

## **Project Description:**

- This project explains operations of a company for its better growth.
- It analysis various aspects of the company to work on areas of improvement.
- As a data analyst, we can deep dive into the database of the company and bring out some insights for the company's future.

## **Approach:**

- I have created the needed database (Project3) and tables in MySQL workbench and imported the entire csv file in the database as tables so that we could analyse the data.
- We write MySQL queries to extract information that is needed for better analysis as the data is huge.
- The tables that are created:
  1. Job\_Data
  2. Events
  3. Email\_Events
  4. Users

## **Insights:**

- We are going to find few insights from the data like-
- Number of jobs reviewed.
- Calculating the percentage share of each language over a period
- Analysing the throughput using window function
- User engagement in a week
- Engagement per device on a weekly basis and more.

## **Result:**

- The project has helped me learn MySQL queries in a very vast manner.
- It expanded my knowledge and improved the ability to think and analyse the data.
- It has helped me to draw conclusions and the requirements for a better understanding of the business growth.
- It has helped me to better understand the MySql Workbench.

## **MySQL queries and output Screenshots:**

## Jobs Reviewed over time-----

```
8 • SELECT
9     ds,
10     COUNT(job_id) AS num_of_jobs,
11     COUNT(*) / 24 AS jobs_per_hr,
12     SUM(time_spent) / 3600 AS hrs_spent
13 FROM
14     job_data
15 WHERE
16     ds BETWEEN '2020-11-01' AND '2020-11-30'
17 GROUP BY ds;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: [IA](#)

	ds	num_of_jobs	jobs_per_hour	hours_spent
▶	2020-11-30	4	0.1667	0.0575
	2020-11-29	3	0.1250	0.0306
	2020-11-28	5	0.2083	0.0469
	2020-11-27	4	0.1667	0.0500
	2020-11-26	6	0.2500	0.0961
	2020-11-25	2	0.0833	0.0367
	2020-11-15	3	0.1250	0.0381
	2020-11-02	6	0.2500	0.0808
	2020-11-20	3	0.1250	0.0542
	2020-11-01	3	0.1250	0.0511
	2020-11-22	7	0.2917	0.1069
	2020-11-16	2	0.0833	0.0342
	2020-11-13	6	0.2500	0.0675
	2020-11-05	1	0.0417	0.0103
	2020-11-09	2	0.0833	0.0350
	2020-11-19	3	0.1250	0.0567

Result 14 x

The query shows number of hours in a day spent on number of jobs in last one month.

**7 day Rolling average:**

SQL Query:

```

SELECT
  ds,
  COUNT(job_id) AS job_count,
  SUM(time_spent) AS total_time_spent,
  CASE
    WHEN SUM(time_spent) <> 0 THEN COUNT(*) / SUM(time_spent)
    ELSE NULL END AS throughput, TRUNCATE(AVG(COUNT(*)) OVER (ORDER BY ds ROWS BETWEEN 6 PRECEDING AND CURRENT ROW) / AVG(SUM(time_spent))
    OVER (ORDER BY ds ROWS BETWEEN 6 PRECEDING AND CURRENT ROW), 3) AS throughput_7dayRolling
FROM job_data
GROUP BY ds;

```

Result Grid:

ds	job_count	total_time_spent	throughput	throughput_7dayRolling
2020-11-01	3	184	0.0163	0.016
2020-11-02	6	291	0.0206	0.018
2020-11-03	5	294	0.0170	0.018
2020-11-04	2	84	0.0238	0.018
2020-11-05	1	37	0.0270	0.019
2020-11-06	3	116	0.0259	0.019
2020-11-07	1	69	0.0145	0.019
2020-11-08	2	159	0.0126	0.019
2020-11-09	2	126	0.0159	0.018
2020-11-10	3	206	0.0146	0.017
2020-11-11	2	126	0.0159	0.016
2020-11-12	6	355	0.0169	0.016
2020-11-13	6	243	0.0247	0.017
2020-11-14	2	87	0.0230	0.017
2020-11-15	3	137	0.0219	0.018
2020-11-16	2	123	0.0163	0.018
2020-11-17	2	97	0.0206	0.019
2020-11-18	3	110	0.0273	0.020
2020-11-19	3	204	0.0147	0.020
2020-11-20	3	195	0.0154	0.018
2020-11-22	7	385	0.0182	0.018

Here we calculate the 7 day rolling average. This is better for huge set of data. Hence this would be preferred more.

## Language Share Analysis-----

Schema: project3

