

56 SQL Problems (Easy to Hard Level)

Problem 1 (Meta Hard Level)	3
Problem 2 (Amazon Hard Level)	5
Problem 3 (Google Medium Level)	7
Problem 4 (Uber Hard Level)	8
Problem 5 (Microsoft Medium Level)	10
Problem 6 (Airbnb Medium Level)	11
Problem 7 (IBM Hard Level)	12
Problem 8 (Tesla Medium Level)	14
Problem 9 (Netflix Hard Problem)	15
Problem 10 (Amazon Medium Level)	16
Problem 11 (Nvidia, Microsoft Medium Level)	17
Problem 12 (LinkedIn, Dropbox Basic Level)	18
Problem 13 (Expedia, Airbnb Basic Level)	18
Problem 14 (Amazon, Salesforce Basic Level)	18
Problem 15 (Google Medium Level)	19
Problem 16 (JP Morgan, Chase, Bloomberg Medium Level)	20
Problem 17 (Uber Medium Level)	21
Problem 18 (Amazon, Doordash Medium Level)	21
Problem 19 (Walmart Hard Level)	22
Problem 20 (Apple, Microsoft, Dell Easy Level)	22
Problem 21 (Microsoft Hard Level)	23
Problem 22 (Walmart, Paypal Medium Level)	24
Problem 23 (Oracle Hard Level)	25
Problem 24 (Amazon Hard Level)	25
Problem 25 (American Express Medium Level)	26
Problem 26 (LinkedIn Hard Level)	26
Problem 27 (Goldman Sachs, Deloitte Hard Level)	27
Problem 28 (Meta, Salesforce Hard Level)	28
Problem 29 (Cisco Hard Level)	28
Problem 30 (Amazon Hard Level)	29
Problem 31 (Spotify Hard Level)	30
Problem 32 (Accenture Medium Level)	30
Problem 33 (Google Hard Level)	31
Problem 34 (Google, Airbnb, Expedia Medium Level)	32
Problem 35 (Twitter Hard Level)	32

Problem 36 (Uber Hard Level)	33
Problem 37 (Netflix Hard Level)	34
Problem 38 (Airbnb Hard Level)	35
Problem 39 (Walmart (Hard Level)	36
Problem 40 (Walmart (Hard Level)	36
Problem 41 (Meta Easy Level)	37
Problem 42 (Microsoft Medium Level)	37
Problem 43 (Apple Hard Level)	38
Problem 44 (Amazon Hard Level)	39
Problem 45 (Visa Medium Level)	39
Problem 46 (EY, TCS, Deloitte Medium Level)	40
Problem 47 (Expedia, Airbnb, Tripadvisor Medium Level)	40
Problem 48 (Amazon, Doordash, Bosch Medium Level)	41
Problem 49 (Goldman Sachs Medium Level)	41
Problem 50 (Meta Hard Level)	42
Problem 51 (Amazon Hard Level)	42
Problem 52 (Tesla Hard Level)	43
Problem 53 (DoorDash Medium Level)	44
Problem 54 (Meta Hard level)	44
Problem 55 (Amazon Hard Level)	45
Problem 56 (ESPN Hard Level)	45

Problem 1 (Meta | Hard Level)

A table named “famous” has two columns called user id and follower id. It represents each user ID has a particular follower ID. These follower IDs are also users of #Facebook / Meta. Then, find the famous percentage of each user.

Famous Percentage = number of followers a user has / total number of users on the platform.

Explanation:

1. distinct_users CTE: Combines user_id and follower_id using UNION to get all unique users on the platform. This helps us determine the total number of users.
2. follower_count CTE: Counts the number of followers for each user_id by grouping the rows in the famous table. This gives a list of users with their follower counts.
3. Final SELECT Statement: Uses the data from follower_count and distinct_users to calculate the famous percentage for each user.

MySQL Solution:

```
1 • WITH distinct_users AS(  
2     SELECT user_id AS users FROM famous  
3     UNION  
4     SELECT follower_id AS users FROM famous  
5 ),  
6 follower_count AS(  
7     SELECT  
8         user_id, COUNT(follower_id) AS followers  
9     FROM  
10        famous  
11    GROUP BY user_id  
12 )  
13 SELECT  
14     f.user_id,  
15     ROUND((f.followers * 100.0) / (SELECT  
16         COUNT(*)  
17     FROM  
18         distinct_users),  
19     2) AS famous_percentage  
20 FROM  
21     follower_count f;
```

MSSQL Solution:

```
1 WITH distinct_users AS(  
2     SELECT user_id AS users FROM famous  
3     UNION  
4     SELECT follower_id AS users FROM famous  
5 ),  
6 follower_count AS(  
7     SELECT  
8         user_id, COUNT(follower_id) AS followers  
9     FROM  
10         famous  
11     GROUP BY user_id  
12 )  
13 SELECT  
14     f.user_id,  
15     CAST((f.followers * 100.0) / (  
16         SELECT  
17             COUNT(*)  
18         FROM  
19             distinct_users) AS DECIMAL(10,2)  
20     ) AS famous_percentage  
21 FROM  
22     follower_count f;
```

Problem 2 (Amazon | Hard Level)

Given a table 'sf_transactions' of purchases by date, calculate the month-over-month percentage change in revenue. The output should include the year-month date (YYYY-MM) and percentage change, rounded to the 2nd decimal point, and sorted from the beginning of the year to the end of the year. The percentage change column will be populated from the 2nd month forward and calculated as

$((\text{this month's revenue} - \text{last month's revenue}) / \text{last month's revenue}) * 100.$

Explanation:

1. MonthlyRevenue CTE: Aggregates the total revenue for each month using FORMAT to convert the created_at date to the format YYYY-MM.
2. RevenueChange CTE: Adds a column previous_revenue using the LAG function, which fetches the total revenue of the previous month for each row.
3. Final SELECT: Calculates the percentage change as $((\text{total_revenue} - \text{previous_revenue}) / \text{previous_revenue}) * 100$. The ROUND function ensures the percentage is rounded to two decimal places. The output is ordered by year_month to display the data chronologically.

MSSQL Server Solution:

```
1 WITH MonthlyRevenue AS (  
2     SELECT  
3         FORMAT(created_at, 'yyyy-MM') AS year_month,  
4         SUM(value) AS total_revenue  
5     FROM  
6         sf_transactions  
7     GROUP BY FORMAT(created_at, 'yyyy-MM')  
8 ),  
9 RevenueChange AS (  
10    SELECT  
11        year_month,  
12        total_revenue,  
13        LAG(total_revenue) OVER (ORDER BY year_month) AS previous_revenue  
14    FROM  
15        MonthlyRevenue  
16 )  
17 SELECT  
18     year_month,  
19     total_revenue,  
20     ROUND(CASE  
21         WHEN previous_revenue IS NULL THEN NULL  
22         ELSE ((total_revenue - previous_revenue) / CAST(previous_revenue AS FLOAT)) * 100  
23     END,  
24     2) AS percentage_change  
25 FROM  
26     RevenueChange  
27 ORDER BY year_month;
```

MySQL Solution:

```
1  WITH MonthlyRevenue AS (  
2      SELECT  
3          DATE_FORMAT(created_at, '%Y-%m') AS yearmonth,  
4          SUM(value) AS total_revenue  
5      FROM  
6          sf_transactions  
7      GROUP BY yearmonth  
8  ),  
9  RevenueChange AS (  
10     SELECT  
11         yearmonth,  
12         total_revenue,  
13         LAG(total_revenue) OVER (ORDER BY yearmonth) AS previous_revenue  
14     FROM  
15         MonthlyRevenue  
16 )  
17 SELECT  
18     yearmonth,  
19     total_revenue,  
20     ROUND(CASE  
21         WHEN previous_revenue IS NULL THEN NULL  
22         ELSE ((total_revenue - previous_revenue) / previous_revenue) * 100  
23     END,  
24     2) AS percentage_change  
25 FROM  
26     RevenueChange  
27 ORDER BY yearmonth;
```

Problem 3 (Google | Medium Level)

You are analyzing a social network dataset at Google. Your task is to find mutual friends between two users, Karl and Hans. There is only one user named Karl and one named Hans in the dataset.

The output should contain 'user_id' and 'user_name' columns.

