# Project: Operation Analytics and Investigating Metric Spike

NAME: MUKESH CHANDRA KAMILA

**Description:**

In this project, I was designated as Data Analyst Lead and provided with different data sets. I was focusing on analysing the data which is provided by the company. So, the task is to derive insights and answer the questions asked by different departments. So that these insights are then used by operation team, support team and marketing team to predict the overall growth or decline of the company. It means better automation, better understanding between cross effective workflows.

In Case Study 1, there is job_data table wherever in Case Study 2 there are users, events and email_event tables.

**Approach:**

To execute the project, I used SQL. So, the first approach was to run SQL queries to create a database using the raw data provided. Once the database was created, I imported the raw data and run various sorting and data extracting queries to get the required insights.

**Tech-Stack:** PostgreSQL 6.19

**Used:**

PostgreSQL 6.19 was used in this project execution. The ease of access and set up with convenient user interface made it a good tool for the project.

**Insights:**

In this project, I learned about advanced SQL. And how to analyse the problem statement and the functions I can use in SQL to solve the problem statement and write the queries to get required output. Following are the questions that has been answered with the help of SQL.

**Result:**

In this project, I have achieved and gain knowledge how to deal with data with help of SQL. And how to interact and run queries to get desired output from database.

## Operation Analytics

**Q1.** Task: Calculate the number of jobs reviewed per hour per day for November 2020?

**Query:**

select ds as Dates, count(job_id)::decimal/sum(time_spent)*3600

as num_jobs_reviewed

from job_data

where ds between '2020-11-01' and '2020-11-30'

group by ds;

**Q2.** Task: Calculate 7 day rolling average of throughput?

For throughput, do you prefer daily metric or 7-day rolling and why?

**Query:**

**--Weekly throughput**

select round(cast(count(event)as decimal)/sum(time_spent),2)

as weekly_throughput

from job_data;

--or

select round(count(event)::decimal /sum(time_spent),2) as weekly_throughput

from job_data;

**--Daily throughput**

select ds as Dates,round(count(event)::decimal/sum(time_spent),2)

as daily_throughput

from job_data

group by ds

order by ds;

File   Object   Tools   Help

Browser   |   Dashboard   Properties   SQL   Statistics   Dependencies   Dependents   Processes   Operation_Analytics.sql
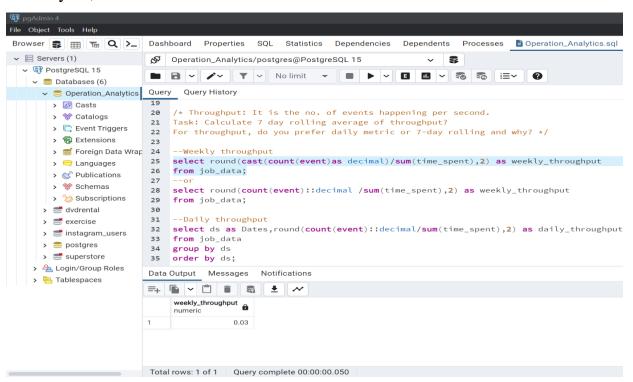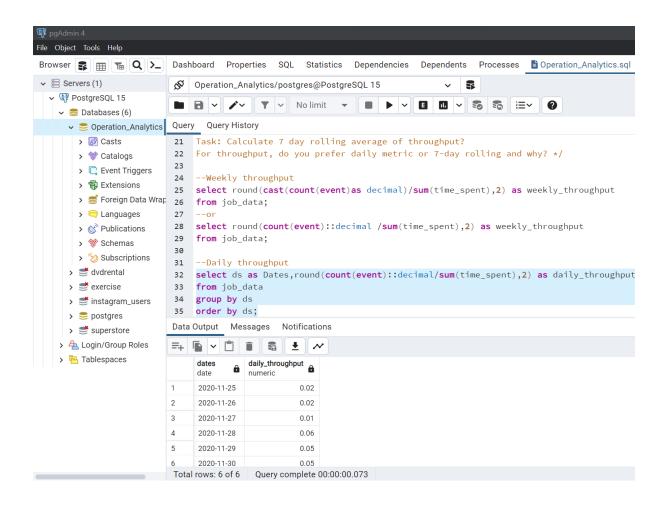
Servers (1)
  PostgreSQL 15
    Databases (6)
      Operation_Analytics
        Casts
        Catalogs
        Event Triggers
        Extensions
        Foreign Data Wrap
        Languages
        Publications
        Schemas
        Subscriptions
      dvdrental
      exercise
      instagram_users
      postgres
      superstore
    Login/Group Roles
    Tablespaces

