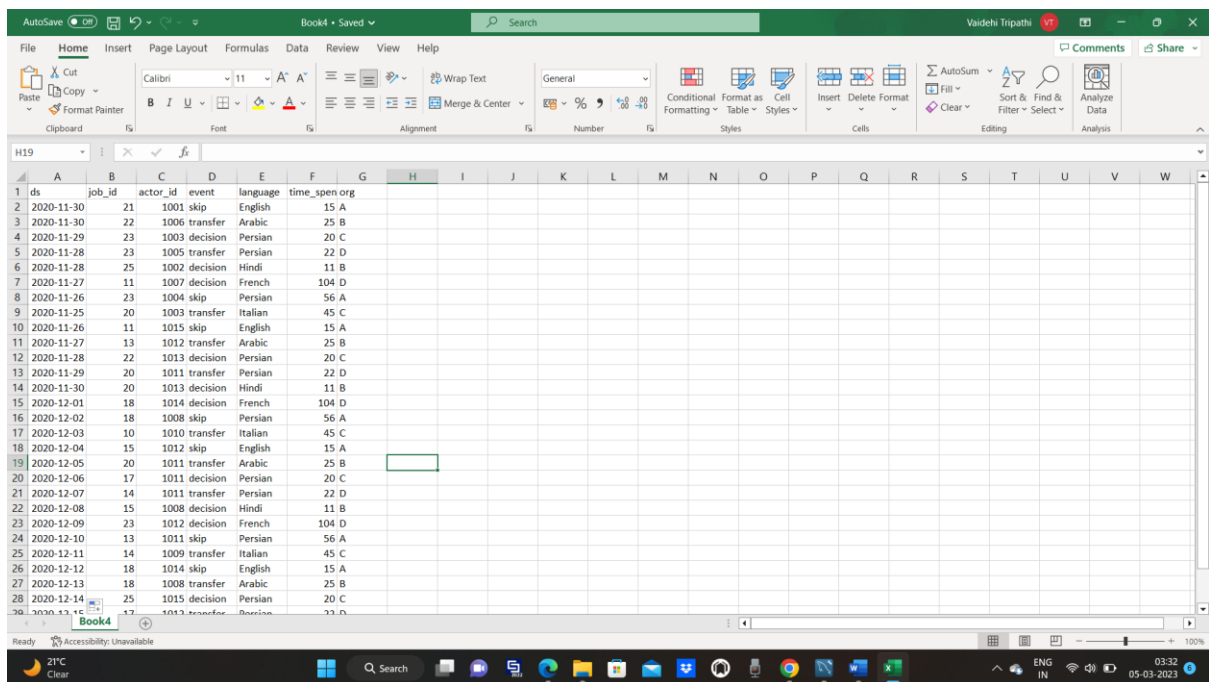**event:** decision/skip/transfer

**language:** language of the content

**time_spent:** time spent to review the job in seconds
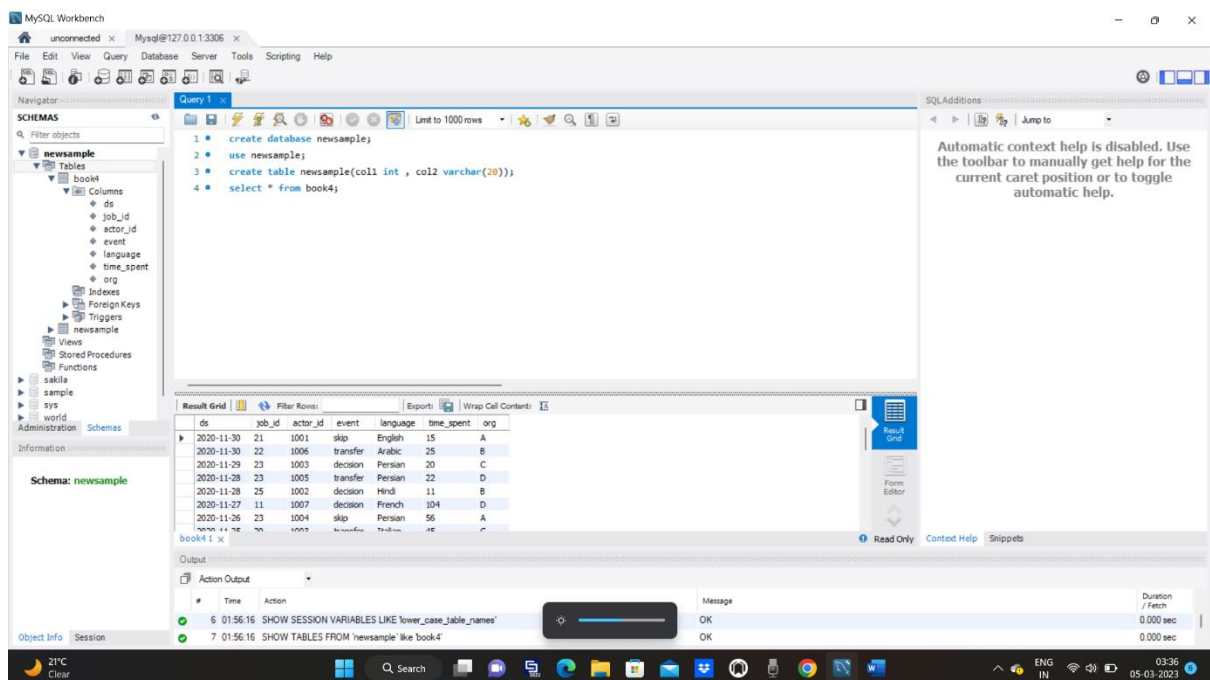
**org:** organization of the actor

**ds:** date in the yyyy/mm/dd format. It is stored in the form of text, and we use presto to run. no need for date function

