

OPERATION ANALYTICS AND INVESTIGATING METRIC SPIKE

TRAINITY PROJECT #3

Submission by :

ANANAYA LAL (Data Analyst Trainee)

PROJECT DESCRIPTION :

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

Our goal is to use advanced SQL skills to analyze the data and provide valuable insights that can help improve the company's operations and understand sudden changes in key metrics.

This project involves conducting a comprehensive analysis of the end-to-end operations of a company based on two case studies. There were 2 case studies given to work upon. A brief description of both case studies and output required from both of them are given further.

Case Studies

Case Study 1: Job Data Analysis

Table name - **job_data**

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 to review the job in seconds.
- **org**: The Organization of the actor
- **ds**: The date in the format yyyy/mm/dd (stored as text).

Case Study 2: Investigating Metric Spike

Table names -

- **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.

Tasks

Case Study 1

- **Jobs Reviewed Over Time**
- **Throughput Analysis**
- **Language Share Analysis**
- **Duplicate Rows Detection**

Case Study 2

- **Weekly User Engagement**
- **User Growth Analysis**
- **Weekly Retention Analysis**
- **Weekly Engagement Per Device**
- **Email Engagement Analysis**

Important Informations

Approach :

- To conduct a comprehensive analysis of a company's end-to-end operations based on case studies, we would typically follow a structured approach that involves several key steps like , define the objective , collect data , analyse data , identify areas of improvement etc .

Tech-Stack Used :

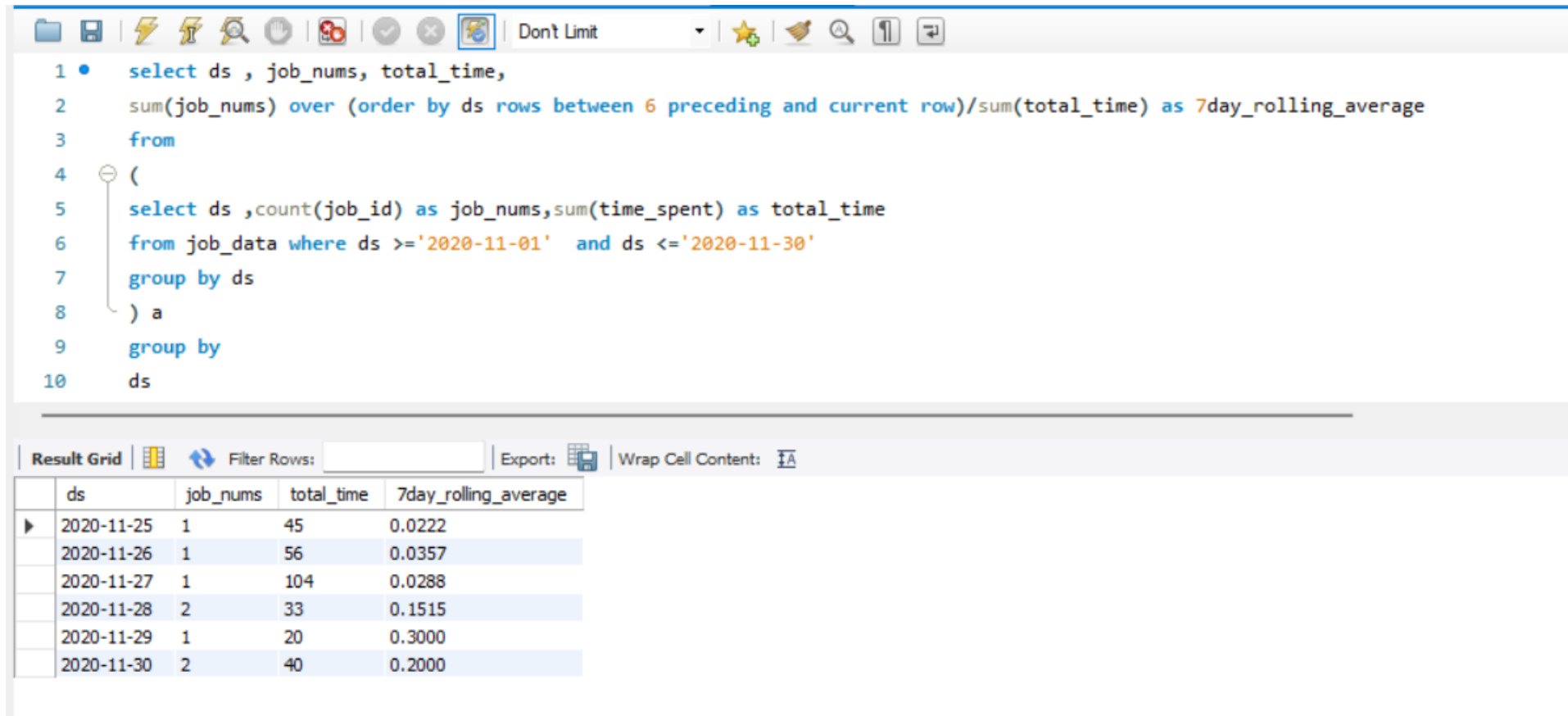
- MySQL
- PostgreSQL
- MS Excel
- MS Word
- Powerpoint Ppt

