

In [8]: [#SQL CheatSheet](https://www.stratascratch.com/blog/sql-cheat-sheet-technical-concepts-for-the)
[#https://www.stratascratch.com/blog/sql-cheat-sheet-technical-concepts-for-the](https://www.stratascratch.com/blog/sql-cheat-sheet-technical-concepts-for-the)

1. Most Profitable Companies

Forbes Medium ID 10354

Find the 3 most profitable companies in the entire world. Output the result along with the corresponding company name. Sort the result based on profits in descending order.

```
#方法1: with
WITH CTE AS(
SELECT
    company,
    profits,
    RANK() OVER(order by profits desc) as rnk
FROM
    forbes_global_2010_2014
)
SELECT
    company,
    profits
FROM CTE
WHERE rnk <= 3
```

```
#方法2: subquery
SELECT
    company,
    profit
FROM
    (SELECT
        *,
        rank() OVER(ORDER BY total_profit DESC) as rank
    FROM
        (SELECT
            company,
            sum(profits) AS total_profit
        FROM forbes_global_2010_2014
        GROUP BY 1) sq) sq2
WHERE rank <= 3
```

```
#方法3: with改写的subquery
WITH CTE1 AS(
SELECT
    company,
    sum(profits) AS total_profit
FROM forbes_global_2010_2014
GROUP BY 1
),
CTE2 AS(
SELECT
    company,
    total_profit,
    RANK() OVER (ORDER BY total_profit DESC) AS rank
```



```

FROM CTE1
)
SELECT
    company,
    total_profit
FROM CTE2
WHERE rank <= 3

```

#方法4: python

```

# Explore Dataset
forbes_global_2010_2014.head()
forbes_global_2010_2014.sample(5)
forbes_global_2010_2014.info()

# Import your libraries
import pandas as pd
import numpy as np

# Group& sort columns
result = forbes_global_2010_2014.groupby('company')
['profits'].sum().reset_index().sort_values(by='profits', ascending = False)

# Rank the companies
result['rank'] = result['profits'].rank(method = 'min', ascending = False)

# Filter the dataset
result = result[result['rank']<=3][['company', 'profits']]

#Optimized Solution
forbes_global_2010_2014.sort_values(by = 'profits', ascending = False)
[['company', 'profits']] head(3)

```

2.Workers With The Highest Salaries

Interview Question Date: July 2021

Amazon DoorDash Easy ID 10353

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.

DataFrames: worker, titleExpected Output Type: pandas.DataFrame

```

-- --方法1: rank
-- WITH CTE AS(
-- SELECT
--     worker_title,
--     RANK() OVER(ORDER BY salary DESC) AS rnk
-- FROM worker w JOIN title t
-- ON w.worker_id = t.worker_ref_id
-- )
-- SELECT
--     worker_title
-- FROM CTE

```

```
-- WHERE rnk = 1
```

```
--方法2: case when
--先求max_salary; case when salary = max_salary then title; select * where
title is not null
```

```
SELECT *
FROM
    (SELECT CASE
        WHEN salary =
            (SELECT max(salary)
             FROM worker) THEN worker_title
        END AS best_paid_title
    FROM worker a
    INNER JOIN title b
    ON b.worker_ref_id = a.worker_id
    ORDER BY best_paid_title
    ) sq
WHERE best_paid_title IS NOT NULL
```

```
--方法3: select title; where salary = (select max(salary) from worker)
```

```
SELECT
    worker_title AS best_paid_title
FROM worker
JOIN title
ON work_id = worker_ref_id
WHERE salary = (SELECT MAX(salary) FROM worker)
```

3. Users By Average Session Time

Interview Question Date: July 2021

Meta/Facebook Medium ID 10352

Calculate each user's average session time. A session is defined as the time difference between a page_load and page_exit. For simplicity, assume a user has only 1 session per day and if there are multiple of the same events on that day, consider only the latest page_load and earliest page_exit, with an obvious restriction that load time event should happen before exit time event. Output the user_id and their average session time.

Table: facebook_web_log

```
--date(ts); ts::timestamp: 只取年月日
--timestamp: 2019-04-25 13:30:15
--date(timestamp): 2019-04-25
--timestamp::date : 2019-04-25
```

```
--方法1: self join
先创建cte表: self join取user_id, date, session(min(t1.time - t2.time)); where
中指定t1.action = 'page_load' t2.action = 'page_exit', t2.time > t1.time;
在表中选择user_id, avg(session)
```

```
WITH all_user_sessions AS(
```

```

SELECT
    t1.user_id,
    --user_id
    t1.timestamp::date as date,
    --date整数日; t1.timestamp::date = date(t1.timestamp)
    min(t2.timestamp) - max(t1.timestamp) as session_duration
    --session = ealiest
page_exit - latest page_load: min(exit) - max(load)
FROM facebook_web_log t1 JOIN facebook_web_log t2
ON t1.user_id = t2.user_id
WHERE t1.action = 'page_load'
    --t1表为load
    AND t2.action = 'page_exit'
    --t2表为exit
    AND t2.timestamp > t1.timestamp
    --load在exit前: exit > load
GROUP BY 1, 2
)
SELECT user_id, avg(session_duration)
    --user_id, avg(session)
FROM all_user_sessions
GROUP BY user_id

```

```

--方法2:
--case when 求出page_load& page_exit timestamp; 再select 两者相减求
avg_session_time; 最后having is not null去空值

--题目条件:
--avg session time: page_load - page_exit
--latest page_load(max); ealiest page_exit(min)
--load在exit前: (exit - load)
--user_id, avg session time

WITH min_max as
(
select
    user_id,
    date(timestamp),
    max(CASE
        WHEN action = 'page_load' then timestamp
    END) as pg_load,
    --latest page_load
    min(CASE
        WHEN action = 'page_exit' then timestamp
    END) as pg_exit
    --ealiest page_exit
FROM facebook_web_log
GROUP BY 1,2
)
SELECT
    user_id,
    avg(pg_exit - pg_load) as avg_session_time
    --avg(exit - load) =
avg_session_time
FROM min_max
GROUP BY 1
HAVING avg(pg_exit - pg_load) is not null
    --去掉结尾空值

```

```

-- 方法3.
-- 创建exit表和load表 (user_id, day, exit_time/load_time) ;
-- 根据user_id和day去join两表; 取user_id, avg(exit - load);

```

-- 每个表在创立时要有day, 代表每天, 否则无法join成功

```
WITH exit as(
SELECT
    user_id,                                --user_id
    date(timestamp) as day,                 --day
    min(timestamp) as exit                  --exit
FROM
    facebook_web_log
WHERE
    action = 'page_exit'
GROUP BY
    1,2
),
load as(
SELECT
    user_id,                                --user_id
    date(timestamp) as day,                 --day

    max(timestamp) as load                  --load
FROM
    facebook_web_log
WHERE
    action = 'page_load'
GROUP BY
    1,2
)
SELECT
    e.user_id,                                --user_id
    avg(exit - load)                          --avg(exit - load)
FROM exit e JOIN load l
ON e.user_id = l.user_id
AND e.day = l.day
group by 1
```

#python:

```
# Import your libraries
import pandas as pd
import numpy as np
facebook_web_log.head()

# Extract page_load and page_exit action
loads = facebook_web_log.loc[facebook_web_log['action'] == 'page_load',
['user_id', 'timestamp']]
exits = facebook_web_log.loc[facebook_web_log['action'] == 'page_exit',
['user_id', 'timestamp']]
#Identify possible sessions of each user
sessions = pd.merge(loads, exits, how = 'inner', on = 'user_id', suffixes =
['_load', '_exit'])
#Filter valid sessions:
#page before page_exit
sessions = sessions[sessions['timestamp_load'] < sessions['timestamp_exit']]
#add a column with the date of a page_load timestamp
sessions['date_load'] = sessions['timestamp_load'].dt.date
#aggregate data by user_id and date, select latest page_load and earliest
page_exit
```

```

sessions = sessions.groupby(['user_id',
                             'date_load']).agg({'timestamp_load': 'max',
                             'timestamp_exit': 'min'}).reset_index()
#calculate the duration of the session
sessions['duration'] = sessions['timestamp_exit'] - sessions['timestamp_load']
sessions
#aggregate to get avg duration by user
result = sessions.groupby('user_id')['duration'].agg(lambda
x: np.mean(x)).reset_index()

```

4.Activity Rank

Interview Question Date: July 2021

Google Medium ID 10351

Find the email activity rank for each user. Email activity rank is defined by the total number of emails sent. The user with the highest number of emails sent will have a rank of 1, and so on. Output the user, total emails, and their activity rank. Order records by the total emails in descending order. Sort users with the same number of emails in alphabetical order. In your rankings, return a unique value (i.e., a unique rank) even if multiple users have the same number of emails. For tie breaker use alphabetical order of the user usernames.

Table: google_gmail_emails

```

--法1: cte
--total emails: from_user, count(to_user), group by 1
--先order by total email再order by 用户字母: row_number() OVER (ORDER BY
total_emails DESC, from_user ASC) as row_number
--最后再order一次

--审题: user, total_email, rank:total_email
--rank total emails desc
--users asc
--unique rank: row_number

--output: from user, total_emails, row_number

WITH CTE AS(
SELECT
    from_user,
    count(to_user) AS total_emails
FROM google_gmail_emails
GROUP BY 1
)
SELECT
    from_user,
    total_emails,
    row_number() OVER (ORDER BY total_emails DESC, from_user ASC) as
row_number
FROM CTE
ORDER BY
    total_emails DESC,

```

```
from_user
```

--法2: count(*)与row_number可以在一个query中实现, group by 1。不用cte

```
SELECT
    from_user,
    count(*) as total_emails,
    row_number() OVER (order by count(*) desc, from_user asc)
FROM
    google_gmail_emails
GROUP BY
    from_user
ORDER BY
    total_emails DESC,
    from_user
```

5.Algorithm Performance

Interview Question Date: July 2021

Meta/Facebook Hard

Meta/Facebook is developing a search algorithm that will allow users to search through their post history. You have been assigned to evaluate the performance of this algorithm.

We have a table with the user's search term, search result positions, and whether or not the user clicked on the search result.

Write a query that assigns ratings to the searches in the following way: • If the search was not clicked for any term, assign the search with rating=1 • If the search was clicked but the top position of clicked terms was outside the top 3 positions, assign the search a rating=2 • If the search was clicked and the top position of a clicked term was in the top 3 positions, assign the search a rating=3

As a search ID can contain more than one search term, select the highest rating for that search ID. Output the search ID and its highest rating.

Example: The search_id 1 was clicked (clicked = 1) and its position is outside of the top 3 positions (search_results_position = 5), therefore its rating is 2.

Table: fb_search_events

```
--法1:
--rating criteria: cte(case when)
-- 1.search not cliked, rating = 1: clicked = 0
-- 2.clicked, top position outsided the top 3, rating = 2: clicked = 1 and
search_results_position > 3
-- 3.clicked, top position in top 3, rating = 3: clicked = 1 and
search_results_position <= 3

--max rating for each user_id: max(rating)

--output: search_id, max_rating
```

```

WITH rating AS(
SELECT
    search_id,
    clicked,
    search_results_position,
CASE
    WHEN clicked = 0 THEN '1'
    WHEN clicked = 1 AND search_results_position > 3 THEN '2'
    WHEN clicked = 1 AND search_results_position <= 3 THEN '3'
    END AS rating
FROM fb_search_events
GROUP BY 1,2,3
)
SELECT
    search_id,
    max(rating) AS max_rating
FROM rating
GROUP BY 1

```

```

In [ ]: --法2:
SELECT
    search_id,
    MAX(CASE WHEN clicked = 0 THEN 1
            WHEN search_results_position > 3 THEN 2
            ELSE 3 END) as max_rating
FROM fb_search_events
GROUP BY search_id

```

```

--#法3:
-- 建立一个矩阵, search_id, 3个case when filter
-- unnest(array[colum_a, colum_b, colum_c]) 是将3列并为一列, 只显示有数字的那列数。
-- max(rating)

WITH CTE AS(
    SELECT
        search_id,
        unnest(array[one, two, three]) AS rating
    FROM
        (SELECT
            search_id,
            CASE
                WHEN clicked = 0 THEN 1
                ELSE 0
            END AS one,
            CASE
                WHEN clicked = 1 AND search_results_position > 3 THEN 2
                ELSE 0
            END AS two,
            CASE
                WHEN clicked = 1 AND search_results_position <= 3 THEN 3
                ELSE 0
            END AS three
        FROM fb_search_events
        ) sq
    )
SELECT
    search_id,

```



```

        max(rating) as max_rating
FROM cte
GROUP BY 1

```

6.Distances Traveled

Interview Question Date: December 2020

Lyft Medium ID 10324

Find the top 10 users that have traveled the greatest distance. Output their id, name and a total distance traveled.

Tables: lyft_rides_log, lyft_users

```

--output: user_id, name, traveled_distance: total_distance_traveled
--top 10; greatest distance

```

```

--方法1: cte
WITH CTE AS(
SELECT
    user_id,
    name,
    sum(distance) as traveled_distance,
    RANK() OVER (ORDER BY sum(distance) DESC) AS rnk
FROM
    lyft_rides_log l
JOIN
    lyft_users u
ON
    l.user_id = u.id
GROUP BY
    1, 2
)
SELECT
    user_id,
    name,
    traveled_distance
FROM
    CTE
WHERE
    rnk <= 10

```

```

--方法2: subquery: from后整体后tab一格, rank后order by另起tab
SELECT
    user_id,
    name,
    traveled_distance
FROM
    (SELECT
        lr.user_id,
        lu.name,
        SUM(lr.distance) AS traveled_distance,
        rank() OVER(
            ORDER BY SUM(lr.distance) DESC) AS rank
        FROM lyft_users lu

```

```

INNER JOIN lyft_rides_log lr ON lu.id = lr.user_id
GROUP BY
    lr.user_id,
    lu.name
ORDER BY
    traveled_distance DESC
) sq
WHERE rank <= 10

```

--方法3:
 --rank() OVER(ORDER BY traveled_distance DESC), rank<=10
 --可被ORDER BY traveled_distance, limit 10替代

```

SELECT
    user_id,
    name,
    sum(distance) as traveled_distance
FROM
    lyft_rides_log l
JOIN
    lyft_users u
ON
    l.user_id = u.id
GROUP BY
    1, 2
ORDER BY
    traveled_distance DESC
LIMIT 10

```

7.Finding User Purchases 【有个video可看】

Interview Question Date: December 2020

Amazon Medium ID 10322

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.

Table: amazon_transactions

#①替换写法

```

b.created_at - a.created_at BETWEEN 0 AND 7
b.created_at - a.created_at <= 7
ABS(DATEDIFF(a.created_at, b.created_at)) <= 7

```

```

a.created_at <= b.created_at
b.created_at >= a.created_at

```

#②self join 要去重自身, 需a.id != b.id

#③lead lag后, 可用where next - current >= a 来filter out:
 lead(X, 1) over(partition by Y order by X asc) as next 【注意是asc】
 lag(X, 1) over(partition by Y order by X ASC) as previous 【注意也是asc】

#4 2nd purchase within 7 days of any other purchase, 包含第二单和第一单在同一天
`b.created_at >= a.created_at`

```
--return active user: 2nd purchase within 7 days of any other pur:
--包含第二次下单和第一次下单是同一天情况, ∴ b.created_at >= a.created_at
--output:user_id

--法1: self join; on user_id=, id<>, where b.created - a.created <= 7 and
a.created_at <=b.created

SELECT
distinct
a.user_id
FROM amazon_transactions a JOIN amazon_transactions b
ON a.user_id = b.user_id
AND a.id <> b.id --排除同一用户的相同购买记录进行比较:只比较不同的购买记录
WHERE
b.created_at - a.created_at BETWEEN 0 AND 7 --可以替换为b.created_at
- a.created_at <=7 ; 也可以替换为ABS(DATEDIFF(a.created_at, b.created_at)) <= 7
--DATEDIFF(a.created_at, b.created_at) <= 7 --为什么写这个不对?
AND a.created_at <= b.created_at --【可以替换为
b.created_at >= a.created_at . 但必须勾括=】
```

```
--法2: self join; on user_id; where .created_at - b.created_at BETWEEN 0 AND 7
AND a.id != b.id
```

```
SELECT
    DISTINCT(a.user_id)
FROM amazon_transactions a
JOIN amazon_transactions b
ON a.user_id = b.user_id
WHERE a.created_at - b.created_at BETWEEN 0 AND 7
AND a.id != b.id
```

```
--法3: Lead建立新的一列, 指前面变量的后一个值, 再用where filter 相减 <=7即可
--lead(X, 1) over(partition by Y order by X asc) as next 【注意是asc】
```

```
WITH next_transaction AS(
SELECT
    user_id,
    created_at,
    LEAD(created_at, 1) OVER(PARTITION BY user_id ORDER BY created_at ASC) AS
next_transaction
FROM amazon_transactions
)
SELECT
    DISTINCT user_id
FROM next_transaction
WHERE next_transaction - created_at <= 7
```

```
--LEAD(created_at, 1)表示获取在当前行之后的下一个行的created_at值。参数1表示获取下一个行 (偏移量为1) 的created_at值
```

```
--法4: lag建立新的一列, 指前面变量的前一个值, 再用where filter 相减 <=7即可
--lag(X, 1) over(partition by Y order by X ASC) as previous 【注意也是asc】
```

```
WITH previous_transaction AS(
```

```

SELECT
    user_id,
    created_at,
    LAG(created_at, 1) OVER(PARTITION BY user_id ORDER BY created_at ASC) AS
previous_transaction
FROM amazon_transactions
)
SELECT
    DISTINCT user_id
FROM previous_transaction
WHERE created_at - previous_transaction <= 7

```

8.Monthly Percentage Difference

Interview Question Date: December 2020

Amazon Hard ID 10319

Given a table 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 can be calculated as $((\text{this month's revenue} - \text{last month's revenue}) / \text{last month's revenue}) * 100$.

Table: sf_transactions

1. Date: timestamp转换为YYYY-MM:
 ②PostgreSQL: `to_char(x, 'YYYY-MM');`; `to_char(date_trunc('month', x), 'YYYY-MM')`
`to_char(created_at::date, 'YYYY-MM')`
`SUBSTRING(created_at::text, 1, 7)`
`LEFT(created_at::text, 7)`
`to_char(DATE_TRUNC('month', created_at), 'YYYY-MM')`
 ①other SQL: `切(left, substring); date_format(date_trunc('month'), x), '%Y-%m'`
`SUBSTRING(created_at, 1, 7)`
`LEFT(created_at, 7)`
`DATE_FORMAT(created_at, '%Y-%m')`
`DATE_FORMAT(DATE_TRUNC('MONTH', created_at), '%Y-%m')`
2. :: 将 created_at 字段从其原始数据类型转换为其他类型, 以便进行后续的操作。
3. 使用 DATE_TRUNC 函数将日期字段截断到月份级别。然后, 使用 TO_CHAR 函数将截断后的日期字段格式化为 "YYYY-MM" 的形式。
4. 除法后保留两位小数: `round(a/b * 100, 2)`
5. 将lag由over后的日期单拿出来在最后写w 可简写`lag(sum(value) - 1) over (w)`

```

--month by month pct change in rev
--output: date:YYYY-MM; % pct change 2nd decimal: (this month rev - last month
rev) / last month rev * 100

--法1: cte
With cte AS(
SELECT

```

```

        to_char(created_at::date, 'YYYY-MM') as year_month,
        sum(value) as revenue
FROM sf_transactions
GROUP BY 1
ORDER BY 1
)
SELECT
    year_month,
    -- revenue,
    -- LAG(revenue, 1) OVER(ORDER BY year_month) as last_month_revenue,
    round((revenue - LAG(revenue, 1) OVER(ORDER BY year_month))/ LAG(revenue,
1) OVER(ORDER BY year_month) * 100, 2) as revenue_diff_pct
FROM CTE

#可互相替代(substring, left, to_char, to_char(date_trunc):
SUBSTRING(created_at::text, 1, 7) as year_month
LEFT(created_at::text, 7) as year_month
to_char(created_at::date, 'YYYY-MM') as year_month
to_char(DATE_TRUNC('month', created_at), 'YYYY-MM') as year_month

```

--法2.0: 大牛写法: windows alias放在最后以保证code易读:

```

SELECT
    to_char(created_at::date, 'YYYY-MM') AS year_month,           --日期
    ROUND((sum(value) - lag(sum(value), 1) OVER (w)) * 100 /      --分子一行
    lag(sum(value), 1) OVER (w),2) as revenue_diff               --分母一行
FROM sf_transactions
GROUP BY 1
Window w as (ORDER BY to_char(created_at::date, 'YYYY-MM'))      --over后较长的
部分写成windows alias放在code最后

```

法2.1: 将lag中over后的日期单拿出来在最后写w

```

SELECT
    to_char(created_at::date, 'YYYY-MM') as year_month,
    round(
        (
            sum(value) - lag(sum(value),1) over (w)
        )
        / lag(sum(value), 1) over (w) * 100
    ,2
    ) as revenue_diff
FROM sf_transactions
GROUP BY year_month
window w as (order by to_char(created_at::date, 'YYYY-MM'))

```

法2.2: 直接将order by写全, 不用写w

```

SELECT
    to_char(created_at::date, 'YYYY-MM') as year_month,
    round(
        (
            sum(value) - lag(sum(value),1) over (order by
to_char(created_at::date, 'YYYY-MM'))
        )
        / lag(sum(value), 1) over (order by to_char(created_at::date, 'YYYY-
MM')) * 100
    ,2
    ) as revenue_diff
FROM sf_transactions
GROUP BY year_month

```

```
--lag(sum(value),1) over (w) 表示对于每一行，它会获取与当前行相同窗口规范 w 中的前
一行的 sum(value) 值。
--窗口规范 w 被定义为 window w as (order by to_char(created_at::date, 'YYYY-
MM'))。
--它指定了按照 created_at 字段的日期部分进行排序，并且窗口函数将在该排序后的顺序中进行计算
```

9.New Products

Interview Question Date: December 2020

Salesforce Tesla Medium ID 10318

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.

