# Operation Analytics and Investigating Metric Spike

—

03. 07. 2023

# Project Description

Operation Analytics is the analysis done for the complete end to end operations of a company. With the help of this, the company then finds the areas on which it must improve upon. You work closely with the ops team, support team, marketing team, etc and help them derive insights out of the data they collect.

Being one of the most important parts of a company, this kind of analysis is further used to predict the overall growth or decline of a company's fortune. It means better automation, better understanding between cross-functional teams, and more effective workflows.

Investigating metric spike is also an important part of operation analytics as being a Data Analyst you must be able to understand or make other teams understand questions like- Why is there a dip in daily engagement? Why have sales taken a dip? Etc. Questions like these must be answered daily and for that its very important to investigate metric spike.

**Below are the two operations mentioning the answers for the related questions:**

Case Study 1 (Job Data)

**Below is the structure of the table with the definition of each column that you must work on:**

- **Table-1:** job_data

  - **job_id:** unique identifier of jobs
  - **actor_id:** unique identifier of actor
  - **event:** decision/skip/transfer
  - **language:** language of the content
  - **time_spent:** time spent to review the job in seconds
  - **org:** organization of the actor
  - **ds:** date in the yyyy/mm/dd format. It is stored in the form of text and we use presto to run. no need for date function

Use the dataset attached in the Dataset section below the project images then answer the questions that follows

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

2. **Throughput:** It is the no. of events happening per second.
   **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?

3. **Percentage share of each language:** Share of each language for different contents.
   **Task:** Calculate the percentage share of each language in the last 30 days?

4. **Duplicate rows:** Rows that have the same value present in them.
   **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)

**The structure of the table with the definition of each column that you must work on is present in the project image**

- **Table-1:** users
  This table includes one row per user, with descriptive information about that user's account.

- **Table-2:** events
  This table includes one row per event, where an event is an action that a user has taken. These events include login events, messaging events, search events, events logged as users progress through a signup funnel, events around received emails.

- **Table-3:** email_events
  This table contains events specific to the sending of emails. It is similar in structure to the events table above.

Use the dataset attached in the Dataset section below the project images then answer the questions that follows

1. **User Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service.
   **Task:** Calculate the weekly user engagement?

2. **User Growth:** Amount of users growing over time for a product.
   **Task:** Calculate the user growth for product?

3. **Weekly Retention:** Users getting retained weekly after signing-up for a product.
   **Task:** Calculate the weekly retention of users-sign up cohort?

4. **Weekly Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly.
   **Task:** Calculate the weekly engagement per device?
5. **Email Engagement:** Users engaging with the email service.
   **Task:** Calculate the email engagement metrics?

# Approach

- First, I created the database and tables in MySQL Workbench, then importing the csv data provided using load data infile methods for tables having large number of rows, further describing tables checking datatypes of columns to modify column if there is different data types and used table import wizard for tables having small number of rows.

- After creating the required tables and data insertion, I spent some time understanding each column in the tables in Microsoft Excel and then I carried my analysis by writing SQL queries in MySQL Workbench and showcasing analysis using Tableau.

# Tech-Stack Used

- **MySQL Workbench v8.0.33.0 community edition** is used for project execution in order to query the database and gather the required results.
- **Microsoft Excel Office 2019 version** was used to understand columns and datasets of different csv files. The Datasets can be referred [here](#).
- **Tableau Desktop v2019.4.1 professional edition** is used to create visual representation of results for creation of graphs and other charts to understand the result better and make better decisions.

# Results and Insights

Case Study-1 (Job Data)

1.1. Number of jobs reviewed per hour per day for November 2020

```
47 •    SELECT
48          ds as date,
49          COUNT(*) AS job_reviewed,
50          sum(time_spent) / 3600 AS time_spent_per_hour
51      FROM
52          jobs
53      WHERE
54          ds >= '2020-11-01' AND ds < '2020-12-01'
55      GROUP BY date;
```

| date | job_reviewed | time_spent_per_hour |
|---|---|---|
| 2020-11-24 | 1 | 0.0158 |
| 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 |

