# **Operation Analytics and Investigating Metric Spike**

## Submitted by:

### Nadendla Dharani

# **Project Description**

The project focuses on Operational Analytics, aiming to analyze end-to-end operations within the company to identify areas for improvement. As a Lead Data Analyst, the goal is to investigate and explain sudden changes in key metrics, addressing questions from various departments. By using advanced SQL skills, the project aims to derive valuable insights that can enhance the company's operations and shed light on metric spikes.

## Approach

In the context of Operational Analytics, the initial step involves a comprehensive data exploration phase. This includes a thorough examination of the structure of provided datasets to identify relevant columns, relationships, and potential key metrics.

Following data exploration, the focus shifts to the crafting of SQL queries tailored to extract, transform, and analyze data to address specific inquiries from different departments. MySQL Workbench is utilized as the primary tool for database design, querying, and visualization, providing an integrated environment for efficient analysis. Statistical analysis tools are incorporated where necessary to delve deeper into insights. The insights generated contribute to data-driven decision-making processes across departments, fostering collaboration and promoting a holistic understanding of company operations.

The overall approach combines meticulous data exploration, strategic SQL query design, and collaborative efforts to derive actionable insights.

#### Tech-Stack Used

The technology stack employed for this project includes MySQL Workbench and SQL queries. MySQL Workbench serves as the primary tool for visualizations and provides a robust platform for interpreting the analyzed data. SQL queries, on the other hand, are instrumental for extracting, manipulating, and analyzing the user data from the given databases efficiently. This combination was chosen for its effectiveness in handling the specific requirements of the project, allowing for seamless data exploration and extraction of meaningful insights from the given data.

## Insights

- Peaks in job reviews during specific hours can be identified, aiding in scheduling resource allocation and workload management.
- Daily throughput and the 7-day rolling average provide different perspectives.
   Daily metrics offer a granular view, while the rolling average smoothens short-term fluctuations, making it preferable for trend analysis.
- This query reveals the percentage share of each language over the last 30 days, assisting in language-related decision-making and content prioritization.
- Identified duplicate rows based on multiple columns, providing a basis for data cleansing and ensuring data integrity.
- By analyzing weekly user engagement, trends in user activity can be identified, potentially aligning with marketing campaigns or product releases.
- Understanding user growth over time is crucial for assessing the product's performance and potential market saturation.
- Weekly retention analysis sheds light on user stickiness, highlighting areas for improvement in onboarding and user engagement strategies.
- Understanding how user engagement varies across devices can inform device-specific optimizations and feature development.
- Analyzing email engagement metrics provides insights into the effectiveness of email campaigns, helping refine communication strategies.

These insights collectively empower data-driven decision-making across departments, enabling proactive responses to operational challenges and driving improvements in overall business performance.

#### Results

### Case Study 1: Job Data Analysis

You will be working with a table named job\_data with the following columns:

- job\_id: Unique identifier of jobs
- actor id: Unique identifier of actor
- event: The type of event (decision/skip/transfer).
- language: The Language of the content
- time spent: Time spent to review the job in seconds.
- org: The Organization of the actor
- **ds:** The date in the format yyyy/mm/dd (stored as text).

## Creating the job\_data table:

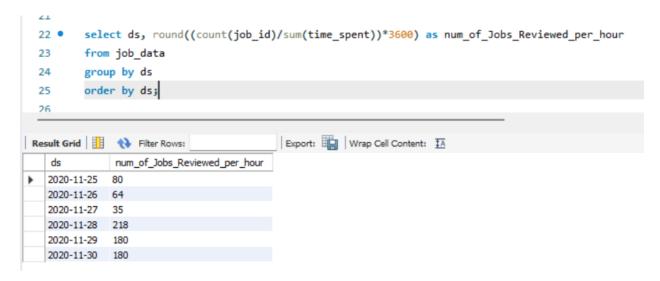
```
create database jobs;
  2 • ⊖ create table job_data (
            ds date,
            job_id int not null,
            actor_id int not null,
  5
            event varchar(50) not null,
  6
  7
            language varchar(50) not null,
  8
             time_spent int not null,
  9
             org char(5)
 10
       );
11 •
       insert into job_data (ds, job_id, actor_id, event, language, time_spent, org)
       values ('2020-11-30', 21, 1001, 'skip', 'English', 15, 'A'),
           ('2020-11-30', 22, 1006, 'transfer', 'Arabic', 25, 'B'),
13
            ('2020-11-29', 23, 1003, 'decision', 'Persian', 20, 'C'),
14
            ('2020-11-28', 23, 1005, 'transfer', 'Persian', 22, 'D'),
           ('2020-11-28', 25, 1002, 'decision', 'Hindi', 11, 'B'),
            ('2020-11-27', 11, 1007, 'decision', 'French', 104, 'D'),
17
            ('2020-11-26', 23, 1004, 'skip', 'Persian', 56, 'A'),
            ('2020-11-25', 20, 1003, 'transfer', 'Italian', 45, 'C');
19
     select * from job data;
```

