# CASE STUDY 1: JOB DATA ANALYSIS

**Project Description:**

The purpose of this project is to perform data analysis on the '**job_data**' table, focusing on four main tasks: Jobs Reviewed Over Time, Throughput Analysis, Language Share Analysis, and Duplicate Rows Detection. The analysis aims to gain insights into the job review patterns, review throughput, language distribution, and data quality issues in the table. By completing these tasks, we can better understand how jobs are reviewed, identify any performance trends, assess language preferences, and ensure data accuracy.

**Approach:**

To accomplish the analysis, I followed these steps:

- Familiarized myself with the '**job_data**' table's structure and columns.
- Formulated SQL queries for each task, considering the appropriate filters, groupings, and calculations needed for each analysis.
- Executed the queries in MySQL Workbench, connected to the database containing the '**job_data**' table.
- Reviewed the results to gain insights into the data patterns and trends.
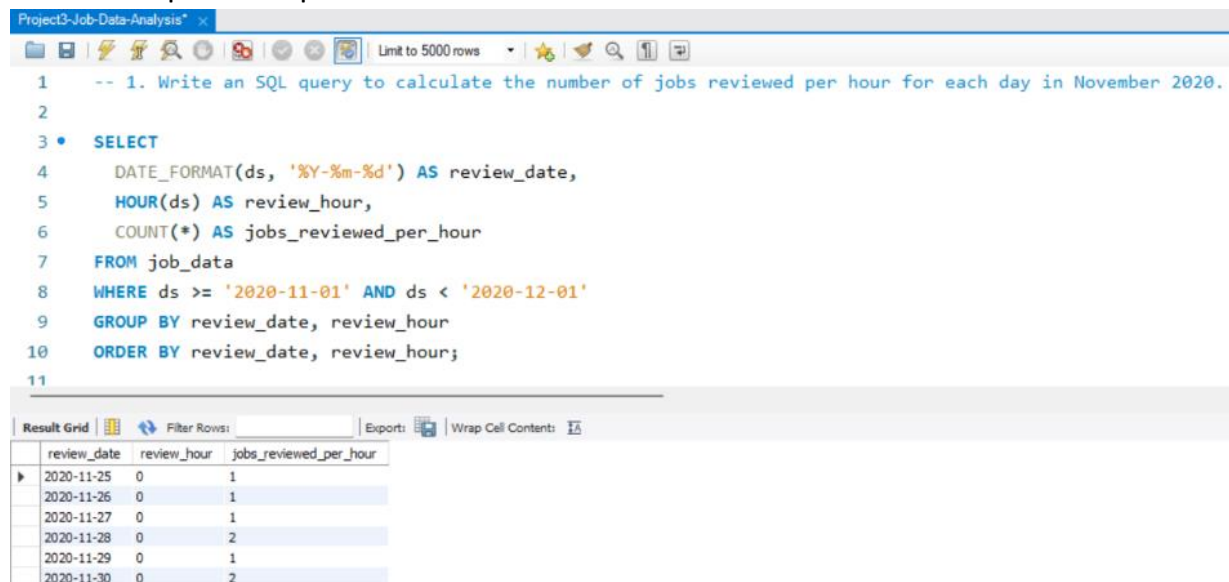- Ensured data quality by identifying and handling duplicate rows.

**Tech-Stack Used:**

For this project, I utilized the following software and tools:

- **MySQL Workbench:** Used to interact with the MySQL database, write and execute SQL queries, and visualize the results.
- **Excel:** Utilized to store and manipulate the dataset.

**Insights:**

a. **Jobs Reviewed Over Time:** The analysis provided a breakdown of the number of jobs reviewed per hour for each day in November 2020. This information revealed peak hours of job reviews, which can be used for resource allocation and process optimization.



The result provides 6 records about the workload distribution throughout November 2020, which can be used to identify peak hours and better manage resources.

b. **Throughput Analysis:** The 7-day rolling average of throughput was calculated to smooth out daily variations and provide a more stable representation of event counts. This approach helped in identifying long-term trends and performance changes, making it preferable over daily metrics for understanding overall system performance.

I prefer using the 7-day rolling average for throughput analysis over the daily metric. Because:

1. **Smoothing Out Variations:** The 7-day rolling average smooths out daily fluctuations, providing a stable trend over time. Daily metrics can be erratic due to events like weekends or holidays, but the rolling average offers a clearer view of the overall trend by reducing short-term impacts.

2. **Identifying Long-Term Patterns:** The rolling average reveals long-term trends and system behaviour effectively. Daily metrics may obscure these patterns with random spikes or fluctuations, making the rolling average a more reliable indicator.

3. **Decision Making:** For making data-driven decisions, a stable and consistent metric is crucial. The rolling average's steadiness makes it a better choice for decision-making, ensuring choices are based on reliable performance data.