1.2. 7-day rolling average of throughput

```
62 •    CREATE TEMPORARY TABLE JOBS_REVIEWED
63  ⊖ (
64      SELECT ds,count(distinct job_id) as jobs_reviewed,CAST(COUNT(DISTINCT JOB_ID)/86400 AS DECIMAL(10,10)) AS THROUGHPUT
65          from jobs
66          group by ds
67          order by ds
68  );
69 •    SELECT ds,jobs_reviewed,throughput
70      ,avg(throughput) OVER(ORDER BY DS ROWS BETWEEN 6 PRECEDING AND CURRENT ROW) AS throughput_7day
71      FROM JOBS_REVIEWED;
72
```
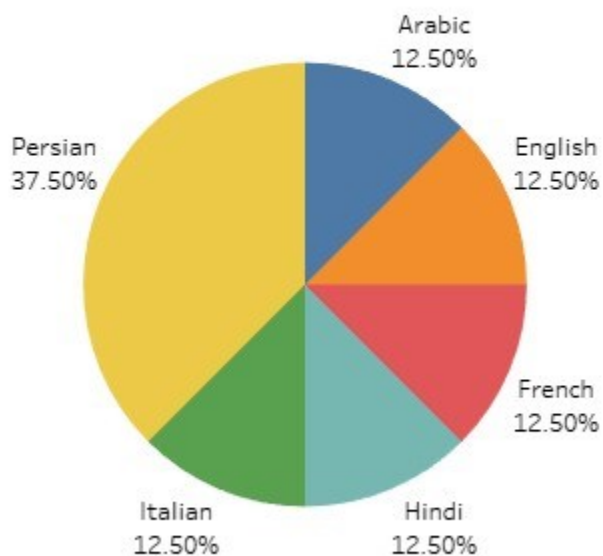
| ds | jobs_reviewed | throughput | throughput_7day |
|---|---|---|---|
| 2020-11-24 | 1 | 0.0000115740 | 0.00001157400000 |
| 2020-11-25 | 1 | 0.0000115740 | 0.00001157400000 |
| 2020-11-26 | 1 | 0.0000115740 | 0.00001157400000 |
| 2020-11-27 | 1 | 0.0000115740 | 0.00001157400000 |
| 2020-11-28 | 2 | 0.0000231480 | 0.00001322742857 |
| 2020-11-29 | 1 | 0.0000115740 | 0.00001322742857 |
| 2020-11-30 | 2 | 0.0000231480 | 0.00001488085714 |

## 1.3. Percentage share of each language in the last 30 days i.e November month.

```
81      -- Percentage share of each language: Share of each language for different contents in last 30 days.
82
83 •    SELECT
84          language, count(job_id) as jobs_applied,
85          COUNT(*) / (SELECT COUNT(*) FROM jobs WHERE ds >= '2020-11-01' AND ds < '2020-12-01') * 100
86          AS percentage_share_of_language
87      FROM
88          jobs
89      WHERE
90          ds >= '2020-11-01' AND ds < '2020-12-01'
91      GROUP BY language;
```

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

| language | jobs_applied | percentage_share_of_language |
|----------|--------------|------------------------------|
| French | 4 | 12.5000 |
| Persian | 12 | 37.5000 |
| Italian | 4 | 12.5000 |
| English | 4 | 12.5000 |
| Arabic | 4 | 12.5000 |
| Hindi | 4 | 12.5000 |

### Percentage of language share



Arabic 12.50%
English 12.50%
French 12.50%
Hindi 12.50%
Italian 12.50%
Persian 37.50%

## 1.4. Display of duplicate rows from the table.

```
95 •    select * from
96  ⊝ (
97      select *,
98      row_number() over (partition by job_id order by ds) as rownum
99      from jobs
100     )as t
101     where rownum>1;
102
```

