# Operation Analytics and Investigating Metric Spike

Advanced SQL

Swet

#### **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. I've worked 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 I've made the 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 its very important to investigate metric spike.

#### **APPROACH**

- 1. Created the Database and Tables: Created a database and then the tables using the structure and links provided.
- 2. Perform Analysis: Used SQL to perform entire analysis answering the questions asked. I went through the given dataset and then created the tables for calculating various queries.
- 3. Additionally, I joined the data bits and structured the tables to derive business insights, fetched the required results and hence, created useful insights for the company to take calculated and planned decisions.

#### **TECH-STACK USED**

#### Software And The Version Used While Making The Project:

- 1. MySQL WorkBench 8.0 (For working, analysing and reporting insight)
- 2. Microsoft Power Point (For presenting the detailed analysis)

The objective of the project is to find out insights about following:

#### → Case 1 : JOB DATA

- i. Number of jobs reviewed: Number of jobs reviewed over time
- ii. Throughput: It is the no. of events happening per second.
- iii. Percentage share of each language: Share of each language for different contents.
- iv. Duplicate rows: Rows that have the same value present in them.

#### → Case 2 : INVESTIGATING METRIC SPIKE

- i. User Engagement: To measure the activeness of a user. Measuring if the user finds quality in a product/service.
- ii. User Growth: Number of users growing over time for a product.
- iii. Weekly Retention: Users getting retained weekly after signingup for a product.
- iv. Weekly Engagement: To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly.
- v. Email Engagement: Users engaging with the email service.

### Case study 1: JOB DATA

#### Creating table job\_data

job_id	actor_id	event	language	time_spent	org	ds
21	1001	skip	English	00:00:15	Α	2020-11-30
22	1006	transfer	Arabic	00:00:25	В	2020-11-30
23	1003	decision	Persian	00:00:20	C	2020-11-29
23	1005	transfer	Persian	00:00:22	D	2020-11-28
25	1002	decsison	Hindi	00:00:11	В	2020-11-28
11	1007	decision	French	00:01:04	D	2020-11-27
23	1004	skip	Persian	00:00:56	Α	2020-11-26
20	1003	transfer	Italian	00:00:45	C	2020-11-25

#### Syntax used for creating database job db and job data table.

```
create database job_db
       use job db
       create table job_data
       job_id int,
 5
       actor id int,
 6
       event varchar(50),
       language varchar(50),
       time spent time,
       org varchar(100),
11
       ds date
12
       insert into job data (job id, actor id, event, language, time spent, org, ds)
13
14
       values
       ('21', '1001', 'skip', 'English', '15', 'A', '2020-11-30'),
15
       ('22', '1006', 'transfer', 'Arabic', '25', 'B', '2020-11-30'),
16
       ('23', '1003', 'decision', 'Persian', '20', 'C', '2020-11-29'),
17
       ('23', '1005', 'transfer', 'Persian', '22', 'D', '2020-11-28'),
18
       ('25', '1002', 'decsison', 'Hindi', '11', 'B', '2020-11-28'),
19
       ('11', '1007', 'decision', 'French', '104', 'D', '2020-11-27'),
20
       ('23', '1004', 'skip', 'Persian', '56', 'A', '2020-11-26'),
21
       ('20', '1003', 'transfer', 'Italian', '45', 'C', '2020-11-25')
22
23
       select * from job data
```

### A. <u>Number of jobs reviewed Objective</u> <u>Calculate the</u> <u>number of jobs reviewed per hour per day for November 2020</u>

```
select count(distinct job_id)/(30*24)

as per_day_jobs

from job_data

per_day_jobs
```

Less than 0.01 jobs were reviewed each hour of the day throughout the month of November.