## NOW FIRST I HAVE FORMED THE DATASET FOR 30 DAYS USING EXCEL

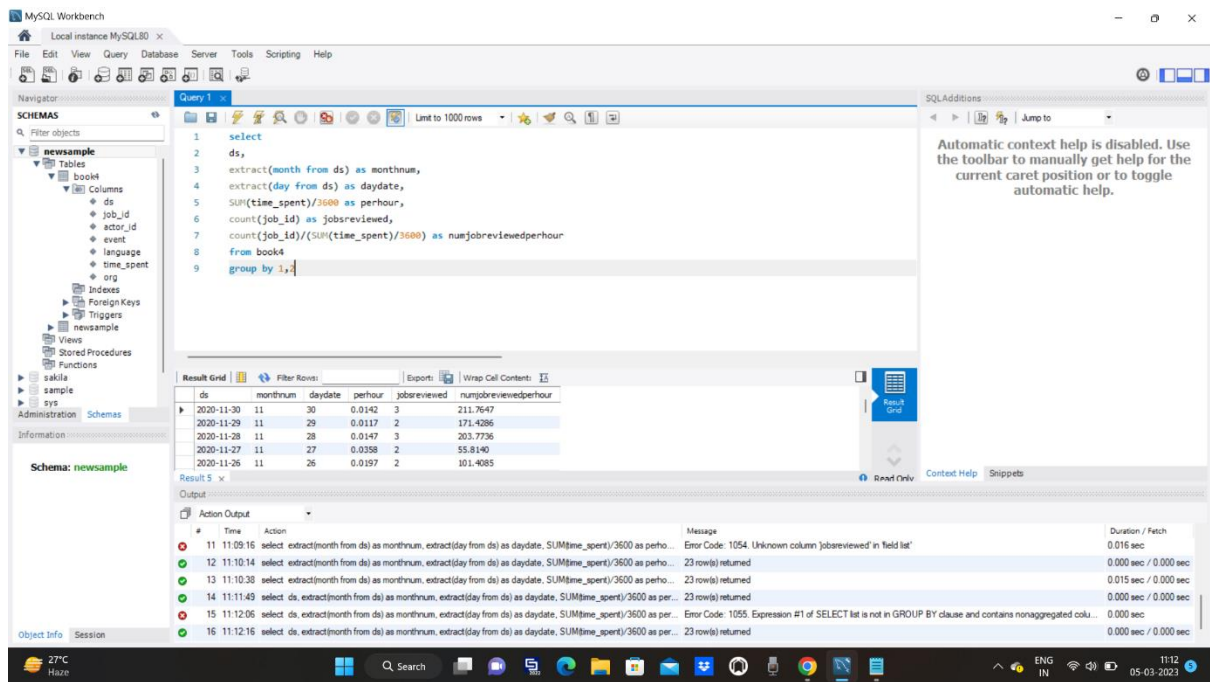## THEN I HAVE IMPORTED THIS .CSV FILE INTO MYSQL PLATFORM



## NOW THE TASKS !

( A ) **Number of jobs reviewed: Amount of jobs reviewed over time.**

**Your task: Calculate   jobs reviewed per hour per day for November 2020?**

**<u>ANSWER:</u>**

select ds,

extract(month from ds) as monthnum,

extract(day from ds) as daydate,

SUM(time_spent)/3600 as perhour,

count(job_id) as jobsreviewed,

count(job_id)/(SUM(time_spent)/3600) as numjobreviewedperhour

from book4

group by 1,2

**INSIGHT:** By using group by and different functions like sum and count, I was able to find out the answer to the question asked. Firstly I extracted day and month from the given dates and then I used sum, count, group by functions.