C1T1 . Amount of jobs reviewed over time.

```
1 • use opr_n_ma ;  
2 • SELECT ((COUNT(job_id))/sum(time_spent)/3600) as job_reviewed  
3 FROM job_data
```

Result Grid		Filter Rows:	Export:	Wrap Cell Content:
	job_reviewed			
▶	0.00000746			

C1T2 . No. of events happening per second – to calculate 7 day rolling average.



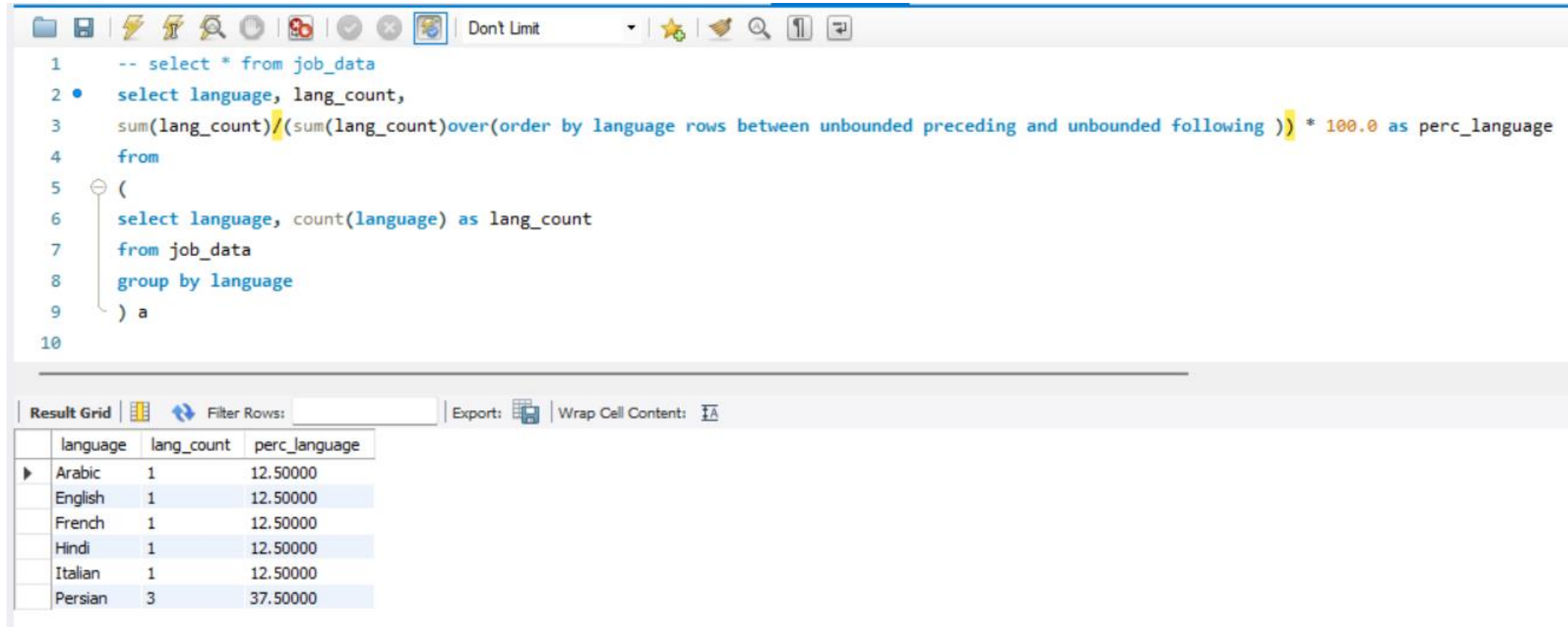
The screenshot shows a SQL IDE interface. The top toolbar includes icons for file operations, execution, and a 'Don't Limit' button. The SQL editor contains a query that calculates a 7-day rolling average of events per second. The query is as follows:

```
1 • select ds , job_nums, total_time,  
2    sum(job_nums) over (order by ds rows between 6 preceding and current row)/sum(total_time) as 7day_rolling_average  
3    from  
4    (  
5      select ds ,count(job_id) as job_nums,sum(time_spent) as total_time  
6      from job_data where ds >='2020-11-01' and ds <='2020-11-30'  
7      group by ds  
8    ) a  
9    group by  
10   ds
```