#### B. <u>Throughput (Number of events happening per</u> <u>second) Objective:</u> <u>Calculate 7 day rolling average of</u> <u>throughput</u>

```
select ds, tot_events,
avg(tot_events) over(order by ds rows between 6 preceding and current row) as 7_day_rolling_average
from
(select ds, count(distinct event) as tot_events
from job_data
group by ds
order by ds)sub;
```

	ds	tot_events	7_day_rolling_average
•	2020-11-25	1	1.0000
	2020-11-26	1	1.0000
	2020-11-27	1	1.0000
	2020-11-28	2	1.2500
	2020-11-29	1	1,2000
	2020-11-30	2	1.3333

Using a 7-day rolling average for throughput can be helpful in understanding trends over time, as it provides a longer-term perspective compared to a daily metric. This can help to smooth out any short-term fluctuations in the data and provide a clearer picture of the overall trend.

# C. <u>Percentage share of each language Objective</u> <u>Calculate the percentage share of each language in the last 30</u> <u>days</u>

```
select language, count(language) as total_language,
count(*)*100/sum(count(*))
over() as percentage
from job_data
group by language
order by language
```

	language	total_language	percentage
١	Arabic	1	12.5000
	English	1	12.5000
	French	1	12,5000
	Hindi	1	12.5000
	Italian	1	12.5000
	Persian	3	37.5000

Persian Language had the highest share among other languages

## D. <u>Duplicate Rows Objective</u> <u>Display duplicate</u> rows if any

```
job_id actor_id event language time_spent org ds row_numb

row_number() over (partition by job_id) as row_numb

from job_data)

select * from cte where row_numb>1

job_id actor_id event language time_spent org ds row_numb

23 1005 transfer Persian 00:00:22 D 2020-11-28 2

23 1004 skip Persian 00:00:56 A 2020-11-26 3
```

The output showed two records as there were two duplicate job id in the dataset

### Case Study 2: Investigating metric spike

#### 1. Creating Users Table



	user_id	created_at	company_id	language	activated_at	state
•	0	2013-01-01 20:59:39	5737	english	2013-01-01 21:01:07	active
	1	2013-01-01 13:07:46	28	english		pending
	2	2013-01-01 10:59:05	51	english		pending
	3	2013-01-01 18:40:36	2800	german	2013-01-01 18:42:02	active
	4	2013-01-01 14:37:51	5110	indian	2013-01-01 14:39:05	active
	5	2013-01-01 13:39:51	2463	spanish		pending
	6	2013-01-01 18:37:27	11699	english	2013-01-01 18:38:45	active
	7	2013-01-01 16:19:01	4765	french	2013-01-01 16:20:28	active
	8	2013-01-01 04:38:30	2698	french	2013-01-01 04:40:10	active
	9	2013-01-01 08:04:17	1	french		pending
	10	2013-01-01 09:36:41	10	arabic		pending
	11	2013-01-01 08:07:45	3745	english	2013-01-01 08:09:17	active
	12	2013-01-01 18:05:05	903	english		pending