**Result Grid** | Filter Rows: [        ] | Export: | Wrap Cell Content: ⫞A

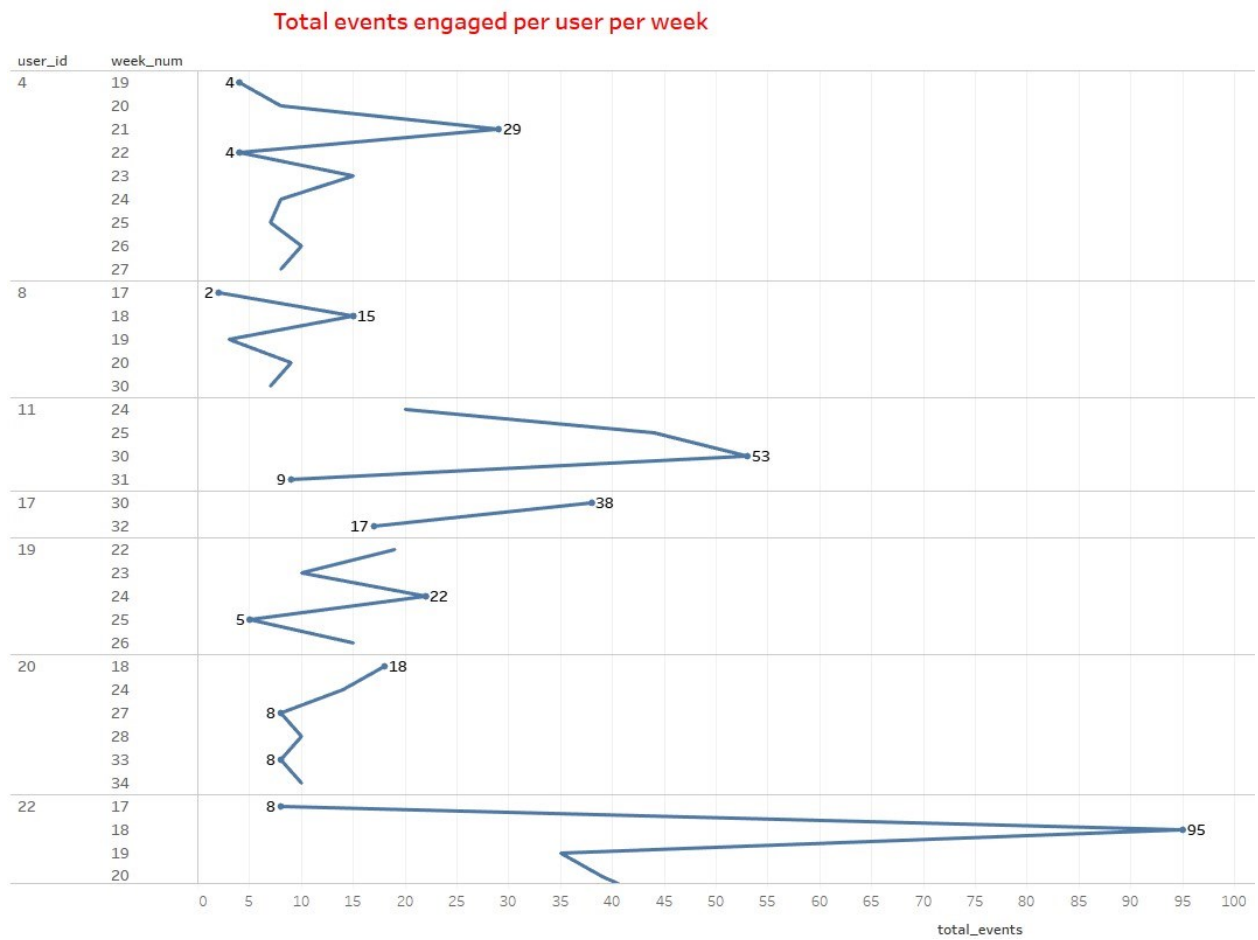| ds | job_id | actor_id | event | language | time_spent | org | rownum |
|----|--------|----------|-------|----------|------------|-----|--------|
| 2020-10-17 | 3 | 1011 | skip | Persian | 63 | A | 2 |
| 2020-11-10 | 3 | 1014 | skip | Persian | 27 | A | 3 |
| 2020-11-07 | 10 | 1020 | decision | Hindi | 14 | B | 2 |
| 2020-09-27 | 11 | 1003 | decision | Persian | 61 | C | 2 |
| 2020-11-27 | 11 | 1007 | decision | French | 104 | D | 3 |
| 2020-09-16 | 12 | 1020 | transfer | Italian | 80 | C | 2 |
| 2020-11-26 | 23 | 1004 | skip | Persian | 56 | A | 2 |
| 2020-11-28 | 23 | 1005 | transfer | Persian | 22 | D | 3 |

Case Study 2 (Investigating metric spike)

2.1. weekly user engagement

```
105 ●   SELECT
106         user_id,
107         EXTRACT(WEEK FROM occurred_at) AS week_num,
108         COUNT(event_type) AS total_events
109     FROM
110         events
111     WHERE
112         event_type = 'engagement'
113     GROUP BY user_id , week_num
114     ORDER BY user_id;
115
```

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

| user_id | week_num | total_events |
|---------|----------|--------------|
| 4 | 19 | 4 |
| 4 | 20 | 8 |
| 4 | 21 | 29 |
| 4 | 22 | 4 |
| 4 | 23 | 15 |
| 4 | 24 | 8 |
| 4 | 25 | 7 |
| 4 | 26 | 10 |