Below the editor is the 'Result Grid' section, which includes a 'Filter Rows' input field, an 'Export' button, and a 'Wrap Cell Content' checkbox. The results are displayed in a table with the following data:

	ds	job_nums	total_time	7day_rolling_average
▶	2020-11-25	1	45	0.0222
	2020-11-26	1	56	0.0357
	2020-11-27	1	104	0.0288
	2020-11-28	2	33	0.1515
	2020-11-29	1	20	0.3000
	2020-11-30	2	40	0.2000

C1T3 . Share of each language for different contents.



```
1  -- select * from job_data
2  • select language, lang_count,
3     sum(lang_count)/(sum(lang_count)over(order by language rows between unbounded preceding and unbounded following )) * 100.0 as perc_language
4  from
5  (
6     select language, count(language) as lang_count
7     from job_data
8     group by language
9  ) a
10
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: |

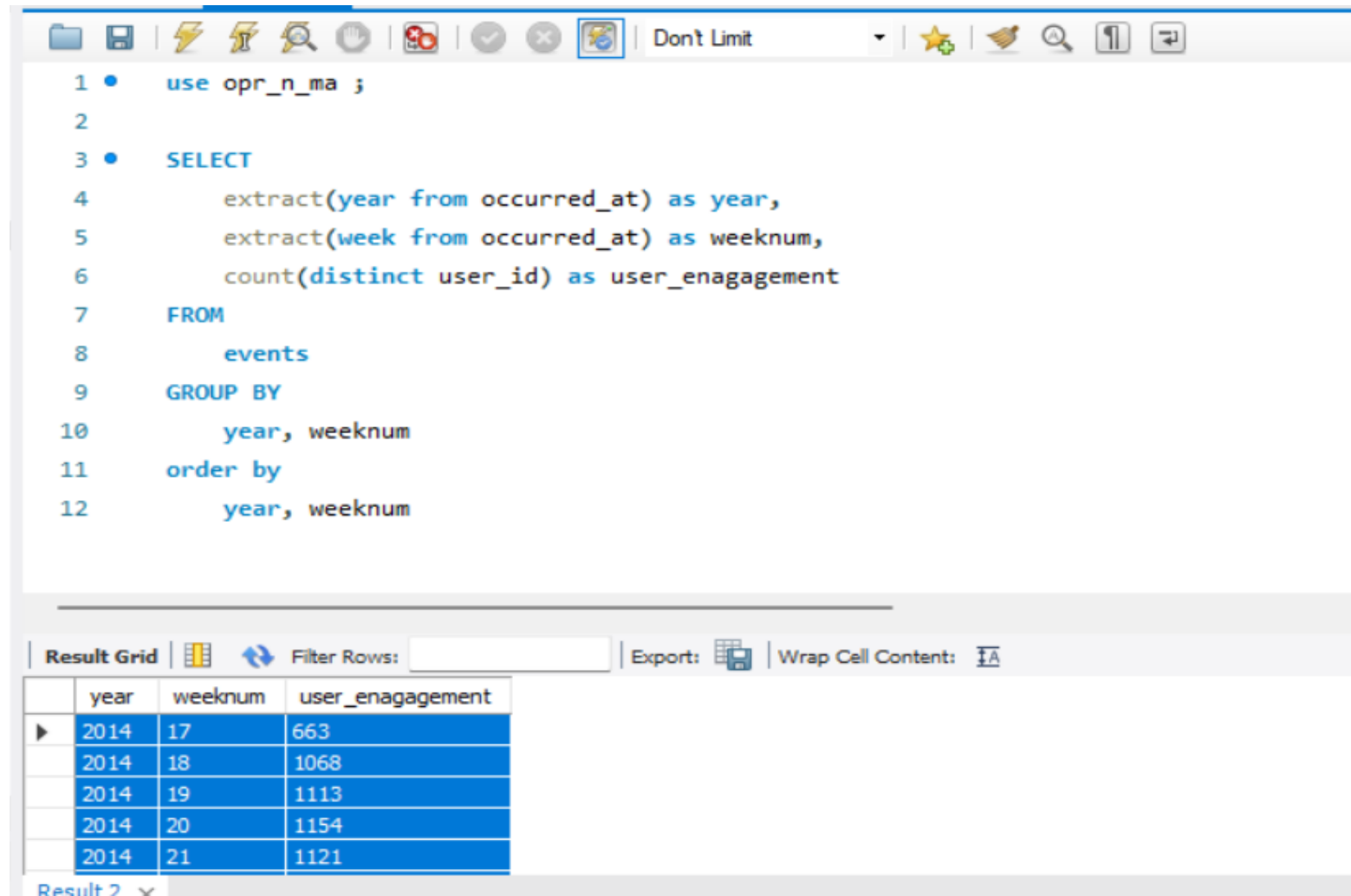
	language	lang_count	perc_language
▶	Arabic	1	12.50000
	English	1	12.50000
	French	1	12.50000
	Hindi	1	12.50000
	Italian	1	12.50000
	Persian	3	37.50000

C1T4 . Rows that have same values present in them.

```
1 • select job_id, count(job_id) as duplicate_count  
2   from job_data  
3   group by job_id  
4   having duplicate_count > 1
```

Result Grid			Filter Rows:	Export:	Wrap Cell Content:
	job_id	duplicate_count			
▶	23	3			

C2T1 . To measure the activeness of a user.



The screenshot displays a SQL query editor interface. The query is as follows:

```
1 • use opr_n_ma ;
2
3 • SELECT
4     extract(year from occurred_at) as year,
5     extract(week from occurred_at) as weeknum,
6     count(distinct user_id) as user_engagement
7 FROM
8     events
9 GROUP BY
10    year, weeknum
11 order by
12    year, weeknum
```

Below the query editor, the 'Result Grid' tab is active, showing the following data:

	year	weeknum	user_engagement
▶	2014	17	663
	2014	18	1068
	2014	19	1113
	2014	20	1154
	2014	21	1121

At the bottom left, the text 'Result 2' is visible with a dropdown arrow.

C2T2 . Amount of users growing over time for a product – number of active users per week.

The screenshot shows a SQL IDE window titled "SQL File 4*" with a tab labeled "events". The query editor contains the following SQL code:

```
1 • use opr_n_ma ;
2
3 • select
4     year,
5     weeknum,
6     new_active_user,
7     sum(new_active_user)over(order by year,weeknum rows between unbounded preceding and current row) as cum_active_user
8 FROM
9 (
10  SELECT
11     extract(year from activated_at) as year,
12     extract(week from activated_at) as weeknum,
13     count(distinct user_id) as new_active_user
14  FROM
15     users
16  where
17     state = 'active'
18  group by
19     year,
20     weeknum
21 ) a
22
```

Below the query editor, the "Result Grid" tab is active, displaying the following data:

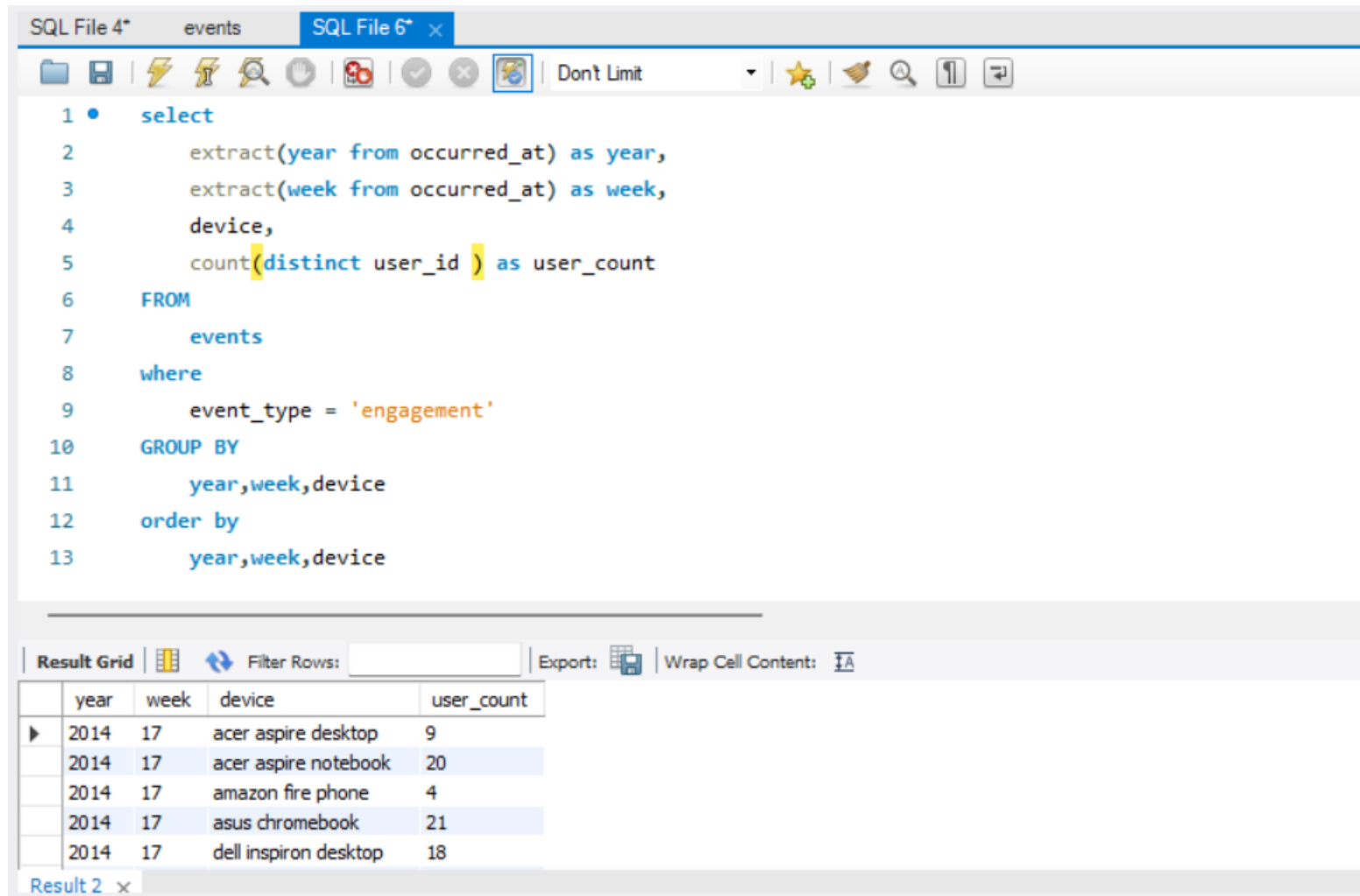
	year	weeknum	new_active_user	cum_active_user
▶	2013	0	23	23
	2013	1	30	53
	2013	2	48	101
	2013	3	36	137
	2013	4	30	167

The IDE interface includes a toolbar with various icons, a "Don't Limit" dropdown, and a "Read Only" status indicator at the bottom right.

C2T3 . Users getting retained weekly after signing up for a product.