### 2. Creating Events table

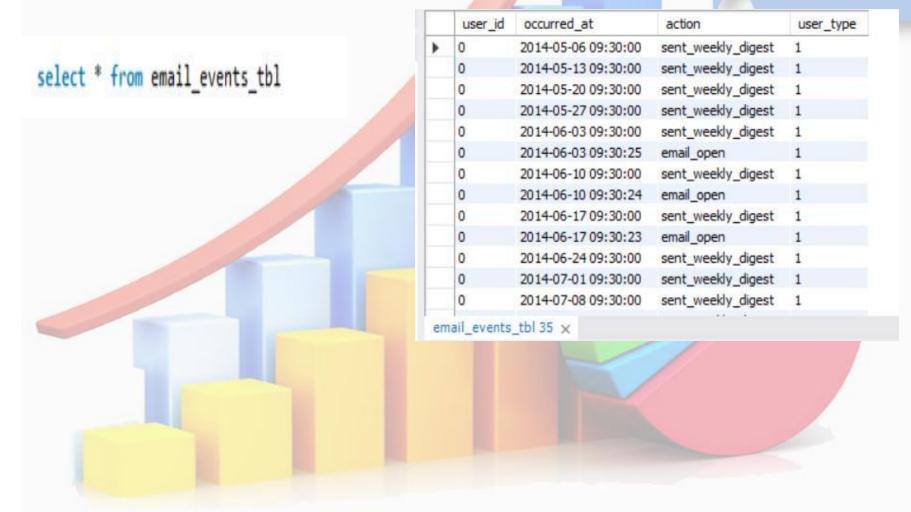
select \* from events\_tbl

0522 0522 0522 0522 0522	2014-05-02 11:02:39 2014-05-02 11:02:53 2014-05-02 11:03:28 2014-05-02 11:04:09 2014-05-02 11:03:16	engagement engagement engagement engagement	login home_page like_message view_inbox	Japan Japan Japan	dell inspiron notebook dell inspiron notebook dell inspiron notebook	3 3 3
0522 0522 0522	2014-05-02 11:03:28 2014-05-02 11:04:09	engagement engagement	like_message	Japan	al little and the second of th	- 7
0522 0522	2014-05-02 11:04:09	engagement	The state of the s	Total States	dell inspiron notebook	3
0522			view_inbox	1		
	2014-05-02 11:03:16	CONTRACTOR OF STREET		Japan	dell inspiron notebook	3
		engagement	search_run	Japan	dell inspiron notebook	3
0522	2014-05-02 11:03:43	engagement	search_run	Japan	dell inspiron notebook	3
0612	2014-05-01 09:59:46	engagement	login	Netherlands	iphone 5	1
0612	2014-05-01 10:00:18	engagement	like_message	Netherlands	iphone 5	1
0612	2014-05-01 10:00:53	engagement	send_message	Netherlands	iphone 5	1
0612	2014-05-01 10:01:24	engagement	home_page	Netherlands	iphone 5	1
0612	2014-05-01 10:01:52	engagement	like_message	Netherlands	iphone 5	1
0612	2014-05-01 10:02:17	engagement	home_page	Netherlands	iphone 5	1
0612	2014-05-01 10:02:51	engagement	view_inbox	Netherlands	iphone 5	1
0	0612 0612 0612 0612	2014-05-01 10:00:53 2014-05-01 10:01:24 2014-05-01 10:01:52 2014-05-01 10:02:17	2014-05-01 10:00:53 engagement 2012 2014-05-01 10:01:24 engagement 2012 2014-05-01 10:01:52 engagement 2012 2014-05-01 10:02:17 engagement	2014-05-01 10:00:53 engagement send_message 2014-05-01 10:01:24 engagement home_page 2014-05-01 10:01:52 engagement like_message 2014-05-01 10:02:17 engagement home_page	2014-05-01 10:00:53         engagement         send_message         Netherlands           2014-05-01 10:01:24         engagement         home_page         Netherlands           2014-05-01 10:01:52         engagement         like_message         Netherlands           2014-05-01 10:02:17         engagement         home_page         Netherlands           2014-05-01 10:02:17         engagement         home_page         Netherlands	0612       2014-05-01 10:00:53       engagement       send_message       Netherlands       iphone 5         0612       2014-05-01 10:01:24       engagement       home_page       Netherlands       iphone 5         0612       2014-05-01 10:01:52       engagement       like_message       Netherlands       iphone 5         0612       2014-05-01 10:02:17       engagement       home_page       Netherlands       iphone 5

events\_tbl 34 ×

### 3. Creating Email Events table

#### " LIVE" DATA



A. User Engagement Objective Calculate the

weekly user engagement

select extract(week from occured\_at) as week\_number,
count(distinct user\_id) as active\_user
from events\_tbl
where event\_type='engagement'
group by week\_number
order by week\_number

TITLE

TITLE

TITLE week\_num num\_users 17 18 1068 19 1113 20 1154 1121 1186 1232 TITLE 1275 25 1264 1302 1372 1365 TITLE TITLE

Week 31 posted the highest user engagement and week 18 posted the minimum user engagement

## B. <u>User Growth:</u> Amount of users growing over time for a product. Your task: Calculate the user growth for product?

```
select year, week_num, num_users, sum(num_users)

over(order by year, week_num) as cum_users

from (
select extract(year from created_at) as year, extract(week from created_at) as week_num, count(distinct user_id) as num_users

from users_tbl

where state='active'

group by year, week_num

order by year, week_num

order by year, week_num)sub
```

