

Project on Operation Analytics and Investigating Metric Spike

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

One of the critical 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, it is our job to answer these questions daily, making it crucial to understand how to investigate these metric spikes.

• Case Study 1: Job Data Analysis

Here, we 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 reviewing the job in seconds.

org: The Organization of the actor

ds: The date in the format yyyy/mm/dd (stored as text)

The problems

A)Jobs Reviewed Over Time:

Objective: Calculate the number of jobs reviewed per hour for each day in November 2020.

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

B)Throughput Analysis:

Objective: Calculate the 7-day rolling average of throughput (number of events per second).

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

C)Language Share Analysis:

Objective: Calculate the percentage share of each language in the last 30 days.

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

D)Duplicate Rows Detection:

Objective: Identify duplicate rows in the data.

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

• Case Study 2: Investigating Metric Spike

Here, we 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.

1. Weekly User engagement:

Objective: Measure the activeness of users on a weekly basis.

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

2. User Growth Analysis:

Objective: Analyze the growth of users over time for a product.

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

3. Weekly Retention Analysis:

Objective: Analyze the retention of users on a weekly basis after signing up for a product.

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

4. Weekly Engagement Per Device:

Objective: Measure the activeness of users on a weekly basis per device.

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

5. Email Engagement Analysis:

Objective: Analyze how users are engaging with the email service.

Task: Write an SQL query to calculate the email engagement metrics.

Database and tools:

- ☐ MySQL
- ☐ MySQL Workbench 8.0 CE

Software used for visualisation:

- ☐ Microsoft Excel

Queries and outputs of Case study-I

- Query to Create database and table:

```
1 • create database project3;
2 • show databases;
3 • use project3;
4 • create table job_data (
5     ds datetime,
6     job_id int,
7     actor_id int,
8     event varchar(50),
9     language varchar(100),
10    time_spent int,
11    org varchar(10)
12 );
13
```

A)Jobs Reviewed Over Time:

```
17 • SELECT
18     ds, COUNT(job_id), sum(time_spent) / 3600 AS time_spent_in_hr
19 FROM
20     job_data
21     group by 1
22     order by 1 asc;
23
24
```

Result Grid	Filter Rows:	Export:	Wrap Cell Content:
ds	COUNT(job_id)	time_spent_in_hr	
11/25/2020	1	0.0125	
11/26/2020	1	0.0156	
11/27/2020	1	0.0289	
11/28/2020	2	0.0092	
11/29/2020	1	0.0056	
11/30/2020	2	0.0111	

B)Throughput Analysis:

```
25     # Calculate the 7-day rolling average of throughput (number of events per second)
26
27 • WITH throughput AS
28     ( SELECT ds, COUNT(job_id) AS total_job, SUM(time_spent) AS total_time
29     FROM job_data
30     GROUP BY 1)
31     SELECT ds, SUM(total_job) OVER (ORDER BY ds ROWS BETWEEN 6 PRECEDING AND CURRENT ROW) /
32     SUM(total_time) OVER (ORDER BY ds ROWS BETWEEN 6 PRECEDING AND CURRENT ROW)
33     AS 7_days_rolling_average FROM throughput;
34
```

Result Grid	Filter Rows:	Export:	Wrap Cell Content:
ds	7_days_rolling_average		
11/25/2020	0.0222		
11/26/2020	0.0198		
11/27/2020	0.0146		
11/28/2020	0.0210		
11/29/2020	0.0233		
11/30/2020	0.0268		

Throughput analysis by 7-day rolling average gives a general idea about the average of consecutive 7 days whereas the daily metric gives information about data on a daily basis.

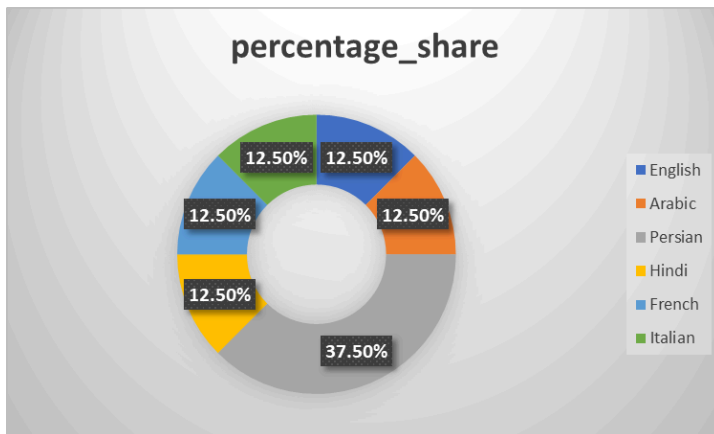
In some cases, where we need an idea about how a parameter is changing over the course of time, it's better to opt for the rolling average than the daily metric.

