

PROJECT III

OPERATION ANALYTICS AND INVESTIGATING METRIC SPIKE

PROJECT DESCRIPTION

The project has two case studies that deal with analyzing two key aspects of operation analytics and investigating metric spikes. This process is very crucial for businesses, using these analyses they can improve their strategies and performance level. As a data analyst, one has to closely monitor the ups and downs in the metrics, derive valuable insights, understand the reason behind them, and report them to the team. Therefore the key purpose of this project is to help these companies so that they can work in fields where they lack improvement and accelerate their growth.

I'm planning to use my advanced SQL skills to perform analysis based on the questions raised by various departments and provide valuable insights that can help the company's operations.

APPROACH

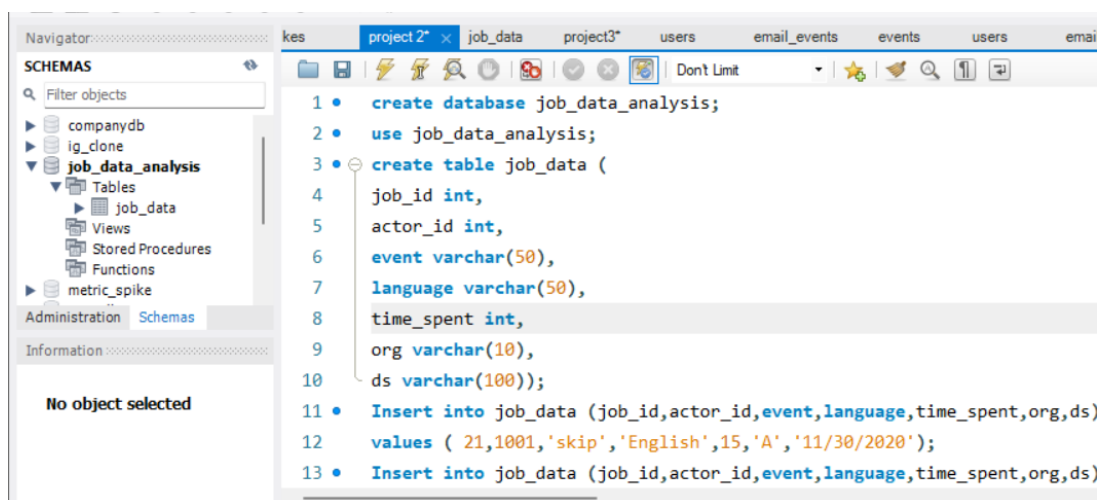
There are two case studies in this project.

Case 1 deals with Job data analysis

Case 2 deals with Investigating Metric Spike

Case Study 1: In Job data analysis there is only one job data table which involves analyzing job data to improve the efficiency in the operations field. Four tasks are to be investigated in the job data table. The tasks are to monitor how many jobs are reviewed, analyze the throughput, and percentage share of languages, and identify whether there are any duplicate rows. These are the tasks to be carried out for the analysis of job_data.

The first thing to be done before running the queries is to create a database. I created a 'job_data_analysis' database and a 'job_data' table.



The screenshot shows a SQL IDE interface. On the left, the 'SCHEMAS' pane displays a tree view with 'companydb', 'ig_clone', 'job_data_analysis' (selected), 'metric_spike', and 'Administration'. Under 'job_data_analysis', there is a 'Tables' folder containing 'job_data'. The main editor area shows the following SQL code:

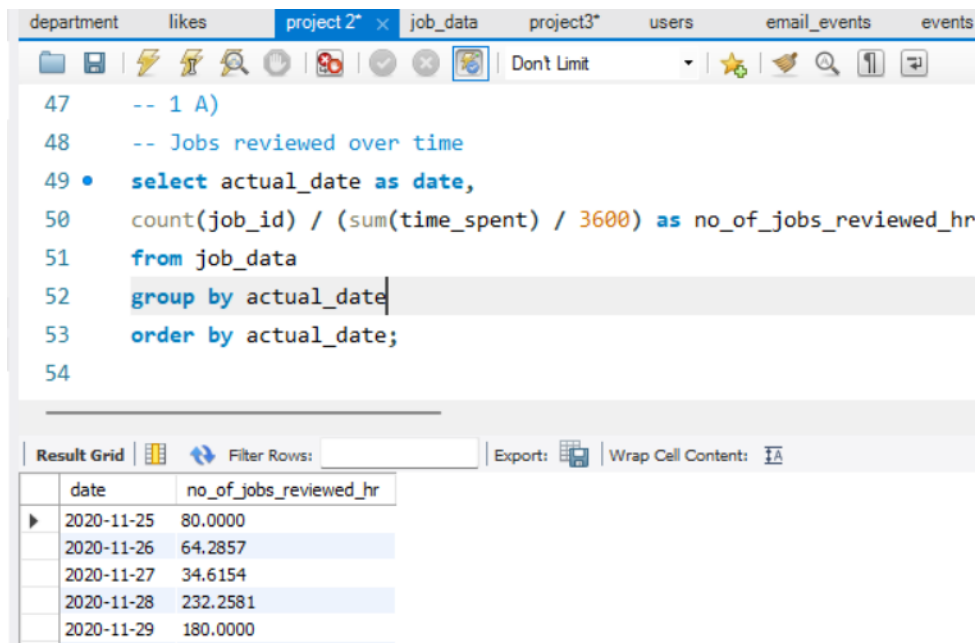
```
1 • create database job_data_analysis;
2 • use job_data_analysis;
3 • create table job_data (
4     job_id int,
5     actor_id int,
6     event varchar(50),
7     language varchar(50),
8     time_spent int,
9     org varchar(10),
10    ds varchar(100));
11 • Insert into job_data (job_id,actor_id,event,language,time_spent,org,ds)
12 values ( 21,1001,'skip','English',15,'A','11/30/2020');
13 • Insert into job_data (job_id,actor_id,event,language,time_spent,org,ds)
```