**(B)Throughput:** It is the no. of events happening per second. **Your task:** Let's say the above metric is called throughput. Calculate 7 day rolling average of throughput? For throughput, do you prefer daily metric or 7-day rolling and why?

**ANSWER:**

select a.*,

avg(throughput_per_day) over(order by ds desc rows between current row and 7 following) as seven_days_rolling_avg_throughput

from

(select

ds,

COUNT(event) AS num_events_per_day,

SUM(time_spent) AS time_per_sec,

COUNT(event)/(SUM(time_spent)) AS throughput_per_day

FROM book4

GROUP BY

1)a



**INSIGHT:** Here I have made use of the concept of subquery and in that subquery, I have calculated throughput per day by using sum, count, group by functions. Then in the outer query I calculated the average for 7 rolling days.

**(C)Percentage share of each language:** Share of each language for different contents.
**Your task:** Calculate the percentage share of each language in the last 30 days?
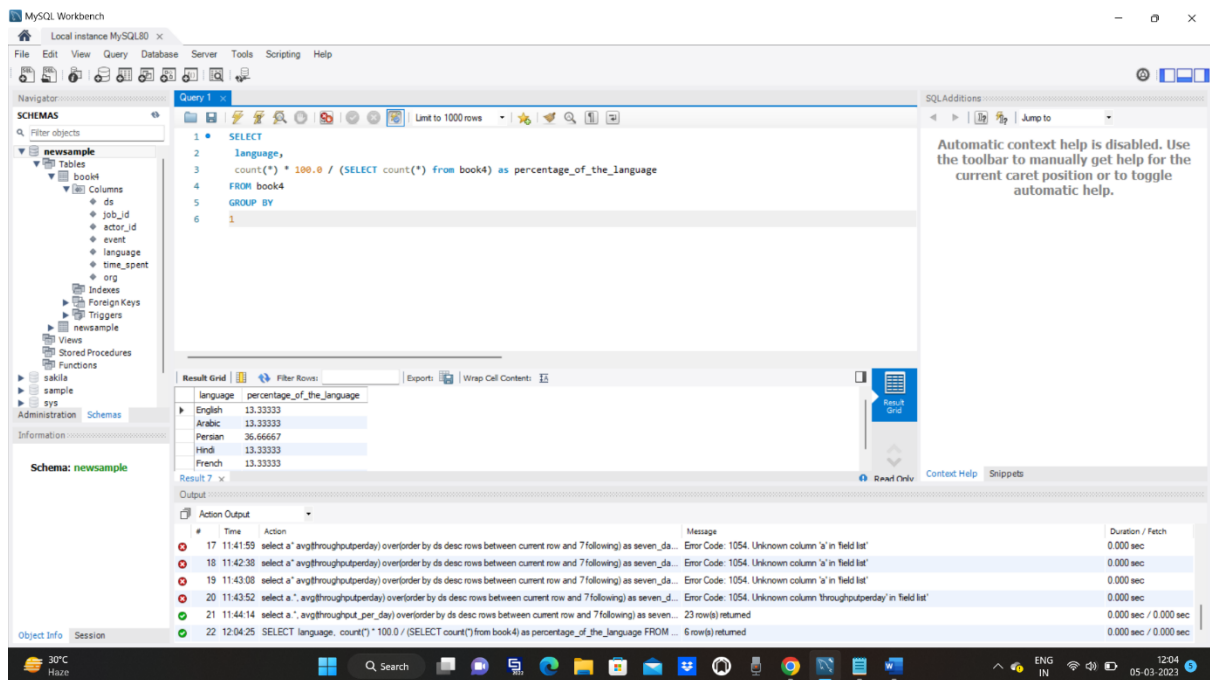
**ANSWER:**