```
1 • SELECT user_id, count(user_id),
2     sum(case when retention_week = 1 then 1 else 0 end) as week_1,
3     sum(case when retention_week = 2 then 1 else 0 end) as week_2,
4     sum(case when retention_week = 3 then 1 else 0 end) as week_3,
5     sum(case when retention_week = 4 then 1 else 0 end) as week_4,
6     sum(case when retention_week = 5 then 1 else 0 end) as week_5,
7     sum(case when retention_week = 6 then 1 else 0 end) as week_6,
8     sum(case when retention_week = 7 then 1 else 0 end) as week_7,
9     sum(case when retention_week = 8 then 1 else 0 end) as week_8,
10    sum(case when retention_week = 9 then 1 else 0 end) as week_9
11  from
12  (
13    SELECT
14      a.user_id,
15      a.signup_week,
16      b.engagement_week,
17      b.engagement_week - a.signup_week as retention_week
18  FROM
```

C2T4 . To measure the activeness of a user weekly -
Calculate the weekly engagement per device.



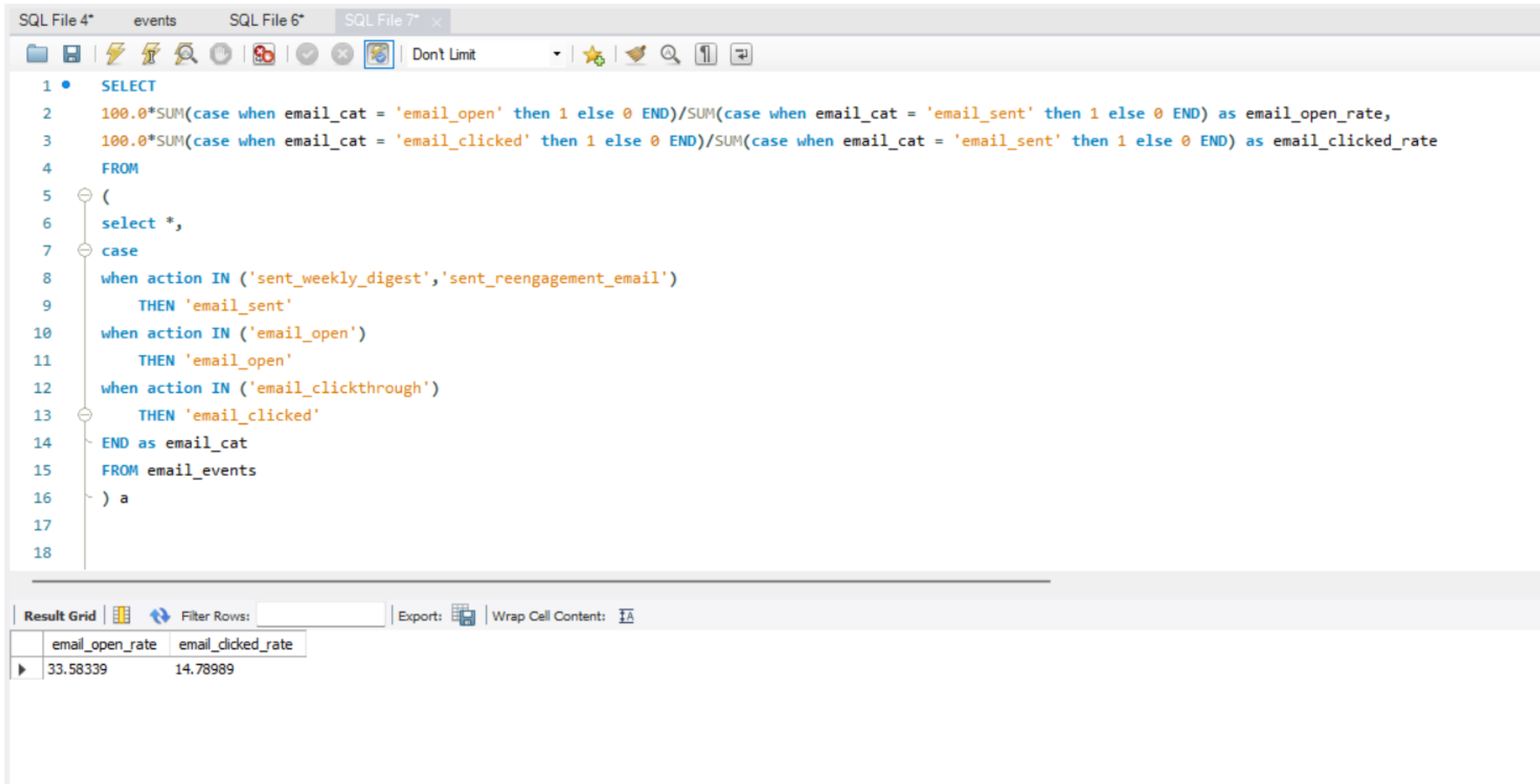
The screenshot shows a SQL IDE interface with a query editor and a result grid. The query is designed to calculate the weekly engagement per device by extracting the year and week from the 'occurred_at' field, counting distinct users, and grouping by year, week, and device.