1 A) Jobs Reviewed Over Time

The first task was to calculate the number of jobs reviewed per hour for each day in November 2020.

In the data given the time_spent was given in seconds. I used a mathematical equation To convert it into hours as mentioned in the task. After running the queries I found that some days in November more jobs are reviewed and some days very few are viewed per hour.

I have attached the SQL queries below:



The screenshot shows a SQL query editor with a query window and a results grid. The query is as follows:

```
-- 1 A)
-- Jobs reviewed over time
select actual_date as date,
count(job_id) / (sum(time_spent) / 3600) as no_of_jobs_reviewed_hr
from job_data
group by actual_date
order by actual_date;
```

The results grid shows the following data:

date	no_of_jobs_reviewed_hr
2020-11-25	80.0000
2020-11-26	64.2857
2020-11-27	34.6154
2020-11-28	232.2581
2020-11-29	180.0000

In this picture, the most number of jobs were viewed on 28th November and only 34 jobs were viewed on 11th November.

1 B) Throughput Analysis:

The second task was to calculate the 7-day rolling average for throughput and explain whether I prefer daily metric or the 7-day rolling average throughput and why.

To perform the second task one must know what a 7-day rolling average is. It is an average taken over the last 7 days i.e. for a given day calculate the average of that particular day and the preceding six days. Using this approach one can observe any trend that occurs over time because it smooths out short-term fluctuations.

First I used a subquery to find the throughput value, it calculates the number of jobs per time for each day. In the main query to find the 7-day rolling average I used a Window function 'average () over()'.The condition was passed in the query and further grouped the results by the 'actual_date column'.

In the daily metric i.e. throughput, it provides a direct day-to-day metric that includes short-term fluctuations. Whereas the 7-day-average-throughput provides a week-long trend of the daily metric. In the end, it depends on the objective, If one wants to analyze immediate trends then the daily-metric approach is the better fit if one wants to observe the trend for a

long time then the rolling average method is the best because it smooths out the short-term fluctuations. I prefer the 7-day rolling average method for this task because it provides more informed decisions and stabilization.

The following queries were used to calculate the throughput and 7-day Rolling average of the throughput.

```

53  order by actual_date;
54
55  -- 1 B)
56  -- Throughput Analysis
57  -- 7day rolling avg
58  • SELECT a.actual_date ,
59        a.throughput,
60        avg(a.throughput) over ( order by a.actual_date rows between 6 preceding and current row ) as
61        7_day_rolling_avg
62  from (
63    select actual_date, count(job_id) / sum(time_spent) as throughput
64    from job_data group by actual_date ) as a;
65
66  -- 1 C)
67  -- Language Share Analysis

```

actual_date	throughput	7_day_rolling_avg
2020-11-25	0.0222	0.02220000
2020-11-26	0.0179	0.02005000
2020-11-27	0.0096	0.01656667
2020-11-28	0.0645	0.02855000
2020-11-29	0.0500	0.03284000
2020-11-30	0.0500	0.03570000

The throughput and 7-day rolling average throughput are calculated as required by the department.

One can see the difference between throughput and the 7-day rolling average.

1 C) Language Share Analysis

The third task of this project was to calculate the percentage share of each language over the last 30 days.

To run this task I used the language column from the job_data table. To find the percentage share divide the count of each language by the total count and multiply by 100. I used the where clause which considers only the last 30 days and then used the group by and order by clauses to sort the query.

After running the queries it was observed that the Persian language had the highest percentage share compared to the other languages

The query below calculates the percentage share of each language in the job_data table.

The screenshot shows a SQL IDE with a query editor and a result grid. The query is as follows:

```

64 group by a.actual_date;
65
66 -- 1 C)
67 -- Language Share Analysis
68 • select language, count(language) as language_count ,
69     round(100*count(*)/(select count(*) from job_data),2) as percentage_share
70 from job_data
71 where actual_date between '2020-11-01' and '2020-11-30'
72 group by language
73 order by language desc;
74
75 -- alternate

```

The result grid displays the following data:

language	language_count	percentage_share
Persian	3	37.50
Italian	1	12.50
Hindi	1	12.50
French	1	12.50
English	1	12.50
Arabic	1	12.50

One can observe that the Persian language has the highest percentage share of 37.50% and other languages are equally balanced out with 12.50%

1D) Duplicate Rows Identification

The final task of this case study was to display duplicate rows from the job_data table.

