

PROJECT NO.3

NAME: - AYUSH RAVINDRA VAIDANDE



# Description:

#### Case Study 1 (Job Data):

A. Number of jobs reviewed: Amount of jobs reviewed over time.

Your task: Calculate the number of jobs reviewed per hour per day for November 2020?

B. Throughput: It is the no. of events happening per second.

Your task: Let's say the above metric is called throughput. Calculate 7 day rolling average of throughput? For throughput, do you prefer daily metric or 7-day rolling and why?

C. Percentage share of each language: Share of each language for different contents.

Your task: Calculate the percentage share of each language in the last 30 days?

D. Duplicate rows: Rows that have the same value present in them.

Your task: Let's say you see some duplicate rows in the data. How will you display duplicates from the table?

#### Case Study 2 (Investigating metric spike):

A. User Engagement: To measure the activeness of a user. Measuring if the user finds quality in a product/service.

Your task: Calculate the weekly user engagement?

B. User Growth: Amount of users growing over time for a product.

Your task: Calculate the user growth for product?

C. Weekly Retention: Users getting retained weekly after signing-up for a product.

Your task: Calculate the weekly retention of users-sign up cohort?

D. Weekly Engagement: To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly.

Your task: Calculate the weekly engagement per device?

E. Email Engagement: Users engaging with the email service.

Your task: Calculate the email engagement metrics?

## SOFTWARE USED

MYSQL WORKBENCH 8.0 CE (COMMUNITY EDITION)

VERSION:-8.0

Number of jobs reviewed: Amount of jobs reviewed over time. Calculate the number of jobs reviewed per hour per day for November 2020?

To find Number of jobs reviewed reviewed per hour per day for november 2020

QUERY:-

SELECT ds AS Dates, ROUND((COUNT(job\_id)/sum(time\_spent))\*3600) AS "jobs reviewed per hour per day" FROM job\_data WHERE ds BETWEEN "2020-11-01" AND "2020-11-30" GROUP BY ds;



Dates	jobs reviewed per hour per day	
2020-11-30	180	
2020-11-29	180	
2020-11-28	218	
2020-11-27	35	
2020-11-26	64	
2020-11-25	80	

- Throughput: It is the no. of events happening per second.
  - Let's say the above metric is called throughput. Calculate 7 day rolling average of throughput? For throughput, do you prefer daily metric or 7-day rolling and why?
- Here if we use weekly thoughput :-
- QUERY:-
- SELECT ROUND(COUNT(event)/SUM(time\_spent),2) AS "Weekly throughput" FROM job\_data
- OUTPUT:-

Weekly throughput

0.03

- But if we use daily throught put
- QUERY:-
- SELECT ds AS Dates, ROUND(COUNT(event)/SUM(time\_spent),2) AS "Daily throughput" FROM job\_data GROUP by ds ORDER by ds
- ORDER:-

Dates	Daily throughput
2020-11-25	0.02
2020-11-26	0.02
2020-11-27	0.01
2020-11-28	0.06
2020-11-29	0.05
2020-11-30	0.05



 Metrics will go up and down on weekly or daily basis, we will get number faster every day or minute if we want, as a result rolling metrics are superb at showing if your metrics are trending up or down on a daily level.

- Percentage share of each language: Share of each language for different contents. Calculate the percentage share of each language in the last 30 days?
- We will first divide the total number of languages (distinct/non-distinct) by the total number of rows presents in the table
- Then we will do the grouping based on the languages
- QUERY:-
- select job\_data.job\_id, job\_data.language, count(distinct job\_data.language) as total\_of\_each\_language, ((count(job\_data.language)/(select count(\*) from job\_data))\*100) as percentage\_share\_of\_each\_distinct\_language from job\_data group by job\_data.language;

## • OUTPUT :-

job_id	language	total_of_each_language	percentage_share_of_each_distinct_language
22	Arabic	1	12.5000
21	English	1	12.5000
11	French	1	12.5000
25	Hindi	1	12.5000
20	Italian	1	12.5000
23	Persian	1	37.5000

- Duplicate rows: Rows that have the same value present in them. Let's say you see some duplicate rows in the data. How will you display duplicates from the table?
- QUERY:-
- SELECT actor\_id ,COUNT(\*) AS Duplicates FROM job\_data GROUP by actor\_id HAVING COUNT(\*) > 1;
- OUTPUT :-

actor_id	Duplicates
1003	2

### Case Study 2 (Investigating metric spike)

- User Engagement: To measure the activeness of a user. Measuring if the user finds quality in a product/service. Calculate the weekly user engagement?
- We will extract the week from the occurred\_at column of the events table using the EXTRACT function and WEEK function
- Then we will be counting the number of distinct user\_id from the events table
- Then we will use the GROUP BY function to group the output w.r.t week from occurred at
- QUERY:-
- SELECT extract (week from occurred\_at) as week\_number, count(distinct user\_id) as number\_of\_users FROM tutorial.yammer\_events group by week\_number;

### • OUTPUT :-

week_number	number_of_users
18	791
19	1244
20	1270
21	1341
22	1293
23	1366
24	1434
25	1462
26	1443
27	1477
28	1556
29	1556
30	1593
31	1685
32	1483
33	1438
34	1412
35	1442

- User Growth: Amount of users growing over time for a product. Calculate the user growth for product?
- To find the user growth (number of active users per week):-
- First we will the extract the year and week for the occurred\_at column of the users table using the extract, year and week functions
- Then we will group the extracted week and year on the basis of year and week number
- Then we ordered the result on the basis of year and week number
- Then we will find the cumm\_active\_users using the SUM, OVER and ROW function between unbounded preceding and current row