Explanation:

1. The CTEs (karl_friends and hans_friends) efficiently find all friends for Karl and Hans, respectively.
2. The main query joins these CTEs with the users table to find the users who are present in both Karl's and Hans's friend lists (mutual friends).

```
25  /*solution 1*/
26  WITH karl_friends AS(
27      SELECT DISTINCT friend_id
28      FROM friends
29      WHERE user_id = (SELECT user_id FROM users WHERE user_name = 'Karl')
30  ),
31  hans_friends AS(
32      SELECT DISTINCT friend_id
33      FROM friends
34      WHERE user_id = (SELECT user_id FROM users WHERE user_name = 'Hans')
35  )
36  SELECT
37      u.user_id, u.user_name
38  FROM
39      users u
40      JOIN
41      karl_friends kf ON u.user_id = kf.friend_id
42      JOIN
43      hans_friends hf ON u.user_id = hf.friend_id;
44
45
46  /*solution 2*/
47  SELECT
48      DISTINCT u.user_id, u.user_name
49  FROM
50      users u
51      JOIN
52      friends f1 ON u.user_id = f1.friend_id
53      JOIN
54      friends f2 ON u.user_id = f2.friend_id
55  WHERE
56      f1.user_id = (SELECT user_id FROM users WHERE user_name = 'Karl')
57      AND f2.user_id = (SELECT user_id FROM users WHERE user_name = 'Hans');
```

Problem 4 (Uber | Hard Level)

Some forecasting methods are extremely simple and surprisingly effective. Naïve forecast is one of them. To create a naïve forecast for "distance per dollar" (defined as $\text{distance_to_travel} / \text{monetary_cost}$), first sum the "distance to travel" and "monetary cost" values monthly. This gives the actual value for the current month. For the forecasted value, use the previous month's value. After obtaining both actual and forecasted values, calculate the root mean squared error (RMSE) using the formula

$$RMSE = \sqrt{\text{mean}(\text{square}(\text{actual} - \text{forecast}))}.$$

Report the RMSE rounded to two decimal places.

✳ Steps to Approach:

1. Group the data by month-year to sum `distance_to_travel` and `monetary_cost`.
2. Calculate distance per dollar as:
$$\text{distance per dollar} = \frac{\text{sum of distance}}{\text{sum of monetary cost}}$$
3. For the forecasted value, use the previous month's value.
4. Calculate the actual and forecasted values for each month from February onward.
5. Compute the Root Mean Squared Error (RMSE) using:

$$RMSE = \sqrt{\text{mean}((\text{actual} - \text{forecast})^2)}$$

MySQL Solution:

```
2  ● WITH monthly_aggregates AS(
3      SELECT
4          DATE_FORMAT(request_date, '%Y-%m') AS yearmonth,
5          SUM(distance_to_travel) AS total_distance,
6          SUM(monetary_cost) AS total_cost
7      FROM
8          uber_request_logs
9      GROUP BY yearmonth
10 ),
11 -- Calculate actual distance per dollar
12 -- actual_distance AS(
13     SELECT
14         yearmonth,
15         total_distance / total_cost AS actual_distance
16     FROM
17         monthly_aggregates
18 -- ),
19 -- Generate Naïve Forecast
20 -- naive_forecast AS(
21     SELECT
22         yearmonth,
23         actual_distance,
24         LAG(actual_distance) OVER(ORDER BY yearmonth) AS forecasted_distance
25     FROM
26         actual_distance
27 -- )
28 -- Calculate RMSE
29 SELECT
30     ROUND(SQRT(AVG(POWER(actual_distance - forecasted_distance, 2))), 2) AS rmse
31 FROM
32     naive_forecast
33 WHERE
34     forecasted_distance IS NOT NULL;
```


MSSQL Solution:

```
35 -- Calculate monthly aggregates
36 WITH monthly_aggregates AS(
37     SELECT
38         FORMAT(request_date, 'yyyy-MM') AS yearmonth,
39         SUM(distance_to_travel) AS total_distance,
40         SUM(monetary_cost) AS total_cost
41     FROM
42         uber_request_logs
43     GROUP BY FORMAT(request_date, 'yyyy-MM')
44 ),
45 -- Calculate actual distance per dollar
46 actual_distance AS(
47     SELECT
48         yearmonth,
49         total_distance / total_cost AS actual_distance
50     FROM
51         monthly_aggregates
52 ),
53 -- Generate Naïve Forecast
54 naive_forecast AS(
55     SELECT
56         yearmonth,
57         actual_distance,
58         LAG(actual_distance) OVER(ORDER BY yearmonth) AS forecasted_distance
59     FROM
60         actual_distance
61 )
62 -- Calculate RMSE
63 SELECT
64     ROUND(SQRT(AVG(POWER(actual_distance - forecasted_distance, 2))), 2) AS rmse
65 FROM
66     naive_forecast
67 WHERE
68     forecasted_distance IS NOT NULL;
```

Problem 5 (Microsoft | Medium Level)

Given a list of projects and employees mapped to each project, calculate by the amount of project budget allocated to each employee. The output should include the project title and the project budget rounded to the closest integer. Order your list by projects with the highest budget per employee first.

Explanation:

1. **Joining Tables:** The initial step involves joining the `ms_projects` and `ms_emp_projects` tables on the project ID to combine project details (including titles and budgets) with employee assignments.
2. **Grouping and Aggregating:** The data is then grouped by project title and budget, allowing for the calculation of budget per employee by dividing the total budget of each project by the count of employees assigned to that project.
3. **Rounding and Ordering:** Finally, the computed budget per employee is rounded to the nearest integer, and the results are ordered in descending order to prioritize projects with the highest budget allocation per employee.

```
39 -- Solution
40 SELECT
41     p.title AS project_title,
42     ROUND(p.budget / COUNT(e.project_id), 0) AS budget_per_employee
43 FROM
44     ms_projects p
45     INNER JOIN
46     ms_emp_projects e ON p.id = e.project_id
47 GROUP BY p.id , p.title , p.budget
48 ORDER BY budget_per_employee DESC;
49
50 -- in real case scenario there will be multiple employees tagged to one project
51 -- this one is correct solution
52 SELECT
53     p.title AS project_title,
54     ROUND(SUM(p.budget) / COUNT(e.project_id), 0) AS budget_per_employee
55 FROM
56     ms_projects p
57     INNER JOIN
58     ms_emp_projects e ON p.id = e.project_id
59 GROUP BY p.id , p.title , p.budget
60 ORDER BY budget_per_employee DESC;
```

Problem 6 (Airbnb | Medium Level)

Find the total number of available beds per hosts' nationality. Output the nationality along with the corresponding total number of available beds. Sort records by the total available beds in descending order.

Explanation

1. **Joining Tables:** The first step involves joining the `airbnb_apartments` and `airbnb_hosts` tables on the `host_id`. This allows us to combine the apartment details (such as the number of beds) with the host's nationality information.
2. **Grouping and Aggregating:** Next, the data is grouped by the host's nationality, so that the total number of beds available for each nationality can be calculated. The `SUM()` function is used to add up the beds (`n_beds`) for all apartments hosted by individuals of the same nationality.
3. **Sorting the Results:** Finally, the results are ordered in descending order based on the total number of available beds.

```
35 -- Solution
36 SELECT
37     h.nationality,
38     SUM(a.n_beds) AS total_available_beds
39 FROM
40     airbnb_apartments a
41     INNER JOIN
42     airbnb_hosts h ON a.host_id = h.host_id
43 GROUP BY h.nationality;
```

Problem 7 (IBM | Hard Level)

IBM is working on a new feature to analyze user purchasing behavior for all Fridays in the first quarter of the year. For each Friday separately, calculate the average amount users have spent per order. The output should contain the week number of that Friday and average amount spent.

To solve this problem, we need to:

1. Identify the Fridays in the first quarter (Q1) of the year.
2. Calculate the week number for each of these Fridays.
3. Group the purchases by week number and calculate the average amount spent per order.

MSSQL Solution

```
29 -- solution with cte
30 WITH q1_friday AS(
31     SELECT
32         amount_spent, DATEPART(WEEK, date) AS week_number
33     FROM
34         user_purchases
35     WHERE
36         day_name = 'Friday'
37 )
38 SELECT
39     week_number, ROUND(AVG(amount_spent), 2) AS avg_amount_spent
40 FROM
41     q1_friday
42 GROUP BY week_number
43 ORDER BY week_number;
44
45 -- without cte
46 SELECT
47     DATEPART(WEEK, date) AS week_number,
48     ROUND(AVG(amount_spent), 2) AS avg_amount_spent
49 FROM
50     user_purchases
51 WHERE
52     day_name = 'Friday'
53 GROUP BY DATEPART(WEEK, date)
54 ORDER BY week_number;
```

MySQL Solution:

```
1  -- solution with cte
2  ● WITH q1_friday AS(
3      SELECT
4          amount_spent, WEEK(date) AS week_number
5      FROM
6          user_purchases
7      WHERE
8          day_name = 'Friday'
9  )
10 SELECT
11     week_number, ROUND(AVG(amount_spent), 2) AS avg_amount_spent
12 FROM
13     q1_friday
14 GROUP BY week_number
15 ORDER BY week_number;
16
17 -- without cte
18 ● SELECT
19     WEEK(date) AS week_number,
20     ROUND(AVG(amount_spent), 2) AS avg_amount_spent
21 FROM
22     user_purchases
23 WHERE
24     day_name = 'Friday'
25 GROUP BY week_number
26 ORDER BY week_number;
```

Problem 8 (Tesla | Medium Level)

You are given a table of product launches by company by year. Write a query to count the net difference between the number of products companies launched in 2020 with the number of products companies launched in the previous year. Output the name of the companies and a net difference of net products released for 2020 compared to the previous year.