**Total events engaged per user per week**



Above visualization, shows the highest and lowest number of events per week for each user.
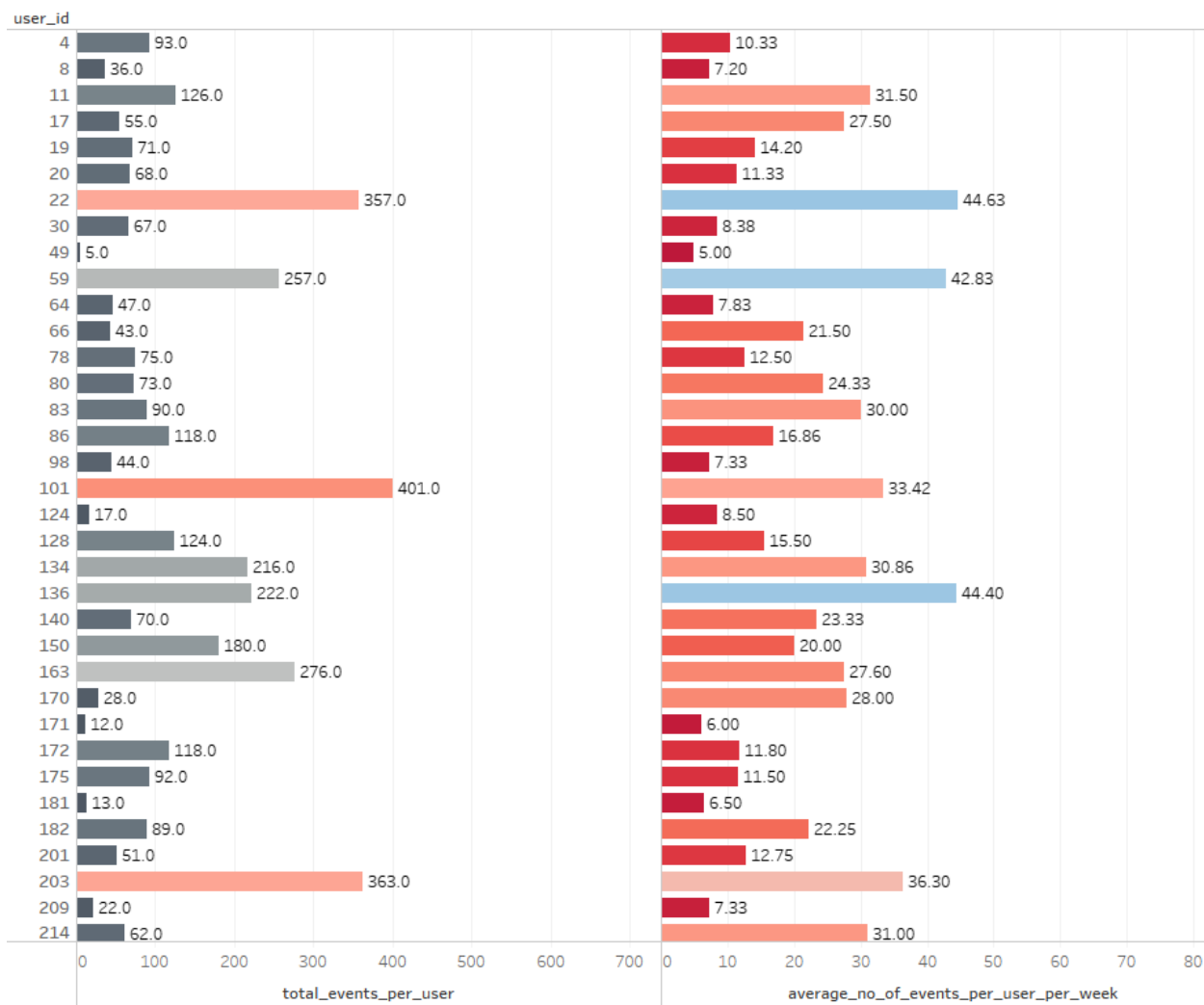
```
---
116 •   SELECT
117         user_id,
118         sum(event_name) as total_events_per_user,
119         AVG(event_name) AS average_no_of_events_per_user_per_week
120   ⊝ FROM (
121         SELECT
122             user_id,
123             COUNT(event_name) AS event_name
124         FROM events
125         WHERE
126             event_type='engagement'
127         GROUP BY user_id, week(occurred_at)
128     ) subquery
129     GROUP BY user_id
130     ORDER BY user_id;
131
```

| user_id | total_events_per_user | average_no_of_events_per_user_per_week |
|---------|----------------------|-----------------------------------------|
| 4 | 93 | 10.3333 |
| 8 | 36 | 7.2000 |
| 11 | 126 | 31.5000 |
| 17 | 55 | 27.5000 |
| 19 | 71 | 14.2000 |
| 20 | 68 | 11.3333 |
| 22 | 357 | 44.6250 |
| 30 | 67 | 8.3750 |