SELECT

 language,

 count(*) * 100.0 / (SELECT count(*) from book4) as percentage_of_the_language

FROM book4

GROUP BY

1

**INSIGHT:** Here I have first grouped the languages using group by function, then I have count the total number of times the particular language is used divided by total number of languages.

**(D)Duplicate rows:** Rows that have the same value present in them. **Your task:** Let's say you see some duplicate rows in the data. How will you display duplicates from the table?

## ANSWER:

SELECT *

FROM book4

GROUP BY ds , job_id, actor_id, event,language, time_spent,org

HAVING COUNT(ds) >1 AND  COUNT(job_id) >1 AND COUNT(actor_id) >1 AND COUNT(event) >1 AND COUNT(language) >1 AND COUNT(time_spent) >1 AND COUNT(org) >1

**INSIGHT**: Here I have used group by to group all the same dates , events, languages, time_spent,actor_id,job_id etc. Rows having all the data same are to be shown, since my dataset does not contains any such row therefore it is showing nothing inside the dataset , which is fine.

# CASE STUDY – 2 (Investigating metric spike)

The structure of the table with the definition of each column that you must work on is present in the project image

**Table-1:** users

This table includes one row per user, with descriptive information about that user's account.

**Table-2:** events

This table includes one row per event, where an event is an action that a user has taken. These events include login events, messaging events, search events, events logged as users progress through a signup funnel, events around received emails.

**Table-3:** email_events

This table contains events specific to the sending of emails. It is similar in structure to the events table above.

Use the dataset then answer the questions that follows:

## TABLES:

https://drive.google.com/drive/folders/1bB-uqONISA6wiI1hw1LzISpe0-kHg0Nx

## NOW THE TASKS!

**(A) User Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service.
**Your task:** Calculate the weekly user engagement?