Re	esult Grid	<b>₹</b> Fil	ter Rows:		Exp	ort: 📳 Wr	ap Cell Co	ntent: TA	
	ds	job_id	actor_id	event	language	time_spent	org		
١	2020-11-30	21	1001	skip	English	15	Α		
	2020-11-30	22	1006	transfer	Arabic	25	В		
	2020-11-29	23	1003	decision	Persian	20	C		
	2020-11-28	23	1005	transfer	Persian	22	D		
	2020-11-28	25	1002	decision	Hindi	11	В		
	2020-11-27	11	1007	decision	French	104	D		
	2020-11-26	23	1004	skip	Persian	56	Α		
	2020-11-25	20	1003	transfer	Italian	45	С		

### Tasks:

### 1. Jobs Reviewed Over Time:

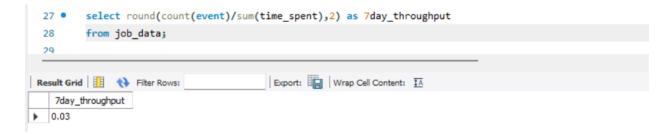
- Objective: Calculate the number of jobs reviewed per hour for each day in November 2020.
- Your Task: Write an SQL query to calculate the number of jobs reviewed per hour for each day in November 2020.

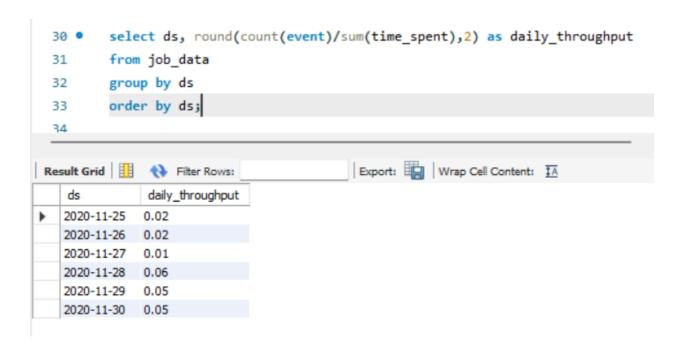


Here I have calculated the number of jobs reviewed per hour by dividing the number of times a user reviewed by total time spent, later this result is multiplied by 3600 because time spent is given in seconds and we need the result in hours. Also round() function is used which returns the only integer value excluding decimal values.

# 2. Throughput Analysis:

- Objective: 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.





Here after finding 7-day rolling throughput and daily throughput, I would say I prefer both, because both offer different perspectives for the trend analysis. So both are crucial for a proper data-driven decision.

# 3. Language Share Analysis:

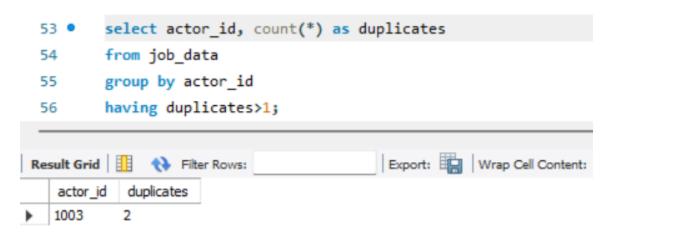
- Objective: Calculate the percentage share of each language in the last 30 days.
- Your Task: Write an SQL query to calculate the percentage share of each language over the last 30 days.

```
61 •
         select language, count(language) as language_count,
         (count(language) / (select count(*) from job_data)) * 100 as percentage
 62
         from job_data
 63
 64
         group by language
 65
         order by language DESC;
Result Grid
                                             Export: Wrap Cell Content: TA
               Filter Rows:
   language
             language_count
                            percentage
  Persian
                            37.5000
   Italian
                            12.5000
  Hindi
                            12.5000
            1
  French
                            12.5000
  English
            1
                            12.5000
  Arabic
                            12.5000
```

Here, there are languages that are repeated multiple times so I displayed the count of each language for better understanding and also found percentage using the basic formula, (number/total)\*100.

## 4. Duplicate Rows Detection:

- Objective: Identify duplicate rows in the data.
- Your Task: Write an SQL query to display duplicate rows from the job\_data table.



actor\_id is a unique attribute, so I counted the number of times each ID repeated and displayed the ID which has counted more than 1.

# Case Study 2: Investigating Metric Spike

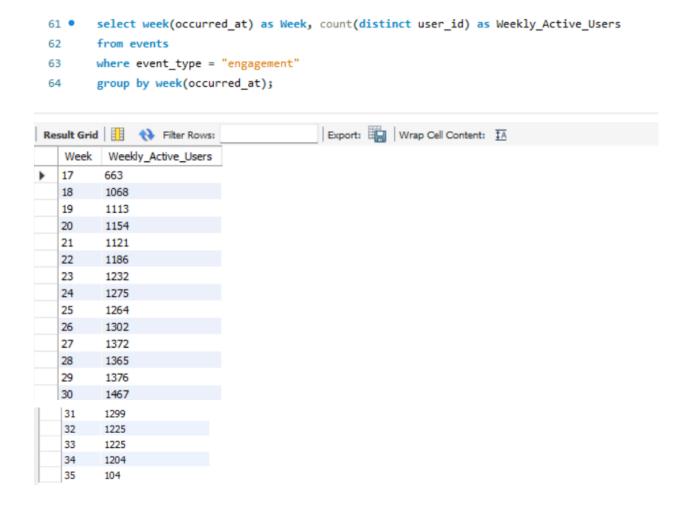
# You will be working with three tables:

- **users**: Contains one row per user, with descriptive information about that user's account.
- **events**: Contains one row per event, where an event is an action that a user has taken (e.g., login, messaging, search).
- email events: Contains events specific to the sending of emails.