It was mentioned in the attribute description that job_id and actor_id uniquely identifies jobs and actor. Therefore I did this task in three parts, First I only considered the actor_id attribute, to find whether there were any duplicate actors with the same actor_id.

This is what I observed when I ran the query to find the duplicates in actor_id.

The screenshot shows a SQL IDE with a query editor and a result grid. The query is as follows:

```

97 -- 1 D)
98 -- Duplicate Row Detection
99 -- If the actor_id attribute alone is considered then there are two duplicates
100 • select * from job_data
101 where (actor_id) in (
102     select actor_id
103     from job_data
104     group by actor_id
105     having count(*) > 1
106 )
107 order by actor_id;
108 -- If the job_id attribute alone is considered then there are three duplicates

```

The result grid displays the following data:

job_id	actor_id	event	language	time_spent	org	actual_date
20	1003	transfer	Italian	45	C	2020-11-25
23	1003	decision	Persian	20	C	2020-11-29
NULL	NULL	NULL	NULL	NULL	NULL	NULL

One can observe that there are two duplicate rows in the actor_id. The duplicate actor_id is 1003.

Then I considered only the `job_id` attribute since that is also unique. In the `job_id` attribute also there are duplicates. I used a subquery to find if there are any duplicate rows in the `job_id` column.

```

106 )
107 order by actor_id;
108 -- If the job_id attribute alone is considered t
109 • SELECT *from job_data
110 where (job_id) in (
111     select job_id
112     from job_data
113     group by job_id
114     having COUNT(*) > 1
115 )
116 order by job_id;
117

```

There were duplicate rows in the `job_id` attribute. Three actors had the same `job_id` (23).

job_id	actor_id	event	language	time_spent	org	actual_date
23	1003	decision	Persian	20	C	2020-11-29
23	1004	skip	Persian	56	A	2020-11-26
23	1005	transfer	Persian	20	C	2020-11-28
NULL	NULL	NULL	NULL	NULL	NULL	NULL

Finally, I considered all the attributes at once to find if there were any duplicates.

There were no duplicate rows when the entire row was considered.

```

112 from job_data
113 group by job_id
114 having COUNT(*) > 1
115 )
116 order by job_id;
117
118 -- If the entire row is considered then there are no duplicates
119 • select * from job_data
120 group by job_id, actor_id, event, language, time_spent, org, actual
121 having count(*) > 1;
122
123

```

There are no duplicate rows when the entire row is considered.

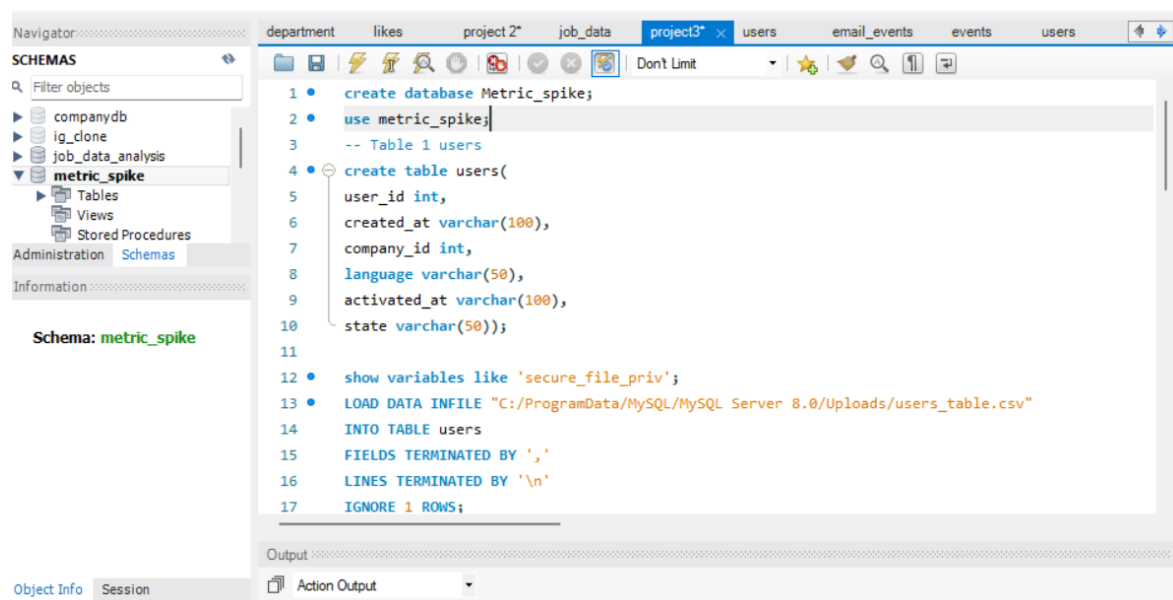
job_id	actor_id	event	language	time_spent	org	actual_date
NULL	NULL	NULL	NULL	NULL	NULL	NULL

To conclude, using various SQL queries the completion of the above tasks was carried out successfully. These insights can be useful for the company to improve its growth by rectifying the changes in fields where it particularly needs to grow.