**ANSWER:**
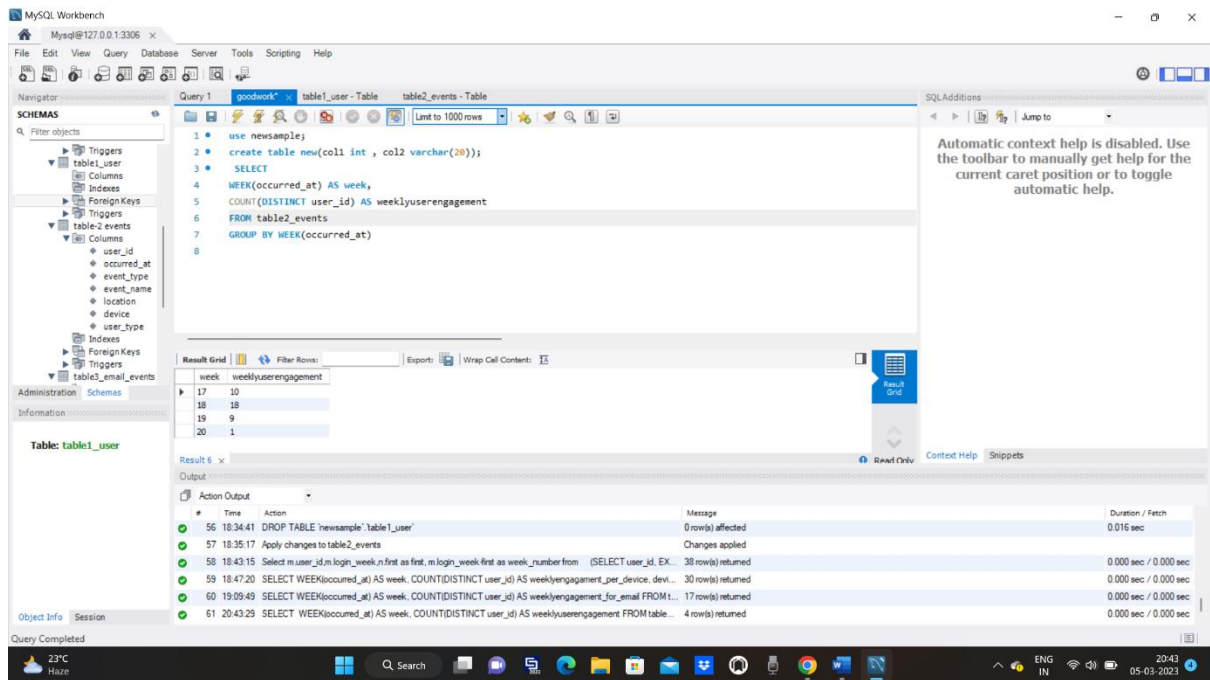
SELECT

WEEK(occurred_at) AS week,

COUNT(DISTINCT user_id) AS weeklyuserengagement

FROM table2_events

GROUP BY WEEK(occurred_at)

**_(B)_User Growth:** Amount of users growing over time for a product.

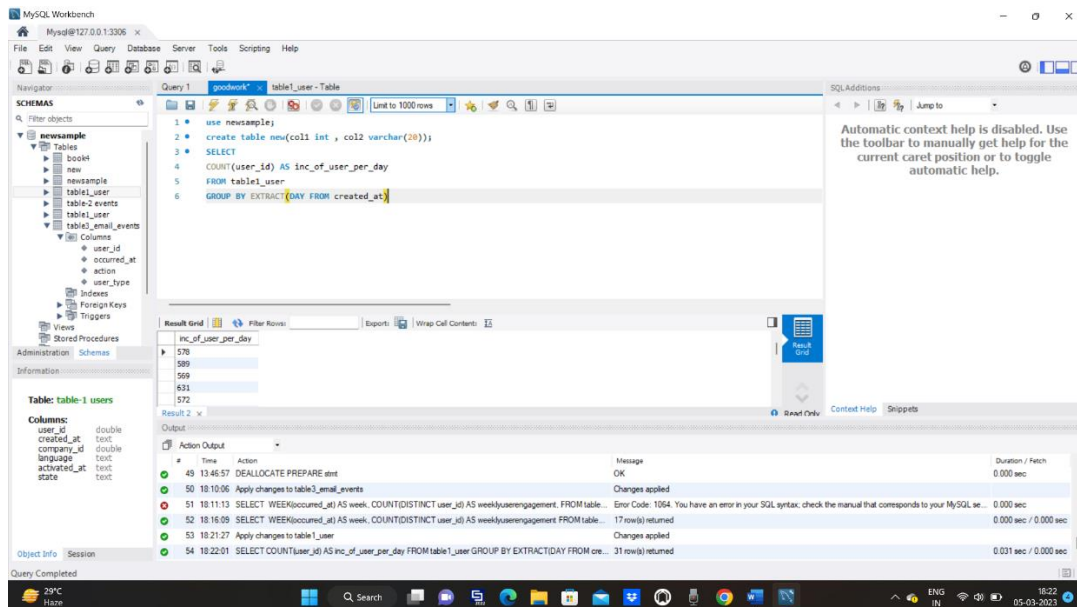**Your task:** Calculate the user growth for product?

## ANSWER:

SELECT

COUNT(user_id) AS inc_of_user_per_day

FROM table1_user

GROUP BY EXTRACT(DAY FROM created_at)

THIS RESULT SHOWS THE NUMBER OF USES THAT JOINED THE PRODUCT ON DAILY BASIS.