- QUERY:-
- select year\_num, week\_num, num\_active\_users, SUM(num\_active\_users)OVER(ORDER BY year\_num, week\_num ROWS BETWEEN
- UNBOUNDED PRECEDING AND CURRENT ROW) AS cum\_active\_users
- from
- (select
- extract (year from a.activated\_at) as year\_num,
- extract (week from a.activated at) as week num,
- count(distinct user\_id) as num\_active\_users
- from
- tutorial.yammer\_users a
- WHERE state = 'active' group by year\_num, week\_num order by year\_num, week\_num) a;

#### • OUTPUT :-

year_num	week_num	num_active_users	cum_active_users	year_num	week_num	num_active_users	cum_active_users
2013				2013	45		
2013				2013	46		
2013				2013	47		
2013				2013	48		
2013				2013	49		
2013				2013	50		
2013				2013	51		
2013				2013	52		
2013				2014	1		
2013				2014	2		
2013		33		2014	3		
2013				2014	4		
2013				2014	5		
2013				2014	6		
2013				2014	7		
2013				2014	8		
2013				2014	9		
2013				2014	10		
2013				2014	11		
2013				2014	12		
2013				2014	13		
2013	22	49	914	2014	14		
2013				2014	15		
2013			1016	2014	16		
2013	25	46	1062	2014	17		
2013	26	57	1119	2014	18		
2013				2014	19		
2013				2014	20		
2013				2014	21		
2013	30	66	1365	2014	22		
2013		69	1434	2014	23		
2013	32	66	1500	2014	24		
2013		73	1573	2014	25		
2013			1643	2014	26		
2013	35	80	1723	2014	27		
2013	36	65	1788	2014	28		
2013	37	71	1859	2014	29		
2013	38	84	1943	2014	30		
2013		92	2035	2014	31		
2013	40	81	2116	2014	32		
2013	41	88		2014	33		
2013		74	2278	2014	34		
2013	43	97	2375	2014	35	266	938:
2013	44	92	2467				

- User Growth = Number of active users per week
- Query :
- select count(\*) from tutorial.yammer\_userswhere state = 'active';
- OUTPUT:-
- there are in total 9381 active users from 1st week of 2013 to the 35<sup>th</sup> week of 2014

count 9381

- Weekly Retention: Users getting retained weekly after signing-up for a product. Calculate the weekly retention
  of users-sign up cohort?
- The weekly retention of users-sign up cohort can be calculated by two means i.e. either by specifying the week number (18 to 35) or for the entire column of occurred\_at of the events table.
- Firstly we will extract the week from occurred\_at column using the extract, week functions
- Then, we will select out those rows in which event\_type = 'signup\_flow' and event\_name = 'complete\_signup'
- If finding for a spectifc week we will spectify the week number using the extract function
- Then using the left join we will join the two tables on the basis of user\_id where event\_type = 'engagement'
- Then we will use the Group By function to group the output table on the basis of
- user\_id
- Then we will use the Order By function to order the result table on the basis of user\_id

- QUERY :-
- SELECT distinct user\_id, COUNT(user\_id), SUM(CASE WHEN retention\_week = 1 Then 1 Else 0 END) as per\_week\_retention FROM ( SELECT a.user\_id, a.signup\_week, b.engagement\_week, b.engagement\_week a.signup\_week as retention\_week FROM ( (SELECT distinct user\_id, extract(week from occurred\_at) as signup\_week from tutorial.yammer\_events WHERE event\_type = 'signup\_flow' and event\_name = 'complete\_signup' )a LEFT JOIN (SELECT distinct user\_id, extract (week from occurred\_at) as engagement\_week FROM tutorial.yammer\_events where event\_type = 'engagement' )b on a.user\_id = b.user\_id ) )d group by user\_id order by user\_id;
- OUTPUT:-
- LINK :- https://drive.google.com/file/d/1mr8dFvW8et3VBKInA9e4XmXoGWvUtgTu/view?usp=share\_link

- Weekly Engagement: To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly. Calculate the weekly engagement per device?
- To find the weekly user engagement per device:-
- Firstly we will extract the year\_num and week\_num from the occurred\_at column of the events table using the extract, year and week function
- Then we will select those rows where event\_type = 'engagement' using the WHERE clause
- Then by using the Group By and Order By function we will group and order the result on the basis of year\_num, week\_num and device

- QUERY :-
- SELECT
- extract(year from occurred\_at) as year\_num,
- extract(week from occurred\_at) as week\_num,
- device, COUNT(distinct user\_id) as no\_of\_users FROM tutorial.yammer\_events where event\_type = 'engagement' GROUP by 1,2,3 order by 1,2,3;
- OUTPUT:-
- LINK :- https://drive.google.com/file/d/1kvv1VpqsZBuUyDCYqdVfC2\_pL4GvHS\_K/view?usp=sharing

- Email Engagement: Users engaging with the email service. Calculate the email engagement metrics?
- To find the email engagement metrics(rate) of users:-
- We will first categorize the action on the basis of email\_sent, email\_opened and email\_clicked using the CASE, WHEN, THEN functions
- Then we select the sum of category of email\_opened divide by the sum of the category of email\_sent and multiply the result by 100.0 and name is as email\_opening\_rate
- Then we select the sum of category of email\_clicked divide by the sum of the category of email\_sent and multiply the result by 100.0 and name is as email\_clicking\_rate
- email\_sent = ('sent\_weekly\_digest', 'sent\_reengagement\_email')
- email\_opened = 'email\_open'
- email\_clicked = 'email\_clickthrough'

- 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 tutorial.yammer\_emails ) a;

### • OUTPUT:-

	A	В
1	email_opening_rate	email_clicking_rate
2	33.58338805	14.78988838