Explanation:

1. Counting Products per Year: Using SUM with CASE statements, we count the number of products launched in 2020 and 2019 separately for each company.
2. Calculating Net Difference: We calculate the difference between 2020 and 2019 product counts to get the net change.
3. Ordering: The results are ordered by net_difference in descending order to show companies with the highest increase first.

```
26 -- solution
27 WITH products_counts AS(
28     SELECT
29         company_name,
30         SUM(CASE WHEN year = 2020 THEN 1 ELSE 0 END) AS products_2020,
31         SUM(CASE WHEN year = 2019 THEN 1 ELSE 0 END) AS products_2019
32     FROM
33         car_launches
34     GROUP BY company_name
35 )
36 SELECT
37     company_name,
38     (products_2020 - products_2019) AS product_difference
39 FROM
40     products_counts
41 ORDER BY product_difference DESC;
```

Problem 9 (Netflix | Hard Problem)

Find the genre of the person with the most number of oscar winnings.

If there are more than one person with the same number of oscar wins, return the first one in alphabetic order based on their name. Use the names as keys when joining the tables.

Explanation:

1. WinnerCount CTE: Calculates the total Oscar wins for each nominee by counting rows where winner = 1.
2. Final Selection: The TOP 1 clause fetches all rows with the highest total_wins, sorted alphabetically by name to handle ties. We join the WinnerCount CTE with nominee_information on the nominee's name to retrieve the top_genre for the top nominee(s) in terms of Oscar wins.

MSSQL Solution

```
41 -- solution
42 WITH winner_count AS(
43     SELECT
44         nominee, COUNT(*) AS total_wins
45     FROM
46         oscar_nominees
47     WHERE
48         winner = 1
49     GROUP BY nominee
50 )
51 SELECT TOP 1
52     ni.top_genre
53 FROM
54     winner_count wc
55     JOIN
56     nominee_information ni ON wc.nominee = ni.name
57 ORDER BY wc.total_wins DESC , ni.name ASC;
```

MySQL Solution

```
1 -- solution
2 WITH winner_count AS(
3     SELECT
4         nominee, COUNT(*) AS total_wins
5     FROM
6         oscar_nominees
7     WHERE
8         winner = 1
9     GROUP BY nominee
10 )
11 SELECT
12     ni.top_genre
13 FROM
14     winner_count wc
15     JOIN
16     nominee_information ni ON wc.nominee = ni.name
17 ORDER BY wc.total_wins DESC , ni.name ASC
18 LIMIT 1;
```

Problem 10 (Amazon | Medium Level)

Write a query that'll identify returning active users. A returning active user is a user that has made a second purchase within 7 days of any other of their purchases. Output a list of user_ids of these returning active users.

MSSQL Solution

```
31 -- solution 1
32 WITH cte_next_purchase AS(
33     SELECT
34         user_id,
35         created_at,
36         LAG(created_at) OVER(PARTITION BY user_id ORDER BY created_at DESC) AS next_purchase
37     FROM
38         amazon_transactions
39 )
40 SELECT
41     user_id
42 FROM
43     cte_next_purchase
44 WHERE
45     DATEDIFF(DAY, created_at, next_purchase) <= 7;
46
47 -- solution 2
48 SELECT DISTINCT
49     a.user_id
50 FROM
51     amazon_transactions a
52     JOIN
53     amazon_transactions b ON a.user_id = b.user_id
54     AND a.created_at < b.created_at
55 WHERE
56     DATEDIFF(DAY, a.created_at, b.created_at) <= 7;
```

MySQL Solution

```
31 -- solution 1
32 WITH cte_next_purchase AS(
33     SELECT
34         user_id,
35         created_at,
36         LAG(created_at) OVER(PARTITION BY user_id ORDER BY created_at DESC) AS next_purchase
37     FROM
38         amazon_transactions
39 )
40 SELECT
41     user_id
42 FROM
43     cte_next_purchase
44 WHERE
45     DATEDIFF(next_purchase, created_at) <= 7;
46
47 -- solution 2
48 SELECT DISTINCT
49     a.user_id
50 FROM
51     amazon_transactions a
52     JOIN
53     amazon_transactions b ON a.user_id = b.user_id
54     AND a.created_at < b.created_at
55 WHERE
56     DATEDIFF(b.created_at, a.created_at) <= 7;
```


Problem 11 (Nvidia, Microsoft | Medium Level)

Find the number of transactions that occurred for each product. Output the product name along with the corresponding number of transactions and order records by the product id in ascending order. You can ignore products without transactions.

Explanation:

1. **Joining Tables:** The INNER JOIN between excel_sql_inventory_data (aliased as inv) and excel_sql_transaction_data (aliased as trans) matches records by product_id. This way, only products with transactions are included.
2. **Counting Transactions:** Using COUNT(trans.transaction_id) counts the number of transactions for each product.
3. **Grouping and Ordering:** GROUP BY inv.product_id, inv.product_name groups by product_id and product_name to get the transaction count per product. ORDER BY inv.product_id ASC sorts the output by product_id in ascending order.

```
52 | -- solution
53 | SELECT
54 |     i.product_name,
55 |     COUNT(t.transaction_id) AS transaction_count
56 | FROM
57 |     excel_sql_inventory_data i
58 |     JOIN
59 |     excel_sql_transaction_data t ON i.product_id = t.product_id
60 | GROUP BY t.product_id , i.product_name
61 | ORDER BY t.product_id ASC;
```

Problem 12 (LinkedIn, Dropbox | Basic Level)

Write a query that calculates the difference between the highest salaries found in the marketing and engineering departments. Output just the absolute difference in salaries.

Explanation

1. CASE is used to selectively get the salary for the "marketing" and "engineering" departments.
2. MAX is applied to retrieve the highest salary in each department.
3. ABS calculates the absolute difference between the two maximum values.

```
35 -- solution
36 SELECT ABS(
37     MAX(CASE WHEN d.department = 'marketing' THEN e.salary END)
38     -
39     MAX(CASE WHEN d.department = 'engineering' THEN e.salary END)
40 ) AS salary_difference
41 FROM
42     db_employee e
43     JOIN
44     db_dept d ON e.department_id = d.id;
```

Problem 13 (Expedia, Airbnb | Basic Level)

Find the number of rows for each review score earned by 'Hotel Arena'. Output the hotel name (which should be 'Hotel Arena'), review score along with the corresponding number of rows with that score for the specified hotel.

```
37 -- solution
38 SELECT
39     hotel_name,
40     reviewer_score,
41     COUNT(*) AS score_count
42 FROM
43     hotel_reviews
44 WHERE
45     hotel_name = 'Hotel Arena'
46 GROUP BY hotel_name, reviewer_score
47 ORDER BY reviewer_score DESC;
```

Problem 14 (Amazon, Salesforce | Basic Level)

What is the total sales revenue of Samantha and Lisa?

```
20 -- solution
21 SELECT
22     SUM(sales_revenue) AS total_sales_revenue
23 FROM
24     sales_performance
25 WHERE
26     salesperson IN ('Samantha', 'Lisa');
```

Problem 15 (Google Medium | Level)

Find all records from days when the number of distinct users receiving emails was greater than the number of distinct users sending emails.

Explanation:

1. The distinct_counts CTE calculates the number of distinct to_user and from_user for each day.
2. The main query joins the original google_gmail_emails table with distinct_counts on the day field, selecting only records where distinct_receivers is greater than distinct_senders.

```
29 -- solution
30 WITH distinct_counts AS(
31     SELECT
32         day,
33         COUNT(DISTINCT from_user) AS distinct_sender,
34         COUNT(DISTINCT to_user) AS distinct_receiver
35     FROM
36         google_gmail_emails
37     GROUP BY day
38 )
39 SELECT
40     g.id, g.from_user, g.to_user, g.day
41 FROM
42     google_gmail_emails g
43     JOIN
44     distinct_counts dc ON g.day = dc.day
45 WHERE
46     dc.distinct_receiver > dc.distinct_sender;
```

Problem 16 (JP Morgan Chase, Bloomberg | Medium Level)

Bank of Ireland has requested that you detect invalid transactions in December 2022. An invalid transaction is one that occurs outside of the bank's normal business hours. The following are the hours of operation for all branches:

Monday - Friday 09:00 - 16:00

Saturday & Sunday Closed

Irish Public Holidays 25th and 26th December

Determine the transaction ids of all invalid transactions.