```

36 • select language, (count(language)/ (SELECT
37         COUNT(language)
38     FROM
39         job_data)) * 100 as language_share from job_data
40     where ds >= '11/01/2020'
41     and ds <= '11/30/2020' group by language;

```

language	language_share
English	3.0303
Arabic	1.0101
Persian	3.0303
Hindi	1.0101
French	1.0101
Italian	2.0202
Romansh	2.0202
Swedish	1.0101

Using this query we can find the share of each language from job\_data

## Duplicate Rows Detection

```

1 • select job_id, ds, count(*) from job_data group by ds, job_id having count(*) > 1;

```

job_id	ds	count(*)
--------	----	----------

There are no duplicate rows which is a good sign of a good data.

## Weekly User Engagement

71	•	SELECT
72		EXTRACT(WEEK FROM occurred_at) AS weeknum,
73		COUNT(DISTINCT user_id) AS weekly_user_engmnt
74		FROM
75		events a
76		GROUP BY weeknum;
77		
78		

  

Result Grid	Filter Rows:	Export:	Wrap Cell Content:
weeknum	weekly_user_engmnt		
17	740		
18	1260		
19	1287		
20	1351		
21	1299		
22	1381		
23	1446		
24	1471		
25	1459		
26	1509		
27	1573		
28	1577		
29	1607		
30	1706		
31	1514		
32	1454		
33	1438		
34	1443		
35	118		

This query shows how much users are engaged in a week.

## User Growth Analysis-----

84	•	SELECT year, weeknum, num_active_user,
85	⊖	SUM(num_active_user)OVER(ORDER BY year,weeknum
86		ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW) AS cum_active_users FROM
87	⊖	(
88		SELECT EXTRACT(year from a.activated_at) AS year,
89		EXTRACT(week from a.activated_at) AS weeknum,
90		COUNT(DISTINCT user_id) AS num_active_user FROM users a
91		GROUP BY year, weeknum ORDER BY year, weeknum ) a;

  

Result Grid	Filter Rows:	Export:	Wrap Cell Content:
year	weeknum	num_active_user	cum_active_users
NULL	NULL	9685	9685
2013	0	23	9708
2013	1	30	9738
2013	2	48	9786
2013	3	36	9822
2013	4	30	9852
2013	5	48	9900
2013	6	38	9938
2013	7	42	9980
2013	8	34	10014
2013	9	43	10057
2013	10	32	10089
2013	11	31	10120
2013	12	33	10153
2013	13	39	10192
2013	14	35	10227
2013	15	43	10270
2013	16	46	10316
2013	17	40	10356

This query finds the growth of number of users for the product.

## Weekly Retention-----

```
1 • select
2   count(user_id), engagement_week, signup_week, retention_week, sum(case when retention_week = 1 then 1 else 0 end) as week_1
3   from
4   (
5     select
6       a.user_id,
7       a.signup_week,
8       b.engagement_week,
9       b.engagement_week - a.signup_week as retention_week
10    from
11    (
12      (select distinct user_id, extract(week from occurred_at) as signup_week from events where event_type = 'signup_flow'
13       and event_name = 'complete_signup' and extract(week from occurred_at) ) a left join( select distinct user_id,
14       extract(week from occurred_at) as engagement_week
15       from events where event_type = 'engagement'
16       ) b on a.user_id = b.user_id) order by
17     a.user_id
18   ) c group by user_id;
```

	count(user_id)	engagement_week	signup_week	retention_week	week_1
1	17	17	17	0	0
1	17	17	17	0	0
2	17	17	17	0	1
3	17	17	17	0	0
5	17	17	17	0	1
2	17	17	17	0	1
1	17	17	17	0	0
3	17	17	17	0	1
2	17	17	17	0	1
6	17	17	17	0	1
2	17	17	17	0	1
6	17	17	17	0	1
10	17	17	17	0	1
2	17	17	17	0	1
2	17	17	17	0	1
1	17	17	17	0	0
1	17	17	17	0	0
2	17	17	17	0	1
6	17	17	17	0	0
3	17	17	17	0	0
2	17	17	17	0	1
4	17	17	17	0	1
3	17	17	17	0	1
7	17	17	17	0	1

This query shows how many users retain after signing up for the cohort that is shown by engagement week.

## Weekly Engagement per Device

Query 1 x SQL File 3\* SQL File 4\* SQL File 6\* SQL File 7\* SQL File 8\*

Limit to 1000 rows