Case study 2: In investigating metric spike there were three tables users, events, and email_events table. In case study 2 there were 5 tasks to be investigated. The tasks included checking the activeness of users, analyzing the growth of users over time, cohort-based analysis, and email engagement analysis. The aim is to analyze the data carefully, draw valuable insights, and understand the sudden change in metrics. These insights are crucial in improving the company's operations.

A database was created to import the metrics data. Then three tables were also created in the databases. Since the dataset was huge "Load Data Infile" command was used to import the data from Excel to the database.

These were the commands used to create the Metric_spike database and three tables.



The screenshot shows a MySQL IDE interface. On the left, the 'SCHEMAS' panel displays a tree view with 'companydb', 'ig_clone', 'job_data_analysis', and 'metric_spike' (selected). Below it, the 'Information' panel shows 'Schema: metric_spike'. The main editor window displays the following SQL commands:

```
1 • create database Metric_spike;
2 • use metric_spike;
3 -- Table 1 users
4 • create table users(
5   user_id int,
6   created_at varchar(100),
7   company_id int,
8   language varchar(50),
9   activated_at varchar(100),
10  state varchar(50));
11
12 • show variables like 'secure_file_priv';
13 • LOAD DATA INFILE "C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/users_table.csv"
14   INTO TABLE users
15   FIELDS TERMINATED BY ','
16   LINES TERMINATED BY '\n'
17   IGNORE 1 ROWS;
```

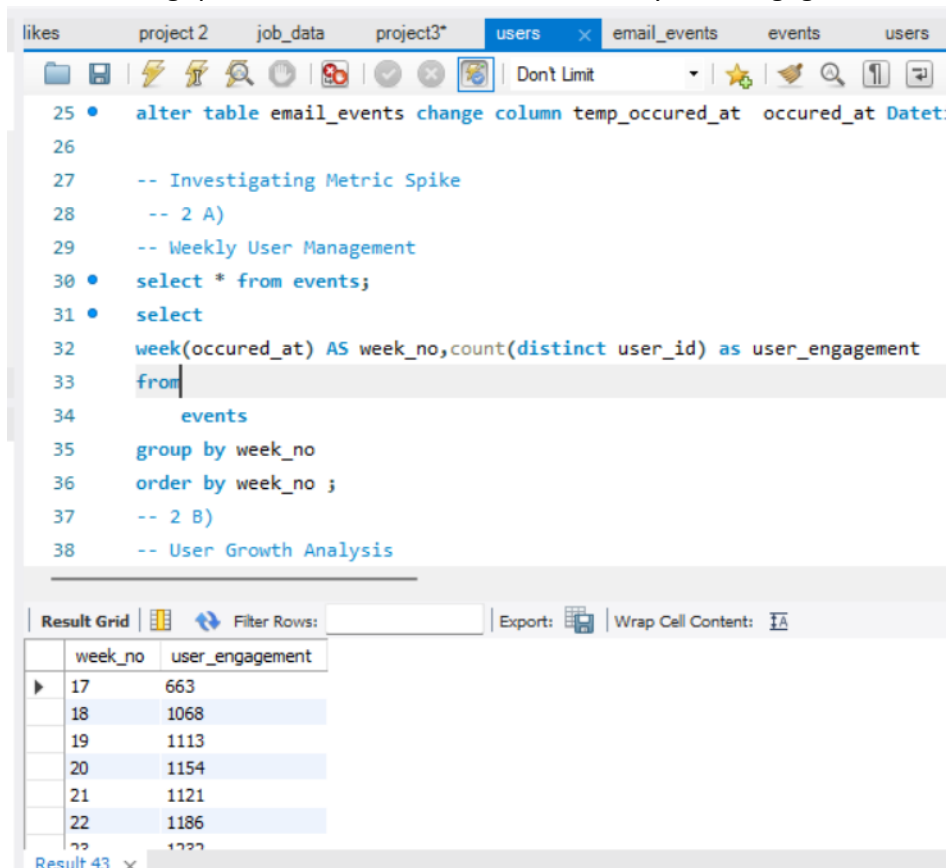
The bottom of the interface shows the 'Output' panel with 'Action Output' selected.

2A) Weekly User Engagement:

The first task to perform in the case study was to calculate the weekly user engagement.

