**Case Study 1: Job Data Analysis**

## Tasks:

A. **Jobs Reviewed Over Time:**
   o Objective: Calculate the number of jobs reviewed per hour for each day in November 2020.
   o Your Task: Write an SQL query to calculate the number of jobs reviewed per hour for each day in November 2020.

Ans:

select ds,count(job_id) as jobs_per_day, sum(time_spent)/3600 as hours_spent
from job_data
where ds >='2020-11-01'  and ds <='2020-11-30'
group by ds ;

B. **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.

Ans:

select ds, count(job_id) as jobs_per_day, sum(time_spent)/3600 as hours_spent
from job_data
where ds >='2020-11-01'  and ds <='2020-11-30'
group by ds;



The daily metric provides accurate values for each day, reflecting the actual throughput for that specific day, but

The rolling average helps smooth out short-term fluctuations, making it easier to identify trends and patterns over a longer period.

If you need a more granular view of **daily fluctuations**, then daily metric might be more appropriate.
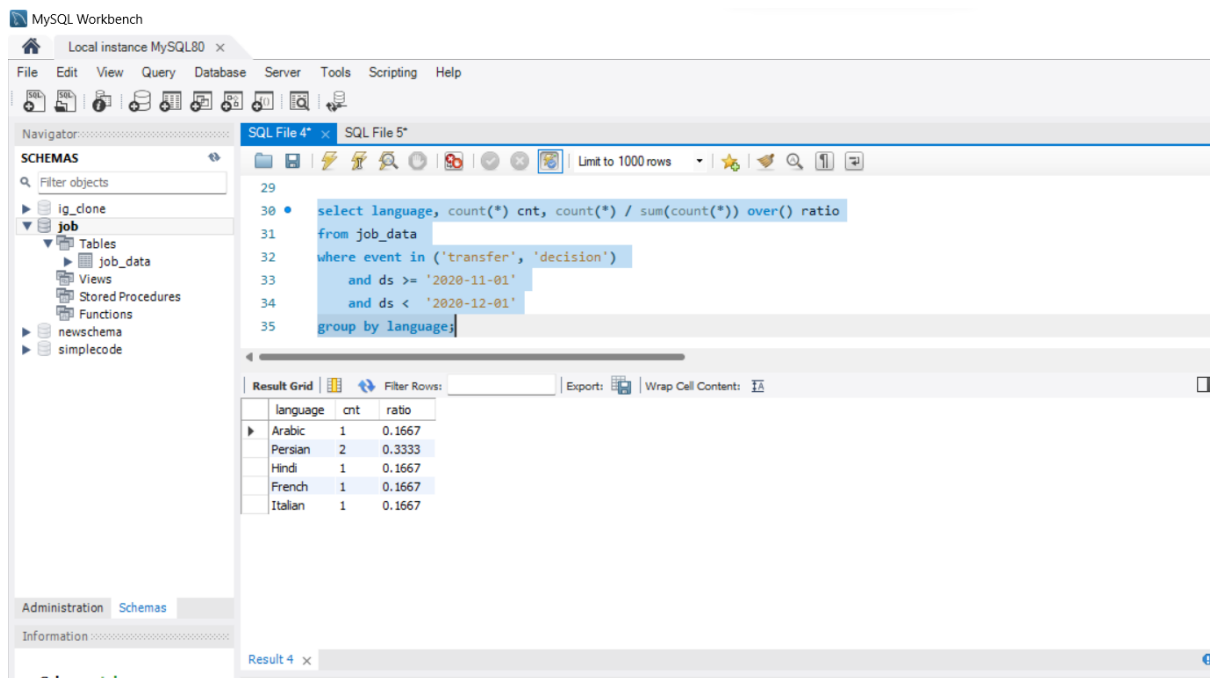
If you're interested in **identifying trends** over a more extended period while smoothing out short-term variations, the 7-day rolling average might be a better choice.

C. **Language Share Analysis:**

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

Ans:

```
select language, count(*) cnt, count(*) / sum(count(*)) over() ratio
from job_data
where event in ('transfer', 'decision')
    and ds >= '2020-11-01'
    and ds <  '2020-12-01'
group by language
```

D. **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.

Ans:

```
SELECT
    actor_id,
    COUNT(DISTINCT language) AS language_count
FROM job_data
GROUP BY actor_id
HAVING COUNT(DISTINCT language) > 1;
```