Explanation:

1. `MONTH(time_stamp) = 12 AND YEAR(time_stamp) = 2022`: This filters transactions to include only those in December 2022.
2. `DATEPART(WEEKDAY, time_stamp) IN (1, 7)`: This checks if the transaction occurred on a Saturday (7) or Sunday (1).
3. `CAST(time_stamp AS TIME) < '09:00:00'`: This checks if the transaction time is before the opening hours.
4. `CAST(time_stamp AS TIME) > '16:00:00'`: This checks if the transaction time is after the closing hours.
5. `(DATEPART(DAY, time_stamp) IN (25, 26) AND MONTH(time_stamp) = 12)`: This checks if the transaction occurred on the public holidays of December 25th or 26th.

MSSQL Solution:

```
30 -- solution
31 SELECT
32     transaction_id
33 FROM
34     boi_transactions
35 WHERE
36     YEAR(time_stamp) = 2022
37     AND MONTH(time_stamp) = 12
38     AND (
39         DATEPART(WEEKDAY, time_stamp) IN (1,7)
40         OR CAST(time_stamp AS TIME) < '09:00:00'
41         OR CAST(time_stamp AS TIME) > '16:00:00'
42         OR (DATEPART(DAY, time_stamp) IN (25,26) AND MONTH(time_stamp) = 12)
43     );
```

MySQL Solution

```
30 -- solution
31 • SELECT
32     transaction_id
33 FROM
34     boi_transactions
35 WHERE
36     YEAR(time_stamp) = 2022
37     AND MONTH(time_stamp) = 12
38     AND (
39         DAYOFWEEK(time_stamp) IN (1, 7)
40         OR CAST(time_stamp AS TIME) < '09:00:00'
41         OR CAST(time_stamp AS TIME) > '16:00:00'
42         OR (DAY(time_stamp) IN (25, 26) AND MONTH(time_stamp) = 12)
43     );
```

Problem 17 (Uber | Medium Level)

You're given a table of Uber rides that contains the mileage and the purpose for the business expense. You're asked to find business purposes that generate the most miles driven for passengers that use Uber for their business transportation. Find the top 3 business purpose categories by total mileage.

```
31 | -- solution
32 | SELECT TOP 3
33 |     purpose,
34 |     SUM(miles) AS total_miles
35 | FROM
36 |     my_uber_drives
37 | GROUP BY purpose
38 | ORDER BY total_miles DESC;
```

Problem 18 (Amazon, Doordash | Medium Level)

You have been asked to find the job titles of the highest-paid employees.

Your output should include the highest-paid title or multiple titles with the same salary.

```
33 | -- solution 1
34 | -- using simple subquery
35 | SELECT
36 |     t.worker_title
37 | FROM
38 |     worker w
39 |     JOIN
40 |     title t ON w.worker_id = t.worker_ref_id
41 | WHERE
42 |     w.salary = (
43 |         SELECT MAX(salary) FROM worker
44 |     );
45 |
46 | -- solution 2
47 | -- using cte and window function
48 | WITH salary_ranks AS(
49 |     SELECT
50 |         t.worker_title,
51 |         w.salary,
52 |         DENSE_RANK() OVER(ORDER BY w.salary DESC) AS worker_rank
53 |     FROM
54 |         worker w
55 |         JOIN
56 |         title t ON w.worker_id = t.worker_ref_id
57 | )
58 | SELECT
59 |     worker_title
60 | FROM
61 |     salary_ranks
62 | WHERE worker_rank = 1;
```

Problem 19 (Walmart | Hard Level)

Identify users who started a session and placed an order on the same day.

For these users, calculate the total number of orders and the total order value for that day.

Your output should include the user, the session date, the total number of orders, and the total order value for that day.

```
37 -- solution
38 SELECT
39     s.user_id,
40     s.session_date,
41     COUNT(os.order_id) AS total_orders,
42     SUM(os.order_value) AS total_order_value
43 FROM
44     sessions s
45     JOIN
46     order_summary os ON s.user_id = os.user_id
47     AND CAST(s.session_date AS DATE) = CAST(os.order_date AS DATE)
48 GROUP BY s.user_id , s.session_date
49 HAVING COUNT(os.order_id) > 0;
```

Problem 20 (Apple, Microsoft, Dell | Easy Level)

Write a query that returns the number of unique users per client per month

MSSQL Solution

```
32 -- solution
33 SELECT
34     client_id,
35     FORMAT(time_id, 'yyyy-MM') AS month_year,
36     COUNT(DISTINCT user_id) AS unique_user_count
37 FROM
38     fact_events
39 GROUP BY client_id, FORMAT(time_id, 'yyyy-MM')
40 ORDER BY client_id;
```

MySQL Solution

```
32 -- solution
33 • SELECT
34     client_id,
35     DATE_FORMAT(time_id, '%Y-%m') AS month_year,
36     COUNT(DISTINCT user_id) AS unique_user_count
37 FROM
38     fact_events
39 GROUP BY client_id , month_year
40 ORDER BY client_id;
```

Problem 21 (Microsoft | Hard Level)

Find the total number of downloads for paying and non-paying users by date. Include only records where non-paying customers have more downloads than paying customers. The output should be sorted by earliest date first and contain 3 columns date, non-paying downloads, paying downloads.

MSSQL Solution:

```
36 -- solution
37 SELECT
38     CONVERT (DATE, d.date) AS [date],
39     SUM(CASE
40         WHEN a.paying_customer = 'NO' THEN d.downloads
41         ELSE 0
42     END) AS non_paying_downloads,
43     SUM(CASE
44         WHEN a.paying_customer = 'YES' THEN d.downloads
45         ELSE 0
46     END) AS paying_downloads
47 FROM
48     ms_user_dimension u
49     JOIN
50     ms_acc_dimension a ON u.acc_id = a.acc_id
51     JOIN
52     ms_download_facts d ON u.user_id = d.user_id
53 GROUP BY CONVERT (DATE, d.date)
54 HAVING
55     SUM(CASE WHEN a.paying_customer = 'NO' THEN d.downloads ELSE 0 END) >
56     SUM(CASE WHEN a.paying_customer = 'YES' THEN d.downloads ELSE 0 END)
57 ORDER BY [date];
```

MySQL Solution:

```

1  -- solution
2  • SELECT
3      DATE(d.date) AS date,
4      SUM(CASE
5          WHEN a.paying_customer = 'NO' THEN d.downloads
6          ELSE 0
7      END) AS non_paying_downloads,
8      SUM(CASE
9          WHEN a.paying_customer = 'YES' THEN d.downloads
10         ELSE 0
11     END) AS paying_downloads
12 FROM
13     ms_user_dimension u
14     JOIN
15     ms_acc_dimension a ON u.acc_id = a.acc_id
16     JOIN
17     ms_download_facts d ON u.user_id = d.user_id
18 GROUP BY date
19 HAVING non_paying_downloads > paying_downloads
20 ORDER BY date;

```

Problem 22 (Walmart, Paypal | Medium Level)

Find managers with at least 7 direct reporting employees. In situations where user is reporting to himself/herself, count that also.

Output first names of managers.

```

38 -- solution
39 SELECT
40     m.first_name
41 FROM
42     employees e
43     JOIN
44     employees m ON e.manager_id = m.id
45 GROUP BY m.first_name
46 HAVING COUNT(e.id) >= 7;

```


Problem 23 (Oracle | Hard Level)

Write a query that compares each employee's salary to their manager's and the average department salary (excluding the manager's salary). Display the department, employee ID, employee's salary, manager's salary, and department average salary. Order by department, then by employee salary (highest to lowest).

```
32 | -- solution
33 | WITH department_avg AS(
34 |     SELECT
35 |         department,
36 |         AVG(salary) AS avg_salary
37 |     FROM
38 |         employee_o
39 |     GROUP BY department
40 | )
41 | SELECT
42 |     e.department,
43 |     e.id,
44 |     e.salary AS employee_salary,
45 |     m.salary AS manager_salary,
46 |     d.avg_salary AS dept_avg_salary
47 | FROM
48 |     employee_o e
49 |     LEFT JOIN
50 |     employee_o m ON e.manager_id = m.id AND e.id != e.manager_id
51 |     JOIN
52 |     department_avg d ON e.department = d.department
53 | ORDER BY e.department, e.salary DESC;
```

Problem 24 (Amazon | Hard Level)

Find products which are exclusive to only Amazon and therefore not sold at Top Shop and Macy's. Your output should include the product name, brand name, price, and rating.