Weekly Engagement per User

## 2.2. User Growth rate
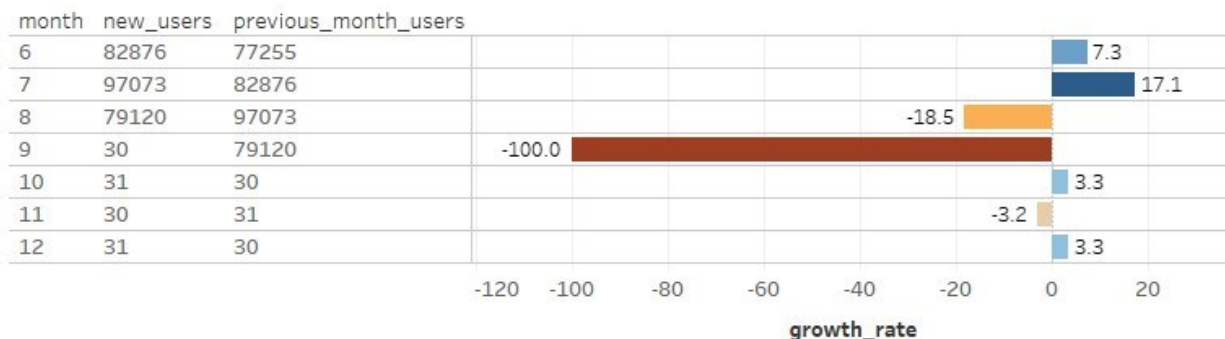
```
134 ●    SELECT
135          extract(month from occurred_at) AS month,
136          COUNT(*) AS new_users,
137          LAG(COUNT(*)) OVER (ORDER BY extract(month from occurred_at)) AS previous_month_users,
138
139          (COUNT(*) - LAG(COUNT(*)) OVER (ORDER BY extract(month from occurred_at))) /
140          LAG(COUNT(*)) OVER (ORDER BY extract(month from occurred_at)) * 100 AS growth_rate
141      FROM
142          events
143      WHERE
144          occurred_at >= '2014-01-01'
145          AND occurred_at< '2015-01-01'
146      GROUP BY
147          extract(month from occurred_at)
148      ORDER BY
149          extract(month from occurred_at);
150
```

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

| month | new_users | previous_month_users | growth_rate |
|---|---|---|---|
| 5 | 77255 | NULL | NULL |
| 6 | 82876 | 77255 | 7.2759 |
| 7 | 97073 | 82876 | 17.1304 |
| 8 | 79120 | 97073 | -18.4943 |
| 9 | 30 | 79120 | -99.9621 |
| 10 | 31 | 30 | 3.3333 |
| 11 | 30 | 31 | -3.2258 |
| 12 | 31 | 30 | 3.3333 |

### User Growth Rate over Month

| month | new_users | previous_month_users | growth_rate |
|---|---|---|---|
| 6 | 82876 | 77255 | 7.3 |
| 7 | 97073 | 82876 | 17.1 |
| 8 | 79120 | 97073 | -18.5 |
| 9 | 30 | 79120 | -100.0 |
| 10 | 31 | 30 | 3.3 |
| 11 | 30 | 31 | -3.2 |
| 12 | 31 | 30 | 3.3 |

growth_rate

July Month has the highest user growth rate.

```
153 •   SELECT
154         extract(month from occurred_at) AS month, extract(week from occurred_at) as week,
155         COUNT(*) AS new_users,
156         LAG(COUNT(*)) OVER (ORDER BY extract(month from occurred_at)) AS previous_new_users,
157         (COUNT(*) - LAG(COUNT(*)) OVER (ORDER BY extract(month from occurred_at))) /
158         LAG(COUNT(*)) OVER (ORDER BY extract(month from occurred_at)) * 100 AS growth_rate
159     FROM
160         events
161     WHERE
162         occurred_at >= '2014-01-01'
163         AND occurred_at< '2015-01-01'
164     GROUP BY
165         extract(month from occurred_at), extract(week from occurred_at)
166     ORDER BY
167         extract(month from occurred_at);
168
```

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

| month | week | new_users | previous_new_users | growth_rate |
|-------|------|-----------|--------------------|-------------|
| 5 | 17 | 7404 | NULL | NULL |
| 5 | 18 | 15850 | 7404 | 114.0735 |
| 5 | 19 | 17155 | 15850 | 8.2334 |
| 5 | 20 | 18775 | 17155 | 9.4433 |
| 5 | 21 | 18071 | 18775 | -3.7497 |
| 6 | 22 | 19444 | 18071 | 7.5978 |
| 6 | 23 | 19317 | 19444 | -0.6532 |
| 6 | 24 | 20217 | 19317 | 4.6591 |