To perform this task I made use of the events table in the metric_spike database. Before that, the date in the 'occurred_at' attribute was in text format. So after importing the dataset I changed the string formatted text into date format and carried out my analysis. Since the task was to calculate user engagement every week, I made use of the 'week' command to extract the week from the occurred_at attribute. Therefore occurred_at and user_id attributes were retrieved from the events table to analyze the weekly user engagement and further sorted them using group by and order by clauses.

The following queries were used to find the weekly user engagement



The screenshot shows a SQL IDE with a query editor and a result grid. The query editor contains the following SQL code:

```
25 • alter table email_events change column temp_occured_at occured_at Datet:
26
27 -- Investigating Metric Spike
28 -- 2 A)
29 -- Weekly User Management
30 • select * from events;
31 • select
32 week(occured_at) AS week_no, count(distinct user_id) as user_engagement
33 from
34     events
35 group by week_no
36 order by week_no ;
37 -- 2 B)
38 -- User Growth Analysis
```

The result grid shows the following data:

week_no	user_engagement
17	663
18	1068
19	1113
20	1154
21	1121
22	1186
23	1222

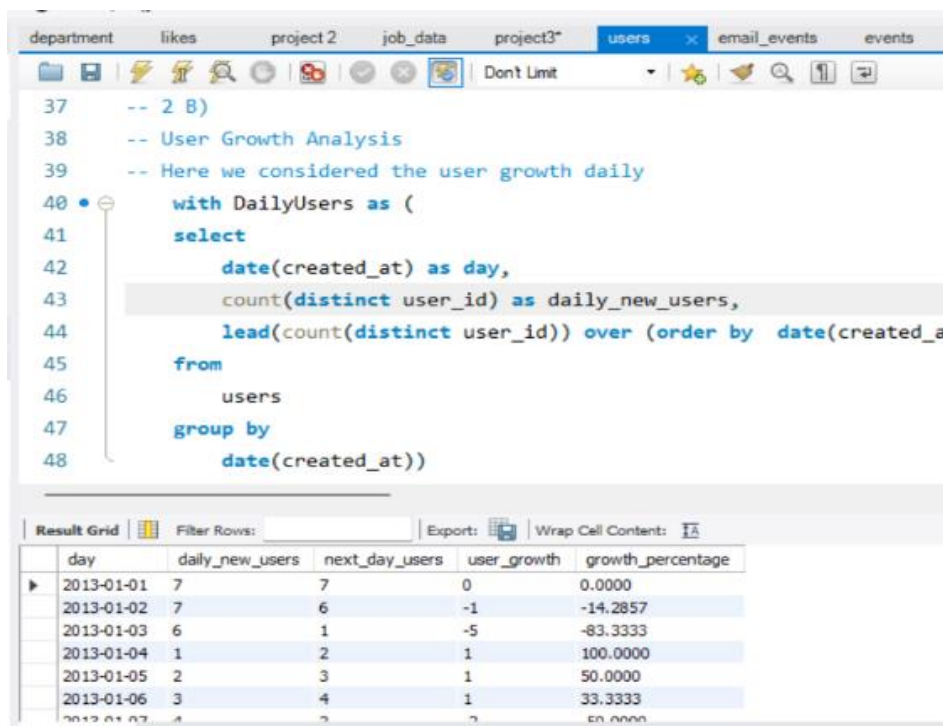
After running the queries it was observed that Week 30 had the highest count of user engagement of 1467, and the lowest count was 105 which was observed in Week 35.

2B) User Growth Analysis

The second task was to calculate the user growth for the product over time.

In this task, we need to find the growth of users for a product over time. Since it's not mentioned to calculate it weekly or monthly or daily. I decided to calculate on a daily as well as on a monthly basis. To calculate it daily a subquery using the CTE (Common Table Expression) was created. In the CTE date was extracted from the 'created_at' attribute, using the count function no of distinct users were also calculated, and then a window function called 'lead' was used to calculate the users of the next day. In that way, we can analyze the user growth daily. Then in the main query, we calculated the user growth by the difference of users on the current day and the next day, and % of user growth was also calculated.

The following queries were used to perform the second task.