**(C)Weekly Retention:** Users getting retained weekly after signing-up for a product.

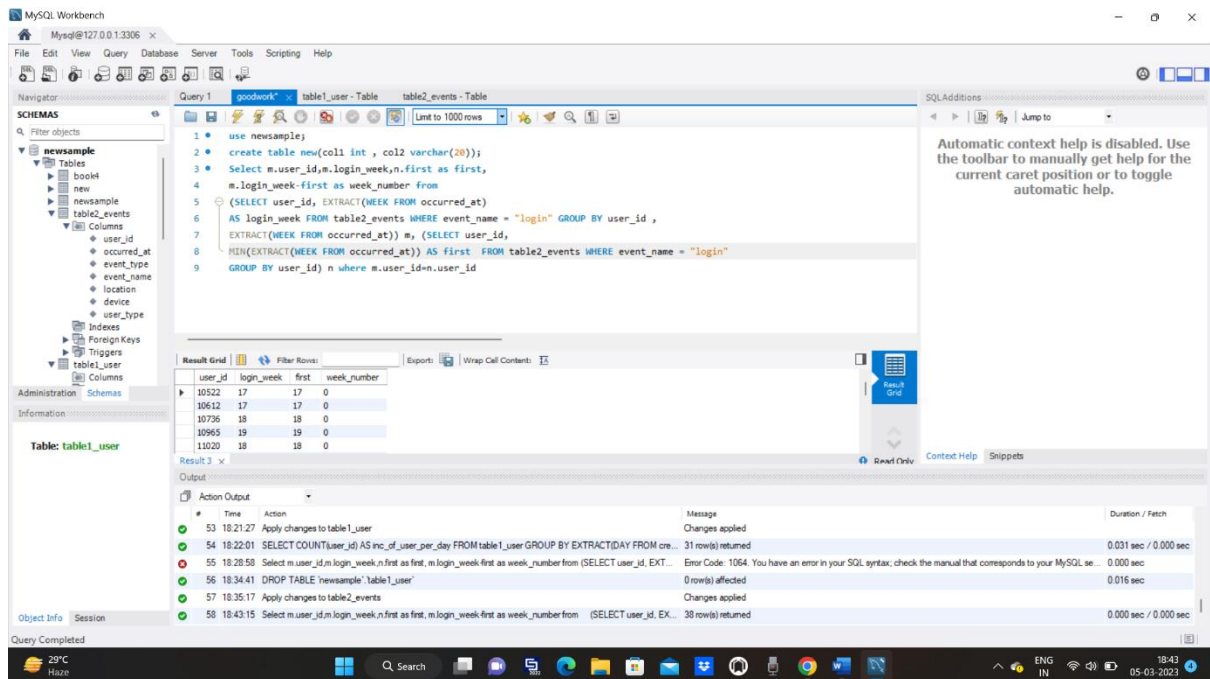**Your task:** Calculate the weekly retention of users-sign up cohort?

## ANSWER:

Select m.user_id,m.login_week,n.first as first,

m.login_week-first as week_number from

(SELECT user_id, EXTRACT(WEEK FROM occurred_at)

AS login_week FROM table2_events WHERE event_name = "login" GROUP BY user_id ,

EXTRACT(WEEK FROM occurred_at)) m, (SELECT user_id,

MIN(EXTRACT(WEEK FROM occurred_at)) AS first  FROM table2_events WHERE event_name = "login"

GROUP BY user_id) n where m.user_id=n.user_id



The results show excellent user retention and that most of the users returned to login in the same week with a very exceptions that returned the next week.