CASE WHEN:

①COUNT和SUM在CASE WHEN中互换:

```
COUNT(CASE WHEN year = 2019 THEN product_name END) AS prev_counts
```

```
SUM(CASE WHEN year = 2019 THEN 1 ELSE 0 END) AS prev_counts
```

②CASE WHEN中else可以省略，end不可省略

```
count(case when year = 2020 then 1 else null end) 等同于 count(case when year =
2020 then product_name end)
```

③count/max(case when X then Y end) as Z: count/max/sum等放在case when括号外
case when不需要group by, count等aggregation function需要group by

④起名如果要起数字开头的名字，要加双引号"2019_counts", 包括运算中也要加，∴最好不要数字开头

⑤去重重复计算 count时看是否要加distinct: count(distinct a_brand_2020)

net diff # product 2020 vs 2019

```
company_name, net_products: 20 - 19: count(product_name) in 20
```

法1:

```
WITH CTE AS(
```

```
SELECT
```

```
    company_name,
```

```
    count(case when year = 2019 then product_name end) as prev_counts,
```

```
    count(case when year = 2020 then product_name end) as current_counts
```

```
FROM
```

```
    car_launches
```

```
GROUP BY 1
```

```
)
```

```
SELECT
```

```
    company_name,
```

```
    current_counts - prev_counts as net_products
```

```
FROM CTE
```

法2: 直接减

```
SELECT
```

```
company_name,
--count(case when year = 2020 then 1 else null end) - count(case when year =
2019 then 1 else null end) as net_diff
count(case when year = 2020 then product_name end) - count(case when year =
2019 then product_name end) as net_diff
FROM car_launches
GROUP BY 1
```

count(case when year = 2020 then 1 else null end) 等同于 count(case when year = 2020 then product_name end)

--法3: 先用where分别选出2019年和2020年product_name; 再outer join on company_name;
--最后取company_name, count(20) - count(19) as net_diff +group by

```
SELECT
    a.company_name,
    (count(distinct a.brand_2020) - count(distinct b.brand_2019)) as
net_products
FROM
    (SELECT
        company_name,
        product_name AS brand_2020
    FROM car_launches
    WHERE YEAR = 2020) a
FULL OUTER JOIN
    (SELECT
        company_name,
        product_name AS brand_2019
    FROM car_launches
    WHERE YEAR = 2019) b ON a.company_name = b.company_name
GROUP BY a.company_name
ORDER BY a.company_name
```

10. Cities With The Most Expensive Homes

Interview Question Date: December 2020

Zillow Medium ID 10315

Write a query that identifies cities with higher than average home prices when compared to the national average. Output the city names.

Table: zillow_transactions

```
output:city
city: > national avg home price
这个城市区域的平均值 大于国家区域的平均值
```

--正确答案: 让城市平均值, select city, group by city, having avg(price) > 国家平均值

```
SELECT DISTINCT city
FROM zillow_transactions
GROUP BY city
HAVING avg(mkt_price) >
    (
        SELECT avg(mkt_price)
```

```
FROM zillow_transactions
)
```

错误答案, 因为没有求city的平均值

```
SELECT DISTINCT city
FROM zillow_transactions
WHERE mkt_price >
      (SELECT avg(mkt_price)
       FROM zillow_transactions
       )
```

11.Revenue Over Time

Interview Question Date: December 2020

Amazon Hard ID 10314

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.

In []: *#本月及前两个月的rolling avg:*
#3 months rolling avg: avg total rev from all purchases for current and preivo
`AVG(t.monthly_revenue) OVER(ORDER BY t.month ROWS BETWEEN 2 PRECEDING AND CURR`

```
--法1:
--AVG(t.monthly_revenue) OVER (ORDER BY t.month ROWS BETWEEN 2 PRECEDING AND
CURRENT ROW) AS avg_revenue

SELECT
  t.month,
  AVG(t.monthly_revenue) OVER (
                                ORDER BY t.month ROWS BETWEEN 2 PRECEDING AND
CURRENT ROW) AS avg_revenue
FROM
  (SELECT
    to_char(created_at::date, 'YYYY-MM') AS month,
    sum(purchase_amt) AS monthly_revenue
  FROM amazon_purchases
  WHERE purchase_amt > 0
  GROUP BY 1
  ORDER BY 1
  ) t
ORDER BY 1
```

法2: i.a表算每个月的rev; ii.b表分别算本月, 上月, 上上月的rev; c表计算每个月有几个rev; 用b和c表join求: avg_rev = sum(revenue) / count(rev)
 最后group by month, avg_count


```

WITH rev AS(
SELECT
    to_char(created_at, 'YYYY-MM') as month,
    sum(purchase_amt) as revenue
FROM amazon_purchases
WHERE purchase_amt >= 1
GROUP BY 1
ORDER BY 1
),
final AS(
SELECT
    month,
    revenue,
    lag(revenue,1) OVER(ORDER BY month) as prev_month,
    lag(revenue,2) OVER(ORDER BY month) as prev_2_month
FROM rev
ORDER BY 1
),
month_count as(
SELECT
    month,
    (count(revenue) + count(prev_month) + count(prev_2_month)) as avg_count
FROM final
GROUP BY 1
ORDER BY 1
)
SELECT
f.month,
sum(COALESCE(f.revenue,0) + COALESCE(f.prev_month,0) +
COALESCE(f.prev_2_month,0)) / mc.avg_count AS avg_revenue
FROM final f
LEFT JOIN month_count mc
ON f.month = mc.month
GROUP BY f.month, mc.avg_count
ORDER BY 1

```

12.Naive Forecasting

Interview Question Date: December 2020

Uber Hard ID 10313

Some forecasting methods are extremely simple and surprisingly effective. Naïve forecast is one of them; we simply set all forecasts to be the value of the last observation. Our goal is to develop a naïve forecast for a new metric called "distance per dollar" defined as the (distance_to_travel/monetary_cost) in our dataset and measure its accuracy.

To develop this forecast, sum "distance to travel" and "monetary cost" values at a monthly level before calculating "distance per dollar". This value becomes your actual value for the current month. The next step is to populate the forecasted value for each month. This can be achieved simply by getting the previous month's value in a separate column. Now, we have actual and forecasted values. This is your naïve forecast. Let's evaluate our model by calculating an error matrix called root mean squared error (RMSE). RMSE is defined as $\sqrt{\text{mean}(\text{square}(\text{actual} - \text{forecast}))}$. Report out the RMSE rounded to the 2nd decimal spot.

Table: uber_request_logs

```
square: power(X,2)
mean: avg(Y)
sqrt: sqrt(Z)
2 decimal: round(M::decimal, 2): 如果M不是小数格式的话
sum(distance_to_travel) / sum(monetary_cost) <=>
avg((distance_to_travel/monetary_cost))
```

#法1.1:

```
-- 将预测值定为上一个月的实际值, 求预测值和实际值的误差MRSE:
-- "distance per dollar" defined as the (distance_to_travel/monetary_cost)
-- month: sum "distance to travel" and "monetary cost" values
-- populate the forecasted value for each month. This can be achieved simply
by getting the previous month's value in a separate column
-- calculating an error matrix called root mean squared error (RMSE).
-- sqrt(mean(square(actual - forecast)))
```

```
WITH CTE AS(
SELECT
    to_char(request_date, 'YYYY-MM') as month,
    sum(distance_to_travel) as distance_to_travel_sum,
    sum(monetary_cost) as monetary_cost_sum
FROM uber_request_logs
GROUP BY 1
),
CTE2 AS(
SELECT
    month,
    distance_to_travel_sum/ monetary_cost_sum as actual_value,
    lag(distance_to_travel_sum/ monetary_cost_sum, 1) over(order by month) as
forecasted_value
FROM CTE
ORDER BY 1
)
SELECT
    ROUND(SQRT(AVG(POWER((actual_value - forecasted_value),2))))::DECIMAL,2) AS
mrse
FROM CTE2
```

#法1.2: 先求sum(A)/sum(B), 和先写出来sum(A), sum(B),下一个cte再相除, 效果是一样的:

```
WITH monthly_actuals AS(
SELECT
    to_char(request_date, 'YYYY-MM') as month,
    sum(distance_to_travel) / sum(monetary_cost) as actual_value
FROM uber_request_logs
GROUP BY 1
),
forecast AS(
SELECT
    *,
    LAG(actual_value, 1) OVER (ORDER BY month) as forecasted_value
FROM monthly_actuals
)
SELECT
```

```

        ROUND(
            Sqrt(
                AVG(
                    POWER((actual_value - forecasted_value),2)
                )
            )::DECIMAL
            ,2) AS mrse
FROM forecast

```

```

#法3:
WITH avg_monthly_dist_per_dollar AS
(SELECT
    to_char(request_date::date, 'YYYY-MM') as request_mnth,
    sum(distance_to_travel) / sum(monetary_cost) AS monthly_dist_per_dollar
FROM uber_request_logs
GROUP BY 1
ORDER BY 1
),
naive_forecast AS
(
SELECT
    request_mnth,
    monthly_dist_per_dollar,
    lag(monthly_dist_per_dollar,1) over (
                                                order by request_mnth)
                                                as previous_monthly_dist_per_dollar
FROM avg_monthly_dist_per_dollar
),
power AS(
SELECT
    request_mnth,
    monthly_dist_per_dollar,
    previous_monthly_dist_per_dollar,
    POWER(previous_monthly_dist_per_dollar - monthly_dist_per_dollar, 2) AS
power
FROM naive_forecast
GROUP BY 1,2,3
ORDER BY 1
)
SELECT round(sqrt(avg(power))::decimal, 2) FROM power

```

13.Class Performance

Interview Question Date: December 2020

Box Medium ID 10310

You are given a table containing assignment scores of students in a class. Write a query that identifies the largest difference in total score of all assignments. Output just the difference in total score (sum of all 3 assignments) between a student with the highest score and a student with the lowest score.

Table: box_scores

```
#法1:
```

```
SELECT
max(assignment1 + assignment2 + assignment3) - min(assignment1 + assignment2 +
assignment3) as difference_in_scores
FROM box_scores
```

```
#法2: subquery
SELECT
    max(score) - min(score) AS difference_in_scores
FROM
    (SELECT
        student,
        sum(assignment1 + assignment2 + assignment3) AS score
    FROM box_scores
    GROUP BY 1
    ) +
```

14. Salaries Differences

Interview Question Date: November 2020

LinkedIn Dropbox Easy ID 10308

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.

Tables: db_employee, db_dept

```
ABS(max(CASE WHEN X THEN Y END) - min(CASE WHEN M THEN N END)) AS diff
SELECT ABS((SELECT A FROM B) - (SELECT C FROM D)) AS diff
WHERE A in ('X', 'Y'): 在A中filter出 X 和 Y
```

max(case when A then B else end): 不能错写成case when a then max(sth)
max_window_function最后的end不能丢

```
#法1: ABS(max(CASE WHEN X THEN Y END) - min(CASE WHEN M THEN N END))

SELECT
-- max(CASE WHEN department = 'marketing' THEN salary END) AS mrt_max,
-- max(CASE WHEN department = 'engineering' THEN salary END) AS eng_max,
    ABS(max(CASE WHEN department = 'marketing' THEN salary END) - max(CASE
WHEN department = 'engineering' THEN salary END)) AS salary_difference
FROM db_employee a
JOIN db_dept b
ON a.department_id = b.id
```

#法2: where 分别filter出mrt和eng的最大值, 在最外面用SELECT ABS(A - B) AS diff 去求, 最外面两者相减时不需要from。

```
SELECT ABS(
    (SELECT
        max(salary)
    FROM db_employee emp
    JOIN db_dept dpt
    ON emp.department_id = dpt.id
    WHERE department = 'marketing'
    )
```

```

-
(
    SELECT
        max(salary)
    FROM db_employee emp
    JOIN db_dept dpt
    ON emp.department_id = dpt.id
    WHERE department = 'engineering'
)
) AS salary_difference

```

#法3: 先单出一列mrt和eng的最大值, 再让两者中的max - min 就得了非负整数diff, 此处没有用abs:

--WHERE A in ('X', 'Y'): 在A中filter出 X 和 Y

```

WITH CTE AS(
SELECT
    department_id,
    max(salary) as highest_salary
FROM db_employee e
JOIN db_dept d
ON d.id = e.department_id
WHERE d.department in ('marketing', 'engineering')
GROUP BY 1
)
SELECT
    max(highest_salary) - min(highest_salary) AS salary_diff
FROM CTE

```

15.Risky Projects

Interview Question Date: November 2020

LinkedIn Medium ID 10304

Identify projects that are at risk for going overbudget. A project is considered to be overbudget if the cost of all employees assigned to the project is greater than the budget of the project.

You'll need to prorate the cost of the employees to the duration of the project. For example, if the budget for a project that takes half a year to complete is 10K, then the total half-year salary of all employees assigned to the project should not exceed \$10K. Salary is defined on a yearly basis, so be careful how to calculate salaries for the projects that last less or more than one year.

Output a list of projects that are overbudget with their project name, project budget, and prorated total employee expense (rounded to the next dollar amount).

HINT: to make it simpler, consider that all years have 365 days. You don't need to think about the leap years.

Tables: linkedin_projects, linkedin_emp_projects, linkedin_employees

1.如果结果只返回0, 不返回0后小数, 可以使分子或分母变为float, 比如365.00, 或者+0.0, 即可返回小数点后位数了

2. ROUND(X,2), CEILING(X), FLOOR(X), ceiling和floor没有,2的设置:

ROUND UP (向上取整): CEILING

CEILING(3.14) 返回 4, CEILING(5.6) 返回 6

ROUND DOWN (向下取整): FLOOR

FLOOR(3.14) 返回 3, FLOOR(5.6) 返回 5

ROUND (四舍五入)

ROUND(3.14) 返回 3, ROUND(5.6) 返回 6

3. 计算天数差值 or 除法: DATE_PART('day', X - Y), 但不能在postgre sql中用

DATE_PART('day', end_date - start_date)

EXTRACT(DAY FROM end_date - start_date) AS duration_days

ROUND(COALESCE(CAST(field1 AS DOUBLE), 0)/field2, 2) FROM TB

4. --SELECT后的部分可以在having中筛选条件使用

5. 如果不写..decimal 或 ..float可能会导致最后一位不准

overbudget: cost of all employees > budget of project

budget: 0.5 yr ~ 10k, then total 0.5 yr salary of all employees < 10k

salary yrly basis

365 days/ yr

output: projects overbudget: name, budget, prorated total employee expense(nt dollar)

title, budget, prorated_employee_expense

prorated_employee_expense = duration in yr * (sum(salary/ yr))

outbudget: budget - pro < 0

#法1: 预计花费 = 每天salary * 天数 = 起终时间相减得天数; 把该项目所有人的salary加起来, 再÷ 365 得每天salary

WITH prorated_expense AS(

SELECT

title,

budget,

(end_date - start_date) AS duration,

SUM(salary) / 365 AS salary_per_day,

(end_date::date - start_date::date) * (SUM(salary) / 365) AS

prorated_employee_expense

FROM linkedin_projects a

JOIN linkedin_emp_projects b

ON a.id = b.project_id

JOIN linkedin_employees c

ON b.emp_id = c.id

GROUP BY 1, 2, 3

ORDER BY 1, 2

)

SELECT

title,

budget,

CEILING(prorated_employee_expense) AS prorated_employee_expense

FROM prorated_expense

WHERE budget < prorated_employee_expense

ORDER BY 1

#法2: 预计花费 = SUM(年salary * 年数) = SUM(年salary * (END - START)/365)

ROUNDUP: CEILING(X)

SELECT后的部分可以在having中筛选条件使用
如果不写::decimal 或 ::float就会导致最后一位不准

```
SELECT
    title,
    budget,
    CEILING(SUM(salary * (end_date - start_date)::decimal / 365 )) AS
prorated_employee_expense
FROM
    linkedin_projects AS p JOIN linkedin_emp_projects AS ep ON p.id =
ep.project_id
                                JOIN linkedin_employees AS e ON ep.emp_id = e.id
GROUP BY
    title,
    budget
HAVING budget < CEILING(SUM(salary * (end_date - start_date)::float / 365 ))
ORDER BY 1
```

16. Top Percentile Fraud

Interview Question Date: November 2020

Google Netflix Hard ID 10303

ABC Corp is a mid-sized insurer in the US and in the recent past their fraudulent claims have increased significantly for their personal auto insurance portfolio. They have developed a ML based predictive model to identify propensity of fraudulent claims. Now, they assign highly experienced claim adjusters for top 5 percentile of claims identified by the model. Your objective is to identify the top 5 percentile of claims from each state. Your output should be policy number, state, claim cost, and fraud score.

Table: fraud_score

```
top 5 percentile:
NTILE(100) OVER (PARTITION BY STATE ORDER BY fraud_score DESC) AS percentile;
percentile <= 5
```

```
SELECT
    policy_num,
    state,
    claim_cost,
    fraud_score
FROM
    (
        SELECT
            *,
            NTILE(100) OVER (PARTITION BY STATE ORDER BY fraud_score DESC) AS
percentile
        FROM fraud_score
    ) t
WHERE percentile <= 5
```

17. Distance Per Dollar

Interview Question Date: November 2020

Uber Hard ID 10302

You're given a dataset of uber rides with the traveling distance ('distance_to_travel') and cost ('monetary_cost') for each ride. For each date, find the difference between the distance-per-dollar for that date and the average distance-per-dollar for that year-month. Distance-per-dollar is defined as the distance traveled divided by the cost of the ride.

The output should include the year-month (YYYY-MM) and the absolute average difference in distance-per-dollar (Absolute value to be rounded to the 2nd decimal). You should also count both success and failed request_status as the distance and cost values are populated for all ride requests. Also, assume that all dates are unique in the dataset. Order your results by earliest request date first.