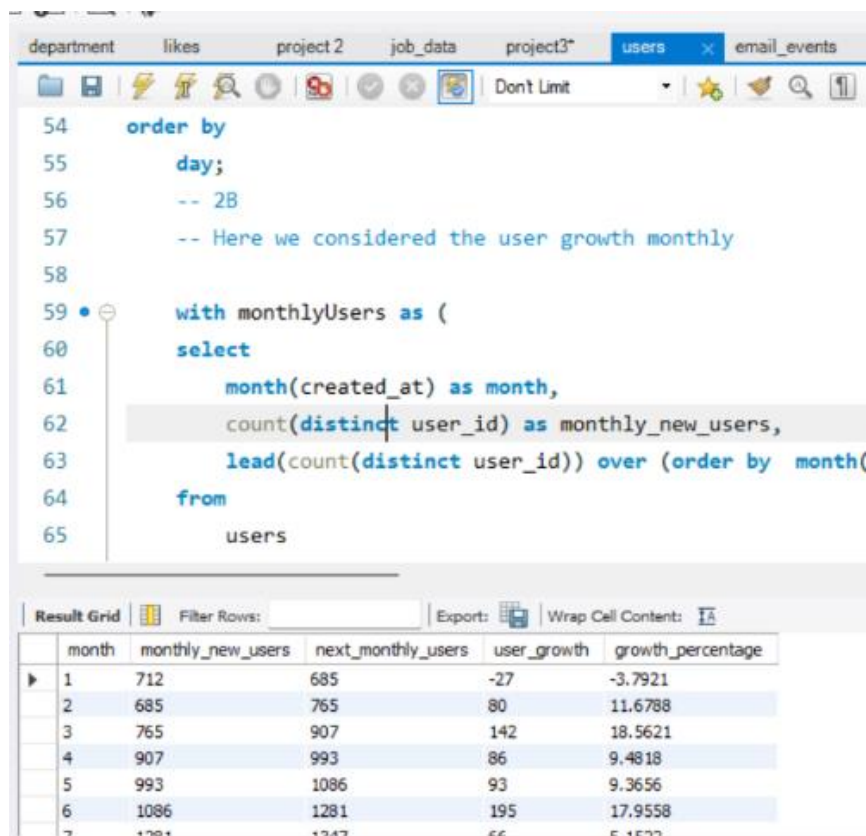


```
37 -- 2 B)
38 -- User Growth Analysis
39 -- Here we considered the user growth daily
40 with DailyUsers as (
41 select
42     date(created_at) as day,
43     count(distinct user_id) as daily_new_users,
44     lead(count(distinct user_id)) over (order by date(created_at))
45 from
46     users
47 group by
48     date(created_at))
```

day	daily_new_users	next_day_users	user_growth	growth_percentage
2013-01-01	7	7	0	0.0000
2013-01-02	7	6	-1	-14.2857
2013-01-03	6	1	-5	-83.3333
2013-01-04	1	2	1	100.0000
2013-01-05	2	3	1	50.0000
2013-01-06	3	4	1	33.3333
2013-01-07	4	5	1	25.0000

The highest user growth was 41 and it was observed on 2014-08-24.

Similarly, if we want to calculate user growth every month then follow the steps except use the month function to extract the month from the created_at attribute.



```
54 order by
55     day;
56 -- 2B
57 -- Here we considered the user growth monthly
58
59 with monthlyUsers as (
60 select
61     month(created_at) as month,
62     count(distinct user_id) as monthly_new_users,
63     lead(count(distinct user_id)) over (order by month)
64 from
65     users
```

month	monthly_new_users	next_monthly_users	user_growth	growth_percentage
1	712	685	-27	-3.7921
2	685	765	80	11.6788
3	765	907	142	18.5621
4	907	993	86	9.4818
5	993	1086	93	9.3656
6	1086	1281	195	17.9558
7	1281	1247	-34	-2.6573

The highest user growth was observed in the 6th month where the user growth spiked to 195.

The lowest growth was observed in the 8th month where user growth plummeted to 1017.

We can also consider it on a weekly basis then the 31st week shows the highest user growth.

2 C) Weekly Retention Analysis

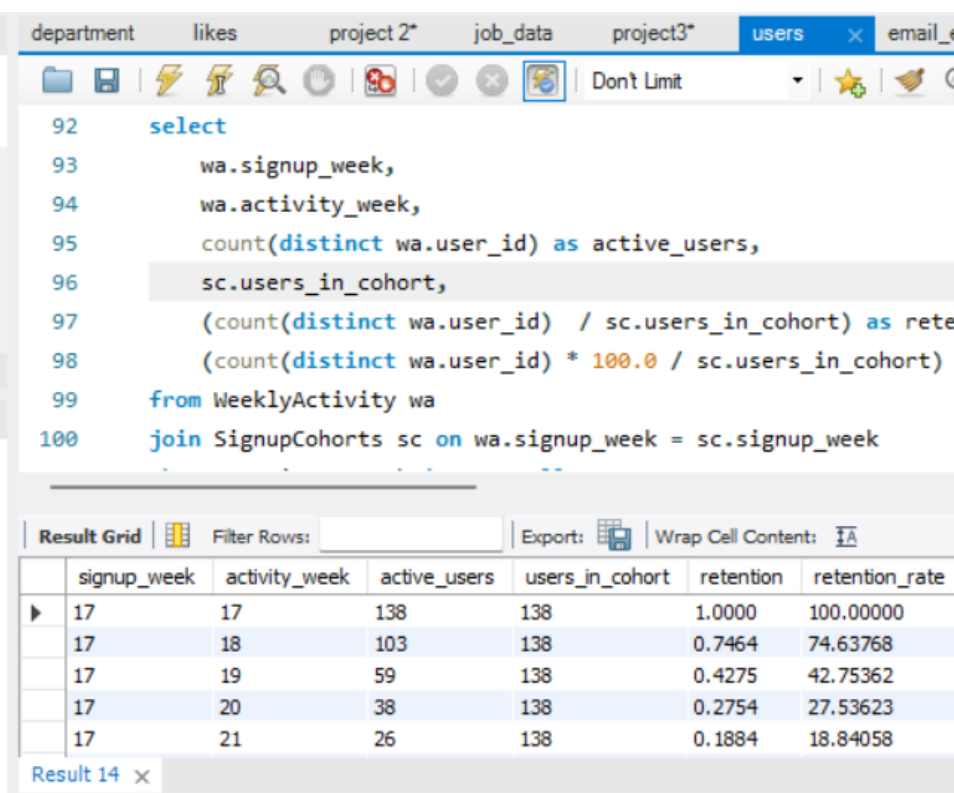
The third task is to calculate the weekly retention of users based on their sign-up cohort.