C)Language Share Analysis:

```
36  #Calculate the percentage share of each language in the last 30 days
37
38  • with count_of_language as
39  (
40  select
41  language,
42  count(language) as languages
43  from job_data
44  group by 1
45  )
46  select language,
47  concat(languages*100/(select sum(languages) as total from count_of_language),'%') as percentage_share
48  from count_of_language
49  group by 1;
```

language	percentage_share
English	12.5000%
Arabic	12.5000%
Persian	37.5000%
Hindi	12.5000%
French	12.5000%
Italian	12.5000%

The percentage share of Persian is the maximum among all the languages and its value is 3 times the percentage share of all other languages.



D)Duplicate Rows Detection:

```
51 # Duplicate rows:
52
53 • with duplicate_value as
54 (
55     select ds, job_id, actor_id,
56     row_number() over(partition by job_id order by job_id) as row_no
57     from job_data
58 )
59 select*
60 from duplicate_value
61 where row_no > 1;
62
```

<				
Result Grid				
Filter Rows:				
Export:				
Wrap Cell Content:				
	ds	job_id	actor_id	row_no
▶	11/28/2020	23	1005	2
	11/26/2020	23	1004	3

Query and Output of Case study-II

- Query to Create database, table and loading data into the table

```
1 • use project3;
2 • create table users
3 • (
4 •     user_id int,
5 •     created_at varchar(100),
6 •     company_id int,
7 •     language varchar(50),
8 •     activated_at varchar(100),
9 •     state varchar(50)
10 • );
11
12
13 • SHOW VARIABLES LIKE 'secure_file_priv';
14
15 • LOAD DATA INFILE "C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/users.csv"
16 INTO TABLE users
17 FIELDS TERMINATED BY ','
18 ENCLOSED BY '"'
19 LINES TERMINATED BY '\n'
20 IGNORE 1 ROWS;
21
22 • alter table users add column temp_created_at datetime;
23 • update users set temp_created_at = str_to_date(created_at, '%d-%m-%Y %H:%i:');
24 • alter table users drop column created_at;
25 • alter table users change temp_created_at created_at datetime;
26 • alter table users add column temp_activated_at datetime;
27 • update users set temp_activated_at = str_to_date(activated_at, '%d-%m-%Y %H:%i:');
28 • alter table users drop column activated_at;
29 • alter table users change temp_activated_at activated_at datetime;
30 • create table events(
31 •     user_id int,
32 •     occurred_at varchar(100),
33 •     event_type varchar(100),
34 •     event_name varchar(50),
35 •     location varchar(50),
36 •     device varchar(50),
37 •     user_type int
38 • );
39 • LOAD DATA INFILE "C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/events.csv"
40 INTO TABLE events
41 FIELDS TERMINATED BY ','
42 ENCLOSED BY '"'
```

```

41     LINES TERMINATED BY '\n'
42     IGNORE 1 ROWS;
43
44
45 •   alter table events add column temp_created_at datetime;
46 •   update events set temp_created_at = str_to_date(occurred_at, '%d-%m-%Y %H:%i:');
47 •   alter table events drop column occurred_at;
48 •   alter table events change temp_created_at occurred_at datetime;
49
50 •   create table email_events
51   (
52     user_id int,
53     occurred_at varchar(100),
54     action varchar(100),
55     user_type int
56   );
57
58 •   LOAD DATA INFILE "C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/email_events.csv"
59   INTO TABLE email_events
60   FIELDS TERMINATED BY ','
61   ENCLOSED BY '"'
62   LINES TERMINATED BY '\n'
63   IGNORE 1 ROWS;
64
65 •   alter table email_events add column temp_created_at datetime;
66 •   update email_events set temp_created_at = str_to_date(occurred_at, '%d-%m-%Y %H:%i:');
67 •   alter table email_events drop column occurred_at;
68 •   alter table email_events change temp_created_at occurred_at datetime;
--

```

1. Weekly User engagement:


```

75     # weekly user engagement
76
77     • select extract(week from occurred_at) as week_no,
78           count(distinct user_id)
79           from events
80           where event_type='engagement'
81           group by 1
82           order by 1;