Operation_Analytics/postgres@PostgreSQL 15

No limit

Query   Query History

```sql
21  Task: Calculate 7 day rolling average of throughput?
22  For throughput, do you prefer daily metric or 7-day rolling and why? */
23
24  --Weekly throughput
25  select round(cast(count(event)as decimal)/sum(time_spent),2) as weekly_throughput
26  from job_data;
27  --or
28  select round(count(event)::decimal /sum(time_spent),2) as weekly_throughput
29  from job_data;
30
31  --Daily throughput
32  select ds as Dates,round(count(event)::decimal/sum(time_spent),2) as daily_throughput
33  from job_data
34  group by ds
35  order by ds;
```

Data Output   Messages   Notifications

| | dates<br>date | daily_throughput<br>numeric |
|---|---|---|
| 1 | 2020-11-25 | 0.02 |
| 2 | 2020-11-26 | 0.02 |
| 3 | 2020-11-27 | 0.01 |
| 4 | 2020-11-28 | 0.06 |
| 5 | 2020-11-29 | 0.05 |
| 6 | 2020-11-30 | 0.05 |

Total rows: 6 of 6      Query complete 00:00:00.073
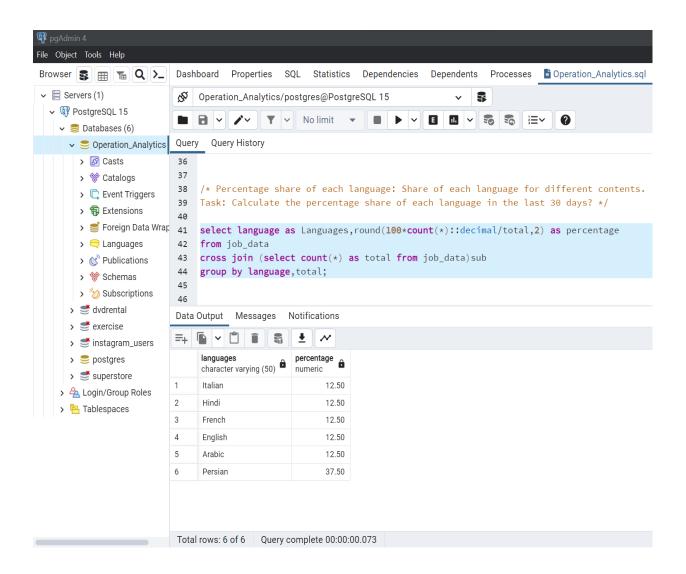
**Page | 5**

**Q3.** Task: Calculate the percentage share of each language in the last 30 days?

**Query:**

select language as Languages,round(100*count(*)::decimal/total,2)

as percentage

from job_data

cross join (select count(*) as total from job_data)sub

group by language,total;

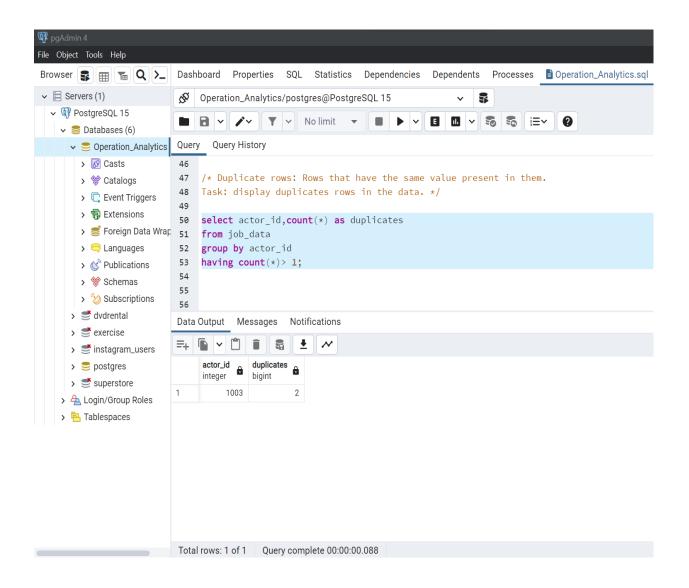**Q4.** Task: display duplicates rows in the data.

**Query:**

select actor_id,count(*) as duplicates

from job_data

group by actor_id

having count(*)> 1;

# Investigating Metric Spike

**Q1.** Task: Calculate the weekly user engagement?

**Query:**

select extract (week from occurred_at) as week_numbers,

    count (distinct user_id) as weekly_active_users

from events

where event_type = 'engagement'

group by 1;

**Q2.** Task: Calculate the user growth for product?

**Query:**

select year, num_week, num_active_users,

sum(num_active_users) over(order by year, num_week rows between unbounded preceding and current row)

as cumm_active_users

from

(select

　　extract(year from a.activated_at) as year,

　　extract(week from a.activated_at)as num_week,

　　count(distinct user_id) as num_active_users

from users a

where state='active'

group by year, num_week

order by year, num_week

)a;