Table: uber_request_logs

求每天ratio 与 月均ratio 的运算:

```
i. 所有, ratio, month
ii. date, month, ratio, avg(ratio) over (partition by month) as month_ratio
iii. distinct month, round(abs(ratio - month_ratio), 2)
```

#法1: 求每天ratio 与 月均ratio 的运算:

```
-- i. 所有, ratio, month
-- ii. date, month, ratio, avg(ratio) over (partition by month) as month_ratio
-- iii. distinct month, round(abs(ratio - month_ratio), 2)
```

```
SELECT
    DISTINCT
    b.request_mnth,
    ROUND(ABS(b.dist_to_cost - b.avg_dist_to_cost)::DECIMAL, 2) AS
mean_deviation
FROM (
    SELECT
        a.request_date,
        a.request_mnth,
        a.dist_to_cost,
        AVG(a.dist_to_cost) OVER (PARTITION BY request_mnth) AS
avg_dist_to_cost
    FROM (
        SELECT
            *,
            (distance_to_travel / monetary_cost) AS dist_to_cost,
            to_char(request_date::date, 'YYYY-MM') AS request_mnth
        FROM uber_request_logs) a
    )b
ORDER BY b.request_mnth
```

#法2:

```
WITH avg_distance_per_dollar AS(
SELECT
    to_char(request_date, 'YYYY-MM') as request_mnth,
    AVG((distance_to_travel) / (monetary_cost)) AS avg_distance_per_dollar
FROM uber_request_logs
GROUP BY 1
),
```



```

distance_per_dollar AS(
SELECT
    request_date,
    to_char(request_date, 'YYYY-MM') as request_mnth,
    distance_to_travel / monetary_cost AS distance_per_dollar
FROM uber_request_logs
),
diff AS(
SELECT
    request_date,
    avg_distance_per_dollar,
    distance_per_dollar,
    ABS(avg_distance_per_dollar - distance_per_dollar) AS difference
FROM distance_per_dollar a
LEFT JOIN avg_distance_per_dollar b
ON a.request_mnth = b.request_mnth
),
final AS(
SELECT
    DISTINCT
    to_char(request_date, 'YYYY-MM') AS year_month,
    ROUND(difference::DECIMAL,2)
FROM diff
)
SELECT * FROM final
ORDER BY 1

```

#法3:

```

SELECT
    request_mnth,
    ROUND(AVG(mean_deviation), 2) AS difference
FROM
    (
    SELECT
        request_mnth,
        ABS(dist_to_cost - monthly_dist_to_cost)::DECIMAL AS mean_deviation
    FROM
        (
        SELECT
            to_char(request_date::date, 'YYYY-MM') AS request_mnth,
            distance_to_travel / monetary_cost AS dist_to_cost,
            SUM(distance_to_travel) OVER (PARTITION BY
to_char(request_date::date, 'YYYY-MM'))
            / SUM(monetary_cost) OVER (PARTITION BY
to_char(request_date::date, 'YYYY-MM')) AS monthly_dist_to_cost
        FROM uber_request_logs
        ) t
        ) t2
GROUP BY 1
ORDER BY 1

```

18. Expensive Projects

Interview Question Date: November 2020

Microsoft Medium ID 10301

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.

Tables: ms_projects, ms_emp_projects

1. 互换:

```
CEILING(budget / COUNT(emp_id)::DECIMAL)
ROUND(budget / COUNT(emp_id)::FLOAT)::numeric, 0)
```

2. 如果最后一位差一点, 试试::float/ ::decimal/ ::numeric

3. ORDER BY 后可直接由 alias name

```
--budget to employee
--output: project, budget_emp_ratio: round to closest int: ceiling
--order by ratio desc

SELECT
    title AS project,
    CEILING(budget / COUNT(emp_id)::DECIMAL) AS budget_emp_ratio
    -- ROUND((budget / COUNT(emp_id)::FLOAT)::numeric, 0) AS budget_emp_ratio
FROM ms_projects a
JOIN ms_emp_projects b
ON a.id = b.project_id
GROUP BY title, budget
ORDER BY budget_emp_ratio DESC
```

19. Premium vs Freemium

Interview Question Date: November 2020

Microsoft Hard ID 10300

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.

Tables: ms_user_dimension, ms_acc_dimension, ms_download_facts

1. CASE WHEN 后的END不要忘记

2. string: yes, no 需要加''

3. having 放在group by 后, 不可以用 alias name

4. case后换行, when接在case后地方, end换行和case对齐

```
--total number of downloads for pay and non-pay by date
--include only download(non-pay) > download(pay)

--output: date, non_paying's downloads, paying's downloads:sum(downloads)
--sort by date asc

SELECT
```

```

date,
SUM(CASE
    WHEN paying_customer = 'no' THEN downloads
END) AS non_paying,
SUM(CASE
    WHEN paying_customer = 'yes' THEN downloads
END) AS paying
FROM ms_user_dimension mud
JOIN ms_acc_dimension mac
ON mud.acc_id = mac.acc_id
JOIN ms_download_facts mdf
ON mud.user_id = mdf.user_id
GROUP BY date
HAVING SUM(CASE WHEN paying_customer = 'no' THEN downloads END) > SUM(CASE
    WHEN paying_customer = 'yes' THEN downloads END)
ORDER BY date

```

20. Finding Updated Records

Interview Question Date: November 2020

Microsoft Easy

We have a table with employees and their salaries, however, some of the records are old and contain outdated salary information. Find the current salary of each employee assuming that salaries increase each year. Output their id, first name, last name, department ID, and current salary. Order your list by employee ID in ascending order.

Table: ms_employee_salary

```

SELECT
    id,
    first_name,
    last_name,
    department_id,
    max(salary) AS max
FROM ms_employee_salary
GROUP BY 1,2,3,4
ORDER BY 1

```

21. Comments Distribution

Interview Question Date: November 2020

Meta/Facebook Hard ID 10297

Write a query to calculate the distribution of comments by the count of users that joined Meta/Facebook between 2018 and 2020, for the month of January 2020.

The output should contain a count of comments and the corresponding number of users that made that number of comments in Jan-2020. For example, you'll be counting how many users made 1 comment, 2 comments, 3 comments, 4 comments, etc in Jan-2020. Your left column in

the output will be the number of comments while your right column in the output will be the number of users. Sort the output from the least number of comments to highest.

To add some complexity, there might be a bug where an user post is dated before the user join date. You'll want to remove these posts from the result.

- ①算 count(comments)~ count(users)
- i. user_id, count(comments) as cmt_count, group by 1
- ii. cmt_count, count(user_id) as user_count, group by 1
- ②COUNT(*) 与 COUNT(X_id) 在一些时候效果相同可替换
- ③BETWEEN 'A' AND 'B'::DATE
- ④考虑是否有等于号=: joined <= created : 考虑join当天发帖的人
- ⑤先join两表 列出最细level的变量 再继续整理计算

```
--comments by the # users joined between 18 and Jan 20; time:for the month of
January 2020.
--output: comment_cnt: #com ; user_cnt#users
--a count of comments and # of users join 18-20 that made that # of comments
in Jan 2020
    -- join date of 18 and 20; comment date of jan 20
--eg. #users~1comment, 2 comments, 3 comments, 4, etc.--
--sort #com asc
--bug: user post date < user join date: remove

--法1: cte
WITH comment_cnt AS(
SELECT
    b.user_id,
    COUNT(b.created_at) AS comment_cnt
COUNT(*)也可以
FROM fb_users a
JOIN fb_comments b
ON a.id = b.user_id
WHERE (b.created_at BETWEEN '2020-01-01' AND '2020-01-31'::DATE)
BETWEEN 'A' AND 'B'::DATE
AND (a.joined_at BETWEEN '2018-01-01' AND '2020-12-31'::DATE)
AND joined_at <= created_at
--joined
<= created : 考虑join当天发帖的人
GROUP BY 1
ORDER BY 2
),
user_cnt AS(
SELECT
    comment_cnt,
    COUNT(user_id) AS user_cnt
FROM comment_cnt
GROUP BY 1
ORDER BY 1
)
SELECT *
```

```
FROM user_cnt
```

法2: subquery

subquery写法: 写新FROM, 在from后(, 在上次code句末), 将上次code整体tab

subquery写法: 先join两表, 把最细level的变量都挑出来, 再count comment, 再count user

```
SELECT
    comment_cnt,
    COUNT(id) AS user_cnt
FROM (
    SELECT
        t.id,
        COUNT(t.user_id) as comment_cnt
    FROM (
        SELECT
            a.id,
            a.joined_at,
            b.user_id,
            b.created_at
        FROM fb_users a
        JOIN fb_comments b
        ON a.id = b.user_id
        WHERE (b.created_at BETWEEN '2020-01-01' AND '2020-01-31')
        AND (a.joined_at BETWEEN '2018-01-01' AND '2020-12-31')
        AND joined_at <= created_at) t
    GROUP BY t.id ) c
GROUP BY comment_cnt
ORDER BY comment_cnt ASC
```

22. Most Active Users On Messenger

Interview Question Date: November 2020

Meta/Facebook Medium ID 10295

Meta/Facebook Messenger stores the number of messages between users in a table named 'fb_messages'. In this table 'user1' is the sender, 'user2' is the receiver, and 'msg_count' is the number of messages exchanged between them. Find the top 10 most active users on Meta/Facebook Messenger by counting their total number of messages sent and received. Your solution should output usernames and the count of the total messages they sent or received

Table: fb_messages

0. union and union all:

此题应该用UNION ALL,

因为如果union user1 and user2 后出现相同记录, 但代表的含义不同的需要全部保留

如union后出现两条记录: UserA, 5, 分别表示sendA发了5条信息和receiverA收了5条信息

union all是全部保留, 用union则distinct记录

1理解: msg_count: 既是user1发的数字, 也是user2收的数字。

2如果想要进行的数值列在名字列之前, 依然可以照常计算: user2, sum(msg_count);

user_name, sum(msg_count)

3union法:

i. 求user1发的数: user1, msg_count

ii. 求user2收的数: user2, msg_count
 iii. 把user1 发的数UNION user2收的数, 是每个user收发的数明细
 iv. user_name, sum(count) 是每个user收发总数
 注: 可以先sum 再union; 也可以先union再sum

4self join法:

i. 先求发过&收过消息的所有人名单表: distinct u1 union distinct u2 as user_name
 ii. 再将原表 a 和新表 b join, on a.user1 = b.user_name OR a.user2 = b.user_name, 得所有人名单收发明细: 加新列命名user_name, 出: user1, user2, user1和user1, user2, user2

	id	date	user1	user2	msg_count	user_name
	1	2020-08-02	kpena	scottmartin	2	kpena
	1	2020-08-02	kpena	scottmartin	2	scottmartin

iii. 取user_name, sum(msg_count), group by 1, 得所有人, 收发总数

5limit和rank互换法:

rank需要多写一个query, 但对于数值相同的记录处理稳妥, 不会少。partition可以省略。
 limit: 适用于没有数值相同的记录

#messages IN users table: fb_messages
 TOP 10 most active user by counting total #messages sent & received
 output: usernames, total_msg_count: count of total mess sent & received

法1:

```
SELECT user_name, SUM(msg_count) AS total_messages
FROM
(
  SELECT user1 AS user_name, msg_count
  FROM fb_messages
  UNION ALL
  SELECT user2 AS user_name, msg_count
  FROM fb_messages
) AS subquery
GROUP BY user_name
ORDER BY total_messages DESC
LIMIT 10
```

法2.1: 先SUM再UNION, limit

```
SELECT *
FROM(
  SELECT
    user1 AS user_name,
    SUM(msg_count) AS s_r_sum
  FROM fb_messages
  GROUP BY user1
  UNION
  SELECT
    user2 AS user_name,
    SUM(msg_count) AS s_r_sum
  FROM fb_messages
  GROUP BY user2
) t
ORDER BY 2 DESC
LIMIT 10
```

法2.2: 先SUM再UNION, rank

```
SELECT
  user_name,
```

```

total_msg_count
FROM(
SELECT
    user_name,
    SUM(s_r_sum) AS total_msg_count,
    RANK() OVER (ORDER BY SUM(s_r_sum) DESC) AS rnk
FROM(
    SELECT
        user1 AS user_name,
        SUM(msg_count) AS s_r_sum
    FROM fb_messages
    GROUP BY user1
    UNION
    SELECT
        user2 AS user_name,
        SUM(msg_count) AS s_r_sum
    FROM fb_messages
    GROUP BY user2
    ) t
GROUP BY user_name,s_r_sum) t2
WHERE rnk <= 10
ORDER BY total_msg_count DESC

```

法3: self join

i. 先求发过&收过消息的所有人名单表: distinct u1 union distinct u2 as user_name
 ii. 再将原表 a 和新表 b join, on a.user1 = b.user_name OR a.user2 = b.user_name,
 得所有人名单收发明细: 加新列命名user_name, 出: user1, user2, user1和user1, user2,
 user2

id	date	user1	user2	msg_count	user_name
1	2020-08-02	kpena	scottmartin	2	kpena
1	2020-08-02	kpena	scottmartin	2	scottmartin

iii. 取user_name, sum(msg_count), group by 1, 得所有人, 收发总数

```

SELECT
    u.user_name, SUM(m.msg_count) AS total_messages
FROM fb_messages m
JOIN (
    SELECT DISTINCT user1 AS user_name FROM fb_messages
    UNION
    SELECT DISTINCT user2 AS user_name FROM fb_messages
    ) u ON m.user1 = u.user_name OR m.user2 = u.user_name
GROUP BY u.user_name
ORDER BY total_messages DESC
LIMIT 10

```

```

WITH user_name AS(
    SELECT DISTINCT user1 AS user_name FROM fb_messages
    UNION
    SELECT DISTINCT user2 AS user_name FROM fb_messages
),
final AS(
    SELECT *
    FROM fb_messages a
    JOIN user_name b
    ON a.user1 = b.user_name OR a.user2 = user_name
    ORDER BY 1
)
SELECT

```

```

        user_name AS username,
        SUM(msg_count) AS total_msg_count
FROM final
GROUP BY 1
ORDER BY 2 DESC
LIMIT 10

```

23.SMS Confirmations From Users

Interview Question Date: November 2020

Meta/Facebook Medium

Meta/Facebook sends SMS texts when users attempt to 2FA (2-factor authenticate) into the platform to log in. In order to successfully 2FA they must confirm they received the SMS text message. Confirmation texts are only valid on the date they were sent.

Unfortunately, there was an ETL problem with the database where friend requests and invalid confirmation records were inserted into the logs, which are stored in the 'fb_sms_sends' table. These message types should not be in the table.

Fortunately, the 'fb_confirmers' table contains valid confirmation records so you can use this table to identify SMS text messages that were confirmed by the user.

Calculate the percentage of confirmed SMS texts for August 4, 2020. Be aware that there are multiple message types, the ones you're interested in are messages with type equal to 'message'.

Tables: fb_sms_sends, fb_confirmers

求部分占整体的比率:

i. 求部分的query:可以多放一些细节level

ii. 求整体的query

iii. 整体 left join 部分 ON 整体1 = 部分1 AND 整体2 = 部分2: 出现以整体为准, 部分与整体不match的表格, 可能出现match部分为null的情况

iv. $\text{count}(\text{部分.id}) / \text{count}(\text{整体.id}) :: \text{float}$: 得出比率 (此处可以填写两个表的任意变量, 结果不变)

```

WITH send AS(
    SELECT
        ds,
        phone_number,
        type
    FROM fb_sms_sends
    WHERE type = 'message'
),
confirmed AS(
    SELECT
        date,
        phone_number
    FROM fb_confirmers
)
SELECT
    COUNT(b.phone_number) / COUNT(a.phone_number)::float * 100 AS perc
FROM send a

```



```
LEFT JOIN confirmed b
ON a.ds = b.date AND a.phone_number = b.phone_number
WHERE a.ds = '08-04-2020'
```

24. Acceptance Rate By Date

Interview Question Date: November 2020

Meta/Facebook Medium ID 10285

What is the overall friend acceptance rate by date? Your output should have the rate of acceptances by the date the request was sent. Order by the earliest date to latest.

Assume that each friend request starts by a user sending (i.e., user_id_sender) a friend request to another user (i.e., user_id_receiver) that's logged in the table with action = 'sent'. If the request is accepted, the table logs action = 'accepted'. If the request is not accepted, no record of action = 'accepted' is logged.

```
--friend acc rate by date
--output: date: request sent date, percentage_acceptance: acc rate
--order by date

--(sent, accepted) / sent

SELECT
a.date,
COUNT(b.user_id_receiver) / COUNT(a.user_id_sender)::float AS acceptance_rate
--b.后写其他b表变量也可; a.后写其他a表变量也可
FROM (
SELECT
    user_id_sender,
    user_id_receiver,
    date,
    action
FROM fb_friend_requests
WHERE action = 'sent') a
LEFT JOIN (
    SELECT
        user_id_sender,
        user_id_receiver,
        date,
        action
    FROM fb_friend_requests
    WHERE action = 'accepted') b
ON b.user_id_sender = a.user_id_sender
-- ① 所有变量:
send left join acc on sender_id = sender id and receiver_id = receiver_id
AND b.user_id_receiver = a.user_id_receiver
GROUP BY 1
-- ② group by
date, count(acc.receiver_id) / count(send.sender_id) = acceptance rate
ORDER BY 1 DESC
```

25. Popularity Percentage

Interview Question Date: November 2020

Meta/Facebook Hard ID 10284

Find the popularity percentage for each user on Meta/Facebook. The popularity percentage is defined as the total number of friends the user has divided by the total number of users on the platform, then converted into a percentage by multiplying by 100. Output each user along with their popularity percentage. Order records in ascending order by user id. The 'user1' and 'user2' column are pairs of friends.

Table: facebook_friends

①用UNION让第一列包含所有用户

第一列是所有用户名，第二列是他的朋友的总表：

```
SELECT u1, u2 FROM facebook_friends
```

```
UNION
```

```
SELECT u2, u1 FROM facebook_friends
```

②在第一列维度下，第二列在一行内实现 部分与整体 的比率

第一列是用户名，第二列是流行率：他的朋友数 / 平台用户总数

```
SELECT
```

```
user1,
```

```
COUNT(user2) / (SELECT COUNT(DISTINCT user1) FROM total)::float * 100
```

```
FROM total
```

user1, count(*) or user1, count(user_2)是在算每个user的朋友数；

distinct user1 from total 是平台所有用户数

```
WITH total AS(
```

```
SELECT
```

```
    DISTINCT
```

```
    user1,
```

```
    user2
```

```
    FROM facebook_friends
```

```
    UNION
```

```
    SELECT
```

```
    user2,
```

```
    user1
```

```
    FROM facebook_friends
```

```
    ORDER BY 1
```

```
),
```

```
friends AS(
```

```
    SELECT
```

```
    user1,
```

```
    COUNT(user2) / (SELECT COUNT(DISTINCT user1) FROM total)::float * 100 AS
```

```
    popularity_percent
```

```
    FROM total
```

```
    GROUP BY 1
```

```
    ORDER BY 1
```

```
)
```

```
SELECT * FROM friends
```

```
WITH total AS(
```

```
    SELECT
```

```
        DISTINCT
```

```
        user1 AS user
```

```
        FROM facebook_friends
```

```
        UNION
```

```
        SELECT
```

```

        user2 AS user
        FROM facebook_friends
        ORDER BY 1
    ),
    total_count AS(
        SELECT *
        FROM total),
    friends AS(
        SELECT
            user1,
            COUNT(user2) AS friend
        FROM facebook_friends
        GROUP BY 1
        UNION
        SELECT
            user2,
            COUNT(user1) AS friend
        FROM facebook_friends
        GROUP BY 1
    ),
    friends_count AS(
        SELECT
            user1,
            SUM(friend) AS friends_count
        FROM friends
        GROUP BY 1
        ORDER BY 1
    ),
    final AS(
        SELECT
            user1,
            friends_count / (SELECT COUNT(*) FROM total_count) * 100 AS
            popularity_percent
        FROM friends_count
    )
SELECT * FROM final

```

-- 求出总数

-- 第一列的为user, count后面的

-- UNION

-- 第二列的为user, count前面的

-- 求出每个user的friend数

-- 每个user的friend数/total数: total可以用select *

26. Find the top-ranked songs for the past 20 years.

Spotify Medium ID 10283

Find all the songs that were top-ranked (at first position) at least once in the past 20 years

Table: billboard_top_100_year_end

过去N年:date_part要加'', extract不加''

```

DATE_PART('year', CURRENT_DATE) - year <= N
EXTRACT(year from CURRENT_DATE) - year <= 20

```

```

SELECT
    DISTINCT

```

```

    song_name
FROM billboard_top_100_year_end
WHERE DATE_PART('year', CURRENT_DATE) - year <= 20
--WHERE extract(year from current_date) - year <= 20
--WHERE year BETWEEN 2003 AND 2023
AND year_rank = 1

```

27. Find all inspections which are part of an inactive program

City of Los Angeles Easy ID 10277

Find all inspections which are part of an inactive program.

```

WHERE filter 条件:
ilike '%inactive'
= 'INACTIVE'

```

```

SELECT *
FROM los_angeles_restaurant_health_inspections
WHERE program_status = 'INACTIVE'

```

28. Find the total number of available beds per hosts' nationality

Airbnb Medium ID 10187

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.