```
1 • select
2     extract(year from occurred_at) as year,
3     extract(week from occurred_at) as week,
4     device,
5     count(distinct user_id ) as user_count
6 FROM
7     events
8 where
9     event_type = 'engagement'
10 GROUP BY
11     year,week,device
12 order by
13     year,week,device
```

The result grid displays the following data:

	year	week	device	user_count
▶	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

C2T5 . Users engaging with the email service.



The screenshot shows a SQL IDE interface with a query editor and a results grid. The query editor contains a SQL query that calculates email engagement rates. The results grid shows the output of the query, which includes two columns: email_open_rate and email_clicked_rate. The values for these rates are 33.58339 and 14.78989, respectively.

```
1 • SELECT
2 100.0*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,
3 100.0*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_clicked_rate
4 FROM
5 (
6 select *,
7 case
8 when action IN ('sent_weekly_digest','sent_reengagement_email')
9 THEN 'email_sent'
10 when action IN ('email_open')
11 THEN 'email_open'
12 when action IN ('email_clickthrough')
13 THEN 'email_clicked'
14 END as email_cat
15 FROM email_events
16 ) a
17
18
```

Result Grid

email_open_rate	email_clicked_rate
33.58339	14.78989

Insights

- 0.00000746 jobs were reviewed over time.
- Highest share of language is of Persian i.e. 37.5 % followed by Arabic, English, French, Hindi and Italian
- Job Id : 23 has 3 duplicate count.
- The highest user engagement was 1154 in week number 20 in 2014.
- The highest weekly engagement as per device has user count 21 by Asus Chromebook in week 17 in 2014.

RESULTS

Results obtained are displayed on each slide respectively.

The project focused on various aspects of advanced SQL and enhanced the skills on the same.