**(D)Weekly Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly. **Your task:** Calculate the weekly engagement per device?
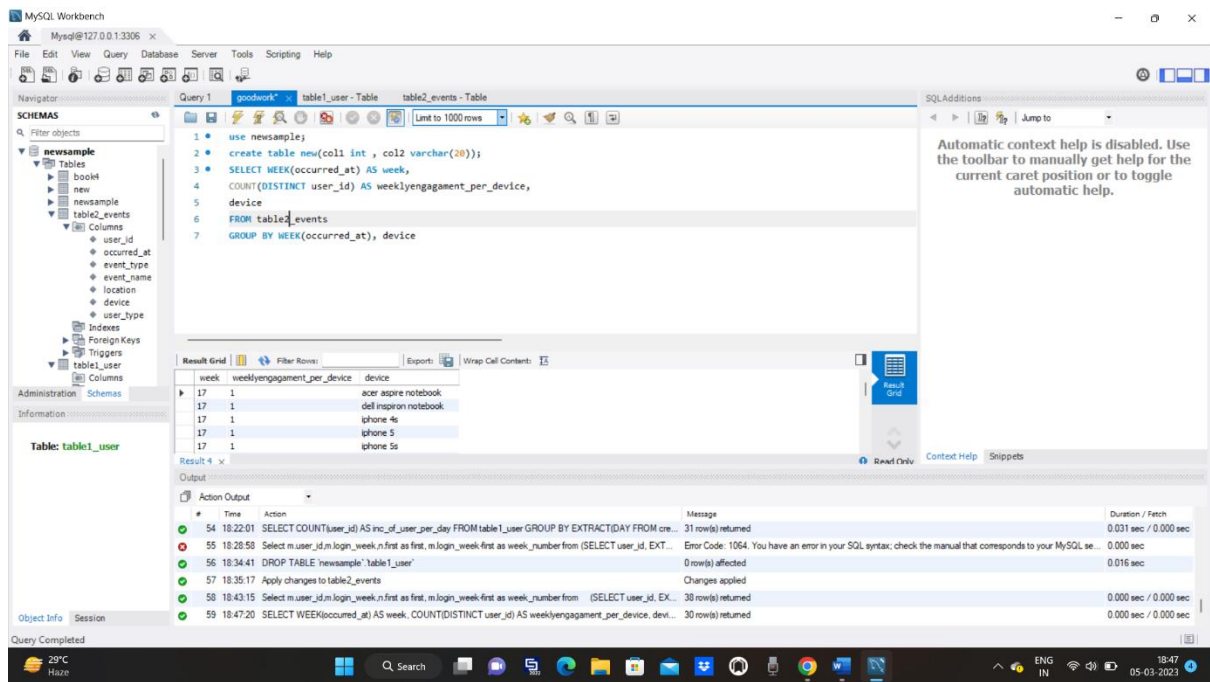
## ANSWER:

SELECT WEEK(occurred_at) AS week,

COUNT(DISTINCT user_id) AS weeklyengagament_per_device,

device

FROM table2_events

GROUP BY WEEK(occurred_at), device

**(E)Email Engagement: Users engaging with the email service.**
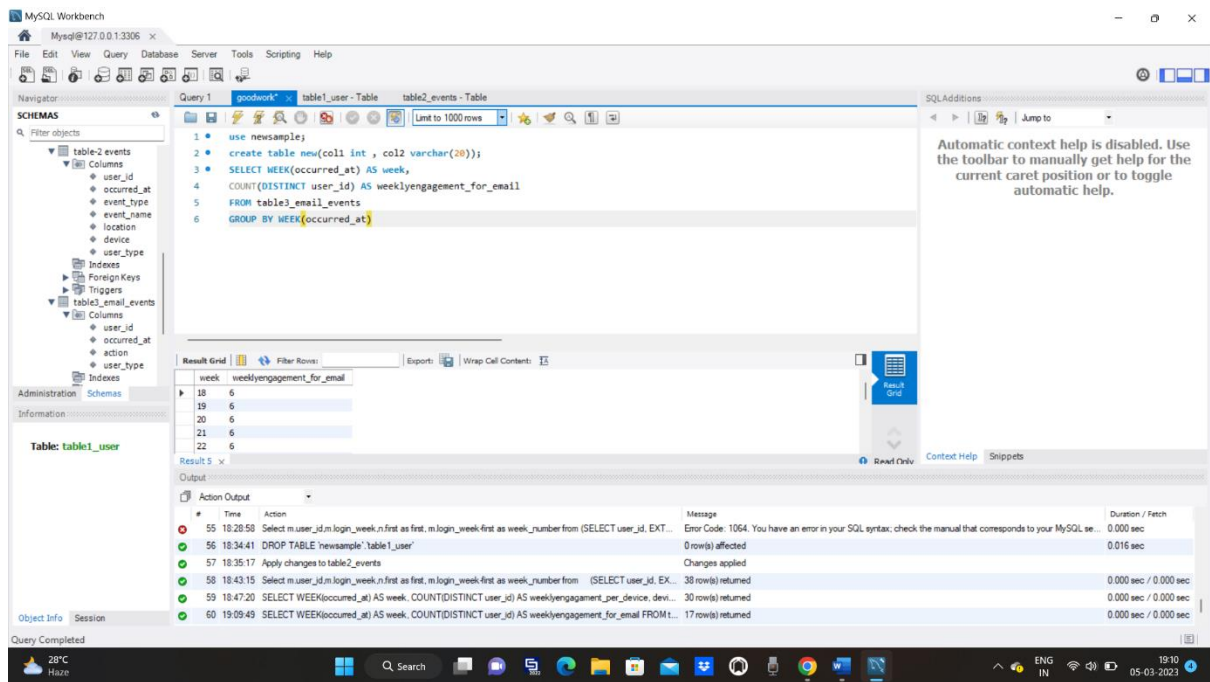**Your task: Calculate the email engagement metrics?**