Tables: `airbnb_apartments`, `airbnb_hosts`

当两张表中都有看起来需要的变量时，多留意下从哪张表取：
第一张表中的country是room所在country；要host nationality要在第二张表中找

```

/*
auto_comment_out_words
*/

```

```

SELECT
    b.nationality,
    SUM(n_beds) AS total_beds_available
    -- country AS nationality, -- 这是room的country, ∴ ×
FROM
    airbnb_apartments a
INNER JOIN
    airbnb_hosts b
ON
    a.host_id = b.host_id
GROUP BY
    nationality
ORDER BY

```

```
total_beds_available DESC
```

29. Order all countries by the year they first participated in the Olympics

ESPN Easy ID 10184

Order all countries by the year they first participated in the Olympics.

Output the National Olympics Committee (NOC) name along with the desired year.

Sort records by the year and the NOC in ascending order.

Table: olympics_athletes_events

求AB为1组，每组的第一个值：

RANK 法: `RANK() OVER (PARTITION BY noc ORDER BY year ASC) AS rnk`

MIN法: `A, MIN(B) + group by 1`

```
/*
year first in Olympic, countries
output: noc, first_time_year
sort year, NOC in asc
*/

-- RANK 法: RANK() OVER (PARTITION BY noc ORDER BY year ASC) AS rnk
WITH rnk AS(
SELECT
    noc,
    year,
    RANK() OVER (PARTITION BY noc ORDER BY year ASC) AS rnk
FROM olympics_athletes_events
)
SELECT
    DISTINCT
    noc,
    year AS first_time_year
FROM rnk
WHERE rnk = 1
```

--MIN法: `A, MIN(B) + group by 1` 求AB为1组，每组的第一个值

```
SELECT
    noc,
    MIN(year) AS first_time_year
FROM olympics_athletes_events
GROUP BY 1
ORDER BY 1, 2
```

30. Total Cost Of Orders

Interview Question Date: July 2020

Amazon Etsy Easy ID 10183

Find the total cost of each customer's orders. Output customer's id, first name, and the total order cost. Order records by customer's first name alphabetically.

Tables: customers, orders

```
--total cost of each customer's orders
--output: id: customer's id; first_name, sum: total order cost
--order by first name asc

SELECT
    c.id,
    c.first_name,
    SUM(o.total_order_cost) AS sum
FROM customers c
JOIN orders o
ON c.id = o.cust_id
GROUP BY c.id,
         c.first_name
ORDER BY c.first_name ASC
```

31. Find the lowest score for each facility in Hollywood Boulevard

Interview Question Date: July 2020

City of Los Angeles City of San Francisco Tripadvisor Medium ID 10180

Find the lowest score per each facility in Hollywood Boulevard. Output the result along with the corresponding facility name. Order the result based on the lowest score in descending order and the facility name in the ascending order.

Table: los_angeles_restaurant_health_inspections

```
前后模糊搜索: ILIKE '%HOLLYWOOD BLVD%'
-- '%hollywood%h%l%v%d%'
```

```
--lowest score per facility in HB
--output: facility_name, min_score
--order lowest score desc, facility name in asc

SELECT
    facility_name,
    MIN(score) AS min_score
FROM los_angeles_restaurant_health_inspections
WHERE facility_address ILIKE '%HOLLYWOOD BLVD%'
GROUP BY facility_name
ORDER BY min_score DESC, facility_name ASC
```

32. Businesses Open On Sunday

Yelp Medium ID 10178

Find the number of businesses that are open on Sundays. Output the slot of operating hours along with the corresponding number of businesses open during those time slots. Order records by total number of businesses opened during those hours in descending order.

```
SELECT
    DISTINCT
    a.Sunday,
    COUNT(is_open) AS total_business      --COUNT(*)也可以
FROM yelp_business_hours a
JOIN yelp_business b
ON a.business_id = b.business_id
WHERE is_open = 1 AND sunday IS NOT NULL
GROUP BY 1
ORDER BY 2 DESC
```

33.Bikes Last Used

Lyft DoorDash Easy ID 10176

Find the last time each bike was in use. Output both the bike number and the date-timestamp of the bike's last use (i.e., the date-time the bike was returned). Order the results by bikes that were most recently used.

Table: dc_bikeshare_q1_2012

```
--last time each bike in use
--output: bike_number: bike number; last_used: date last use
--order by most recent used

--1.RANK法:
WITH rnk AS(
SELECT
    bike_number,
    end_time AS last_used,
    RANK() OVER (PARTITION BY bike_number ORDER BY end_time DESC) AS rnk
FROM dc_bikeshare_q1_2012
GROUP BY bike_number, end_time      --rank()前group by所有变量
)
SELECT
    bike_number,
    last_used
FROM rnk
WHERE rnk = 1
```

```
2.MAX法:
SELECT
    bike_number,
    MAX(end_time) AS last_used
FROM dc_bikeshare_q1_2012
GROUP BY bike_number
ORDER BY 2 DESC
```

34.Days At Number One

Spotify Hard ID 10173

Find the number of days a US track has stayed in the 1st position for both the US and worldwide rankings. 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.

Tables: spotify_daily_rankings_2017_us, spotify_worldwide_daily_song_ranking

法1: 整体left join部分

世界; US; 世界left join US: date = date and name = name; filter 世界榜1US榜1交集: a.date is not null and b.date is not null; name, count(*)

法2.1: inner join; count(*)

us inner join world; date = date name = name: 既在US也在世界榜上出现的记录 where filter us.position = 1 AND world.position = 1, 就是既在US也在世界榜1了, 直接COUNT(*)即可

法2.2: 计算既在总表也在部分表出现的次数: 总表inner join部分表; 两个表中选变量完成 sum() over (partition by) + case when name max()

```
-- days a US track 1st for both US and World
-- region 'US' appear -> include in world ranking
-- output: trackname, n_days_on_n1_position
-- order name asc
```

--法1: 整体left join部分

--世界; US; 世界left join US: date = date and name = name; filter 世界榜1US榜1交集: a.date is not null and b.date is not null; name, count(*)

```
WITH world AS(
    SELECT
        date,
        position,
        trackname
    FROM spotify_worldwide_daily_song_ranking
    WHERE position = 1
),
us AS(
    SELECT
        date,
        position,
        trackname
    FROM spotify_daily_rankings_2017_us
    WHERE position = 1
)
SELECT
    a.trackname,
    COUNT(*) AS n_days_on_n1_position
FROM world a LEFT JOIN us b
ON a.date = b.date AND a.trackname = b.trackname
WHERE a.date IS NOT NULL and b.date IS NOT NULL
GROUP BY 1
```

法2.1: inner join; count(*)

us inner join world; date = date name = name: 既在US也在世界榜上出现的记录

where filter us.position = 1 AND world.position = 1, 就是既在US也在世界榜1了, 直接COUNT(*)即可

```
SELECT
    us.trackname,
    COUNT(*) AS n_days_on_n1_position
FROM spotify_daily_rankings_2017_us us
INNER JOIN spotify_worldwide_daily_song_ranking world
ON world.trackname = us.trackname AND world.date = us.date
WHERE us.position = 1 AND world.position = 1
GROUP BY 1
```

法2.2: 计算既在总表也在部分表出现的次数: 总表inner join部分表; 两个表中选变量完成 sum() over (partition by) + case when; name, max()

① US榜1: US与world join: date = date and name = name; us.position = 1

② US世界榜1天数: us.name, SUM(case when world.position = 1 then 1 else 0 end) OVER (PARTITION BY us.name) AS days 【COUNT也可以】

此处用us.name, count(*) 不行, 因为也许存在us第一但不是world第一的情况。但本题没有出现特例。

③ distinct选出歌名和天数: name, max(days)

us INNER JOIN world ; date = date AND name = name: 同一天既在us榜单也在world榜单出现, US榜1的歌

SUM(CASE WHEN world.position = 1 THEN 1 ELSE 0 END) OVER (PARTITION BY us.trackname) AS n_days:

对于us榜1的歌来说, 如果也在世界榜1就计算为1次, sum/count有多少次, 就是both在us和世界榜1的天数。

```
SELECT
    trackname,
    MAX(n_days_on_n1_position) AS n_days_on_n1_position
FROM
    (
        SELECT
            us.trackname,
            SUM(CASE
                WHEN world.position = 1 THEN 1
                ELSE 0
            END) OVER (PARTITION BY us.trackname) AS n_days_on_n1_position
        FROM spotify_daily_rankings_2017_us us
        INNER JOIN spotify_worldwide_daily_song_ranking world
        ON world.trackname = us.trackname AND world.date = us.date
        WHERE us.position = 1
    ) tmp
GROUP BY trackname
ORDER BY trackname
```

35.Best Selling Item

Interview Question Date: July 2020

Amazon Ebay Best Buy Hard ID 10172

Find the best selling item for each month (no need to separate months by year) where the biggest total invoice was paid. The best selling item is calculated using the formula (unitprice * quantity). Output the description of the item along with the amount paid.

Table: online_retail

1. SUM(unitprice * quantity)与unitprice * quantity的区别:
sum能一个月内购买相同产品的的不同记录加总, 不加sum只能算出一条记录
sum时只需要group by sum前的变量, 不加sum时, 需要group by sum前的变量和sum内的变量

2. 将SUM(P*N)与 RANK放在一个Query里更有效率

法1:

```
WITH total AS(
SELECT
EXTRACT(month from invoicedate) AS month,
description,
SUM(unitprice * quantity) AS total_paid
FROM online_retail
GROUP BY month, description
--, unitprice, quantity          - 不加sum需要Group by所有变量
),
rank AS(
SELECT
    month,
    description,
    total_paid,
    DENSE_RANK() OVER (PARTITION BY month ORDER BY total_paid DESC) AS rnk
FROM total
)
SELECT
    month,
    description,
    total_paid
FROM rank
WHERE rnk = 1
```

法2: 将SUM(P*N)与 RANK放在一个Query里更有效率

```
SELECT
    MONTH,
    description,
    total_paid
FROM
    (SELECT
        date_part('month', invoicedate) AS MONTH,
        description,
        SUM(unitprice * quantity) AS total_paid,
        RANK() OVER (PARTITION BY date_part('month', invoicedate)
                     ORDER BY SUM(unitprice * quantity) DESC) AS rnk
    FROM online_retail
    GROUP BY MONTH, description ) tmp
WHERE rnk = 1
```

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

Netflix Hard ID 10171

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.

Tables: oscar nominees. nominee information

In []: 求user的count, count最大值记录对应的属性:
 i. 求每个演员得奖几次, 对应的排名: name, count(*), rank() over (order by COUNT(*))
 ii. 求genre让名字从刚才的query中取rank1: name IN (select nominee from wins where

```
nominee, COUNT(*) OVER (PARTITION BY nominee) = nominee, COUNT(*), R + GROUP
```

```
-- person with most # of oscar 's genre
-- >1 person same # of oscar wins, return 1st based on name asc
-- names as keys join table
--output: top_genre
```

法1:

i. 求每个演员得奖几次, 对应的排名: name, count(*), rank() over (order by COUNT(*) DESC);
 ii. 求genre让名字从刚才的query中取rank1: name IN (select nominee from wins where
 rnk = 1)

```
WITH wins AS(
SELECT
    a.nominee,
    COUNT(*) AS wins,
    RANK() OVER (ORDER BY COUNT(*) DESC) AS rnk
FROM oscar_nominees a
JOIN nominee_information b
ON a.nominee = b.name
WHERE winner = TRUE
GROUP BY 1
)
SELECT DISTINCT top_genre
FROM nominee_information
WHERE name IN (SELECT nominee FROM wins WHERE rnk = 1)
```

法2:

```
nominee, COUNT(*) OVER (PARTITION BY nominee) = nominee, COUNT(*), R + GROUP
BY 1
```

```
SELECT
    DISTINCT top_genre
FROM nominee_information info
INNER JOIN (
    SELECT
        nominee,
        n_winnings
    FROM (
        SELECT
            nominee,
            COUNT(*) OVER (PARTITION BY nominee) AS n_winnings
        FROM oscar_nominees
        WHERE winner = true
    ) tmp
    WHERE n_winnings = 2
```

```
ORDER BY 2 DESC, 1 ASC
) tmp
ON tmp.nominee = info.name
```

37. Gender With Most Doctor Appointments

HealthTap Natera Easy ID 10170

Find the gender that has made the most number of doctor appointments. Output the gender along with the corresponding number of appointments.

Table: medical_appointments

```
--gender most # doc app
--output:gender, n_appointments

--MOST: order by desc
--aggreagte: group by

SELECT
    gender,
    COUNT(appointmentid) AS n_appointments
FROM medical_appointments
GROUP BY gender
ORDER BY n_appointments DESC
LIMIT 1
```

38. Highest Total Miles

Interview Question Date: July 2020

Uber Medium ID 10169

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.

Table: my_uber_drives

```
-- RANK() OVER(ORDER BY X DESC) AS rnk + where rnk <=3; LIMIT 3 可以互换
```

```
--Uber rides: mileage and the purpose for the biz expense
-- most miles
-- top 3 biz purpose categories by total mileage.
--output: purpose, miles_sum

SELECT
    purpose,
    miles_sum
FROM
    (SELECT
        DISTINCT
```

```

        purpose,
        SUM(miles) AS miles_sum,
        RANK() OVER (ORDER BY SUM(miles) DESC) AS rnk
    FROM my_uber_drives
    WHERE category = 'Business'
    GROUP BY purpose
    ORDER BY miles_sum DESC) t
WHERE rnk <= 3

```

39. Number Of Records By Variety

Microsoft Linux Easy ID 10168

Find the total number of records that belong to each variety in the dataset. Output the variety along with the corresponding number of records. Order records by the variety in ascending order.

Table: iris

```

-- # total belong to each variety
--output: variety, n_total_varieties: variety along with the corresponding
number of records
--order records asc
SELECT
    variety,
    COUNT(*) AS n_total_varieties
FROM iris
GROUP BY 1
ORDER BY 2 ASC

```

40. Total Number Of Housing Units

Airbnb Zillow Easy ID 10167

Find the total number of housing units completed for each year. Output the year along with the total number of housings. Order the result by year in ascending order.

Note: Number of housing units in thousands.

Table: housing_units_completed_us

```

-- each year total # housing units completed
--output: year, n_units: total number of housings.
--order by year asc

--Number of housing units in thousands

SELECT
    year,
    SUM(south + west + midwest + northeast) AS n_units
FROM housing_units_completed_us
GROUP BY year

```

ORDER BY year

41.Reviews of Hotel Arena

Airbnb Expedia Easy ID 10166

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.

Table: hotel_reviews

```
--# of rows for each review score ~ Hotel Arena
--output: hotel_name: Hotel Arena; reviewer_score; count

SELECT
    hotel_name,
    reviewer_score,
    COUNT(*) AS count
FROM hotel_reviews
WHERE hotel_name = 'Hotel Arena'
GROUP BY hotel_name, reviewer_score
```

42.Total AdWords Earnings

Interview Question Date: July 2020

Google Easy ID 10164

Find the total AdWords earnings for each business type. Output the business types along with the total earnings.

Table: google_adwords_earnings

```
--total Adwords earnings for each biz type
--output: business_type, earnings: total earnings.

SELECT
    business_type,
    SUM(adwords_earnings) AS earnings
FROM google_adwords_earnings
GROUP BY business_type
```

43.Product Transaction Count

Microsoft Nvidia Medium ID 10163

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.

Tables: excel_sql_inventory_data, excel_sql_transaction_data

order by的变量, 如果没有在select中出现过, 则需要在group by中出现一次

```
--# transactions for each product
--output: product_name, count: number of transactions
--order records by the product id in ascending order
--ignore products without transactions.

SELECT
    i.product_name,
    COUNT(*) AS count
FROM excel_sql_transaction_data t
LEFT JOIN excel_sql_inventory_data i
ON t.product_id = i.product_id
GROUP BY product_name, t.product_id
ORDER BY t.product_id
```

44. Number Of Acquisitions

Crunchbase Easy ID 10162

Find the number of acquisitions that occurred in each quarter of each year. Output the acquired quarter in YYYY-Qq format along with the number of acquisitions and order results by the quarters with the highest number of acquisitions first.

Table: crunchbase_acquisitions

```
--# acquisition occ in each quarter of each year

--output: acquired_quarter: 2013-Q4; cnt_acq: the number of acquisitions
--order by cntacq DESC: quarters with highest # acquisition first

SELECT
    acquired_quarter,
    COUNT(*) AS cnt_acq
FROM crunchbase_acquisitions
GROUP BY acquired_quarter
ORDER BY cnt_acq DESC
```

45. Ranking Hosts By Beds

Interview Question Date: July 2020

Airbnb Medium ID 10161

Rank each host based on the number of beds they have listed. The host with the most beds should be ranked 1 and the host with the least number of beds should be ranked last. Hosts that have the same number of beds should have the same rank but there should be no gaps between ranking values. A host can also own multiple properties. Output the host ID, number of beds, and rank from highest rank to lowest.

```
--have the same number of beds should have the same rank but there should be no gaps between ranking values.
```

DENSE_RANK: 同一数量的赋值相同的rank数; rank每等之间没有gap

```
SELECT
    host_id,
    SUM(n_beds) AS number_of_beds,
    DENSE_RANK() OVER (ORDER BY SUM(n_beds) DESC) AS rnk
FROM airbnb_apartments
GROUP BY host_id
```

46. Rank guests based on their ages

Airbnb Easy ID 10160

Rank guests based on their ages. Output the guest id along with the corresponding rank. Order records by the age in descending order.

Table: airbnb_guests

RANK() OVER 前的变量, 不需要group by

```
--rank guest ages
--output: guest id, rank
--order records by age desc

SELECT
    guest_id,
    RANK() OVER (ORDER BY age DESC) AS rank
FROM airbnb_guests
```

47. Ranking Most Active Guests

Airbnb Medium ID 10159

Rank guests based on the number of messages they've exchanged with the hosts. Guests with the same number of messages as other guests should have the same rank. Do not skip rankings if the preceding rankings are identical. Output the rank, guest id, and number of total messages they've sent. Order by the highest number of total messages first.

Table: airbnb_contacts

rank也可以放在第一列

```
--rank guest ; # messages eachanged with hosts
--dense_rank: guests same # messages ~ same rank; not skip rank if preceding ranks are identical
--output: ranking, id_guest, sum_n_messages

SELECT
```



```

DENSE_RANK() OVER(ORDER BY SUM(n_messages) DESC) AS ranking,
id_guest,
SUM(n_messages) AS sum_n_messages
FROM
    airbnb_contacts
GROUP BY
    id_guest

```

48. Number Of Units Per Nationality

Airbnb Medium ID 10156

Find the number of apartments per nationality that are owned by people under 30 years old.

Output the nationality along with the number of apartments.

Sort records by the apartments count in descending order.

Tables: airbnb_hosts, airbnb_units

apartment count是count(distinct unit_id)
因为airbnb_hosts表中有很多重复的host_id = 1的记录, 所以和airbnb_units JOIN之后, 有很多重复记录 此时不能直接COUNT(*) 要COUNT(DISTINCT unit_id)

```

--# apt per nationality owned by people age <30
--output: nationality, apartment_count
--sort by apt count desc

SELECT
    h.nationality,
    COUNT(DISTINCT u.unit_id) AS apartment_count
FROM airbnb_hosts h
JOIN airbnb_units u
ON h.host_id = u.host_id
WHERE h.age < 30
AND unit_type = 'Apartment'
GROUP BY h.nationality
ORDER BY apartment_count DESC

```

49. Find the number of Yelp businesses that sell pizza

Yelp Easy ID 10153

Find the number of Yelp businesses that sell pizza.

Table: yelp_business

```
LOWER(X) LIKE '%A%' = X ILIKE '%A%'
```

```
lower(categories) like '%pizza%' = categories ILIKE '%Pizza%'
```

```

SELECT
    COUNT(*) AS count

```

```
FROM yelp_business
WHERE categories LIKE '%Pizza%'
```

50. Workers With The Highest And Lowest Salaries

Amazon Siemens Medium ID 10152

You have been asked to find the employees with the highest and lowest salary.