The 33th week of 2014 saw the greatest number of users actively engaging with the product or service, while the 35th week of 2014 had the lowest number of active users.

	year	week_num	num_users	cum_users
•	2013	0	23	23
	2013	1	30	53
	2013	2	48	101
	2013	3	36	137
	2013	4	30	167
	2013	5	48	215
	2013	6	38	253
	2013	7	42	295
	2013	8	34	329
	2013	9	43	372
	2013	10	32	404
	2013	11	31	435
	2013	12	33	468

## C. Weekly Retention Objective Calculate the weekly retention of the users sign-up cohort

```
with ctel as (
select distinct user id,
Extract(week from occurred_at) as signup_week
from events tbl
where event type = 'signup flow'
and event_name = 'complete_signup' and extract (week from occurred_at) = 18 ),
cte2 as (select distinct user id,
Extract(week from occurred at) as engagement week
from events tbl
where event type = 'engagement')
select count(user id) total engaged users,
sum(case when retention week > 0 then 1 else 0 end) as retained users
from (select a.user_id, a.signup_week,
  b.engagement week, b.engagement week-a.signup week as retention week
  from ctel a
  LEFT JOIN cte2 b
  on a.user_id = b.user_id
  order by a.user id ) sub
```

total_engaged_users	retained_users
317	236

- 30% of the users retained in week 18 were retained only for the next 7 days.
- User 11816 was retained for the longest duration of 17 weeks

## D. Weekly Engagement ObjectiveCalculate the the weekly engagement per device

with cte as (select extract(year from occurred\_at)||'-'||extract(week from occurred\_at) as weeknum
device, count(distinct user\_id) as usercnt

from events_tbl	weeknum	device	usercnt
where event_type = 'engagement'	2014-18	acer aspire desktop	10
group by weeknum, device	2014-18	acer aspire notebook	21
order by weeknum)	2014-18	amazon fire phone	4
select weeknum, device, usercnt	2014-18	asus chromebook	23
from cte	2014-18	dell inspiron desktop	21

Weeks 31 & 32 of the year 2014 had the highest user engagement of 317 users each week for the product and the device being used was 'MacBook Pro' for both the weeks

## E. <u>Email Engagement Objective Calculate the</u> email engagement metrics

```
select
100 * sum(case when email_cat = 'email_open' then 1 else 0 end)/
    sum(case when email_cat = 'email_sent' then 1 else 0 end) as email_open_rate,
100 * sum(case when email_cat = 'email_clicked' then 1 else 0 end)/
    sum(case when email_cat = 'email_sent' then 1 else 0 end) as email_click_rate

from (select *,
    Case
    When action in ('sent_weekly_digest', 'sent_reengagement_email') then 'email_sent'
    when action in ('email_open') then 'email_open'
    when action in ('email_clickthrough') then 'email_clicked'
    end as email_cat
    from email_events ) sub
```

	email_open_rate	email_dick_rate
١	31.1921	10.4745

 Out of the total emails sent, around 35.73% of them were opened and only 15.74% of those emails were clicked

### **Insights**

- Less than 0.01 jobs were reviewed each hour of the day throughout the month of November.
- 7 day rolling average is best for throughput.
- The Persian Language had the highest share among other languages.
- Out of the total emails sent, around 35.73% of them were opened and only 15.74% of those emails were clicked
- The weekly user engagement is highest in 31<sup>t</sup> week.
- 33<sup>rd</sup> and 35<sup>th</sup> week of 2014 were the highest and lowest of user activity engagement respectively.
- Maximum retained users were only retained for a week, the retention rates dropped weekyweek.
- Users who had the highest engagement with the product were operating or Macbook Pro'.
- During the month of August, users received the highest number of weekly digest emails.

### Result

#### "LIVE" DATA

• The project's key results included the identification of reviewed jobs and their distribution across languages, the calculation of retention rates, and the identification of retained users through an in-depth analysis that relied on predefined assumptions. SQL is one of the most crucial skills for anyone in a data driven position. Additionally, this project helped me to gain insight of various factors which are crucially important for the business to run for a long period and grow as well. Brainstorming is the key to run successful business.