To perform the above task one has to break it down into bits. I used a subquery to store the weekly activity of events. In this query, we can find how each event is associated with the week. In the weekly activity subquery, I used a few window functions such as partition by and case to correctly identify the user's signup on the specific event. A min function was used because if all more events were associated with the users then only one sign-up activity would be considered. I joined the events and the users table using the join clause. Another subquery was used for the signup cohort. The main use of this CTE is to find the size of each weekly sign_up_cohort. In the main query, we use the retention formula to find the retention rate and we also join the CTE to find the weekly retention analysis and further it grouped by signup_week.

The formula used to calculate the retention rate is =

$$(\text{No of users in the end} - \text{No of users during the event}) / (\text{No of users in the starting period of the event}) * 100$$

The following queries were used to perform the third task.



```
92 select
93     wa.signup_week,
94     wa.activity_week,
95     count(distinct wa.user_id) as active_users,
96     sc.users_in_cohort,
97     (count(distinct wa.user_id) / sc.users_in_cohort) as rete
98     (count(distinct wa.user_id) * 100.0 / sc.users_in_cohort)
99 from WeeklyActivity wa
100 join SignupCohorts sc on wa.signup_week = sc.signup_week
```

	signup_week	activity_week	active_users	users_in_cohort	retention	retention_rate
▶	17	17	138	138	1.0000	100.00000
	17	18	103	138	0.7464	74.63768
	17	19	59	138	0.4275	42.75362
	17	20	38	138	0.2754	27.53623
	17	21	26	138	0.1884	18.84058

Result 14 x

One can observe the weekly user retention for each group of users.

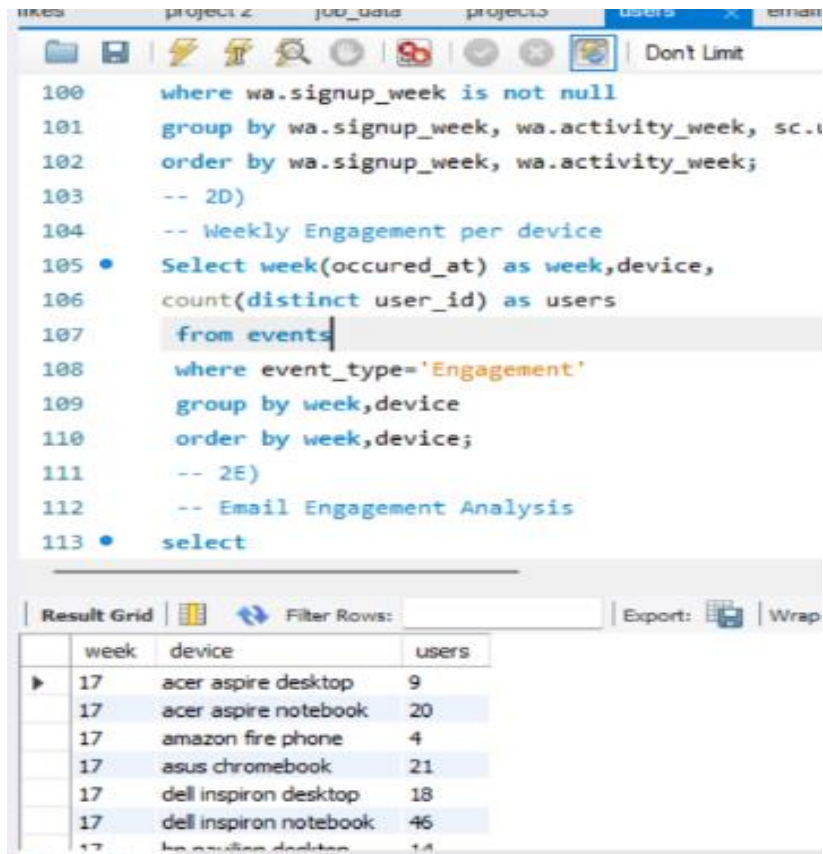
One can observe various ranges of retention rates in the subsequent weeks. In some weeks the rates are really low and in some weeks it's high. This can be due to various factors.

2 D) Weekly Engagement per Device

The fourth task was to calculate the weekly engagement per device.