Your output should include the employee's ID, salary, and department, as well as a column salary_type that categorizes the output by:

'Highest Salary' represents the highest salary

'Lowest Salary' represents the lowest salary

Tables: worker, title

1. case when可以全命名
2. case when也可以只命名最高最低
3. 正反rank法:
 - i. 赋值正反rank: lowest_sal, highest_sal
 - ii. case when highest_sal = 1 THEN 'A' ELSE 'B': 赋值最高rank为A, 其他rank为B
 - iii. where highest_sal = 1 OR lowest_sal = 1: 只挑出最高和最低
4. '' AS X 命名 + where select max/min, 分别求出最高和最低UNION ALL

5. case when写法:

```
CASE
  WHEN salary = (SELECT MAX(salary) FROM worker) THEN 'Highest Salary'
  WHEN salary = (SELECT MIN(salary) FROM worker) THEN 'Lowest Salary'
ELSE 'N/A'
END AS salary_type
```

6. 两种case方法可以替换:

①

```
CASE
  WHEN salary = (SELECT max(salary) FROM total) THEN 'Highest Salary'
  WHEN salary = (SELECT min(salary) FROM total) THEN 'Lowest Salary'
ELSE NULL
END
```

②

```
CASE
  WHEN salary = (SELECT max(salary) FROM total THEN 'Highest Salary')
ELSE 'Lowest Salary'
END
```

```
--find employees with highest and lowest salary
--salary_type : Highest Salary, Lowest Salary
--output: worker_id, salary, department, salary_type
```

```
--法1: 只命名最高最低, where 取出 rnk 在 最高和最低 中的 rnk值(union)
WITH total AS(
```

```

SELECT
    worker_id,
    salary,
    department,
    RANK() OVER (ORDER BY salary DESC) AS rnk
FROM worker
)
SELECT
    worker_id,
    salary,
    department,
    -- CASE
    --     WHEN salary = (SELECT max(salary) FROM total) THEN 'Highest Salary'
    --     WHEN salary = (SELECT min(salary) FROM total) THEN 'Lowest Salary'
    -- ELSE NULL
    -- END AS salary_type
    CASE
        WHEN salary = (SELECT max(salary) FROM total) THEN 'Highest Salary'
        ELSE 'Lowest Salary'
    END AS salary_type
FROM total
WHERE rnk IN (SELECT min(rnk) FROM total UNION SELECT max(rnk) FROM total)

```

--法2: 全命名, where只取出<> 'N/A'的记录, 即最大和最小的记录

```

WITH total AS(
SELECT
    worker_id,
    salary,
    department,
    CASE
        WHEN salary = (SELECT MAX(salary) FROM worker) THEN 'Highest Salary'
        WHEN salary = (SELECT MIN(salary) FROM worker) THEN 'Lowest Salary'
        ELSE 'N/A'
    END AS salary_type
FROM worker
)
SELECT * FROM total
WHERE salary_type <> 'N/A'

```

--法3: 正反rank法:

-- i. 赋值正反rank: lowest_sal, highest_sal
 -- ii. case when highest_sal = 1 THEN 'A' ELSE 'B': 赋值最高rank为A, 其他rank为B
 -- iii. where highest_sal = 1 OR lowest_sal = 1: 只挑出最高和最低

```

WITH CTE AS(
SELECT
    *,
    RANK() OVER (ORDER BY salary) AS lowest_sal,
    RANK() OVER (ORDER BY salary DESC) AS highest_sal
FROM worker
)
SELECT
    worker_id,
    salary,
    department,
    -- lowest_sal,
    -- highest_sal,
    CASE

```

```

        WHEN highest_sal = 1 THEN 'Highest Salary'
        ELSE 'Lowest Salary'
        END AS salary_type
FROM CTE
WHERE highest_sal = 1 OR lowest_sal = 1

```

```

--法4: '' AS X 命名 + where select max/min, 分别求出最高和最低UNION ALL
with highest_salary as(
select
    worker_id,
    salary,
    department,
    'Highest Salary' as salary_type
from worker
where salary = (select max(salary) from worker)
),
lowest_salary as(
select
    worker_id,
    salary,
    department,
    'Lowest Salary' as salary_type
from worker
where salary = (select min(salary) from worker)
)
select * from highest_salary
union all
select * from lowest_salary

```

51. Gender With Generous Reviews

Interview Question Date: June 2020

Airbnb Easy ID 10149

Write a query to find which gender gives a higher average review score when writing reviews as guests. Use the from_type column to identify guest reviews. Output the gender and their average review score.

Tables: airbnb_reviews, airbnb_guests

注意: airbnb_reviews中的from_user同时包含guest_id, user_id
 ①需要先where挑选出from_type = 'guest', 才能得出guest_id的记录, 再与airbnb_guests join
 ②也可以两个表先join一下, 再只where filter出 from_type = 'guest'的记录就可以了

```

-- find which gender gives a higher avg review score when as guests
-- output: gender, avg_score

SELECT
    gender,
    AVG(review_score) AS avg_score
FROM
    (
        SELECT
            from_user,

```

```

        review_score
    FROM airbnb_reviews r
    WHERE from_type = 'guest'
    ) t
    JOIN airbnb_guests g
    ON t.from_user = g.guest_id
GROUP BY gender
ORDER BY avg_score DESC
LIMIT 1

```

```

SELECT
    g.gender,
    AVG(review_score) AS avg_score
FROM airbnb_reviews r
JOIN airbnb_guests g
ON r.from_user = g.guest_id
WHERE r.from_type = 'guest'
GROUP BY g.gender
ORDER BY avg_score DESC
LIMIT 1

```

52. Find the top 5 cities with the most 5 star businesses

Yelp Medium ID 10148

Find the top 5 cities with the most 5-star businesses. Output the city name along with the number of 5-star businesses. In the case of multiple cities having the same number of 5-star businesses, use the ranking function returning the lowest rank in the group and output cities with a rank smaller than or equal to 5.

Table: yelp_business

limit 5: 只展示5个记录
rank <=5: 将rank = 5个以内的所有记录挑出

```

-- top 5 cities most 5-star biz
-- DENSE_RANK: rank return lowest rank in the group; rank <=5
-- output: city, count_of_5_stars

-- 'In the case of multiple cities having the same number of 5-star
businesses, use the ranking function returning the lowest rank in the group'
-- It means that you should use:

-- method = 'min' if you're using Python and Pandas
-- rank() if you use PostgreSQL (minimum method is default to it)
-- It's just a way of making it clear what rank function should be used.

SELECT
    city,
    count_of_5_stars
FROM (
    SELECT
        city,

```

```

COUNT(*) AS count_of_5_stars,
RANK() OVER (ORDER BY COUNT(*) DESC) as rnk
FROM yelp_business
WHERE stars = 5
GROUP BY city) t
WHERE rnk <= 5

```

53. Find countries that are in winemag_p1 dataset but not in winemag_p2

Interview Question Date: June 2020

Wine Magazine Medium ID 10147

Find countries that are in winemag_p1 dataset but not in winemag_p2. Output distinct country names. Order records by the country in ascending order.

Tables: winemag_p1, winemag_p2

```

# countries in p1 but not in p2
# output country: distinct country
# order by country asc

SELECT country FROM winemag_p1
WHERE country NOT IN (SELECT country FROM winemag_p2)
ORDER BY country ASC

```

54. Make a pivot table to find the highest payment in each year for each employee

City of San Francisco Hard ID 10145

Make a pivot table to find the highest payment in each year for each employee. Find payment details for 2011, 2012, 2013, and 2014. Output payment details along with the corresponding employee name. Order records by the employee name in ascending order

Table: sf_public_salaries

题中提到pivot <=> case when: eg. MAX(case when) + group by
 else 0 把其余值赋为0
 MAX(pay_year) 去找到这年最高的totalpay 因为可能一年有两个total pay

```

-- select * from sf_public_salaries;

-- PV to find each yr, each employee, highest payment
--2011 - 2014
--output: payment, employee name
--employee name, pay_2011, pay_2012, pay_2013, pay_2014
--order name asc

SELECT
    employeename,
    -- year,

```

```
--totalpay,
MAX(CASE WHEN year = 2011 THEN totalpay ELSE 0 END) AS pay_2011,
MAX(CASE WHEN year = 2012 THEN totalpay ELSE 0 END) AS pay_2012,
MAX(CASE WHEN year = 2013 THEN totalpay ELSE 0 END) AS pay_2013,
MAX(CASE WHEN year = 2014 THEN totalpay ELSE 0 END) AS pay_2014
-- totalpaybenefits
FROM sf_public_salaries
GROUP BY employeename
ORDER BY employeename
```

55.Average Weight of Medal-Winning Judo

ESPN Medium ID 10144

Find the average weight of medal-winning Judo players of each team with a minimum age of 20 and a maximum age of 30. Consider players at the age of 20 and 30 too. Output the team along with the average player weight.

Table: olympics_athletes_events

```
SELECT
    team,
    AVG(weight) AS average_player_weight
FROM olympics_athletes_events
WHERE sport = 'Judo' AND age BETWEEN 20 AND 30 AND medal IS NOT NULL
GROUP BY team
```

```
SELECT
    team,
    avg(weight) AS average_player_weight
FROM olympics_athletes_events
WHERE sport = 'Judo'
AND medal IS NOT NULL
GROUP BY 1
HAVING min(age) >= 20 and max(age) <= 30
ORDER BY 1
```

56.Find players who participated in the Olympics representing more than one team

ESPN Easy ID 10143

Find players who participated in the Olympics representing more than one team. Output the player name, team, games, sport, and the medal.

Table: olympics_athletes_events

Find players who participated in the Olympics representing more than one team

ESPN
Easy
ID 10143

Find players who participated in the Olympics representing more than one team. Output the player name, team, games, sport, and the medal.

Table: olympics_athletes_events

X column has A: X ILIKE '%A%'

每次代表多于一个team: 在team中 包含/的值: X ILIKE '%/%'

```
--players:OLYMPIC, more than 1 team
--output: name, team, games, sport, medal

--X column has A: X ILIKE '%A%'
--每次代表多于一个team: 在team中, 包含/的值: X ILIKE '%/%'

SELECT
    name,
    team,
    games,
    sport,
    medal
FROM olympics_athletes_events
WHERE team ILIKE '%/%'
```

57. Apple Product Counts

Google Apple Medium ID 10141

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.

Tables: playbook_events, playbook_users

```
1. ILIKE '%macbook pro%' AND '%iphone 5s%' AND 'ipad air')
2. order by 可以用产生的alias
3. both count distinct user_id
```

```
SELECT
    language,
    COUNT(DISTINCT CASE
        WHEN device IN('macbook pro',
                        'iphone 5s',
                        'ipad air') THEN a.user_id
        END) AS n_apple_users,
    COUNT(DISTINCT a.user_id) AS n_total_users
FROM playbook_events a
JOIN playbook_users b
ON a.user_id = b.user_id
WHERE device IS NOT NULL
GROUP BY language
ORDER BY n_total_users DESC
```

58. MacBook Pro Events

Google Apple Medium ID 10140

Find how many events happened on MacBook-Pro per company in Argentina from users that do not speak Spanish. Output the company id, language of users, and the number of events performed by users.

Tables: `playbook_events`, `playbook_users`

Spanish Not English

```
-- events on macbook-pro in argentina from users not speak Spanish
-- output: company_id, language: language of users, n_macbook_pro_events: the
number of events performed by users.
```

```
SELECT
    u.company_id,
    u.language,
    COUNT(*) AS n_macbook_pro_events
FROM playbook_events e
JOIN playbook_users u
ON e.user_id = u.user_id
WHERE e.location = 'Argentina'
AND e.device = 'macbook pro'
AND u.language != 'spanish'
GROUP BY u.company_id, u.language
```

59. Number of Speakers By Language

Google Apple Medium ID 10139

Find the number of speakers of each language by country. Output the country, language, and the corresponding number of speakers. Output the result based on the country in ascending order.

Tables: `playbook_events`, `playbook_users`

```
--country, language, #speaker
--output: location, language, n_speakers
--order country asc
```

```
SELECT
    e.location,
    u.language,
    COUNT(DISTINCT u.user_id)
FROM playbook_events e
JOIN playbook_users u
ON e.user_id = u.user_id
GROUP BY e.location, u.language
ORDER BY e.location
```

60. Even-numbered IDs Hired in June

Amazon Bosch Easy ID 10137

Find employees who started in June and have even-numbered employee IDs.

Table: worker

1. 奇偶数表达

偶数 even:

$\text{num} \% 2 = 0$

$\text{mod}(\text{num}, 2) = 0$

奇数 odd:

$\text{num} \% 2 \neq 0$

$\text{mod}(\text{num}, 2) \neq 0$

2. 月份表达

$\text{EXTRACT}(\text{MONTH FROM X}): \text{EXTRACT}(\text{MONTH FROM joining_date}) \text{ AS month}$

$\text{DATE_PART}(\text{'MONTH'}, X): \text{DATE_PART}(\text{'MONTH'}, \text{joining_date}) \text{ AS month}$

3. 没在变量中出现的, 也可以直接where filter出

$\text{WHERE } \text{EXTRACT}(\text{MONTH FROM joining_date}) = 6$

--#法1:

$\text{SELECT } * \text{ FROM worker}$

$\text{WHERE } \text{mod}(\text{worker_id}, 2) = 0$

$\text{AND } \text{EXTRACT}(\text{MONTH FROM joining_date}) = 6$

--#法2:

$\text{WITH odd_num_worker AS}(\text{$

SELECT

$*$,

$\text{EXTRACT}(\text{MONTH FROM joining_date}) \text{ AS month}$

$--\text{DATE_PART}(\text{'MONTH'}, \text{joining_date}) \text{ AS month}$

FROM worker

$\text{WHERE } \text{worker_id} \% 2 = 0$

$)$

SELECT

$\text{worker_id},$

$\text{first_name},$

$\text{last_name},$

$\text{salary},$

$\text{joining_date},$

department

$\text{FROM odd_num_worker}$

$\text{WHERE } \text{month} = 6$

61. Odd-numbered ID's Hired in February

Amazon Bosch Easy

Find employees who started in February and have odd-numbered employee IDs.

Table: worker

$\text{SELECT } * \text{ FROM worker}$

$\text{WHERE } \text{mod}(\text{worker_id}, 2) \neq 0$

```
AND EXTRACT(MONTH FROM joining_date) = 2
```

62.Spam Posts

Interview Question Date: June 2020

Meta/Facebook Medium ID 10134

Calculate the percentage of spam posts in all viewed posts by day. A post is considered a spam if a string "spam" is inside keywords of the post. Note that the facebook_posts table stores all posts posted by users. The facebook_post_views table is an action table denoting if a user has viewed a post.

Tables: facebook_posts, facebook_post_views

求部分与整体比值:

1. $\text{COUNT}(\text{CASE WHEN X THEN ID END}) / \text{COUNT}(\text{ID}) + \text{JOIN}$
2. $\text{SUM}(\text{CASE WHEN X THEN 1 ELSE 0 END}) * 100 / \text{COUNT}(\text{*}) + \text{JOIN}$
3. 整体 left join 部分: $\text{SELECT 部分} / \text{整体求比值}$

```
--percentage spam posts in all viewed posts by day
--spam: post = 'spam'
--facebook_posts stores all posts by user; view: view post

--output: post_date, spam_share

--#法1:
WITH total AS(
SELECT
    p.post_date,
    COUNT(v.post_id) AS total_viewed_posts,
    COUNT(CASE
        WHEN p.post_keywords ILIKE '%spam%' THEN v.post_id
    END) AS spam_posts
FROM facebook_posts p
JOIN facebook_post_views v
ON p.post_id = v.post_id
GROUP BY post_date
)
SELECT
    post_date,
    spam_posts / total_viewed_posts::float * 100 AS spam_share
FROM total
```

```
--#法2: 部分与整体比值:  $\text{SUM}(\text{CASE WHEN X THEN 1 ELSE 0 END}) * 100 / \text{COUNT}(\text{*}) + \text{JOIN}$ 
SELECT
    post_date,
    SUM(CASE WHEN post_keywords ILIKE '%spam%' THEN 1 ELSE 0 END) * 100 /
COUNT(*)
FROM facebook_posts p
JOIN facebook_post_views v
ON p.post_id = v.post_id
```

```
GROUP BY post_date
```

```
--#法3: 整体 left join 部分; SELECT 部分/整体求比率

SELECT
    spam_summary.post_date,
    n_spam / n_posts::FLOAT * 100 AS spam_share
FROM
    (SELECT
        post_date,
        SUM(CASE
            WHEN v.viewer_id IS NOT NULL THEN 1
            ELSE 0
        END) AS n_posts
    FROM facebook_posts p
    JOIN facebook_post_views v
    ON p.post_id = v.post_id
    GROUP BY post_date) posts_summary
LEFT JOIN
    (SELECT
        post_date,
        SUM(CASE
            WHEN v.viewer_id IS NOT NULL THEN 1
            ELSE 0
        END) AS n_spam
    FROM facebook_posts p
    JOIN facebook_post_views v
    ON p.post_id = v.post_id
    WHERE post_keywords ILIKE '%spam%'
    GROUP BY post_date) spam_summary
ON posts_summary.post_date = spam_summary.post_date
```

63.Requests Acceptance Rate

Airbnb Medium ID 10133

Find the acceptance rate of requests which is defined as the ratio of accepted contacts vs all contacts. Multiply the ratio by 100 to get the rate.

Table: airbnb_contacts

In []: 整体/部分比值:

```
1. COUNT(ts.部分) / COUNT(整体)
2. SUM(CASE WHEN X THEN 1 ELSE 0 END) / COUNT(*)
3. COUNT(CASE WHEN Y THEN 1) / COUNT(*)
```

```
--#法1: COUNT(ts.部分) / COUNT(整体)
SELECT
    COUNT(ts_accepted_at) / COUNT(ts_contact_at)::FLOAT * 100 AS
    acceptance_rate
FROM airbnb_contacts
```

```
--#法2: SUM(CASE WHEN X THEN 1 ELSE 0 END) / COUNT(*)
SELECT
    100 * SUM(CASE
        WHEN ts_accepted_at IS NOT NULL THEN 1
```

```

        ELSE 0
      END) / COUNT(*) AS acceptance_rate
FROM aishbh_contacts

```

```

--#法3: COUNT(CASE WHEN X THEN A) / COUNT(*)
SELECT
  COUNT(CASE
    WHEN ts_accepted_at IS NOT NULL THEN ts_accepted_at
  END) / COUNT(*)::FLOAT * 100 AS acceptance_rate
FROM aishbh_contacts

```

64. Highest Crime Rate

City of San Francisco Easy ID 10132

Find the number of crime occurrences for each day of the week. Output the day alongside the corresponding crime count.

Table: sf_crime_incidents_2014_01

```

--# crime occurrences for each day of the week

--day_of_week; n_occurences

SELECT
  day_of_week,
  COUNT(*) AS n_occurences
FROM sf_crime_incidents_2014_01
GROUP BY
  day_of_week
ORDER BY
  n_occurences DESC

```

65. Business Name Lengths

Interview Question Date: June 2020

City of San Francisco Hard ID 10131

Find the number of words in each business name. Avoid counting special symbols as words (e.g. &). Output the business name and its count of words.

Table: sf_restaurant_health_violations

1. `regexp_replace(business_name, '[^a-zA-Z0-9]', '', 'g')`: 去除除字母以外的字符
将源文本中所有匹配上述正则表达式的字符都移除，从而得到一个只包含字母、数字和空格的新字符串。

'`^[a-zA-Z0-9]`' 是一个正则表达式模式，用于匹配任何不是字母、数字或空格的字符。

'`''`' 是一个空字符串，表示将匹配到的文本部分替换为空，即移除匹配到的字符。

'`g`' 是一个标志，表示全局匹配，会替换所有匹配到的字符，而不仅仅是第一个匹配项。

2. `array_length(regexp_split_to_array(b_name, '\s+'), 1)` : 计算词组数量
将会计算 `b_name` 经过空格字符拆分后的数组的长度，即计算 `b_name` 中包含的单词数量。

`array_length()` 是一个函数, 它用于获取数组的长度 (即数组中元素的数量)
`'\s+'` 是一个正则表达式模式 表示匹配一个或多个连续的空格字符

```
SELECT
    DISTINCT business_name,
    array_length(regexp_split_to_array(b_name, '\s+'), 1) AS word_count
FROM
    (SELECT
        business_name,
        regexp_replace(business_name, '^[a-zA-Z0-9 ]', '', 'g') AS b_name
    FROM sf_restaurant_health_violations) AS sfr
```

66. Find the number of inspections for each risk category by inspection type

Interview Question Date: June 2020

City of San Francisco Medium ID 10130 25

Data Engineer Data Scientist BI Analyst Data Analyst Find the number of inspections that resulted in each risk category per each inspection type. Consider the records with no risk category value belongs to a separate category. Output the result along with the corresponding inspection type and the corresponding total number of inspections per that type. The output should be pivoted, meaning that each risk category + total number should be a separate column. Order the result based on the number of inspections per inspection type in descending order.