Two products are considered equal if they have the same product name and same maximum retail price (mrp column).

```
86 | -- solution
87 | SELECT
88 |     a.product_name,
89 |     a.brand_name,
90 |     a.price,
91 |     a.rating
92 | FROM
93 |     innerwear_amazon_com a
94 |     LEFT JOIN
95 |     innerwear_macys_com m ON a.product_name = m.product_name AND a.mrp = m.mrp
96 |     LEFT JOIN
97 |     innerwear_topshop_com t ON a.product_name = t.product_name AND a.mrp = t.mrp
98 | WHERE
99 |     m.product_name IS NULL
100 |     AND t.product_name IS NULL
101 | ORDER BY a.product_name;
```

Problem 25 (American Express | Medium Level)

American Express is reviewing their customers' transactions, and you have been tasked with locating the customer who has the third highest total transaction amount. The output should include the customer's id, as well as their first name and last name. For ranking the customers, use type of ranking with no gaps between subsequent ranks.

```
42 | -- solution
43 | WITH customer_ranking AS(
44 |     SELECT
45 |         c.id,
46 |         c.first_name,
47 |         c.last_name,
48 |         SUM(o.total_order_cost) AS total_cost,
49 |         DENSE_RANK() OVER(ORDER BY SUM(o.total_order_cost) DESC) AS customer_rank
50 |     FROM
51 |         customers c
52 |         JOIN
53 |         card_orders o ON c.id = o.cust_id
54 |     GROUP BY c.id, c.first_name, c.last_name
55 | )
56 | SELECT
57 |     id, first_name, last_name
58 | FROM
59 |     customer_ranking
60 | WHERE
61 |     customer_rank = 3;
```

Problem 26 (LinkedIn | Hard Level)

Consider all LinkedIn users who, at some point, worked at Microsoft. For how many of them was Google their next employer right after Microsoft (no employers in between)?

```
27 | -- solution
28 | WITH employer_details AS(
29 |     SELECT
30 |         user_id,
31 |         employer,
32 |         position,
33 |         start_date,
34 |         end_date,
35 |         LAG(employer) OVER(PARTITION BY user_id ORDER BY start_date) AS previous_employer,
36 |         LEAD(employer) OVER(PARTITION BY user_id ORDER BY start_date) AS next_employer
37 |     FROM
38 |         linkedin_users
39 | )
40 | SELECT
41 |     user_id,
42 |     employer,
43 |     position,
44 |     start_date,
45 |     end_date
46 | FROM
47 |     employer_details
48 | WHERE
49 |     employer = 'Microsoft' AND next_employer = 'Google';
```

Problem 27 (Goldman Sachs, Deloitte | Hard Level)

You are given a day worth of scheduled departure and arrival times of trains at one train station. One platform can only accommodate one train from the beginning of the minute it's scheduled to arrive until the end of the minute it's scheduled to depart. Find the minimum number of platforms necessary to accommodate the entire scheduled traffic.

Explanation:

1. TrainTimes CTE: We combine both the arrival and departure times into one unified dataset. For arrivals, we use 1 as the event_type to indicate the need for a platform. For departures, we use -1 as the event_type to indicate the freeing of a platform.
2. PlatformCount Subquery: We use a SUM with a window function (OVER clause) to maintain a running count of the platforms needed. The event_type for arrival adds one platform and the event_type for departure subtracts one.
3. Max(platforms_needed): Finally, we get the maximum value from the platforms_needed, which represents the maximum number of platforms required at any point in time.

```
34 -- solution
35 WITH TrainTimes AS(
36     SELECT
37         arrival_time AS event_time,
38         1 AS event_type
39     FROM
40         train_arrivals
41     UNION ALL
42     SELECT
43         departure_time AS event_time,
44         - 1 AS event_type
45     FROM
46         train_departures
47 ),
48 PlatformCount AS(
49     SELECT
50         event_time,
51         SUM(event_type) OVER(ORDER BY event_time) AS platform_needed
52     FROM
53         TrainTimes
54 )
55 SELECT
56     MAX(platform_needed) AS min_platform
57 FROM
58     PlatformCount;
```

Problem 28 (Meta, Salesforce | Hard Level)

Find the highest salary among salaries that appears only once.

```
35 | -- solution
36 | WITH salarycount AS(
37 |     SELECT
38 |         salary,
39 |         COUNT(salary) AS salary_count
40 |     FROM employee_meta
41 |     GROUP BY salary
42 | )
43 | SELECT
44 |     MAX(salary) AS highest_unique_salary
45 | FROM
46 |     salarycount
47 | WHERE salary_count = 1;
```

Problem 29 (Cisco | Hard Level)

Convert the first letter of each word found in content_text to uppercase, while keeping the rest of the letters lowercase. Your output should include the original text in one column and the modified text in another column.

Explanation:

1. STRING_SPLIT(content_text, ' ') splits content_text by spaces, breaking it down into words.
2. CROSS APPLY allows the function to be applied to each row, splitting content_text into individual words as rows.
3. UPPER(LEFT(value, 1)) converts the first letter of each word to uppercase.
4. LOWER(SUBSTRING(value, 2, LEN(value))) converts the rest of each word to lowercase.
5. STRING_AGG(..., ' ') aggregates the words back into a single string, with each word separated by a space.
6. The result is grouped by content_id and content_text, displaying both the original and modified text.

MSSQL Solution:

```
25 | -- solution
26 | SELECT
27 |     content_id,
28 |     content_text AS original_text,
29 |     STRING_AGG(UPPER(LEFT(value, 1)) + LOWER(SUBSTRING(value, 2, LEN(value))), ' ') AS modified_text
30 | FROM
31 |     user_content
32 | CROSS APPLY
33 |     string_split(content_text, ' ')
34 | GROUP BY
35 |     content_id, content_text;
```

Problem 30 (Amazon | Hard Level)

Given the users' sessions logs on a particular day, calculate how many hours each user was active that day.

Note: The session starts when state=1 and ends when state=0.

Explanation:

1. SessionDurations CTE: Use the LAG() function to get the timestamp of the previous state for each cust_id. Filter rows where state=0 because this indicates the end of a session. The session_start is derived from the LAG() value, and session_end is the current timestamp.
2. ActiveHours CTE: Calculate the active duration for each session in minutes using DATEDIFF(MINUTE, session_start, session_end).
3. Final SELECT: Sum the active minutes for each cust_id. Divide the total minutes by 60 to convert them to hours.

```
23 -- solution
24 WITH SessionDuration AS(
25     SELECT
26         cust_id,
27         CAST(LAG(timestamp) OVER (PARTITION BY cust_id ORDER BY timestamp) AS TIME) AS session_start,
28         CAST(timestamp AS TIME) AS session_end
29     FROM
30         customer_state_log
31     WHERE state = 0
32 ),
33 ActiveHours AS(
34     SELECT
35         cust_id,
36         DATEDIFF(HOUR, session_start, session_end) AS active_hours
37     FROM
38         SessionDuration
39     WHERE session_start IS NOT NULL
40 )
41 SELECT
42     cust_id,
43     SUM(active_hours) AS total_active_hours
44 FROM
45     ActiveHours
46 GROUP BY cust_id;
```

```
1 -- solution
2 WITH SessionDuration AS (
3     SELECT
4         cust_id,
5         CAST(LAG(timestamp) OVER (PARTITION BY cust_id ORDER BY timestamp) AS TIME) AS session_start,
6         CAST(timestamp AS TIME) AS session_end
7     FROM
8         customer_state_log
9     WHERE state = 0
10 ),
11 ActiveHours AS (
12     SELECT
13         cust_id,
14         TIMESTAMPDIFF(HOUR, session_start, session_end) AS active_hours
15     FROM
16         SessionDuration
17     WHERE session_start IS NOT NULL
18 )
19 SELECT
20     cust_id,
21     SUM(active_hours) AS total_active_hours
22 FROM
23     ActiveHours
24 GROUP BY cust_id;
```

Problem 31 (Spotify | Hard Level)

Find the number of days a US track has stayed in the 1st position for both the US and worldwide rankings on the same day. Output the track name and the number of days in the 1st position. Order your output alphabetically by track name. If the region 'US' appears in dataset, it should be included in the worldwide ranking

```
44 -- solution
45 SELECT
46     us.trackname,
47     COUNT(*) AS days_in_first_position
48 FROM
49     spotify_daily_rankings_2017_us us
50     JOIN
51     spotify_worldwide_daily_song_ranking w ON us.trackname = w.trackname
52     AND us.date = w.date
53 WHERE
54     us.position = 1
55     AND w.position = 1
56     AND w.region = 'US'
57 GROUP BY us.trackname
58 ORDER BY us.trackname;
```

Problem 32 (Accenture | Medium Level)

Following a recent advertising campaign, the marketing department wishes to classify its efforts based on the total number of units sold for each product.

