

Case Study 1: Job Data Analysis

- Creating database and table

```
1 • create database casestudy_1;
2 • CREATE TABLE job_data
3 (
4     ds DATE,
5     job_id INT NOT NULL,
6     actor_id INT NOT NULL,
7     event VARCHAR(15) NOT NULL,
8     language VARCHAR(15) NOT NULL,
9     time_spent INT NOT NULL,
10    org CHAR(2)
11 );
12
13 • INSERT INTO job_data (ds, job_id, actor_id, event, language, time_spent, org)
14 VALUES ('2020-11-30', 21, 1001, 'skip', 'English', 15, 'A'),
15         ('2020-11-30', 22, 1006, 'transfer', 'Arabic', 25, 'B'),
16         ('2020-11-29', 23, 1003, 'decision', 'Persian', 20, 'C'),
17         ('2020-11-28', 23, 1005, 'transfer', 'Persian', 22, 'D'),
18         ('2020-11-28', 25, 1002, 'decision', 'Hindi', 11, 'B'),
19         ('2020-11-27', 11, 1007, 'decision', 'French', 104, 'D'),
20         ('2020-11-26', 23, 1004, 'skip', 'Persian', 56, 'A'),
21         ('2020-11-25', 20, 1003, 'transfer', 'Italian', 45, 'C');
```

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.

```
• SELECT DATE(ds) AS review_date,
        COUNT(*) AS jobs_reviewed_per_day, sum(time_spent)/3600 as jobs_reviewed_per_hr
FROM job_data
WHERE ds between '2020-11-01' and '20-11-30'
GROUP BY review_date
ORDER BY review_date;
```

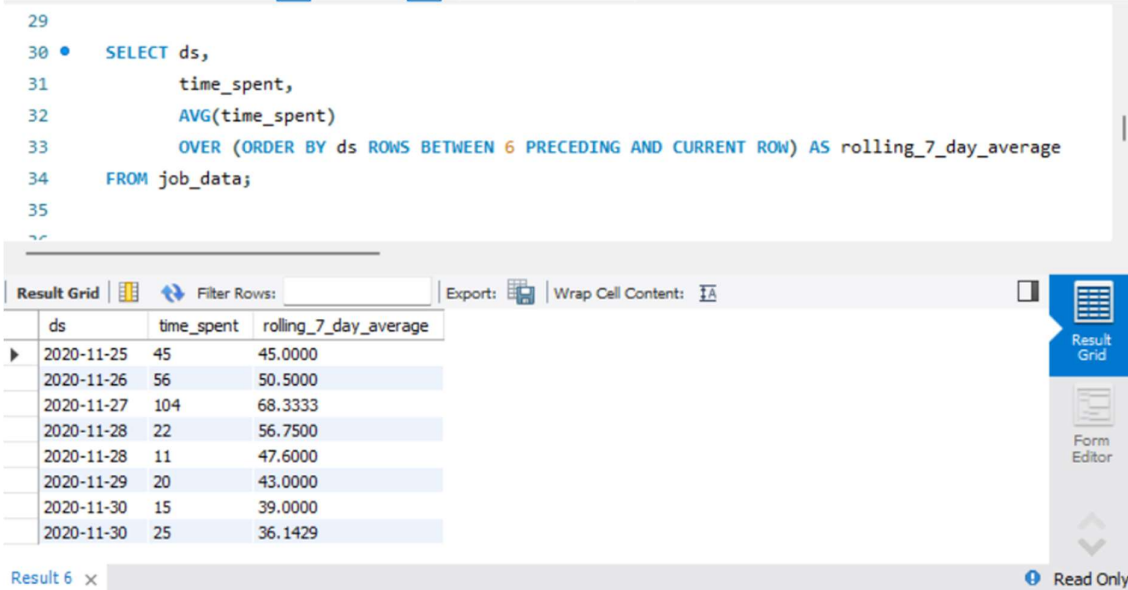
Result Grid	Filter Rows:	Export:	Wrap Cell Content:
review_date	jobs_reviewed_per_day	jobs_reviewed_per_hr	
2020-11-25	1	0.0125	
2020-11-26	1	0.0156	
2020-11-27	1	0.0289	
2020-11-28	2	0.0092	
2020-11-29	1	0.0056	
2020-11-30	2	0.0111	

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.



The screenshot shows a SQL query editor with the following query:

```
29
30 • SELECT ds,
31       time_spent,
32       AVG(time_spent)
33       OVER (ORDER BY ds ROWS BETWEEN 6 PRECEDING AND CURRENT ROW) AS rolling_7_day_average
34 FROM job_data;
35
```

Below the query editor is a 'Result Grid' showing the results of the query. The grid has three columns: 'ds', 'time_spent', and 'rolling_7_day_average'. The data is as follows:

ds	time_spent	rolling_7_day_average
2020-11-25	45	45.0000
2020-11-26	56	50.5000
2020-11-27	104	68.3333
2020-11-28	22	56.7500
2020-11-28	11	47.6000
2020-11-29	20	43.0000
2020-11-30	15	39.0000
2020-11-30	25	36.1429

The interface also includes a 'Filter Rows' field, an 'Export' button, a 'Wrap Cell Content' checkbox, and a 'Form Editor' button. The status bar at the bottom indicates 'Result 6' and 'Read Only'.

Reason:

Daily Metric: This provides the raw, day-to-day values of throughput without any smoothing or averaging. It might be useful when examining short-term trends or fluctuations. Daily metrics are sensitive to sudden changes and can provide insights into specific days' performances.

7-Day Rolling Average: This smoothed metric helps in identifying overall trends and patterns by averaging the data over a 7-day period, reducing the impact of daily fluctuations or irregularities. It's particularly useful for understanding long-term trends, identifying broader patterns, and eliminating noise in the data.

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.

```

36 • select language,
37       round(100*count(*)/(select count(*) from job_data),2) as percentage_share
38 from job_data
39 group by language
40 order by language desc;

```

language	percentage_share
Persian	37.50
Italian	12.50
Hindi	12.50
French	12.50
English	12.50
Arabic	12.50

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.

```

42 • select *
43 from
44 (select *,
45  row_number() over(partition by ds,actor_id,job_id) as row_num
46  from job_data) a
47 where row_num>1;

```

ds	job_id	actor_id	event	language	time_spent	org	row_num
----	--------	----------	-------	----------	------------	-----	---------