## ANSWER:

SELECT WEEK(occurred_at) AS week,

COUNT(DISTINCT user_id) AS weeklyengagement_for_email

FROM table3_email_events

GROUP BY WEEK(occurred_at)

## calculating email service retention:

Select m.user_id,m.login_week,n.first as first,
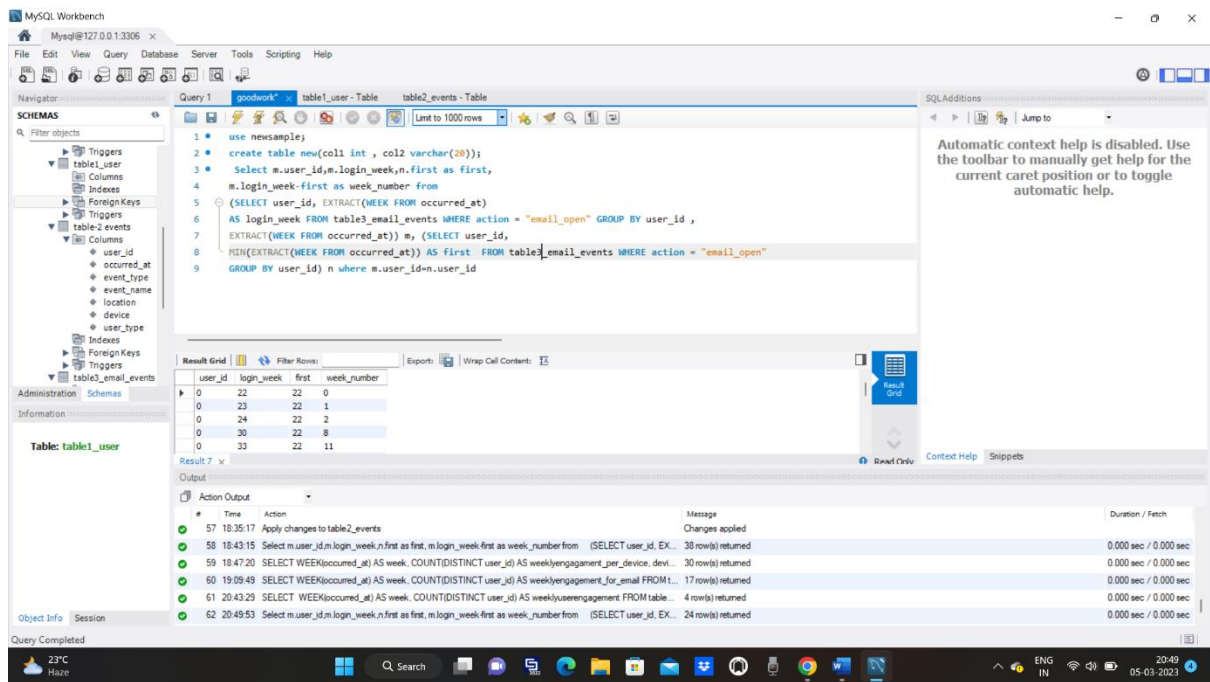
m.login_week-first as week_number from

(SELECT user_id, EXTRACT(WEEK FROM occurred_at)

AS login_week FROM table3_email_events WHERE action = "email_open" GROUP BY user_id ,

EXTRACT(WEEK FROM occurred_at)) m, (SELECT user_id,

MIN(EXTRACT(WEEK FROM occurred_at)) AS first  FROM table3_email_events WHERE action = "email_open"

GROUP BY user_id) n where m.user_id=n.user_id

## RESULT:

This project gave me deep knowledge about the topics that I learned from the mentors and while doing this project I had to search few things and concepts which again helped me in gaining more knowledge. I am happy to complete this project at my level.

**THANK YOU**