## 2.3. Weekly Retention rate: Users engaging product weekly after signing-up for a product
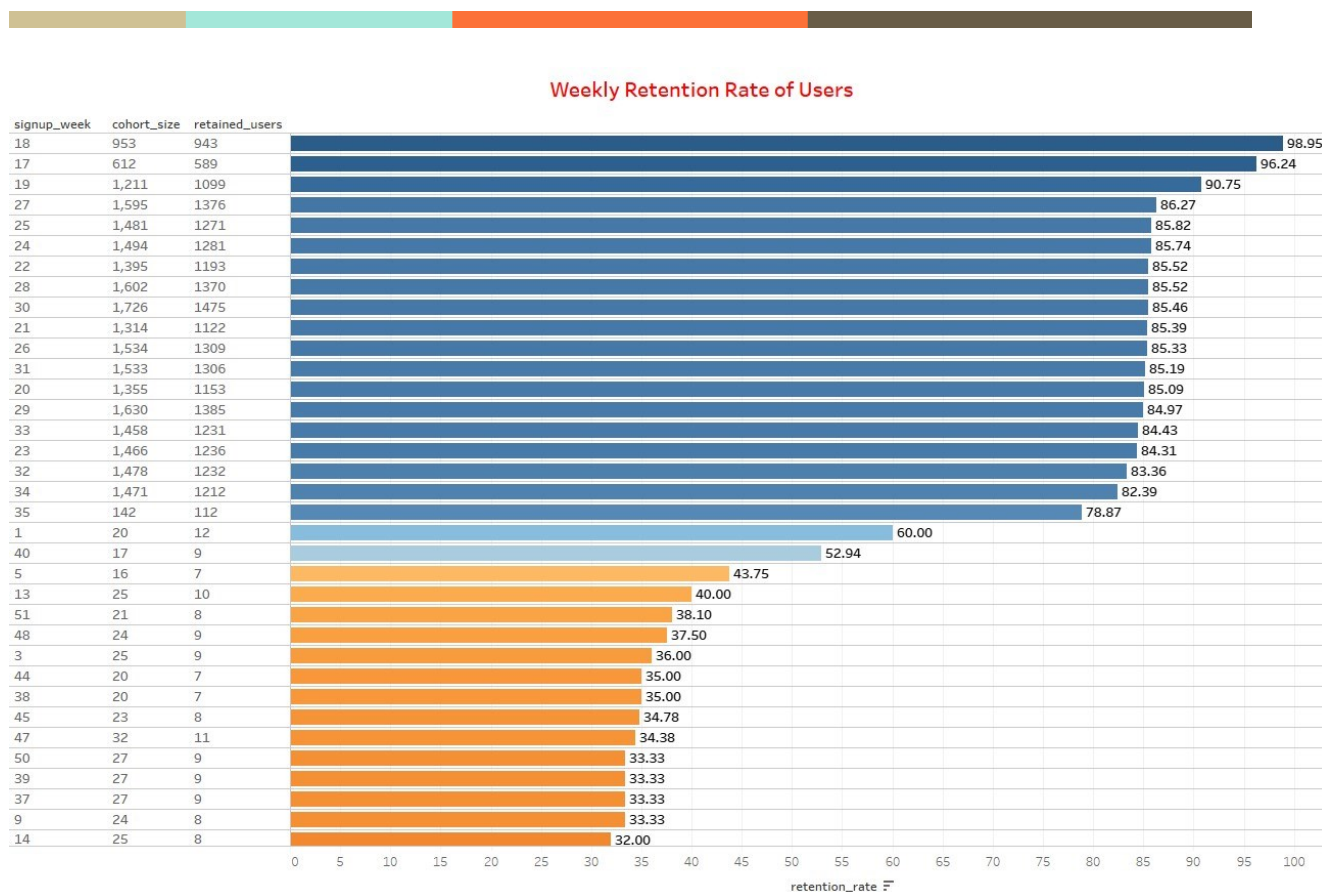
```sql
171    SELECT
172        signup_week,
173        cohort_size,
174        retained_users,
175        (retained_users / cohort_size) * 100 AS retention_rate
176    FROM (
177        SELECT
178            signup_week,
179            COUNT(DISTINCT user_id) AS cohort_size,
180            count(DISTINCT CASE WHEN signup_week < activity_week AND event_type = 'signup_flow'
181            and event_name='complete_signup' then user_id END) AS signup__users,
182            count(DISTINCT CASE WHEN activity_week >= signup_week AND event_type = 'engagement'
183            and event_name='login' then user_id END) AS retained_users
184        FROM (
185            SELECT
186                user_id,
187                event_type, event_name,
188                extract(week from occurred_at) AS signup_week,
189                extract(week from occurred_at) AS activity_week
190            FROM
191                events
192        ) sub1
193        GROUP BY
194            signup_week
195    ) sub2
196    ORDER BY
197        signup_week;
```

| signup_week | cohort_size | retained_users | retention_rate |
|---|---|---|---|
| 0 | 17 | 4 | 23.5294 |
| 1 | 20 | 12 | 60.0000 |
| 2 | 27 | 5 | 18.5185 |
| 3 | 25 | 9 | 36.0000 |
| 4 | 27 | 6 | 22.2222 |
| 5 | 16 | 7 | 43.7500 |
| 6 | 19 | 6 | 31.5789 |
| 7 | 21 | 6 | 28.5714 |