```

Result Grid		
Filter Rows:		
Export: Wrap Cell Content:		
	week_no	count(distinct user_id)
▶	17	663
	18	1068
	19	1113
	20	1154
	21	1121
	22	1186
	23	1232
	24	1275
	25	1264
	26	1302
	27	1372
	28	1365
	29	1376

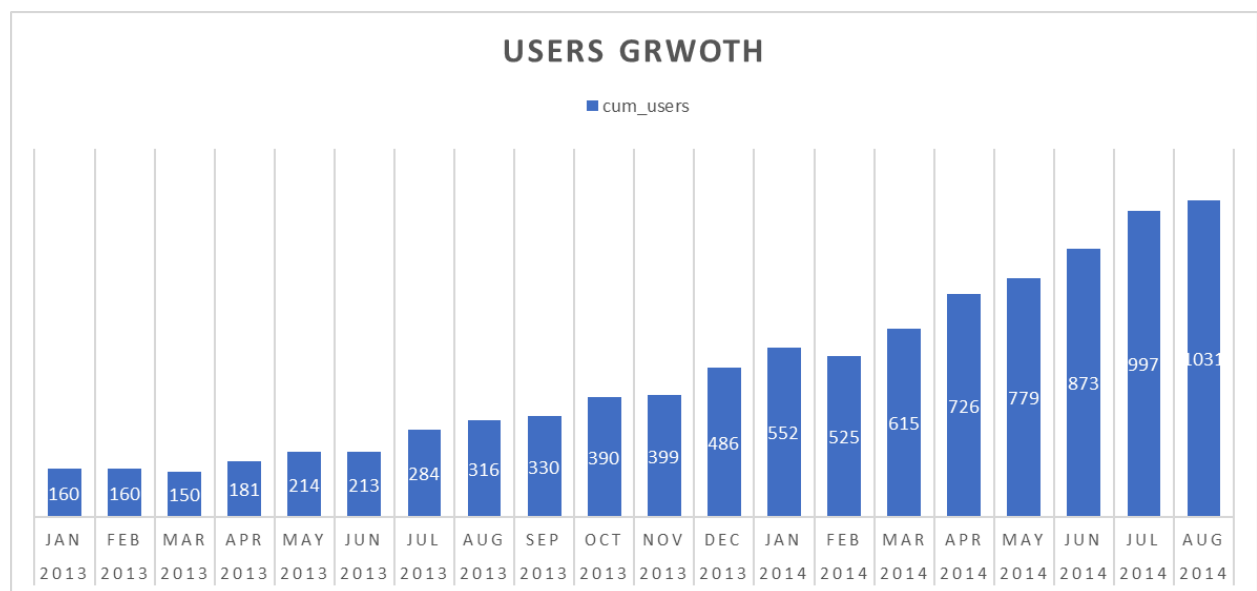
2. User Growth Analysis:

```

82     # USER GROWTH ANALYSIS
83
84     • with growth as
85     (
86     select
87     extract(year from activated_at) as year_no,
88     extract(month from activated_at) as month_no,
89     count(distinct user_id) as cum_users
90     from users
91     group by 1,2
92     )
93     select*, sum(cum_users) over(order by year_no, month_no rows between unbounded preceding and current row) as cumulative_users
94     from growth;

```

year_no	month_no	cum_users	cumulative_users
2013	1	160	160
2013	2	160	320
2013	3	150	470
2013	4	181	651
2013	5	214	865
2013	6	213	1078
2013	7	284	1362
2013	8	316	1678
2013	9	330	2008
2013	10	390	2398
2013	11	399	2797
2013	12	486	3283
2014	1	552	3835
2014	2	525	4360
2014	3	615	4975
2014	4	726	5701
2014	5	779	6480
2014	6	873	7353
2014	7	997	8350
2014	8	1031	9381






3. Weekly Retention Analysis:

```

96     # WEEKLY USER RETENTION
97
98     • SELECT
99     a.cohort_week,
100     b.engagement_week,
101     count(distinct a.user_id) as engagement
102     FROM
103     ((SELECT distinct user_id, extract(week from occurred_at) as cohort_week from events
104     WHERE event_type = 'signup_flow'
105     and event_name = 'complete_signup') as a
106     LEFT JOIN
107     (SELECT distinct user_id, extract(week from occurred_at) as engagement_week FROM events
108     where event_type = 'engagement') as b
109     on a.user_id = b.user_id)
110     group by 1,2
111     order by 1,2