Table: `sf_restaurant_health_violations`

法1: `COUNT(CASE WHEN X THEN id END)`
 法2: `SUM(CASE WHEN X THEN 1 ELSE 0 END)`

```
#法1: COUNT(CASE WHEN X THEN id END)
SELECT
    inspection_type,
    COUNT(CASE WHEN risk_category IS NULL THEN inspection_id END) AS
no_risk_results,
    COUNT(CASE WHEN risk_category = 'Low Risk' THEN inspection_id END) AS
low_risk_results,
    COUNT(CASE WHEN risk_category = 'Moderate Risk' THEN inspection_id END) AS
medium_risk_results,
    COUNT(CASE WHEN risk_category = 'High Risk' THEN inspection_id END) AS
high_risk_results,
    COUNT(inspection_id) AS total_inspections
FROM sf_restaurant_health_violations
GROUP BY inspection_type
ORDER BY total_inspections DESC
```

```
#法2: SUM(CASE WHEN X THEN 1 ELSE 0 END)
SELECT
    inspection_type,

SUM(CASE
    WHEN risk_category IS NULL THEN 1 ELSE 0
END) AS no_risk_results,
```

```
SUM(CASE
    WHEN risk_category = 'Low Risk' THEN 1 ELSE 0
END) AS low_risk_results,

SUM(CASE
    WHEN risk_category = 'Moderate Risk' THEN 1 ELSE 0
END) AS medium_risk_results,

SUM(CASE
    WHEN risk_category = 'High Risk' THEN 1 ELSE 0
END) AS high_risk_results,

COUNT(*) AS total_inspections

FROM
    sf_restaurant_health_violations

GROUP BY
    inspection_type

ORDER BY
    total_inspections DESC
```

67. Count the number of movies that Abigail Breslin nominated for oscar

Google Netflix Easy ID 10128

Count the number of movies that Abigail Breslin was nominated for an oscar.

Table: oscar_nominees

去重重复的记录: COUNT(DISTINCT X)

```
SELECT
    COUNT(DISTINCT movie) AS n_movies_by_abi
FROM oscar_nominees
WHERE nominee = 'Abigail Breslin'
```

68. Calculate Samantha's and Lisa's total sales revenue

Amazon Groupon Salesforce Easy ID 10127

What is the total sales revenue of Samantha and Lisa?

Table: sales_performance

去AB两个对应的值:
X IN ('A', 'B')
X = 'A' OR 'B'

```
WHERE salesperson IN ('Samantha', 'Lisa')
WHERE salesperson = 'Samantha' OR salesperson = 'Lisa'
```

```
SELECT
    SUM(sales_revenue) AS total_revenue
FROM sales_performance
WHERE salesperson IN ('Samantha', 'Lisa')
--WHERE salesperson = 'Samantha' OR salesperson = 'Lisa'
```

69.Bookings vs Non-Bookings

Interview Question Date: May 2020

Airbnb Medium ID 10124

Display the average number of times a user performed a search which led to a successful booking and the average number of times a user performed a search but did not lead to a booking. The output should have a column named action with values 'does not book' and 'books' as well as a 2nd column named average_searches with the average number of searches per action. Consider that the booking did not happen if the booking date is null. Be aware that search is connected to the booking only if their check-in dates match.

Tables: airbnb_contacts, airbnb_searches

整体 (有book的和没有book的) LEFT JOIN 部分 (book的) ON guest book的条件
 整体: 做了case when: ts_book不为空且ds.check-in相等 为'book', 否则'does not book', AVG(n_searches)
 部分: 所有ts_book不为空的记录
 ON: id_user = id_guest AND ds_checkin = ds_checkin: 把部分 (book的记录) 与整体 (book&unbook的记录) match上

```
--avg # times a user search ~ succ booking & avg # times a user ~ not booking
--output: action: books, does not book , average_searches: average number of
searches per action.
-- book not happen if book date is null; book happen only if check-in dates
match
```

法1:

```
SELECT
    CASE
        WHEN c.ts_booking_at IS NOT NULL AND c.ds_checkin = s.ds_checkin
    THEN 'books'
        ELSE 'does not book'
    END AS action,
    AVG(n_searches) AS average_searches
FROM airbnb_searches s
LEFT JOIN
(SELECT *
    FROM airbnb_contacts
WHERE ts_booking_at IS NOT NULL) c
ON s.id_user = c.id_guest
--ON s.id_user = c.id_guest AND s.ds_checkin = c.ds_checkin 效果与ON s.id_user
= c.id_guest 一样
```



```
GROUP BY 1
```

```
法2: 全 left join 部分
WITH c as
(select
    *
FROM airbnb_contacts
WHERE ts_booking_at is not null    --book的
)
SELECT                                --book+not book的
CASE WHEN c.ts_booking_at IS NOT NULL AND c.ds_checkin = s.ds_checkin THEN
'books'  --search ~ booking: ts_checkin date matches
      ELSE 'does not book'
END AS action,
AVG(n_searches) AS average_searches
FROM airbnb_searches s LEFT JOIN c --用全的left join部分的
ON s.id_user = c.id_guest          --顾客id = 用户id
GROUP BY 1

-- ※此处s是有book的, 有没book的; c是全book的
-- s left join c
```

70. Find the total number of searches for houses Westlake neighborhood with a TV

Airbnb Easy ID 10122

Find the total number of searches for houses in Westlake neighborhood with a TV among the amenities.

Table: airbnb_search_details

```
SELECT
    COUNT(*)
FROM airbnb_search_details
WHERE neighbourhood = 'Westlake'
AND amenities ILIKE '%TV%'
AND property_type = 'House'
```

71. Number Of Custom Email Labels

Google Medium ID 10120

Find the number of occurrences of custom email labels for each user receiving an email. Output the receiver user id, label, and the corresponding number of occurrences.

Tables: google_gmail_emails, google_gmail_labels

```
--# custom eail labels for each user receive email
--output:  user_id: receiver user id, label, n_occurences: # occurence

SELECT
    to_user AS user_id,
```

```

    label,
    COUNT(*) AS n_occurrences
FROM google_gmail_emails e
JOIN google_gmail_labels l
ON e.id = l.email_id
WHERE label ILIKE '%Custom%'
GROUP BY to_user, label

```

72. User Exile

Meta/Facebook Easy ID 10091

Find the number of relationships that user with id == 1 is not part of.

Table: facebook_friends

```

user1!=1 AND user2!=1
1 NOT IN (user1, user2)

```

```

SELECT
    COUNT(*) AS user1_not_in_relationship
FROM facebook_friends
WHERE 1 NOT IN (user1, user2)

```

```

SELECT
    COUNT(*) AS user1_not_in_relationship
FROM facebook_friends
WHERE user1!=1 AND user2!=1

```

73. Find the percentage of shipable orders

Google Amazon Medium ID 10090

Find the percentage of shipable orders. Consider an order is shipable if the customer's address is known.

Tables: orders, customers

```

1.COUNT(CASE WHEN X THEN Y END) / COUNT(*)::FLOAT AS pct
2.CASE WHEN X THEN FALSE ELSE TRUE END AS is_shipable 建立T/F列
SUM(CASE WHEN is_shipable THEN 1 ELSE 0 END) / COUNT(*)::NUMERIC AS pct
3.COUNT(A) * 100 / COUNT(*)

```

法1:

```

COUNT(CASE WHEN X THEN Y END) / COUNT(*)::FLOAT AS pct
SELECT
    COUNT(CASE WHEN c.address IS NOT NULL THEN o.id END) / COUNT(*)::FLOAT *
    100 AS percent_shipable
FROM orders o
JOIN customers c
ON o.cust_id = c.id

```

法2:

```

CASE WHEN X THEN FALSE ELSE TRUE END AS is_shipable 建立T/F列
SUM(CASE WHEN is_shipable THEN 1 ELSE 0 END) / COUNT(*)::NUMERIC AS pct

WITH is_shipable AS(
SELECT
    o.id,
    CASE WHEN address IS NULL THEN FALSE ELSE TRUE END AS is_shipable
FROM orders o
JOIN customers c
ON o.cust_id = c.id
)
SELECT
    SUM(CASE WHEN is_shipable THEN 1 ELSE 0 END)::NUMERIC / COUNT(*) * 100 AS
percent_shipable
FROM is_shipable

```

```

--法3: shippable = #orders with address/#orders

SELECT
    COUNT(c.address) * 100 / COUNT(*) AS percent_shipable
FROM orders o
JOIN customers c
ON o.cust_id = c.id

```

74. Find the number of customers without an order

Google Amazon Medium ID 10089

Find the number of customers without an order.

Tables: orders, customers

```

韦恩图不在求法:
X NOT IN SELECT (B FROM C)
LEFT JOIN + X IS NULL

```

```

SELECT
    COUNT(c.id) AS n_customers_without_orders
FROM customers c
LEFT JOIN orders o
ON o.cust_id = c.id
--WHERE c.id NOT IN (SELECT cust_id FROM orders)
WHERE o.cust_id IS NULL

```

75. Liked' Posts

Meta/Facebook Medium ID 10088

Find the number of posts which were reacted to with a like.

Tables: facebook_reactions, facebook_posts

```
--法1: 整体 left join部分 + 创建like列 T/F, SUM(CASE WHEN like_post = 'TRUE'
THEN 1 ELSE 0 END)
--法2: 算每个post有多少like: SUM(CASE WHEN X THE 1 ELSE 0 END); COUNT(*) +
num_like_each_post !=0
--法3: COUNT(DISTINCT ID) + WHERE reaction = 'like'
```

```
--法1: 整体 left join部分 + 创建like列 T/F, SUM(CASE WHEN like_post = 'TRUE'
THEN 1 ELSE 0 END)
WITH like_post AS(
SELECT
    DISTINCT
    p.post_id,
    CASE WHEN r.reaction = 'like' THEN TRUE ELSE FALSE END AS like_posts
FROM facebook_posts p
LEFT JOIN facebook_reactions r
ON p.post_id = r.post_id
)
SELECT
    SUM(CASE WHEN like_posts = 'TRUE' THEN 1 ELSE 0 END) AS
n_posts_with_a_like
FROM like_post
```

```
--法2: 算每个post有多少like: SUM(CASE WHEN X THE 1 ELSE 0 END); COUNT(*) +
num_like_each_post !=0
WITH num_like_each_post AS(
SELECT
    DISTINCT p.post_id,
    SUM(CASE WHEN reaction = 'like' THEN 1 ELSE 0 END) AS num_like_each_post
FROM facebook_posts p
LEFT JOIN facebook_reactions r
ON p.post_id = r.post_id
GROUP BY 1
)
SELECT
    COUNT(*)
FROM num_like_each_post
WHERE num_like_each_post != 0
```

```
--法3: COUNT(DISTINCT ID) + WHERE reaction = 'like'
SELECT
    COUNT(DISTINCT p.post_id) AS n_posts_with_a_like
FROM facebook_posts p
LEFT JOIN facebook_reactions r
ON p.post_id = r.post_id
WHERE r.reaction = 'like'
```

76.Find all posts which were reacted to with a heart

Meta/Facebook Easy ID 10087

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

Tables: facebook_reactions, facebook_posts

去除重复记录: `SELECT DISTINCT table.*`

```
--# posts reacted with a heart
--output all columns

SELECT
    DISTINCT p.*
FROM
    facebook_posts p
JOIN
    facebook_reactions r
ON
    p.post_id = r.post_id
WHERE
    r.reaction = 'heart'
```

77.Email Details Based On Sends

Google Medium ID 10086

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

Table: google_gmail_emails

法1: `COUNT(A), COUNT(B) + HAVING COUNT(A) < COUNT(B)`

法2: 整体 JOIN 部分 `(COUNT(A) / COUNT(B))`: 等于在原表后加了一列ratio; filter ratio < 1

法1: `COUNT(A), COUNT(B) + HAVING COUNT(A) < COUNT(B)`

```
WITH base AS(
SELECT
    day,
    COUNT(DISTINCT from_user) AS users_sending_emails,
    COUNT(DISTINCT to_user) AS users_receiving_emails
FROM google_gmail_emails
GROUP BY 1
HAVING COUNT(DISTINCT from_user) < COUNT(DISTINCT to_user)
)
SELECT *
FROM google_gmail_emails
WHERE day IN (SELECT day FROM base)
```

法2: 整体 JOIN 部分 `(COUNT(A) / COUNT(B))`: 等于在原表后加了一列ratio; filter ratio < 1

```
SELECT
    g.*
FROM google_gmail_emails g
JOIN
    (SELECT
        day,
        COUNT(DISTINCT from_user)::NUMERIC /
        COUNT(DISTINCT to_user) AS sent_received_ratio
    FROM
        google_gmail_emails
```

```

GROUP BY
    day) base
ON g.day = base.day
AND base.sent_received_ratio < 1

```

78. Meta/Facebook Matching Users Pairs

Meta/Facebook Medium ID 10085

Find matching pairs of Meta/Facebook employees such that they are both of the same nation, different age, same gender, and at different seniority levels. Output ids of paired employees.

Table: facebook_employees

法1: self join on e1.id < e2.id for non-repeat pairs

法2: self join: 将所有条件做在join里面 + id is not null

```

--法1: self join on e1.id < e2.id for non-repeat pairs
SELECT
    a.id AS employee_1,
    b.id AS employee_2
FROM facebook_employees a
JOIN facebook_employees b
ON a.id != b.id
WHERE a.location = b.location
AND a.age != b.age
AND a.gender = b.gender
AND a.is_senior != b.is_senior

```

--法2: self join: 将所有条件做在join里面 + id is not null

```

SELECT
    e1.id AS employee_1,
    e2.id AS employee_2
FROM
    facebook_employees e1
JOIN
    facebook_employees e2
ON
    e1.location = e2.location AND
    e1.age != e2.age AND
    e1.gender = e2.gender AND
    e1.is_senior != e2.is_senior
WHERE
    e1.id IS NOT NULL AND
    e2.id IS NOT NULL

```

79. Cum Sum Energy Consumption

Interview Question Date: April 2020

Meta/Facebook Hard ID 10084

Calculate the running total (i.e., cumulative sum) energy consumption of the Meta/Facebook data centers in all 3 continents by the date. Output the date, running total energy consumption, and running total percentage rounded to the nearest whole number.

1. Running Total 和 All Days Total
 i. total表: 先把三个地区的表 union all一起 (如有同一天的相同数字记录会保留)
 ii. by date表: date, sum(c) AS day_total求每天的消费明细
 iii. running total: date, sum(day_total) over (order by date) 求截止该日为止的消费总和和明细
 iv. all days total: sum(day_total)是所有日的消费总和
 v. pct_of_total: running_toal / all days total = sum(day_total) over (order by date) / sum(day_total))
 2. 保留整数: round(X / Y, 0)

```
-- running total (cum sum) energy consump in all 3 continents by date
--output date, cumulative_total_energy: running total,
percentage_of_total_energy: running total pct rounded to nearest whole number

WITH total AS(
    SELECT * FROM fb_eu_energy
    UNION ALL
    SELECT * FROM fb_na_energy
    UNION ALL
    SELECT * FROM fb_asia_energy
    ORDER BY 1
),
energy_by_date AS(
    SELECT
        date,
        SUM(consumption) AS total_energy
    FROM total
    GROUP BY date
    ORDER BY date
)
SELECT
    date,
    SUM(total_energy) OVER (ORDER BY DATE) AS cumulative_total_energy,
    -- (SELECT SUM(total_energy) FROM energy_by_date) AS
all_days_total_energy,
    ROUND(SUM(total_energy) OVER(ORDER BY DATE) * 100 / (SELECT
SUM(total_energy) FROM energy_by_date), 0) AS percentage_of_total_energy --cum
/ all_days_total
FROM energy_by_date
```

80.Start Dates Of Top Drivers

Lyft Medium ID 10083

Find contract starting dates of the top 5 most paid Lyft drivers. Consider only drivers who are still working with Lyft.

Table: lyft_drivers

```
--contract starting dates of top 5 most paid drivers
```

```
--only drivers still work with lyft: end_date IS NULL: end_date是空值

--output: start_date

SELECT start_date
FROM
    (SELECT
        *,
        RANK() OVER (
            ORDER BY yearly_salary DESC) AS rnk
        FROM lyft_drivers
        WHERE end_date IS NULL) t
WHERE rnk <= 5
```

81. Find the number of employees who received the bonus and who didn't

Microsoft Dell Hard ID 10081

Find the number of employees who received the bonus and who didn't. Bonus values in employee table are corrupted so you should use values from the bonus table. Be aware of the fact that employee can receive more than bonus. Output value inside has_bonus column (1 if they had bonus, 0 if not) along with the corresponding number of employees for each.

Tables: employee, bonus

```
-- # employees receive the bonus and who did not
-- bonus values in employee table x; use bonus values from the bonus table
-- can receive bonus > 1
-- output: has_bonus: has_bonus (1,0) , n_employees: # employess

SELECT
CASE WHEN b.bonus_amount IS NOT NULL THEN 1
      ELSE 0
      END AS has_bonus,
COUNT(DISTINCT e.id)
FROM
    employee e
LEFT JOIN
    bonus b
ON
    b.worker_ref_id = e.id
GROUP BY 1
```

82. Find matching hosts and guests in a way that they are both of the same gender and nationality

Airbnb Medium ID 10078

Find matching hosts and guests pairs in a way that they are both of the same gender and nationality. Output the host id and the guest id of matched pair.

DISTINCT 去除ID PAIR中重复的记录

```
SELECT
    DISTINCT
    h.host_id,
    g.guest_id
FROM airbnb_hosts h
JOIN airbnb_guests g
ON h.nationality = g.nationality AND h.gender = g.gender
```

83.Income By Title and Gender

City of San Francisco Medium ID 10077

Find the average total compensation based on employee titles and gender. Total compensation is calculated by adding both the salary and bonus of each employee. However, not every employee receives a bonus so disregard employees without bonuses in your calculation. Employee can receive more than one bonus. Output the employee title, gender (i.e., sex), along with the average total compensation.

Tables: sf_employee, sf_bonus

```
--存在同一个人同一个salary对应多个bonus的情况:
--需要先将第二个表, 将一个worker_id, 对应一个total_bonus
--再将两表join起来, salary + total_bonus作为total_comp, 再avg()

with cte as(
select
    worker_ref_id,
    sum(bonus) as total_bonus
from sf_bonus
group by 1
)
select
employee_title,
sex,
avg(a.salary + b.total_bonus) as avg_total_com
from sf_employee a join cte b
on a.id = b.worker_ref_id
group by 1, 2
```

84.Find the average age of guests reviewed by each host

Airbnb Medium ID 10074

Find the average age of guests reviewed by each host. Output the user along with the average age.

Tables: airbnb_reviews, airbnb_guests

```
-- avg age of guests reviewed by each host
--output: from_user: HOST, average_age: guests

--1.题意是每个host 去评价的 客户的平均年龄: from_type = 'host'; from_user AS
host_id; avg(age) AS avg_guest_age

--2.在应该保留全部记录的题目中, 不要使用distinct。
--一个host可以给同一个guest多次打分, 而且存在多次打分中仍然打一样的分。
--这样导致如果用total使用distinct, 则有不该被省略的记录被省略
```

```
WITH total AS(
SELECT
    *
FROM airbnb_reviews r
JOIN airbnb_guests g
ON r.to_user = g.guest_id
WHERE from_type = 'host'
ORDER BY 1,2
)
SELECT
    from_user,
    AVG(age) AS average_age
FROM total
GROUP BY 1
ORDER BY 1
```

85. Favorite Host Nationality

Interview Question Date: April 2020

Airbnb Medium ID 10073

For each guest reviewer, find the nationality of the reviewer's favorite host based on the guest's highest review score given to a host. Output the user ID of the guest along with their favorite host's nationality. In case there is more than one favorite host from the same country, list that country only once (remove duplicates).