## Weekly Retention Rate of Users

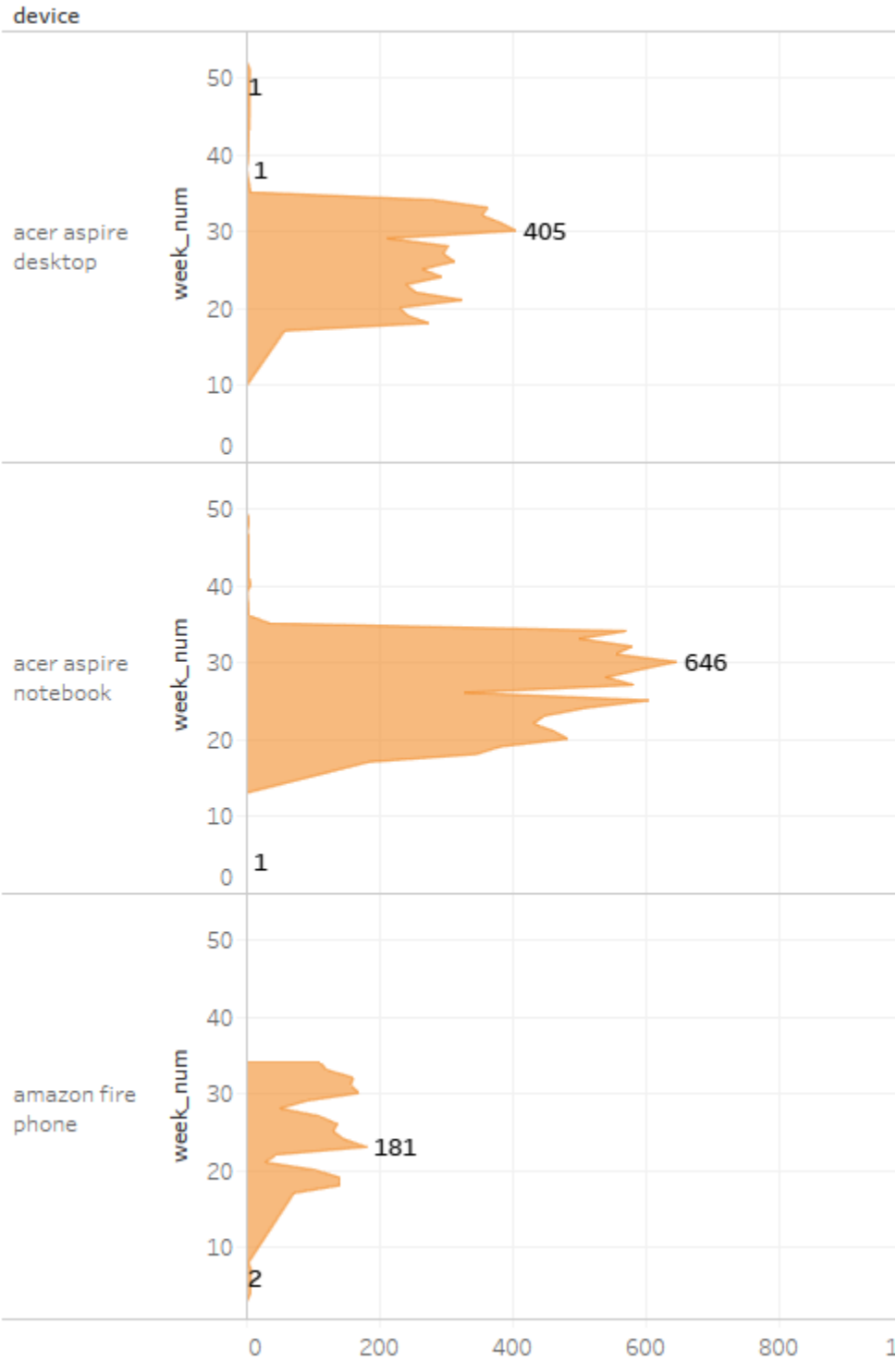| signup_week | cohort_size | retained_users | retention_rate |
|---|---|---|---|
| 18 | 953 | 943 | 98.95 |
| 17 | 612 | 589 | 96.24 |
| 19 | 1,211 | 1099 | 90.75 |
| 27 | 1,595 | 1376 | 86.27 |
| 25 | 1,481 | 1271 | 85.82 |
| 24 | 1,494 | 1281 | 85.74 |
| 22 | 1,395 | 1193 | 85.52 |
| 28 | 1,602 | 1370 | 85.52 |
| 30 | 1,726 | 1475 | 85.46 |
| 21 | 1,314 | 1122 | 85.39 |
| 26 | 1,534 | 1309 | 85.33 |
| 31 | 1,533 | 1306 | 85.19 |
| 20 | 1,355 | 1153 | 85.09 |
| 29 | 1,630 | 1385 | 84.97 |
| 33 | 1,458 | 1231 | 84.43 |
| 23 | 1,466 | 1236 | 84.31 |
| 32 | 1,478 | 1232 | 83.36 |
| 34 | 1,471 | 1212 | 82.39 |
| 35 | 142 | 112 | 78.87 |
| 1 | 20 | 12 | 60.00 |
| 40 | 17 | 9 | 52.94 |
| 5 | 16 | 7 | 43.75 |
| 13 | 25 | 10 | 40.00 |
| 51 | 21 | 8 | 38.10 |
| 48 | 24 | 9 | 37.50 |
| 3 | 25 | 9 | 36.00 |
| 44 | 20 | 7 | 35.00 |
| 38 | 20 | 7 | 35.00 |
| 45 | 23 | 8 | 34.78 |
| 47 | 32 | 11 | 34.38 |
| 50 | 27 | 9 | 33.33 |
| 39 | 27 | 9 | 33.33 |
| 37 | 27 | 9 | 33.33 |
| 9 | 24 | 8 | 33.33 |
| 14 | 25 | 8 | 32.00 |

Week-18 has the highest retention rate of signup users.