To find the weekly engagement per device, I used the events table. I extracted the week from the occurred_at attribute and found the distinct users and devices from the events table. I used the 'where' condition to select only the events with "Engagement" and then sorted them by group by and order by clauses.

The following queries were used to perform the fourth task.



```
100 where wa.signup_week is not null
101 group by wa.signup_week, wa.activity_week, sc.i
102 order by wa.signup_week, wa.activity_week;
103 -- 2D)
104 -- Weekly Engagement per device
105 • Select week(occured_at) as week,device,
106 count(distinct user_id) as users
107 from events
108 where event_type='Engagement'
109 group by week,device
110 order by week,device;
111 -- 2E)
112 -- Email Engagement Analysis
113 • select
```

	week	device	users
▶	17	acer aspire desktop	9
	17	acer aspire notebook	20
	17	amazon fire phone	4
	17	asus chromebook	21
	17	dell inspiron desktop	18
	17	dell inspiron notebook	46
	17	hp pavilion desktop	14

Mac Book Pro had the highest weekly engagement per device with 252 users using it on the 18th week.

2 E) Email Engagement per Device

The final task of this case study is to calculate the email engagement metrics.

To perform this task I used the email_events table. In the email_events table, the action attribute has four categories. These four actions' weekly engagement can be calculated using the case function and then grouped by using the group by and order by clauses.

The following queries were used to find the email engagement action per device

```

106 count(distinct user_id) as users
107 from events
108 where event_type='Engagement'
109 group by week,device
110 order by week,device;
111 -- 2E)
112 -- Email Engagement Analysis
113 • select
114     WEEK(occured_at) as week,
115     count(distinct case when action = "sent_weekly_digest" then user_id end) as Weekly_dige
116     count(distinct case when action = "sent_reengagement_email" then user_id end) as reenga
117     count(distinct case when action = "email_open" then user_id end) as opened_email,
118     count(distinct case when action = "email_clickthrough" then user_id end) as email_click
119 from email_events

```

week	Weekly_digest	reengagement_mail	opened_email	email_clickthrough
17	908	73	310	166
18	2602	157	900	425
19	2665	173	961	476
20	2733	191	989	501
21	2822	164	996	436
22	2911	192	965	478

This query gives the Weekly engagement actions per device every week.

```

110 group by week,device
111 order by week,device;
112 -- 2E)
113 -- Email Engagement Analysis
114 • Select avg(weekly_digest) as AVG_WEEKLY_DIGEST,avg(reengagement_email) as AVG_REENGAGEMENT_EM
115 (select
116     WEEK(occured_at) as week,
117     count(distinct case when action = "sent_weekly_digest" then user_id end) as Weekly_digest,
118     count(distinct case when action = "sent_reengagement_email" then user_id end) as reengagem

```

AVG_WEEKLY_DIGEST	AVG_REENGAGEMENT_EMAIL	AVG_OPENED_MAIL	AVG_EMAIL_CLICK_THROUGH
3014.0526	192.2632	1061.1579	469.3684

From this query, one can notice that Weekly Digest has the highest weekly engagement with an average of 3014

TECH-STACK USED

I used MySQL workbench version 8.0.34 build 3263449 CE(64 bits) for this project. The reason I used MySQL workbench for my project is because I was more comfortable using this software. I practiced my SQL queries mostly in this software. MySQL workbench is user-friendly and is also open-source.

INSIGHTS

This project is the most challenging one I have done so far. To do this project one should possess a strong knowledge of SQL concepts. Once you dive into the data and get to

understand more about it, then it gets more and more interesting. Each task was challenging and nerve-racking in its way. I'm glad for this opportunity wherein I got to apply my SQL knowledge in a Live project.

The key insights I found while working on projects are listed below

In Case Study 1

- The average number of jobs reviewed per hour per day in November alone is 128 and most number of jobs were reviewed on 28th Nov.
- The 7-day rolling average of throughput is better than the daily metric if you want to observe any long-term trend.
- Persian language is the most used language with a share of 37.5%
- If we consider the overall rows to find duplicates then there are none, but there are duplicate rows in the actor_id and job_id attribute if we consider them individually.

In Case Study 2

- There's a spike in the initial weeks but towards the end, it takes a dip. There are many possible reasons for such occurrences. One such reason can be due to the lack of satisfaction with the product.
- The sixth month has the highest user growth
- The weekly retention varies each week
- Most of the users use MacBook Pro on a weekly basis
- Weekly Digest emails have the highest weekly engagement.

RESULT

It was indeed a great experience for me as I got hands-on experience in working with the Live projects. This is my second project working in SQL, I had to work extra on so many concepts to tackle the challenges faced during this project. This in turn helped me to gain more in-depth knowledge of SQL concepts and I believe this project helped me to hone my SQL skills to a great extent. I gained a lot of knowledge and understanding while working on this Operation analytic and investigating metric spike project

END

**Submitted by,
S Nandana.**