Both the from_user and to_user columns are user IDs.

Tables: airbnb_reviews, airbnb_hosts

稳妥的方法:

i. join两表开新表total时, 保留全部记录*, rnk
ii. 在计算时 取出所需变量 视情况看是否需要distinct

```
--each guest reviewer ~ nationality of the reviewer's favorite host ~ guest's
highest review score

--output: from_user: user ID of the guest; nationality: favorite host

--in case > 1 favorite host from the same country, list country only
once(remove duplicate)

WITH total AS(
```

```

SELECT
DISTINCT
    *,
    -- from_user,
    -- to_user,
    -- nationality,
    -- review_score,
    RANK() OVER (PARTITION BY from_user ORDER BY review_score DESC) as rnk
FROM airbnb_reviews r
JOIN airbnb_hosts h
ON r.to_user = h.host_id
WHERE from_type = 'guest'
ORDER BY 1
)
SELECT
    DISTINCT
    from_user,
    nationality
FROM total
WHERE rnk = 1

```

86. Guest Or Host Kindness

Interview Question Date: April 2020

Airbnb Easy ID 10072

Find whether hosts or guests give higher review scores based on their average review scores. Output the higher of the average review score rounded to the 2nd decimal spot (e.g., 5.11).

Table: airbnb_reviews

order by中可以有新产生的变量名

--是H or G 给了高分, 依据他们的平均分 ==> 留一个最高平均分的记录

```

SELECT from_type,
       ROUND(AVG(review_score), 2) winner
FROM airbnb_reviews
GROUP BY from_type
ORDER BY winner DESC
LIMIT 1

```

```

SELECT
    from_type,
    ROUND(av,2)
FROM(
    SELECT
        from_type,
        AVG(review_score) AS av,
        DENSE_RANK() OVER (ORDER BY AVG(review_score) DESC) AS rank
    FROM airbnb_reviews
    GROUP BY from_type) m

```

```
WHERE rank = 1
```

87.Hosts' Abroad Apartments

Airbnb Medium ID 10071

Find the number of hosts that have accommodations in countries of which they are not citizens.

Tables: airbnb_hosts, airbnb_apartments

```
--# of hosts that have accomadations in countries of which host are not
citizens
--output: count

SELECT
    COUNT(DISTINCT a.host_id)
FROM airbnb_hosts h
JOIN airbnb_apartments a
ON h.host_id = a.host_id
WHERE h.nationality != a.country
```

88.DeepMind employment competition

Google Medium ID 10070 14

Data Engineer Data Scientist BI Analyst Data Analyst Find the winning teams of DeepMind employment competition. Output the team along with the average team score. Sort records by the team score in descending order.

Tables: google_competition_participants, google_competition_scores

```
sum(b.member_score) / count(a.member_id) = avg(b.member_score)
```

```
SELECT
    team_id,
    AVG(member_score) AS team_score
FROM google_competition_participants p
JOIN google_competition_scores s
ON p.member_id = s.member_id
GROUP BY 1
ORDER BY 2 DESC
```

89.Correlation Between E-mails And Activity Time

Interview Question Date: April 2020

Google Hard ID 10069

There are two tables with user activities. The `google_gmail_emails` table contains information about emails being sent to users. Each row in that table represents a message with a unique identifier in the `id` field. The `google_fit_location` table contains user activity logs from the Google Fit app. Find the correlation between the number of emails received and the total exercise per day. The total exercise per day is calculated by counting the number of user sessions per day.

i. 每个用户每天收#邮件: 取to_user, day, count(*), 从google_gmail_emails
 ii. 每个用户每天exercise: 取user_id, day, count(distinct session_id)
 【为什么distinct?】: 在google_fit_location中, 因为step_id不同, 会出现同一user_id同一天同一session_id的不同记录, 所以要distinct(session_id)
 iii. join两个表: email和session, on两个条件day and user_id 【为什么也on user_id?】
 vi. COALESCE (A,0) : A值缺少, 则取0
 v. CORR(X,Y): XY两变量间correlation

```
-- google_gmail_emails: emails sent to users; each row ~ unique identifier ~ a message
-- google_fit_location: user activity logs from google fit app

--relationship between # emails received and total exercise per day:counting the number of user sessions per day
--corr

--写法1: CTE
WITH SESSIONS AS(
SELECT
    user_id,
    user_id,
    day,
    COUNT(DISTINCT session_id) AS session
FROM google_fit_location
GROUP BY 1,2
),
EMAIL AS(
SELECT
    to_user,
    user_id,
    day,
    COUNT(id) AS email_received
FROM google_gmail_emails
GROUP BY 1,2
)
SELECT
    CORR(COALESCE(email_received,0), COALESCE(session,0)) AS corr
FROM EMAIL e
FULL JOIN SESSIONS s
ON e.day = s.day AND user_id = to_user
```

写法2: subquery

```
SELECT
    corr(COALESCE(n_emails:: NUMERIC, 0), COALESCE(total_exercise, 0))
FROM
```

```

(
    SELECT
    用户每天收#邮件
        to_user,
        day,
        COUNT(*) AS n_emails
    FROM google_gmail_emails
    GROUP BY 1,2) mail_base
FULL OUTER JOIN
(
    用户每天exercise
    SELECT
        user_id,
        day,
        COUNT(DISTINCT session_id) AS total_exercise
    FROM google_fit_location
    GROUP BY 1,2
    ) total_exercise
ON mail_base.to_user = total_exercise.user_id
AND mail_base.day = total_exercise.day
to_user = user_id AND day = day
JOIN google_fit_location f
ON g.to_user = f.user_id
GROUP BY 1
--每1
--每1
--ON

```

90. User Email Labels

Interview Question Date: April 2020

Google Medium ID 10068

Find the number of emails received by each user under each built-in email label. The email labels are: 'Promotion', 'Social', and 'Shopping'. Output the user along with the number of promotion, social, and shopping mails count,.

Tables: google_gmail_emails, google_gmail_labels Hints Expected Output All required columns and the first 5 rows of the solution are shown

```

COUNT(case when A then 1 else null end) = SUM(case when A then 1 else 0 end)
case when 后的end别忘记

```

```

-- # emails received by each user ~ each built-in email label
-- labels: 'Promotion', 'Social', and 'Shopping'.
-- output: user ~ # of p, s, s mails count
-- output: to_user, promotion_count, social_count, shopping_count

SELECT
    to_user,
    COUNT(CASE WHEN label = 'Promotion' THEN id END) AS promotion_count,
    COUNT(CASE WHEN label = 'Social' THEN id END) AS social_count,
    COUNT(CASE WHEN label = 'Shopping' THEN id END) AS shopping_count
FROM google_gmail_emails e
JOIN google_gmail_labels l
ON e.id = l.email_id

```

GROUP BY 1

91. Google Fit User Tracking

92. Fans vs Opposition

Interview Question Date: March 2020

Meta/Facebook Hard ID 10062

Meta/Facebook is quite keen on pushing their new programming language Hack to all their offices. They ran a survey to quantify the popularity of the language and send it to their employees. To promote Hack they have decided to pair developers which love Hack with the ones who hate it so the fans can convert the opposition. Their pair criteria is to match the biggest fan with biggest opposition, second biggest fan with second biggest opposition, and so on. Write a query which returns this pairing. Output employee ids of paired employees. Sort users with the same popularity value by id in ascending order.

Duplicates in pairings can be left in the solution. For example, (2, 3) and (3, 2) should both be in the solution.

Table: facebook_hack_survey

求一高一低PAIR ID:

用row_number保持唯一性

法1:

i.love_rnk, hate_rnk

ii. ID PAIR: id1, id2 + self join on love_rnk = hate_rnk

法2:

i. 分别建立两个表: 各取id, love/hate rnk

ii. ID PAIR: id1, id2 + join on love_rnk = hate_rnk

法1:

i.love_rnk, hate_rnk

ii. ID PAIR: id1, id2 + self join on love_rnk = hate_rnk

WITH total AS(

SELECT

*,

RANK() OVER (ORDER BY popularity DESC, employee_id ASC) AS love_rnk,

RANK() OVER (ORDER BY popularity, employee_id ASC) AS hate_rnk

FROM facebook_hack_survey

)

SELECT

t1.employee_id,

t2.employee_id

FROM total t1

JOIN total t2

ON t1.love_rnk = t2.hate_rnk

法2:

i. 分别建立两个表: 各取id, love/hate rnk

ii.ID PAIR: id1, id2 + join on love_rnk = hate_rnk

```

WITH love_rnk AS(
    SELECT
        employee_id,
        ROW_NUMBER() OVER (ORDER BY popularity DESC, employee_id ASC) AS
love_rnk
    FROM facebook_hack_survey),
hate_rnk AS(
    SELECT
        employee_id,
        ROW_NUMBER() OVER (ORDER BY popularity, employee_id ASC) AS hate_rnk
    FROM facebook_hack_survey
)
SELECT
    l.employee_id AS employee_fan_id,
    h.employee_id AS employee_opposition_id
FROM love_rnk l
JOIN hate_rnk h
ON l.love_rnk = h.hate_rnk

```

93. Find the number of reviews received by Lo-Lo's Chicken & Waffles for each star

Yelp Easy ID 10058

Find the number of reviews received by Lo-Lo's Chicken & Waffles for each star. Output the number of stars along with the corresponding number of reviews. Sort records by stars in ascending order.

Table: yelp_reviews

```

--ILIKE 'Lo-Lo_s Chicken & Waffles'
--like 'Lo-Lo%s Chicken & Waffles'
--如果名字中有撇号，可以把'写成或者% 注意LIKE

```

```

--# reviews received for each star
--output:stars, n_reviews: # stars and # reviews
--order by stars asc

SELECT
    stars,
    COUNT(*) AS n_reviews
FROM yelp_reviews
WHERE business_name ILIKE 'Lo-Lo_s Chicken & Waffles'
GROUP BY 1
ORDER BY 1

```

94.Popularity of Hack

Interview Question Date: March 2020

Meta/Facebook Easy ID 10061

Meta/Facebook has developed a new programming language called Hack. To measure the popularity of Hack they ran a survey with their employees. The survey included data on previous programming familiarity as well as the number of years of experience, age, gender and most importantly satisfaction with Hack. Due to an error location data was not collected, but your supervisor demands a report showing average popularity of Hack by office location. Luckily the user IDs of employees completing the surveys were stored. Based on the above, find the average popularity of the Hack per office location. Output the location along with the average popularity.

95. Popularity of Hack

Interview Question Date: March 2020

Meta/Facebook Easy ID 10061

Meta/Facebook has developed a new programming language called Hack. To measure the popularity of Hack they ran a survey with their employees. The survey included data on previous programming familiarity as well as the number of years of experience, age, gender and most importantly satisfaction with Hack. Due to an error location data was not collected, but your supervisor demands a report showing average popularity of Hack by office location. Luckily the user IDs of employees completing the surveys were stored. Based on the above, find the average popularity of the Hack per office location. Output the location along with the average popularity.

Tables: facebook_employees, facebook_hack_survey

```
--error location data not collect
--avg pop of hack by office
--avg pop of the hack per office location
--output: location, avg_popularity
```

```
SELECT
    location,
    AVG(popularity) AS avg_popularity
FROM facebook_employees e
JOIN facebook_hack_survey s
ON e.id = s.employee_id
GROUP BY 1
```

Top Cool Votes

Interview Question Date: March 2020

Yelp Medium

Find the review_text that received the highest number of 'cool' votes. Output the business name along with the review text with the highest number of 'cool' votes.

Table: yelp_reviews

```
--review_text receive highest # cool votes
--output: business_name, review_text
```

```
--法1: rank cte做法:
```

```
WITH rnk AS(
SELECT
    business_name,
    review_text,
    cool,
    RANK() OVER (ORDER BY cool DESC) AS rnk
FROM yelp_reviews
)
SELECT
    business_name,
    review_text
FROM rnk
WHERE rnk = 1
```

```
--法2: max subquery做法:
```

```
SELECT
    business_name,
    review_text
FROM yelp_reviews
WHERE cool =
    (
    SELECT
    MAX(cool) AS max_cool
    FROM yelp_reviews
    )
```

Find the number of reviews received by Lo-Lo's Chicken & Waffles for each star

Yelp Easy ID 10058

Find the number of reviews received by Lo-Lo's Chicken & Waffles for each star. Output the number of stars along with the corresponding number of reviews. Sort records by stars in ascending order.

Table: yelp_reviews

```
ILIKE 'Lo-Lo's Chicken & Waffles'
like 'Lo-Lo%s Chicken & Waffles'
如果名字中有撇号，可以把撇号写成或者%来代替LIKE
```

```
--# reviews received for each star
--output:stars, n_reviews: # stars and # reviews
--order by stars asc
```

```
SELECT
    stars,
    COUNT(*) AS n_reviews
FROM yelp_reviews
```

```
WHERE business_name ILIKE 'Lo-Lo_s Chicken & Waffles'
GROUP BY 1
ORDER BY 1
```

Most Checkins

Yelp Medium ID 10053

Find the top 5 businesses with the most check-ins. Output the business id along with the number of check-ins.

Table: yelp_checkin

```
--top 5 biz with most check-ins
--output biz id with # check-ins
--output: business_id, n_checkins

--rank法求top 5:
WITH checkins_rank AS(
SELECT
    business_id,
    SUM(checkins) AS n_checkins,
    RANK() OVER (ORDER BY SUM(checkins) DESC) AS rnk
FROM yelp_checkin
GROUP BY 1
)
SELECT
    business_id,
    n_checkins
FROM checkins_rank
WHERE rnk <= 5
```

```
--limit法求top 5:
SELECT
    business_id,
    SUM(checkins) AS counts
FROM yelp_checkin
GROUP BY 1
ORDER BY 2 DESC
LIMIT 5
```

Find the average number of stars for each state

Yelp Easy ID 10052

Find the average number of stars for each state. Output the state name along with the corresponding average number of stars.

Table: yelp_business

```
--select * from yelp_business;
```

```
SELECT
    state,
    AVG(stars) AS average_stars
FROM yelp_business
GROUP BY 1
```

Find the number of open businesses

Yelp Easy ID 10051

Find the number of open businesses.

Table: yelp_business

```
--select * from yelp_business;

SELECT
    COUNT(*) AS business_open
FROM yelp_business
WHERE is_open = 1
```

Find the review count for one-star businesses from yelp

Yelp Easy ID 10050

Find the review count for one-star businesses from yelp. Output the name along with the corresponding review count.

Table: yelp_business

```
--select * from yelp_business;

SELECT
    name,
    review_count
FROM yelp_business
WHERE stars = 1
```

Reviews of Categories

Interview Question Date: March 2020

Yelp Medium ID 10049

Find the top business categories based on the total number of reviews. Output the category along with the total number of reviews. Order by total reviews in descending order.

Table: yelp_business

```
--top biz categories ~ total # reviews
--output: category, review_cnt:
```

```
--order by total reviews desc

--UNNEST(STRING_TO_ARRAY(A, ';')): 把A中用; 分隔开的每个词 分解取出

SELECT
    UNNEST(STRING_TO_ARRAY(categories, ';')) AS category,
    SUM(review_count) AS review_cnt
FROM yelp_business
GROUP BY 1
ORDER BY 2 DESC
```

Top Businesses With Most Reviews

Yelp Medium ID 10048

Find the top 5 businesses with most reviews. Assume that each row has a unique business_id such that the total reviews for each business is listed on each row. Output the business name along with the total number of reviews and order your results by the total reviews in descending order.

Table: yelp_business

```
In [ ]: --top 5 biz with most reviews
--unique biz_id ~ total reviews for each biz
--output: name, review_count: business name along with the total number of re
--order by total reviews desc

--LIMIT 5 不好, rank 求top 5才最精准:
WITH business_rank AS(
SELECT
    name,
    review_count,
    RANK() OVER (ORDER BY review_count DESC) AS rnk
FROM yelp_business
)
SELECT
    name,
    review_count
FROM business_rank
WHERE rnk <= 5
```

Top 5 States With 5 Star Businesses

Interview Question Date: March 2020

Yelp Hard ID 10046

Find the top 5 states with the most 5 star businesses. Output the state name along with the number of 5-star businesses and order records by the number of 5-star businesses in descending order. In case there are ties in the number of businesses, return all the unique states. If two states have the same result, sort them in alphabetical order.

Table: yelp_business

```

-- top 5 states with most 5 star biz
-- output: state, n_businesses:
-- order by # 5 star biz in desc order, sort state asc
-- in case ties in # biz, return all unique states

--求top 5, 用limit 5会不精准-因为会漏项, 还是rank精准

WITH biz_rank AS(
SELECT
    state,
    COUNT(business_id) AS n_businesses,
    RANK() OVER (ORDER BY COUNT(business_id) DESC) AS rnk
FROM yelp_business
WHERE stars = 5
GROUP BY 1
ORDER BY 2 DESC, 1 ASC
)
SELECT
    state,
    n_businesses
FROM biz_rank
WHERE rnk <= 5

```

Highest Priced Wine In The US

Interview Question Date: March 2020

Wine Magazine Medium ID 10044 13

Data Engineer Data Scientist BI Analyst Data Analyst Find the highest price in US country for each variety produced in English speaking regions, but not in Spanish speaking regions, with taking into consideration varieties that have earned a minimum of 90 points for every country they're produced in. Output both the variety and the corresponding highest price.

Let's assume the US is the only English speaking region in the dataset, and Spain, Argentina are the only Spanish speaking regions in the dataset. Let's also assume that the same variety might be listed under several countries so you'll need to remove varieties that show up in both the US and in Spanish speaking countries.

Table: winemag_p1

```

--highest price in US ~ each variety produced in English regions but not
spanish regions
--√ varieties earned mini 90 for every country produced
--output: variety, max: variety and highest price
--US ~ only English speaking region; Spain, Argentina ~ Spanish only
--remove varieties show up in both US & SPain

--问题: 一个variety种类, 可以有多个points; 一个variety种类, 可以有多个价格

WITH remove AS(

```

```

SELECT
    variety
FROM winemag_p1
WHERE country IN ('Spain', 'Argentina')
),
residue AS(
SELECT
    variety
FROM winemag_p1
WHERE variety NOT IN (SELECT variety FROM remove)
GROUP BY 1
HAVING MIN(points) >= 90
)
SELECT
    variety,
    MAX(price)
FROM winemag_p1
WHERE variety IN (SELECT variety FROM residue)
AND country = 'US'
GROUP BY 1

```

Median Price Of Wines

Interview Question Date: March 2020

Wine Magazine Hard ID 10043

Find the median price for each wine variety across both datasets. Output distinct varieties along with the corresponding median price.