#### Tasks:

# 1. Weekly User Engagement:

- Objective: Measure the activeness of users on a weekly basis.
- Your Task: Write an SQL query to calculate the weekly user engagement.



Here I have extracted week from the occured\_at column, and in the event\_type column I have selected rows having "engagement" and printed the output resulting number of active users in each particular week.

# 2. User Growth Analysis:

- Objective: Analyze the growth of users over time for a product.
- Your Task: Write an SQL query to calculate the user growth for the product.

	months	total_users	Growth_inc
•	1	712	NULL
	2	685	-3.79
	3	765	11.68
	4	907	18.56
	5	993	9.48
	6	1086	9.37
	7	1281	17.96
	8	1347	5.15
	9	330	-75.50
	10	390	18.18
	11	399	2.31
	12	486	21.80

Here I have used the 'Lag' function, which is used to access previous rows. It is a useful function in comparing the current row value from the previous row value.

# 3. Weekly Retention Analysis:

- Objective: 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.

```
with users_start as (
    select user_id, extract(week from created_at) as start_week
    from users
),

user_activity as (
    select user_id, extract(week from occurred_at) as activity_week
    from events
)

select users_start.start_week as cohort_week,
    user_activity.activity_week as retention_week,
    count(distinct user_activity.user_id) as total_users

from users_start

left join user_activity
on users_start.user_id = user_activity.user_id and user_activity.activity_week >= users_start.start_week
group by users_start.start_week, user_activity.activity_week
order by users_start.start_week, user_activity.activity_week;
```

Re	sult Grid	Filter Rows:		Export:
	cohort_week	retention_week	total_users	
•	0	NULL	0	-
	0	17	5	
	0	18	11	
	0	19	15	
	0	20	12	
	0	21	12	
	0	22	14	
	0	23	12	
	0	24	19	
	0	25	12	
	0	26	13	
	0	27	11	
	0	28	8	
	0	29	6	

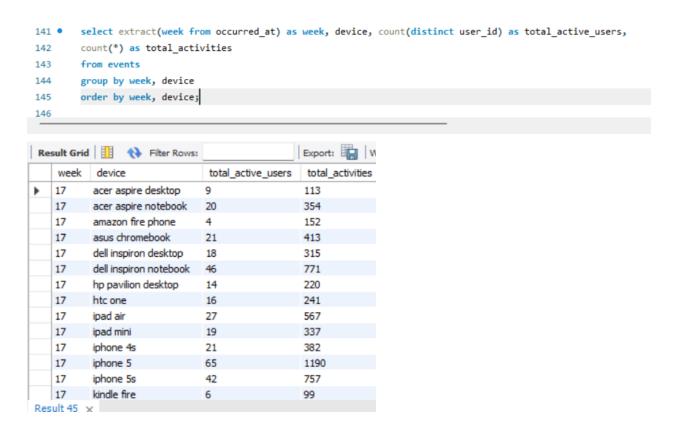
I solved the question using 'WITH' clause subquery and created two subqueries, users\_start and user\_activity. Also in the main subquery I performed a LEFT JOIN operation on users\_start and user\_activity.

Also, the output for this question is very large. So I have exported the output to an excel sheet and I am attaching the drive link to the sheet here. For referral, I have provided a small screenshot of the total output.

https://drive.google.com/file/d/1jLvo0geYZno3iAkMc9CcXb5pWVATKYWm/view?usp=s haring

## 4. Weekly Engagement Per Device:

- Objective: Measure the activeness of users on a weekly basis per device.
- Your Task: Write an SQL query to calculate the weekly engagement per device.



For this question the output is very large, so I have exported the output to an excel sheet and I am attaching the drive link to the sheet here. For referral, I have provided a small screenshot of the total output

https://drive.google.com/file/d/1Lubn0VJBwAsUjwyw1Vv4KJ67uj-mOJna/view?usp=sharing

# 5. Email Engagement Analysis:

- Objective: Analyze how users are engaging with the email service.
- Your Task: Write an SQL query to calculate the email engagement metrics.

```
select action, count(distinct user_id) as distinct_users, count(*) as total_actions
148 •
        from email events
149
150
        group by action
        order by action;
151
152
Export: Wrap Cell Content: IA
   action
                        distinct_users
                                    total_actions
  email_dickthrough
                        5277
                                    9010
  email_open
                        5927
                                    20459
  sent_reengagement_email
                       3653
                                    3653
  sent_weekly_digest
                       4111
                                    57267
```

Through the GROUP BY clause, I have provided the output of a unique number of users and total actions taken place in every action respectively.