```



Result Grid  Filter Rows: <input type="text"/> Export:  Wrap Cell Content: 			
	cohort_week	engagement_week	engagement
▶	17	17	72
	17	18	59
	17	19	24
	17	20	16
	17	21	11
	17	22	16
	17	23	11
	17	24	9
	17	25	6
	17	26	8
	17	27	8
	17	28	8

4. Weekly Engagement Per Device:

```

111 # weekly engagement per device
112 • select
113     extract(year from occurred_at) as year_no,
114     extract(week from occurred_at) as week_no,
115     device,
116     count(distinct user_id)
117 from events
118 where event_type='engagement'
119 group by 1,2,3
120 order by 1;

```

Result Grid				
Filter Rows: <input type="text"/>				
Export:  Wrap Cell Content: 				
	year_no	week_no	device	count(distinct user_id)
▶	2014	17	acer aspire desktop	9
	2014	17	acer aspire notebook	20
	2014	17	amazon fire phone	4
	2014	17	asus chromebook	21
	2014	17	dell inspiron desktop	18
	2014	17	dell inspiron notebook	46
	2014	17	hp pavilion desktop	14
	2014	17	htc one	16
	2014	17	ipad air	27
	2014	17	ipad mini	19
	2014	17	iphone 4s	21
	2014	17	iphone 5	65
	2014	17	iphone 5s	42
	2014	17	kindle fire	6
	2014	17	lenovo thinkpad	86
	2014	17	mac mini	6
	2014	17	macbook air	54

Knowing the engagement per device can help us make inferences like what device most of the users are using and if a certain device has fewer users then it can be an indication that we need to improve our services for those devices.

5. Email Engagement Analysis:

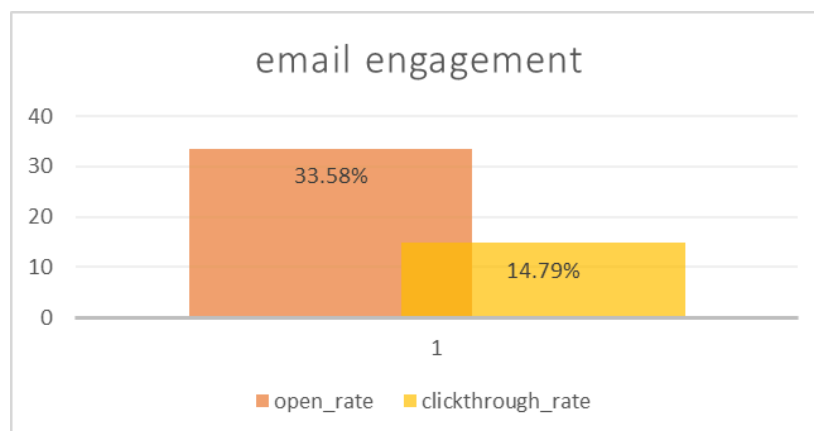
```

125  #email engagement
126
127  select
128    100.0*sum(case when emails='opened' then 1 else 0 end)/ sum(case when emails='sent' then 1 else 0 end ) as open_rate,
129    100.0*sum(case when emails='clicked' then 1 else 0 end)/ sum(case when emails='sent' then 1 else 0 end ) as clickthrough_rate
130  from
131    (select
132      case
133        when action='email_clickthrough' then 'clicked'
134        when action='email_open' then 'opened'
135        when action in ( 'sent_reengagement_email' , 'sent_weekly_digest' ) then 'sent'
136      end as emails
137    from email_events)x;

```

Result Grid	Filter Rows:	Export:	Wrap Cell Content:
open_rate	clickthrough_rate		
33.58339	14.78989		

Here, it shows that the opening rate of emails is 33.58%, i.e. 33 out of 100 users have opened the email. The Click-through rate (CTR) is 14.78%, which is the percentage of users who click on the link in the email. This CTR here, indicates that we need to have a more catchy subject line and captivating body of our email personalized for each user to make our users click the link and boost our CTR.



Conclusion

We can conclude that analyzing operation analytics and investigating metric spikes is crucial for understanding business growth so that we can work on our areas of improvement to provide a better experience to users/clients.

So, it is very important that any organization analyze it on a periodic basis(weekly or monthly or quarterly or yearly) based on their needs for making better decisions for their company/business growth.