You have been tasked with calculating the total number of units sold for each product and categorizing ad performance based on the following criteria for items sold:

- Outstanding: 30+
- Satisfactory: 20 - 29
- Unsatisfactory: 10 - 19
- Poor: 1 - 9

Your output should contain the product ID, total units sold in descending order, and its categorized ad performance.

```
31 -- solution
32 SELECT
33     product_id,
34     SUM(quantity) AS total_units_sold,
35     CASE
36         WHEN SUM(quantity) >= 30 THEN 'Outstanding'
37         WHEN SUM(quantity) BETWEEN 20 AND 29 THEN 'Satisfactory'
38         WHEN SUM(quantity) BETWEEN 10 AND 19 THEN 'Unsatisfactory'
39         WHEN SUM(quantity) BETWEEN 1 AND 9 THEN 'Poor'
40         ELSE 'No Sales'
41     END AS categorized_ad_performance
42 FROM
43     marketing_campaign
44 GROUP BY product_id
45 ORDER BY total_units_sold DESC;
```

Problem 33 (Google | Hard Level)

Calculate the average session distance traveled by Google Fit users using GPS data for two scenarios:
Considering Earth's curvature (Haversine formula).

Assuming a flat surface.

For each session, use the distance between the highest and lowest step IDs, and ignore sessions with only one step.

Calculate and output the average distance for both scenarios and the difference between them.

Formulas:

1. *Curved Earth*: $d = 6371 \times \arccos(\sin(\phi_1) \times \sin(\phi_2) + \cos(\phi_1) \times \cos(\phi_2) \times \cos(\lambda_2 - \lambda_1))$

2. *Flat Surface*: $d = 111 \times (\text{lat}_2 - \text{lat}_1)^2 + (\text{lon}_2 - \text{lon}_1)^2$

```
35 -- solution
36 WITH SessionDetails AS (
37     SELECT
38         session_id,
39         MIN(step_id) AS min_step_id,
40         MAX(step_id) AS max_step_id
41     FROM
42         google_fit_location
43     GROUP BY
44         session_id
45     HAVING
46         COUNT(step_id) > 1 -- Ignore sessions with only one step
47 ),
48 SessionCoordinates AS (
49     SELECT
50         sd.session_id,
51         gfl_min.latitude AS lat1,
52         gfl_min.longitude AS lon1,
53         gfl_max.latitude AS lat2,
54         gfl_max.longitude AS lon2
55     FROM
56         SessionDetails sd
57     JOIN google_fit_location gfl_min
58         ON sd.session_id = gfl_min.session_id AND sd.min_step_id = gfl_min.step_id
59     JOIN google_fit_location gfl_max
60         ON sd.session_id = gfl_max.session_id AND sd.max_step_id = gfl_max.step_id
61 ),
62 Distances AS (
63     SELECT
64         session_id,
65         -- Haversine formula for curved Earth distance
66         6371 * ACOS(
67             COS(RADIANS(lat1)) * COS(RADIANS(lat2)) * COS(RADIANS(lon2) - RADIANS(lon1))
68             + SIN(RADIANS(lat1)) * SIN(RADIANS(lat2))
69         ) AS curved_distance,
70         -- Flat surface formula
71         111 * SQRT(POWER(lat2 - lat1, 2) + POWER(lon2 - lon1, 2)) AS flat_distance
72     FROM
73         SessionCoordinates
74 )
75 SELECT
76     AVG(curved_distance) AS avg_curved_distance,
77     AVG(flat_distance) AS avg_flat_distance,
78     ABS(AVG(curved_distance) - AVG(flat_distance)) AS distance_difference
79 FROM
80     Distances;
```


Problem 34 (Google, Airbnb, Expedia | Medium Level)

Find the three ten hotels with the highest ratings. Output the hotel name along with the corresponding average score.

Sort records based on the average score in descending order.

```
36 | -- solution
37 | SELECT TOP 3
38 |     hotel_name,
39 |     average_score
40 | FROM
41 |     hotel_address
42 | ORDER BY average_score DESC;

36 | -- solution
37 | • SELECT
38 |     hotel_name,
39 |     average_score
40 | FROM
41 |     hotel_address
42 | ORDER BY average_score DESC
43 | LIMIT 3 OFFSET 0;
```

Problem 35 (Twitter | Medium Level)

Find the top three distinct salaries for each department. Output the department name and the top 3 distinct salaries by each department.

Order your results alphabetically by department and then by highest salary to lowest.

```
37 | -- solution
38 | WITH DeptSalaryRank AS (
39 |     SELECT
40 |         department,
41 |         salary,
42 |         DENSE_RANK() OVER(PARTITION BY department ORDER BY salary DESC) AS salaryrank
43 |     FROM
44 |         employees_x
45 |     WHERE salary IS NOT NULL
46 |     GROUP BY department, salary
47 | )
48 | SELECT
49 |     department,
50 |     salary
51 | FROM
52 |     DeptSalaryRank
53 | WHERE salaryrank <= 3
54 | ORDER BY department ASC, salary DESC;
```


Problem 36 (Uber | Hard Level)

Find the most profitable location.

Write a query that calculates the average signup duration and average transaction amount for each location, and then compare these two measures together by taking the ratio of the average transaction amount and average duration for each location.

```
38 -- solution
39 WITH SignupDuration AS (
40     SELECT
41         signup_id,
42         location,
43         DATEDIFF(MINUTE, signup_start_date, signup_stop_date) AS signup_duration_minute
44     FROM
45         signups
46 ),
47 TransactionAmount AS(
48     SELECT
49         signup_id,
50         AVG(amt) AS avg_transaction_amt
51     FROM
52         transactions
53     GROUP BY signup_id
54 )
55 SELECT
56     s.location,
57     AVG(s.signup_duration_minute) AS avg_signup_duration,
58     AVG(t.avg_transaction_amt) AS avg_transaction_amount,
59     CASE
60         WHEN AVG(s.signup_duration_minute) = 0 THEN 0
61         ELSE AVG(t.avg_transaction_amt) / AVG(s.signup_duration_minute)
62     END AS ratio
63 FROM
64     SignupDuration s
65     JOIN
66     TransactionAmount t ON s.signup_id = t.signup_id
67 GROUP BY s.location
68 ORDER BY ratio DESC;
```

Problem 37 (Netflix | Hard Level)

Find all the users who were active for 3 consecutive days or more.

MSSQL Solution:

```
19 -- solution
20 WITH activity_history AS(
21     SELECT
22         user_id,
23         LAG(date) OVER(PARTITION BY user_id ORDER BY date ASC) AS previous_date,
24         date AS present_date,
25         LEAD(date) OVER(PARTITION BY user_id ORDER BY date ASC) AS next_date
26     FROM
27         sf_events
28 )
29 SELECT
30     [user_id]
31 FROM
32     activity_history
33 WHERE
34     DATEDIFF(DAY, previous_date, present_date) = 1
35     AND DATEDIFF(DAY, present_date, next_date) = 1;
```

MySQL Solution:

```
19 -- solution
20 WITH activity_history AS(
21     SELECT
22         user_id,
23         LAG(date) OVER(PARTITION BY user_id ORDER BY date ASC) AS previous_date,
24         date AS present_date,
25         LEAD(date) OVER(PARTITION BY user_id ORDER BY date ASC) AS next_date
26     FROM
27         sf_events
28 )
29 SELECT
30     user_id
31 FROM
32     activity_history
33 WHERE
34     DATEDIFF(present_date, previous_date) = 1
35     AND DATEDIFF(next_date, present_date) = 1;
```

Problem 38 (Airbnb | Hard Level)

Estimate the growth of Airbnb each year using the number of hosts registered as the growth metric. The rate of growth is calculated by taking

$$\left(\frac{\text{(number of hosts registered in the current year} \right. \\ \left. - \text{number of hosts registered in the previous year)} \right. \\ \left. / \text{the number of hosts registered in the previous year} \right) * 100$$

Output the year, number of hosts in the current year, number of hosts in the previous year, and the rate of growth. Round the rate of growth to the nearest percent and order the result in the ascending order based on the year.

```
43 -- Solution
44 WITH YearlyHostCount AS (
45     SELECT
46         YEAR(host_since) AS year,
47         COUNT(DISTINCT id) AS hosts_in_current_year
48     FROM
49         airbnb_search_details
50     WHERE
51         host_since IS NOT NULL
52     GROUP BY
53         YEAR(host_since)
54 ),
55 GrowthCalculation AS (
56     SELECT
57         yc1.year AS current_year,
58         yc1.hosts_in_current_year,
59         yc2.hosts_in_current_year AS hosts_in_previous_year,
60         ROUND(
61             CASE
62                 WHEN yc2.hosts_in_current_year = 0 THEN NULL
63                 ELSE ((yc1.hosts_in_current_year - yc2.hosts_in_current_year) * 100.0) / yc2.hosts_in_current_year
64             END, 0
65         ) AS growth_rate
66     FROM
67         YearlyHostCount yc1
68     LEFT JOIN
69         YearlyHostCount yc2 ON yc1.year = yc2.year + 1
70 )
71 SELECT
72     current_year,
73     hosts_in_current_year,
74     hosts_in_previous_year,
75     growth_rate
76 FROM
77     GrowthCalculation
78 ORDER BY
79     current_year;
```