**Q3.** Task: Calculate the weekly retention of users-sign up cohort?

**Query:**

```
select count(user_id),
     sum(case when retention_week = 1 then 1 else 0 end) as per_week_retention
from
(
select a.user_id,
     a.sign_up_week,
     b.engagement_week,
     b.engagement_week - a.sign_up_week as retention_week
from
(
(select distinct user_id, extract(week from occurred_at) as sign_up_week
from events
where event_type = 'signup_flow'
and event_name = 'complete_signup'
and extract(week from occurred_at)=18)a
left join
(select distinct user_id, extract(week from occurred_at) as engagement_week
from events
where event_type = 'engagement')b
on a.user_id = b.user_id
))sub
group by user_id
order by user_id;
```

```
37  select count(user_id),
38         sum(case when retention_week = 1 then 1 else 0 end) as per_week_retention
39  from
40  (
41  select a.user_id,
42         a.sign_up_week,
43         b.engagement_week,
44         b.engagement_week - a.sign_up_week as retention_week
45  from
46  (
47  (select distinct user_id, extract(week from occurred_at) as sign_up_week
48  from events
49  where event_type = 'signup_flow'
50  and event_name = 'complete_signup'
51  and extract(week from occurred_at)=18)a
52  left join
53  (select distinct user_id, extract(week from occurred_at) as engagement_week
54  from events
55  where event_type = 'engagement')b
56  on a.user_id = b.user_id
57  ))sub
58  group by user_id
59  order by user_id;
```

Data Output | Messages | Notifications

| | count bigint | per_week_retention bigint |
|---|---|---|
| 1 | 1 | 0 |
| 2 | 1 | 0 |
| 3 | 2 | 1 |
| 4 | 3 | 0 |
| 5 | 5 | 1 |
| 6 | 2 | 1 |
| 7 | 1 | 0 |
| 8 | 3 | 1 |
| 9 | 2 | 1 |
| 10 | 6 | 1 |
| 11 | 2 | 1 |

Total rows: 81 of 81    Query complete 00:00:00.270

**Q4.** Task: Calculate the weekly engagement per device?

**Query:**

select

    extract(year from occurred_at) as num_years,

    extract(week from occurred_at) as num_weeks,

    device,

    count(distinct user_id) as num_of_users

from events

where event_type = 'engagement'

group by 1,2,3

order by1,2,3;

**Q5.** Task: Calculate the email engagement metrics?

**Query:**

```
select
100.0 * sum(case when email_cat = 'email_opened' then 1 else 0 end)
     /sum(case when email_cat = 'email_sent' then 1 else 0 end)
as email_opening_rate,
100.0 * 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_clicking_rate
from
(
select *,
case when action in ('sent_weekly_digest', 'sent_reengagement_email')
    then 'email_sent'
    when action in ('email_open')
    then 'email_opened'
    when action in ('email_clickthrough')
    then 'email_clicked'
end as email_cat
from email
)a;
```

File  Object  Tools  Help

Browser  Dashboard  Properties  SQL  Statistics  Dependencies  Dependents  Processes  📄 Investigating Metric Spike.sql*

Operation_Analytics/postgres@PostgreSQL 15

No limit

Query    Query History

```
76   Task: Calculate the email engagement metrics? */
77   select
78   100.0 * sum(case when email_cat = 'email_opened' then 1 else 0 end)
79         /sum(case when email_cat = 'email_sent' then 1 else 0 end)
80   as email_opening_rate,
81   100.0 * sum(case when email_cat = 'email_clicked' then 1 else 0 end)
82         /sum(case when email_cat = 'email_sent' then 1 else 0 end)
83   as email_clicking_rate
84   from
85   (
86   select *,
87   case when action in ('sent_weekly_digest', 'sent_reengagement_email')
88       then 'email_sent'
89       when action in ('email_open')
90       then 'email_opened'
91       when action in ('email_clickthrough')
92       then 'email_clicked'
93   end as email_cat
94   from email
95   )a;
```

**Servers (1)**
- **PostgreSQL 15**
  - **Databases (6)**
    - **Operation_Analytics**
      - Casts
      - Catalogs
      - Event Triggers
      - Extensions
      - Foreign Data Wrap
      - Languages
      - Publications
      - Schemas
      - Subscriptions
    - dvdrental
    - exercise
    - instagram_users
    - postgres
    - superstore
  - Login/Group Roles
  - Tablespaces

Data Output    Messages    Notifications

| | email_opening_rate 🔒 numeric | email_clicking_rate 🔒 numeric |
|---|---|---|
| 1 | 33.5833880499015102 | 14.7898883782009192 |

Total rows: 1 of 1    Query complete 00:00:00.105