```

79 -- Weekly Engagement Per Device:
80 • SELECT
81   EXTRACT(year FROM occurred_at) AS year,
82   EXTRACT(week FROM occurred_at) AS week,
83   device,
84   COUNT(distinct user_id) as weekly_engmnt_count
85 FROM events WHERE event_type = 'engagement' GROUP BY 1,2,3 ORDER BY 1,2,3;

```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: [IA](#)

year	week	device	weekly_engmnt_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
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
2014	17	macbook pro	143
2014	17	nexus 10	16
2014	17	nexus 5	40
2014	17	nexus 7	18
2014	17	nokia lumia 635	17
2014	17	samsung galaxy tablet	8

Result 55 x

year	week	device	weekly_engmnt_count
2014	18	acer aspire notebook	33
2014	18	amazon fire phone	9
2014	18	asus chromebook	42
2014	18	dell inspiron desktop	58
2014	18	dell inspiron notebook	77
2014	18	hp pavilion desktop	37
2014	18	htc one	19
2014	18	ipad air	52
2014	18	ipad mini	30
2014	18	iphone 4s	46
2014	18	iphone 5	113
2014	18	iphone 5s	73
2014	18	kindle fire	27
2014	18	lenovo thinkpad	153
2014	18	mac mini	13
2014	18	macbook air	121
2014	18	macbook pro	252
2014	18	nexus 10	30
2014	18	nexus 5	73
2014	18	nexus 7	30
2014	18	nokia lumia 635	33
2014	18	samsung galaxy tablet	11
2014	18	samsung galaxy note	15

Result 55 x

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

year	week	device	weekly_engmnt_count
2014	19	dell inspiron notebook	83
2014	19	hp pavilion desktop	40
2014	19	htc one	30
2014	19	ipad air	55
2014	19	ipad mini	36
2014	19	iphone 4s	44
2014	19	iphone 5	115
2014	19	iphone 5s	79
2014	19	kindle fire	21
2014	19	lenovo thinkpad	178
2014	19	mac mini	18
2014	19	macbook air	112
2014	19	macbook pro	266
2014	19	nexus 10	25
2014	19	nexus 5	87
2014	19	nexus 7	41
2014	19	nokia lumia 635	23
2014	19	samsung galaxy tablet	6
2014	19	samsung galaxy note	11
2014	19	samsung galaxy s4	91
2014	19	windows surface	16
2014	20	acer aspire desktop	23
2014	20	acer aspire notebook	40

In this query, weekly engagement per device is calculated, if its more the users are satisfied with the product.

# Email Engagement Analysis

The screenshot shows a database IDE interface. On the left, the 'SCHEMAS' pane displays a tree view of the database structure, including 'project3' with tables 'email\_events', 'events', 'job\_data', and 'users'. Below this, the 'Table: job\_data' is detailed with its columns: 'ds' (text), 'job\_id' (int), 'actor\_id' (int), 'event' (text), 'language' (text), and 'time\_spent' (int). The main editor displays two SQL queries. The first query (lines 77-81) counts the number of emails for each action. The second query (lines 83-92) calculates the email open and click rates for specific actions, using a CASE statement to map actions to categories like 'email\_sent', 'email\_open', and 'email\_clicked'. The 'Result Grid' at the bottom shows the results of the first query:

action	num_email
sent_weekly_digest	57267
email_open	20459
email_clickthrough	9010
sent_reengagement_email	3653

The screenshot shows a database IDE interface with a SQL query editor. The query (lines 79-98) calculates the email open and click rates for specific actions, using a CASE statement to map actions to categories like 'email\_sent', 'email\_open', and 'email\_clicked'. The 'Result Grid' at the bottom shows the results of the query:

email_open_rate	email_clicked_rate
33.58339	14.78989



The screenshot displays a database interface with a left-hand navigation pane showing a schema named 'sys' containing tables like 'events', 'job\_data', 'users', 'Views', 'Stored Procedures', and 'Functions'. The main area shows a SQL query being executed:

```

95
96 SELECT *,
97 CASE
98 WHEN action IN ('sent_weekly_digest', 'sent_reengagement_email')
99 THEN 'email_sent'
100 WHEN action IN ('email_open') THEN 'email_open'
101 WHEN action in ('email_clickthrough') THEN 'email_clicked'
102 END AS email_metrics
103 FROM
104 email_events;

```

Below the query, the 'Result Grid' shows the following data:

user_id	occurred_at	action	user_type	email_metrics
0	06-05-2014 09:30	sent_weekly_digest	1	email_sent
0	13-05-2014 09:30	sent_weekly_digest	1	email_sent
0	20-05-2014 09:30	sent_weekly_digest	1	email_sent
0	27-05-2014 09:30	sent_weekly_digest	1	email_sent
0	03-06-2014 09:30	sent_weekly_digest	1	email_sent
0	03-06-2014 09:30	email_open	1	email_open
0	10-06-2014 09:30	sent_weekly_digest	1	email_sent
0	10-06-2014 09:30	email_open	1	email_open
0	17-06-2014 09:30	sent_weekly_digest	1	email_sent
0	17-06-2014 09:30	email_open	1	email_open
0	24-06-2014 09:30	sent_weekly_digest	1	email_sent
0	01-07-2014 09:30	sent_weekly_digest	1	email_sent
0	08-07-2014 09:30	sent_weekly_digest	1	email_sent

On the left, the 'Table: job\_data' is detailed with the following columns:

- ds: text
- job\_id: int
- actor\_id: int
- event: text
- language: text
- time\_spent: int
- org: text

This result gives us an overview of the engagement of the users with emails.

## CODE FOR CREATING THE TABLE AND IMPORTING DATA

```

CREATE database Project3;
Use project3;

```

```

Create Table job_data(
ds DATE,
job_id INT NOT NULL,
actor_id INT NOT NULL,
event Varchar(10) NOT NULL,
language Varchar(15) NOT NULL,
time_spent INT NOT NULL,
org CHAR(3)
);

```

```

CREATE TABLE users
(
    user_id int null,
    created_at datetime null,
    company_id int null,
    language VARCHAR(10) null,
    activated_at datetime null,
    state varchar(25) null
);

```

```
CREATE TABLE events
(
    user_id int null,
    occurred_at datetime null,
    event_type VARCHAR(100) null,
    event_name VARCHAR(100) null,
    location varchar(215) null,
    device varchar(215) null,
    user_type int null
);
```

```
CREATE TABLE email_events
(
    user_id int null,
    occurred_at datetime null,
    action VARCHAR(225) null,
    user_type int null
);
```

```
LOAD DATA LOCAL INFILE 'C:/ProgramData/MySQL/MySQL Server
8.0/Uploads/events.csv'
INTO TABLE events
FIELDS TERMINATED BY ','
ENCLOSED BY '"'
LINES TERMINATED BY '\n'
IGNORE 1 ROWS;
```

```
show variables like "secure_file_priv";
```

Load data result:

15	13:20:57	LOAD DATA INFILE 'C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/email_events.csv' INTO TABLE e...	90389 row(s) affected Records: 90389 Deleted: 0 Skipped: 0 Warnings: 0
16	13:35:08	LOAD DATA INFILE 'C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/users.csv' INTO TABLE users FIE...	Error Code: 1146. Table 'project3.users' doesn't exist
17	13:35:23	CREATE TABLE users (user_id int null, created_at date null, company_id int null, language VARCHAR...	0 row(s) affected
18	13:35:28	LOAD DATA INFILE 'C:/ProgramData/MySQL/MySQL Server 8.0/Uploads/users.csv' INTO TABLE users FIE...	9381 row(s) affected Records: 9381 Deleted: 0 Skipped: 0 Warnings: 0