## 2.4.  Weekly Engagement per Device

```
198
199 •     SELECT
200           device,
201           EXTRACT(WEEK FROM occurred_at) AS week_num,
202           COUNT(event_name) AS total_events
203       FROM
204           events
205       WHERE
206           event_type = 'engagement'
207       GROUP BY device , week_num;
208
```

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

| device | week_num | total_events |
| --- | --- | --- |
| dell inspiron notebook | 17 | 485 |
| dell inspiron notebook | 18 | 877 |
| iphone 5 | 18 | 1273 |
| iphone 5 | 19 | 1166 |
| iphone 5 | 20 | 1281 |
| iphone 4s | 20 | 607 |
| windows surface | 20 | 189 |
| macbook air | 21 | 1279 |

## Total events engaged per device per week
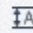
device



Above visualization, shows highest and lowest point of device engagement per week.

```
209 •    SELECT
210          device,
211          sum(event_name) as total_events_per_device,
212          AVG(event_name) AS average_no_of_events_per_device_per_week
213    ⊖ FROM (
214        SELECT
215            device,
216            COUNT(event_name) AS event_name
217        FROM events
218        where
219            event_type='engagement'
220        GROUP BY device, week(occurred_at)
221    ) subquery
222    GROUP BY device;
223
```

| device | total_events_per_device | average_no_of_events_per_device_per_week |
|---|---|---|
| dell inspiron notebook | 19470 | 499.2308 |
| iphone 5 | 25576 | 710.4444 |
| iphone 4s | 9506 | 352.0741 |
| windows surface | 3399 | 147.7826 |
| macbook air | 26482 | 615.8605 |
| iphone 5s | 15744 | 414.3158 |
| macbook pro | 56670 | 1069.2453 |
| kindle fire | 4040 | 202.0000 |

## Weekly Engagement per Device



| device | total_events_per_device | average_no_of_events_per_device_per_week |
|---|---|---|
| macbook pro | 56,670 | 1,069.2 |
| lenovo thinkpad | 36,583 | 762.1 |
| macbook air | 26,482 | 615.9 |
| iphone 5 | 25,576 | 710.4 |
| dell inspiron notebook | 19,470 | 499.2 |
| samsung galaxy s4 | 18,436 | 409.7 |
| nexus 5 | 16,332 | 453.7 |
| iphone 5s | 15,744 | 414.3 |
| dell inspiron desktop | 10,037 | 418.2 |
| iphone 4s | 9,506 | 352.1 |
| asus chromebook | 9,428 | 349.2 |
| ipad air | 9,341 | 322.1 |
| acer aspire notebook | 8,824 | 284.6 |
| hp pavilion desktop | 8,788 | 283.5 |
| nexus 7 | 6,456 | 230.6 |
| nokia lumia 635 | 5,549 | 213.4 |
| ipad mini | 5,525 | 157.9 |
| acer aspire desktop | 5,120 | 204.8 |
| nexus 10 | 5,067 | 220.3 |
| mac mini | 4,410 | 183.8 |
| htc one | 4,228 | 151.0 |
| kindle fire | 4,040 | 202.0 |
| windows surface | 3,399 | 147.8 |
| samsung galaxy note | 2,642 | 101.6 |
| amazon fire phone | 2,136 | 89.0 |
| samsumg galaxy tablet | 1,786 | 94.0 |

## 2.5. Email-engagement metrics

```sql
226 •   SELECT
227         ee.action,
228         COUNT(ee.action) AS event_count,
229         COUNT(*) / (select count(*) from email_events) * 100 AS percentage_share
230     FROM
231         email_events ee
232     GROUP BY
233         ee.action;
234
```

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

| action | event_count | percentage_share |
|---|---|---|
| sent_weekly_digest | 57267 | 63.3562 |
| email_open | 20459 | 22.6344 |
| email_clickthrough | 9010 | 9.9680 |
| sent_reengagement_email | 3653 | 4.0414 |

## Email Engagement Metrics



email_clickthrough
9,010
9.97

email_open
20,459
22.63

sent_reengagement_email
3,653
4.04

sent_weekly_digest
57,267
63.36