4. **Forecasting:** By analysing the historical rolling average, we can predict future throughput trends with greater accuracy. This helps in planning resources and capacity more effectively for the system's future needs.

c. **Language Share Analysis:** The percentage share of each language in the last 30 days was calculated to assess language preferences. This analysis helped identify which languages were most frequently reviewed, enabling better content localization and user satisfaction.



The result is helpful in understanding which languages are in high demand or need more attention. It can also be valuable for resource allocation, language-specific optimizations, and catering to the preferences of the majority of users.
**(We don't have any such language over last 30 days records).**

**d. Duplicate Rows Detection:** Identifying duplicate rows in the '**job_data**' table was crucial for data quality assurance. This process allowed for the detection of potential data entry errors or system issues that could impact the accuracy of analysis results.



The result shows that **we don't have any duplicate records in our dataset.**

## Result:

Through the project, we successfully conducted a comprehensive analysis of the '**job_data**' table. The insights gained from the analysis provided valuable information for optimizing job review processes, understanding throughput trends, and tailoring language support for users. Additionally, detecting and handling duplicate rows contributed to data integrity, ensuring reliable and accurate analyses in the future. The project's outcomes have enhanced decision-making capabilities, improved data quality, and facilitated a better understanding of job review patterns and user preferences within the organization.

# CASE STUDY 2: INVESTIGATING METRIC SPIKE

**Project Description:**
The purpose of this project is to analyse user engagement and behaviour in a product using data from three tables: **users**, **events and email_events**. The main objectives are to measure weekly user engagement, analyse user growth over time, conduct weekly retention analysis, measure weekly engagement per device, and analyse email engagement metrics. By performing this analysis, we aim to gain valuable insights into user activity, identify trends, and make data-driven decisions to enhance the product's user experience and engagement.

**Approach:**
To accomplish the project objectives, we will use SQL queries to extract and analyse data from the three tables.

- For weekly user engagement, we will group events by week and user, calculating the number of engagements per user on a weekly basis.
- For user growth analysis, we will focus on the created_at column in the users table, grouping users by their registration week to track growth over time.
- For weekly retention analysis, we will identify sign-up cohorts by finding the minimum occurred_at date for each user in the events table. We will then join this cohort information with user engagements in different weeks to calculate retention rates.
- To measure weekly engagement per device, we will group events by week and device, providing the count of engagements for each device type on a weekly basis.
- For email engagement analysis, we will utilize the email_events table to count the number of engaged users and total engagements for each email action.

**Tech-Stack Used:**
For this project, I utilized the following software and tools:
- **MySQL Workbench:** Used to interact with the MySQL database, write and execute SQL queries, and visualize the results.
- **Excel:** Utilized to store and manipulate the dataset.

**Insights:**
a. **Weekly User Engagement:** The task calculates the weekly user engagement for each user by counting the number of engagement events they performed in each week.

The result shows 22015 records of the user_id, extracts the week number from the occurred_at column, and counts the number of engagement events for each user during that particular week. It also filters the events table for event_type 'engagement' to focus on general product usage after the user has signed up.

b. **User Growth Analysis:** This task calculates the user growth over time by counting the number of new users who signed up on each date.

The result shows 605 records of new users signed up for an account on the platform or service on a particular registration date.

c. **Weekly Retention Analysis:** The task calculates the weekly retention of users based on their sign-up cohort. It tracks users who signed up in a specific week (cohort_week) and then calculates the number of users from that cohort who remained engaged in each subsequent week.



The result of 1002 records provide insights for users who signed up in cohort creation week, the retention is being measured. Out of all the users who signed up in cohort week, number of users were retained and performed an 'engagement' event in retention week.

d. **Weekly Engagement Per Device:** This task calculates the weekly user engagement per device type by counting the number of engagement events performed on each device in each week.

The result shows 491 records of the weekly engagement for users on a specific device during calculated week number. The engagement count represents the total number of engagement events logged by users on their particular devices during that particular week.

e. **Email Engagement Analysis:** This task calculates the email engagement metrics by counting the occurrences of each action in the email_events table.

The result shows 4 records of total engagements where users interacted with emails, and out of those engagements, counts unique users clicked through the links within the emails. This implies that a significant number of users found the content of the emails compelling or relevant enough to take action and click on the links provided.

**Result:**

The project provides valuable insights into user behaviour and engagement patterns within the product. This information facilitates informed decision-making, leading to improvements in user experience, and optimization of the product for better user retention and growth. By identifying areas for enhancement, such as onboarding optimization, email content improvement, and personalized feature implementation, the project contributes significantly to overall product development and marketing strategies. The achieved insights will guide future strategies to create a more engaging and user-centric product experience.