Problem 39 (Walmart (Hard Level))

Identify users who started a session and placed an order on the same day. For these users, calculate the total number of orders and the total order value for that day. Your output should include the user, the session date, the total number of orders, and the total order value for that day.)

```
35 -- solution
36 SELECT
37     s.user_id,
38     CAST(s.session_date AS DATE) AS session_date,
39     COUNT(o.order_id) AS total_orders,
40     SUM(o.order_value) AS total_order_value
41 FROM
42     sessions s
43     JOIN
44     order_summary o ON s.user_id = o.user_id
45 WHERE
46     CAST(s.session_date AS DATE) = CAST(o.order_date AS DATE)
47 GROUP BY s.user_id , CAST(s.session_date AS DATE)
48 ORDER BY s.user_id , session_date;
```

Problem 40 (Walmart | Hard Level)

Identify users who started a session and placed an order on the same day. For these users, calculate the total number of orders and the total order value for that day. Your output should include the user, the session date, the total number of orders, and the total order value for that day.)

```
35 -- solution
36 SELECT
37     s.user_id,
38     CAST(s.session_date AS DATE) AS session_date,
39     COUNT(o.order_id) AS total_orders,
40     SUM(o.order_value) AS total_order_value
41 FROM
42     sessions s
43     JOIN
44     order_summary o ON s.user_id = o.user_id
45 WHERE
46     CAST(s.session_date AS DATE) = CAST(o.order_date AS DATE)
47 GROUP BY s.user_id , CAST(s.session_date AS DATE)
48 ORDER BY s.user_id , session_date;
```

Problem 41 (Meta | Easy Level)

Meta/Facebook (Easy Level)

Find all posts which were reacted to with a heart. For such posts output all columns from facebook_posts table.

```
36 -- solution
37 SELECT
38     fp.post_id,
39     fp.poster,
40     fp.post_text,
41     fp.post_keywords,
42     fp.post_date
43 FROM
44     facebook_reactions fr
45     JOIN
46     facebook_posts fp ON fr.post_id = fp.post_id
47 WHERE
48     fr.reaction = 'heart';
```

Problem 42 (Microsoft | Medium Level)

Write a query that returns the company (customer id column) with highest number of users that use desktop only.

```
26 -- solution
27 WITH DesktopOnlyUser AS(
28     SELECT
29         customer_id,
30         COUNT(CASE WHEN client_id = 'desktop' THEN user_id END) AS desktop_only_user_count,
31         client_id,
32         DENSE_RANK() OVER (ORDER BY COUNT(CASE WHEN client_id = 'desktop' THEN user_id END) DESC) AS rank
33     FROM
34         microsoft_fact_events
35     GROUP BY customer_id, user_id, client_id
36 )
37 SELECT
38     customer_id,
39     desktop_only_user_count
40 FROM
41     DesktopOnlyUser
42 WHERE rank = 1;
```

Problem 43 (Apple | Hard Level)

Find the number of Apple product users and the number of total users with a device and group the counts by language. Assume Apple products are only MacBook-Pro, iPhone 5s, and iPad-air. Output the language along with the total number of Apple users and users with any device. Order your results based on the number of total users in descending order.

```
46 -- solution
47 WITH AppleUsers AS(
48     SELECT
49         pu.language,
50         COUNT(DISTINCT pe.user_id) AS apple_user_count
51     FROM
52         playbook_users pu
53         JOIN
54         playbook_events pe ON pu.user_id = pe.user_id
55     WHERE
56         pe.device IN ('MacBook-Pro', 'iPhone 5s', 'iPad-air')
57     GROUP BY pu.language
58 ),
59 TotalUsers AS(
60     SELECT
61         pu.language,
62         COUNT(DISTINCT pe.user_id) AS total_user_count
63     FROM
64         playbook_users pu
65         JOIN
66         playbook_events pe ON pu.user_id = pe.user_id
67     WHERE
68         pe.device IS NOT NULL
69     GROUP BY pu.language
70 )
71 SELECT
72     au.language,
73     COALESCE(au.apple_user_count, 0) AS apple_user_count,
74     tu.total_user_count
75 FROM
76     AppleUsers au
77     JOIN
78     TotalUsers tu ON au.language = tu.language
79 ORDER BY total_user_count DESC;
```

Problem 44 (Amazon | Hard Level)

You are given the table with titles of recipes from a cookbook and their page numbers. You are asked to represent how the recipes will be distributed in the book. Produce a table consisting of three columns: left_page_number, left_title and right_title.

The k-th row (counting from 0), should contain the number and the title of the page with the number $2 \times k$ in the first and second columns respectively, and the title of the page with the number $2 \times k + 1$ in the third column.

```
23 -- solution
24 WITH LeftPage AS(
25     SELECT
26         page_number,
27         title
28     FROM
29         cookbook_titles
30     WHERE page_number % 2 = 0
31 ),
32 RightPage AS(
33     SELECT
34         page_number,
35         title
36     FROM
37         cookbook_titles
38     WHERE page_number % 2 = 1
39 )
40 SELECT
41     lp.page_number,
42     lp.title AS left_page_title,
43     rp.title AS right_page_title
44 FROM
45     LeftPage lp
46     LEFT JOIN
47     RightPage rp ON lp.page_number + 1 = rp.page_number
48 ORDER BY lp.page_number;
```

Problem 45 (Visa | Medium Level)

Identify the top 3 areas with the highest customer density.

Customer density = (total number of unique customers in the area / area size)

Your output should include the area name and its calculated customer density.

```
38 SELECT TOP 3
39     s.area_name,
40     (COUNT(DISTINCT tr.customer_id) * 1.0 / s.area_size) AS customer_density
41 FROM
42     transaction_records tr
43     JOIN
44     stores s ON tr.store_id = s.store_id
45 GROUP BY s.area_name, s.area_size
46 ORDER BY customer_density DESC;
```

Problem 46 (EY, TCS, Deloitte | Medium Level)

In a marathon, gun time is counted from the moment of the formal start of the race while net time is counted from the moment a runner crosses a starting line. Both variables are in seconds.

You are asked to check if the interval between the two times is different for male and female runners. First, calculate the average absolute difference between the gun time and net time. Group the results by available genders (male and female). Output the absolute difference between those two values.

```
44 -- solution
45 WITH AvgTimeDiff AS (
46     SELECT
47         'Male' AS gender,
48         AVG(ABS(gun_time - net_time)) AS avg_time_diff
49     FROM
50         marathon_male
51     UNION ALL
52     SELECT
53         'Female' AS gender,
54         AVG(ABS(gun_time - net_time)) AS avg_time_diff
55     FROM
56         marathon_female
57 )
58 SELECT
59     ABS(
60         MAX(CASE WHEN gender = 'Male' THEN avg_time_diff END) -
61         MAX(CASE WHEN gender = 'Female' THEN avg_time_diff END)
62     ) AS absolute_diff
63 FROM
64     AvgTimeDiff;
```

Problem 47 (Expedia, Airbnb, Tripadvisor | Medium Level)

Find the top two hotels with the most negative reviews.

Output the hotel name along with the corresponding number of negative reviews.

Negative reviews are all the reviews with text under negative review different than "No Negative".

Sort records based on the number of negative reviews in descending order.

```
41 -- solution
42 SELECT TOP 2
43     hotel_name,
44     COUNT(negative_review) AS count_of_negative_reviews
45 FROM
46     airbnb_hotel_reviews
47 WHERE
48     negative_review != 'No Negative'
49 GROUP BY hotel_name
50 ORDER BY count_of_negative_reviews DESC;
```


Problem 48 (Amazon, Doordash, | Bosch Medium Level)

Find all employees who have or had a job title that includes manager.

Output the first name along with the corresponding title.

```
37 -- solution
38 SELECT
39     w.first_name,
40     t.worker_title
41 FROM
42     workers w
43     JOIN
44     titles t ON w.worker_id = t.worker_ref_id
45 WHERE
46     LOWER(worker_title) LIKE '%manager%';
```

Problem 49 (Goldman Sachs | Medium Level)

You work for a multinational company that wants to calculate total sales across all their countries they do business in.

You have 2 tables, one is a record of sales for all countries and currencies the company deals with, and the other holds currency exchange rate information.

Calculate the total sales, per quarter, for the first 2 quarters in 2020, and report the sales in USD currency.

```
32 -- solution
33 SELECT
34     DATEPART(QUARTER, sa.sales_date) AS sales_quarter,
35     SUM(sa.sales_amount * er.exchange_rate) AS total_sales_usd
36 FROM
37     sf_exchange_rate er
38     JOIN
39     sf_sales_amount sa ON er.source_currency = sa.currency
40     AND CAST(sa.sales_date AS DATE) = CAST(er.date AS DATE)
41 WHERE
42     sa.sales_date >= '2020-01-01' AND sa.sales_date <= '2020-07-01'
43     AND er.target_currency = 'USD'
44 GROUP BY DATEPART(QUARTER, sa.sales_date);
```

Problem 50 (Meta | Hard Level)

Market penetration is an important metric for Spotify's growth in different regions. As part of the analytics team, calculate the active user penetration rate in specific countries. Active Users must meet these criteria:

Interacted with Spotify within the last 30 days (`last_active_date >= 2024-01-01`). At least 5 sessions. At least 10 listening hours.

Formula: Active User Penetration Rate

= (Number of Active Spotify Users in the Country / Total Users in the Country)

Output: country, active_user_penetration_rate (rounded to 2 decimals).

```
29 -- solution
30 SELECT
31     country,
32     ROUND(
33         CAST(
34             SUM(
35                 CASE WHEN last_active_date >= '2024-01-01'
36                     AND sessions >= 5
37                     AND listening_hours >= 10
38                     THEN 1
39                     ELSE 0
40             ) AS FLOAT) / COUNT(*) * 100
41     , 2) AS active_user_penetration_rate
42 FROM
43     penetration_analysis
44 GROUP BY country;
```

Problem 51 (Amazon | Medium Level)

You have been asked to find the fifth highest salary without using TOP or LIMIT. Note: Duplicate salaries should not be removed.

```
24 -- solution
25 SELECT DISTINCT
26     salary
27 FROM (
28     SELECT
29         salary,
30         DENSE_RANK() OVER (ORDER BY salary DESC) AS rank
31     FROM
32         com_worker
33 ) AS SalaryRanks
34 WHERE rank = 5;
```

Problem 52 (Tesla | Hard Level)

The company you are working for wants to anticipate their staffing needs by identifying their top two busiest times of the week. To find this, each day should be segmented into different parts using following criteria:

- Morning: Before 12 p.m. (not inclusive)
- Early afternoon: 12 -15 p.m.
- Late afternoon: after 15 p.m. (not inclusive)

Your output should include the day and time of day combination for the two busiest times, i.e. the combinations with the most orders, along with the number of orders (e.g. top two results could be Friday Late afternoon with 12 orders and Sunday Morning with 10 orders).

The company has also requested that the day be displayed in text format (i.e. Monday).

```
28 | --solution
29 | WITH TimePeriodOrders AS (
30 |     SELECT
31 |         DATENAME(WEEKDAY, timestamp) AS day_of_week, -- Get the day name (e.g., Monday, Tuesday)
32 |         CASE
33 |             WHEN DATEPART(HOUR, timestamp) < 12 THEN 'Morning'
34 |             WHEN DATEPART(HOUR, timestamp) >= 12 AND DATEPART(HOUR, timestamp) < 15 THEN 'Early afternoon'
35 |             ELSE 'Late afternoon'
36 |         END AS time_of_day,
37 |         COUNT(order_id) AS order_count
38 |     FROM sales_log
39 |     GROUP BY DATENAME(WEEKDAY, timestamp),
40 |         CASE
41 |             WHEN DATEPART(HOUR, timestamp) < 12 THEN 'Morning'
42 |             WHEN DATEPART(HOUR, timestamp) >= 12 AND DATEPART(HOUR, timestamp) < 15 THEN 'Early afternoon'
43 |             ELSE 'Late afternoon'
44 |         END
45 | ),
46 | RankedOrders AS (
47 |     SELECT
48 |         day_of_week,
49 |         time_of_day,
50 |         order_count,
51 |         RANK() OVER (ORDER BY order_count DESC) AS rank
52 |     FROM TimePeriodOrders
53 | )
54 | SELECT
55 |     day_of_week,
56 |     time_of_day,
57 |     order_count
58 | FROM RankedOrders
59 | WHERE rank <= 2
60 | ORDER BY rank, order_count DESC;
```

Problem 53 (DoorDash | Medium Level)

Calculate the average net earnings per order grouped by weekday (in text format, e.g., Monday) and hour from customer_placed_order_datetime.

The net earnings are computed as: order_total + tip_amount - discount_amount - refunded_amount. Round the result to 2 decimals.

```
36 -- solution
37 SELECT
38     DATENAME(WEEKDAY, customer_placed_order_datetime) AS weekday,
39     DATEPART(HOUR, customer_placed_order_datetime) AS order_hour,
40     ROUND(AVG(order_total + tip_amount - discount_amount - refunded_amount), 2) AS average_earning
41 FROM
42     doordash_delivery
43 GROUP BY
44     DATENAME(WEEKDAY, customer_placed_order_datetime),
45     DATEPART(HOUR, customer_placed_order_datetime)
46 ORDER BY
47     weekday,
48     order_hour;
```

Problem 54 (Meta | Hard level)

The sales department has given you the sales figures for the first two months of 2023. You've been tasked with determining the percentage of weekly sales on the first and last day of every week. Consider Sunday as last day of week and Monday as first day of week.

In your output, include the week number, percentage sales for the first day of the week, and percentage sales for the last day of the week. Both proportions should be rounded to the nearest whole number.

```
27 -- solution
28 WITH WeeklySales AS (
29     SELECT
30         DATEPART(WEEK, invoicedate) AS week_number,
31         SUM(quantity * unitprice) AS total_weekly_sales
32     FROM early_sales
33     WHERE invoicedate BETWEEN '2023-01-01' AND '2023-02-28'
34     GROUP BY DATEPART(WEEK, invoicedate)
35 ),
36 SalesByDay AS (
37     SELECT
38         DATEPART(WEEK, invoicedate) AS week_number,
39         DATEPART(WEEKDAY, invoicedate) AS day_of_week,
40         SUM(quantity * unitprice) AS daily_sales
41     FROM early_sales
42     WHERE invoicedate BETWEEN '2023-01-01' AND '2023-02-28'
43     GROUP BY DATEPART(WEEK, invoicedate), DATEPART(WEEKDAY, invoicedate)
44 ),
45 FirstAndLastDaySales AS (
46     SELECT
47         s.week_number,
48         COALESCE(SUM(CASE WHEN s.day_of_week = 2 THEN s.daily_sales END), 0) AS monday_sales,
49         COALESCE(SUM(CASE WHEN s.day_of_week = 7 THEN s.daily_sales END), 0) AS sunday_sales
50     FROM SalesByDay s
51     GROUP BY s.week_number
52 )
53 SELECT
54     ws.week_number,
55     ROUND(100.0 * fl.monday_sales / ws.total_weekly_sales, 0) AS monday_sales_percentage,
56     ROUND(100.0 * fl.sunday_sales / ws.total_weekly_sales, 0) AS sunday_sales_percentage
57 FROM WeeklySales ws
58 JOIN FirstAndLastDaySales fl ON ws.week_number = fl.week_number
59 ORDER BY ws.week_number;
```

Problem 55 (Amazon | Hard Level)

Find the 3-month rolling average of total revenue from purchases given a table with users, their purchase amount, and date purchased. Do not include returns which are represented by negative purchase values. Output the year-month (YYYY-MM) and 3-month rolling average of revenue, sorted from earliest month to latest month.

A 3-month rolling average is defined by calculating the average total revenue from all user purchases for the current month and previous two months.

The first two months will not be a true 3-month rolling average since we are not given data from last year. Assume each month has at least one purchase.

```
24 -- solution
25 WITH MonthlyRevenue AS (
26     SELECT
27         FORMAT(created_at, 'yyyy-MM') AS YearMonth,
28         SUM(CASE WHEN purchase_amt > 0 THEN purchase_amt ELSE 0 END) AS TotalRevenue
29     FROM amazon_purchases
30     GROUP BY FORMAT(created_at, 'yyyy-MM')
31 ),
32 RollingAverage AS (
33     SELECT
34         mr1.YearMonth,
35         AVG(mr2.TotalRevenue * 1.0) AS RollingAvgRevenue
36     FROM MonthlyRevenue mr1
37     JOIN MonthlyRevenue mr2
38         ON DATEDIFF(MONTH, mr2.YearMonth + '-01', mr1.YearMonth + '-01') BETWEEN 0 AND 2
39     GROUP BY mr1.YearMonth
40 )
41 SELECT
42     YearMonth,
43     ROUND(RollingAvgRevenue, 2) AS RollingAvgRevenue
44 FROM RollingAverage
45 ORDER BY YearMonth;
```

Problem 56 (ESPN | Medium Level)

Find the quarterback who threw the longest throw in 2016. Output the quarterback name along with their corresponding longest throw.

The 'lg' column contains the longest completion by the quarterback.

```
35 -- solution
36 SELECT
37     qb AS Quarterback,
38     MAX(CAST(lg AS INT)) AS Longest_Throw
39 FROM
40     qbstats_2015_2016
41 WHERE
42     year = 2016
43 GROUP BY
44     qb
45 ORDER BY
46     Longest_Throw DESC;
```