Tables: winemag_p1, winemag_p2

①当两个dataset列数不一致时，可以选出需要的列进行union操作；如果列数一直，可以直接select * 进行union操作
 ②此处应该UNION ALL而不是UNION，选择两个dataset的全部
 ③median:
 i. percentile_cont(0.5) within GROUP (ORDER BY A) AS median_price
 ii. avg(A) + WHERE row_number_asc IN (row_number_desc, row_number_desc - 1, row_number_desc + 1);
 如果是奇数个记录，则选两个row_num相等的是中位数(eg. n = 5, median = asc = desc = 3)的记录选出来求平均值，即这个记录的值
 如果是偶数个记录，则选asc = desc - 1 和asc = desc + 1的两个记录选出来求平均值

```

--法1: percentile_cont(0.5) within GROUP (ORDER BY A)
WITH total AS(
    SELECT
        variety,
        price::float
    FROM winemag_p1
    UNION ALL
    SELECT
        variety,
        price::float
    FROM winemag_p2
)
SELECT

```

```

        variety,
        percentile_cont(0.5) within GROUP (ORDER BY price) AS median_price
FROM total
GROUP BY 1

```

```

--法2: row_number: avg(A) + WHERE row_number_asc IN (row_number_desc,
row_number_desc - 1, row_number_desc + 1 )
WITH total AS(
    SELECT
        variety,
        price
    FROM winemag_p1
    WHERE price IS NOT NULL
    UNION ALL
    SELECT
        variety,
        price
    FROM winemag_p2
    WHERE price IS NOT NULL
),
w1 AS(
    SELECT
        variety,
        price,
        ROW_NUMBER() OVER (PARTITION BY variety ORDER BY price ASC) AS rnk_asc,
        ROW_NUMBER() OVER (PARTITION BY variety ORDER BY price DESC) AS rnk_desc
    FROM total
)
SELECT
    variety,
    AVG(price)
FROM w1
WHERE rnk_asc IN (rnk_desc, rnk_desc - 1, rnk_desc + 1)
GROUP BY 1

```

Top 3 Wineries In The World

Interview Question Date: March 2020

Wine Magazine Hard ID 10042

Find the top 3 wineries in each country based on the average points earned. In case there is a tie, order the wineries by winery name in ascending order. Output the country along with the best, second best, and third best wineries. If there is no second winery (NULL value) output 'No second winery' and if there is no third winery output 'No third winery'. For outputting wineries format them like this: "winery (avg_points)"

Table: winemag_p1

```

①连接括号写成(C)的形式: CASE WHEN A THEN B || ' (' || C || ')' ELSE NULL END AS D
②ROW_NUMBER() OVER (PARTITION BY A ORDER BY B, C) AS row_num
③max(A) 取一个值

```


④如果没有A, 就写成B: COALESCE(max(A), 'B')

```
-- top 3 wineries in each country ~ avg points earned
-- order wineries by winery name in sac
-- output:country, top_winery, second_winery, third_winery: winery (avg_points
-- If there is no second winery (NULL value) output 'No second winery' and if
there is no third winery output 'No third winery'

SELECT country,
       max(top_winery) AS top_winery,
       COALESCE(max(second_winery), 'No second winery') AS second_winery,
       COALESCE(max(third_winery), 'No third winery') AS third_winery
FROM
  (SELECT country,
         CASE
           WHEN POSITION = 1 THEN winery || ' (' || round(avg_points) ||
' )'
           ELSE NULL
         END AS top_winery,
         CASE
           WHEN POSITION = 2 THEN winery || ' (' || round(avg_points) ||
' )'
           ELSE NULL
         END AS second_winery,
         CASE
           WHEN POSITION = 3 THEN winery || ' (' || round(avg_points) ||
' )'
           ELSE NULL
         END AS third_winery
   FROM
     (SELECT country,
            winery,
            ROW_NUMBER() OVER (PARTITION BY country
                               ORDER BY avg_points DESC, winery ASC) AS
POSITION,
            avg_points
      FROM
        (SELECT country,
               winery,
               avg(points) AS avg_points
         FROM winemag_p1
         WHERE country IS NOT NULL
         GROUP BY country,
                winery) tmp1) tmp2
   WHERE POSITION <= 3) tmp3
GROUP BY country
```

Most Expensive And Cheapest Wine

Interview Question Date: March 2020

Wine Magazine Hard ID 10041

Find the cheapest and the most expensive variety in each region. Output the region along with the corresponding most expensive and the cheapest variety. Be aware that there are 2 region columns, the price from that row applies to both of them.

①total: SELECT * FROM A UNION ALL SELECT * FROM B
 ②计算A中, C在B最大/小时的值:
 cte1: A, ROW_NUMBER() OVER (ORDER BY B DESC) AS rnk1, ROW_NUMBER() OVER (ORDER BY B ASC) AS rnk2,
 cte2: A, MAX(CASE WHEN rnk1 = 1 THEN C END) AS max_price_B; MIN(CASE WHEN rnk1 = 1 THEN C END) AS min_price_B + GROUP BY A
 ③当遇到case when创建后, 本应放在一行的记录, 放成了几行, 中间可能有空格。解决方法是A, MAX() + group by 1 可以将其并为一行

```
WITH total AS(
    SELECT
        region_1 AS region,
        price,
        variety
    FROM winemag_p1
    WHERE region_1 IS NOT NULL
    UNION ALL
    SELECT
        region_2 AS region,
        price,
        variety
    FROM winemag_p1
    WHERE region_2 IS NOT NULL
),
result AS (
    SELECT
        *,
        ROW_NUMBER() OVER (PARTITION BY region ORDER BY price DESC) AS
        expensive_rank,
        ROW_NUMBER() OVER (PARTITION BY region ORDER BY price ASC) AS
        cheap_rank
    FROM total
    WHERE price IS NOT NULL
    AND region IS NOT NULL
    ORDER BY region, price DESC
)
SELECT
    region,
    MAX(CASE
        WHEN expensive_rank = 1 THEN variety
        ELSE NULL
    END) AS most_expensive_variety,
    MIN(CASE
        WHEN cheap_rank = 1 THEN variety
        ELSE NULL
    END) AS cheapest_variety
FROM result
GROUP BY 1
```

Find all wines from the winemag_p2 dataset which are produced in countries that have the highest sum of points in the winemag_p1 dataset

Wine Magazine Hard ID 10040

Find all wines from the winemag_p2 dataset which are produced in the country that have the highest sum of points in the winemag_p1 dataset.

Tables: winemag_p1, winemag_p2

```
--all wines from winemag_p2 produced in the country ~ highest sum of points in 1_p1

--select * from winemag_p1;
WITH sum_points_country AS(
SELECT
    country,
    SUM(points) AS sum_points,
    RANK() OVER (ORDER BY SUM(points) DESC) AS rnk
FROM winemag_p1
GROUP BY country
)
SELECT
    *
FROM winemag_p2
WHERE country IN (SELECT country FROM sum_points_country WHERE rnk = 1)
```

Macedonian Vintages

Wine Magazine Medium ID 10039

Find the vintage years of all wines from the country of Macedonia. The year can be found in the 'title' column. Output the wine (i.e., the 'title') along with the year. The year should be a numeric or int data type.

Table: winemag_p2

```
--select * from winemag_p2;

--vintage years of all wines from Macedonia
--year ~ title
--output: title, year

SELECT
    title,
    substring(title, '[0-9]{4}'):: integer AS year
--    NULLIF(regex_replace(title, '\D','', 'g'), ''):: NUMERIC AS year
--    split_part(title, ' ', 2) :: NUMERIC AS year
FROM winemag_p2
```

```
WHERE country = 'Macedonia'
```

Find all provinces which produced more wines in 'winemag_p1' than they did in 'winemag_p2'

Interview Question Date: March 2020

Wine Magazine Medium ID 10038

Find all provinces which produced more wines in 'winemag_p1' than they did in 'winemag_p2'. Output the province and the corresponding wine count. Order records by the wine count in descending order.

Tables: winemag_p1, winemag_p2

```
-- all provinces produce more wines in p1 > p2
-- output: province, cnt_1
-- order by cnt desc

SELECT
    a.province,
    a.cnt_1
FROM
    (SELECT
        province,
        COUNT(winery) AS cnt_1      --COUNT(*)也可
    FROM winemag_p1
    GROUP BY province) a
    JOIN
    (SELECT
        province,
        COUNT(winery) AS cnt_2
    FROM winemag_p2
    GROUP BY province) b
ON a.province = b.province
WHERE cnt_1 - cnt_2 > 0            --ON两个条件也可
ORDER BY cnt_1 DESC
```

Find Favourite Wine Variety

Interview Question Date: March 2020

Wine Magazine Hard ID 10037

Find each taster's favorite wine variety. Consider that favorite variety means the variety that has been tasted by most of the time. Output the taster's name along with the wine variety.

Table: winemag_p2

```
--each taster's fa wine variety
--tasted by most of the time
```

```
--output: taster_name, variety:
-- select * from winemag_p2
-- WHERE taster_name = 'Joe Czerwinski'
--AND variety = 'Champagne Blend'

--#算favarite 是算一个taster_name, variety 有多少个记录: count(*)
--#修好: 有一个没有年份的漏算记录; 有重复年份的按多个记录算

WITH total AS(
SELECT
    taster_name,
    variety,
    substring(title, '[0-9]{4}')::integer AS year
FROM winemag_p2
WHERE taster_name IS NOT NULL
--AND substring(title, '[0-9]{4}') IS NOT NULL
),
year_rank AS(
SELECT
    taster_name,
    variety,
    COUNT(*) AS years,
    RANK() OVER (PARTITION BY taster_name ORDER BY COUNT(*) DESC) AS rnk
FROM total
GROUP BY 1, 2
)
SELECT
    taster_name,
    variety
FROM year_rank
WHERE rnk = 1
```

Find the number of wines with and without designations per country

Interview Question Date: March 2020

Wine Magazine Medium ID 10035

Find the number of wines with and without designations per country. Output the country along with the total without designations, total with designations, and the final total of both.

Table: winemag_p2

```
--# wines with and without designations per country
--output: country, total_without_designation, total_with_designation,
grand_total

--该值为空 IS NULL
--当计算该值有多少个时: SUM(CASE WHEN A THEN 1 ELSE 0 END) ; 不能将SUM改为COUNT,
因为如用COUNT(CASE WHEN A THEN 1 ELSE 0 END)代表共计多少个 = COUNT(*)

SELECT
    country,
```

```

SUM(CASE WHEN designation IS NULL THEN 1 ELSE 0 END) AS
total_without_designation,
SUM(CASE WHEN designation IS NOT NULL THEN 1 ELSE 0 END) AS
total_with_designation,
COUNT(*) AS grand_total
FROM winemag_p2
GROUP BY country

```

```

--# wines with and without designations per country
--output: country, total_without_designation, total_with_designation,
grand_total

--该值为空 IS NULL
--当计算该值有多少个时: SUM(CASE WHEN A THEN 1 ELSE 0 END) ; 不能将SUM改为COUNT,
因为如用COUNT(CASE WHEN A THEN 1 ELSE 0 END)代表共计多少个 = COUNT(*)

SELECT
    country,
    SUM(CASE WHEN designation IS NULL THEN 1 ELSE 0 END) AS
total_without_designation,
    SUM(CASE WHEN designation IS NOT NULL THEN 1 ELSE 0 END) AS
total_with_designation,
    COUNT(*) AS grand_total
FROM winemag_p2
GROUP BY country

```

Wine Variety Revenues

Interview Question Date: February 2020

Wine Magazine Medium ID 10033

Find the total revenue made by each region from each variety of wine in that region. Output the region, variety, and total revenue.

Take into calculation both region_1 and region_2. Remove the duplicated rows where region, price and variety are exactly the same.

Table: winemag_p1

```

-- total rev made by each rgn ~ each variety of wine in region
-- output: region, variety, sum: total revenue
-- remove duplicated rows where region, price, variety are the same

--一个区域内, 一个variety, 可以对应多个price, 把所有price相加一起, 就是revenue

WITH total AS(
    (SELECT
        region_1 AS region,
        price,
        variety
    FROM
        winemag_p1
    )

```

```

        UNION
        (SELECT
            region_2 AS region,
            price,
            variety
        FROM
            winemag_p1
        )
    )
    SELECT
        DISTINCT
        region,
        variety,
        SUM(price) AS sum
    FROM total
    WHERE region IS NOT NULL AND price IS NOT NULL      -- 此句不可少
    GROUP BY region, variety
    ORDER BY SUM(price) DESC

```

Best Wines By Points-To-Price

Interview Question Date: February 2020

Wine Magazine Medium ID 10032

Find the wine with the highest points to price ratio. Output the title, points, price, and the corresponding points-to-price ratio.

Table: winemag_p2

```

--①如果项不是数字格式，要换为数字： ::NUMERIC
--②找最大值的两种写法：
-- i. rank = 1
-- ii. A = SELECT MAX(A) FROM B

```

```

WITH ratio AS(
    SELECT
        title,
        points,
        price,
        points:: NUMERIC / price:: NUMERIC AS points_price_ratio,
        RANK() OVER (ORDER BY points / price DESC) AS rnk
    FROM winemag_p2
    WHERE points / price IS NOT NULL
)
SELECT
    title,
    points,
    price,
    points_price_ratio
FROM ratio
--WHERE rnk = 1
WHERE points_price_ratio = (SELECT max(points_price_ratio) FROM ratio)

```

Find the number of Bodegas outside of Spain that produce wines with the blackberry taste

Interview Question Date: February 2020

Wine Magazine Medium ID 10031

Find the number of Bodegas (wineries with "bodega" pattern inside the name) outside of Spain that produce wines with the blackberry taste (description contains blackberry string). Group the count by country and region. Output the country, region along with the number of bodegas. Order records by the number of bodegas in descending order.

Table: winemag_p1



①一个地区可能有多个相同名称的winery, 此时应该用count(distinct winery) 去重:
②将region_1和region_2合并union格式, 并进行去重:

```
SELECT
    *
FROM
    (SELECT a FROM A
    UNION
    SELECT b FROM B
    WHERE a != b) t
```

```
--# of B(wineries) outside of Spain that produce wines with balackberry
taste(description)
--group the count by country and region
--output: country, region, n_bodegas
```

```
SELECT
    country,
    region,
    COUNT(DISTINCT winery) AS n_bodegas
FROM
```

```
    (SELECT
        country,
        region_1 AS region,
        winery,
        description
    FROM winemag_p1
    UNION ALL
    SELECT
```

```
        country,
        region_2 AS region,
        winery,
        description
    FROM winemag_p1
    WHERE region_1 != region_2) t
```

```
WHERE winery ILIKE '%bodega%'
AND country != 'Spain'
AND description ILIKE '%blackberry%'
GROUP BY 1,2
```

--去除掉Region为空的部分, 以及两个
region名称相同的部分

ORDER BY 3 DESC

Price Of Wines In Each Country

Interview Question Date: February 2020

Wine Magazine Medium ID 10029

Find the minimum, average, and maximum price of all wines per country. Assume all wines listed across both datasets are unique. Output the country name along with the corresponding minimum, maximum, and average prices.

Tables: winemag_p1, winemag_p2

①看清是否需要合并两个数据集，如果数据集的记录没有重叠，无需去重，直接使用union all；如果有需要去重的需要使用union
②转换格式：::NUMERIC, ::TEXT
③求最大平均最小值：min, max, avg

```
--min, avg, max price of all wines per country
--all wines unique
--output: country, min_price, avg_price, max_price

SELECT
    country,
    MIN(price) AS min_price,
    AVG(price) AS avg_price,
    MAX(price) AS max_price
FROM
    (SELECT
        country:: TEXT,
        price:: NUMERIC
        FROM winemag_p1
        UNION ALL
        SELECT
            country:: TEXT,
            price:: NUMERIC
            FROM winemag_p2) t
GROUP BY 1
ORDER BY 1
```

Find the number of wines each taster tasted within the variation

Interview Question Date: February 2020

Wine Magazine Medium ID 10028

Find the number of wines each taster tasted within the variation. Output the tester's name, variety, and the number of tastings. Order records by taster name and the variety in ascending order and by the number of tasting in descending order.

```
--# wines ~ taster ~ winthin variation
--output: taster_name, variety, n_tastings
--order by taster_name, variety ASCC, n_tastings DESC

SELECT
    taster_name,
    variety,
    COUNT(*) AS n_tastings
FROM winemag_p2
WHERE taster_name IS NOT NULL
GROUP BY
    taster_name,
    variety
ORDER BY
    taster_name, variety ASC,
    COUNT(*) DESC
```

Find the number of US-based wineries that have expensive wines (price >= 200)

Wine Magazine Easy ID 10027

Find the number of US-based wineries that have expensive wines. A wine is considered to be expensive if its price is \$200 or more.

Table: winemag_p1

当记录中，有重复名字的A，使用COUNT(*)无法去重，应该使用COUNT(DISTINCT A)

```
--# US based wineries ~ expensive wines
--wine expensive if price >= 200
--output: n_wineries

SELECT
    COUNT(DISTINCT winery) AS n_wineries
FROM winemag_p1
WHERE price >= 200
AND country = 'US'
```

Find all wineries which produce wines by possessing aromas of plum, cherry, rose, or hazelnut

Wine Magazine Medium ID 10026

Find all wineries which produce wines by possessing aromas of plum, cherry, rose, or hazelnut. To make it more simple, look only for singular form of the mentioned aromas. HINT: if one of the specified words is just a substring of another word, this should not be a hit, but a miss.

Example Description: Hot, tannic and simple, with cherry jam and currant flavors accompanied by high, tart acidity and chile-pepper alcohol heat. Therefore the winery Bella Piazza is expected in the results.

`lower(a) ~ '\y(A|B|C|D)\y'`: a中仅包含ABCD的记录, 不包含ABCD的衍生词
 因为description中包含的不是关键词, 而是含有关键词的词, 应该被去除: eg. prosecco (rose), plump (plum), plummy (plum)
 用ILIKE无法实现去除功能, `'\y(A|B|C|D)\y'`: 确保我们匹配的词汇是一个完整的单词, 不是部分匹配

```
SELECT DISTINCT winery
FROM winemag_p1
WHERE lower(description) ~ '\y(plum|cherry|rose|hazelnut)\y'
```

--ILIKE错误方法, 无法排除 specified words is just a substring of another word的情况: Carpena Malvolti, Finca El Origen, La Fiammenga

```
-- SELECT
--     DISTINCT winery
--     -- description
-- FROM winemag_p1
-- WHERE description LIKE '%plum%'
-- or description LIKE '%cherry%'
-- or description LIKE '%rose%'
-- or description LIKE '%hazelnut%'
```

--错误原因:

```
-- Carpena Malvolti: Prosecco
-- Finca El Origen: plump
-- La Fiammenga: plummy
-- select * from winemag_p1
```

```
-- where winery IN ('Carpena Malvolti', 'Finca El Origen', 'La Fiammenga')
```

Find all possible varieties which occur in either of the winemag datasets

Wine Magazine Medium ID 10025

Find all possible varieties which occur in either of the winemag datasets. Output unique variety values only. Sort records based on the variety in ascending order.

Tables: winemag_p1, winemag_p2

```
--all possible varieties ~ in either winemag
--output: variety: unique variety values
--sort by variety asc
--将variety从每个表中选出, 再distinct, 效果等同于distinct variety再union
```

```
SELECT
    DISTINCT variety
FROM
    (SELECT
        variety
    FROM winemag_p1
    UNION ALL
```

```
SELECT
    variety
FROM winemag_p2) t
ORDER BY variety
```

Wine varieties tasted by 'Roger Voss'

Wine Magazine Easy ID 10024

Find wine varieties tasted by 'Roger Voss' and with a value in the 'region_1' column of the dataset. Output unique variety names only.

Table: winemag_p2

```
SELECT
    DISTINCT variety
FROM winemag_p2
WHERE taster_name = 'Roger Voss'
AND region_1 IS NOT NULL
```

Find all wine varieties which can be considered cheap based on the price

Wine Magazine Easy ID 10022

Find all wine varieties which can be considered cheap based on the price. A variety is considered cheap if the price of a bottle lies between 5 to 20 USD. Output unique variety names only.

Table: winemag_p1

```
--all wine varieties ~ cheap ~ price
--a vareity cheap: 5 < price < 20
--output: unique variety names
```

```
SELECT
    DISTINCT variety
FROM winemag_p1
WHERE price BETWEEN 5 AND 20
```

Find all top-rated wineries based on points

Wine Magazine Easy ID 10021

Find all top-rated wineries based on points. Consider a top-rated winery has been awarded points more or equal than 95.

Table: winemag_p1

```
--all top rated wineries ~ points
```

```
--top rated winery awarded points more or equal than 95
```

```
SELECT
    DISTINCT winery
FROM winemag_p1
WHERE points >= 95
```

Find prices for Spanish, Italian, and French wines

Wine Magazine Easy ID 10020

Find prices for Spanish, Italian, and French wines. Output the price.

Table: winemag_p1

```
SELECT
    price
FROM winemag_p1
WHERE country IN ('Spain', 'Italy', 'France')
```

Find the fraction of rides for each weather and the hour

Interview Question Date: February 2020

Lyft Hard ID 10019

Find the fraction (percentage divided by 100) of rides each weather-hour combination constitutes among all weather-hour combinations. Output the weather, hour along with the corresponding fraction.

Table: lyft_rides

①计算部分与整数的比值:

法1: i. 共计数cte ii. part数cte, part left join total ON TRUE, 即可以直接使用所有变量

法2: i. 共计数cte, ii. part数cte, 在part cte中如用total数, 可以直接写select A from total

②管probability遇到都是0的情况 ...NUMERIC 或 ...DECIMAL 可以转换为小数

```
WITH total AS(
SELECT
    COUNT(*) AS all_counts
FROM lyft_rides
),
part AS (
SELECT
    weather,
    hour,
    COUNT(*) AS part_counts
FROM lyft_rides
GROUP BY 1,2
```

```
)
--法1:
SELECT
    weather,
    hour,
    part_counts / all_counts :: DECIMAL
FROM part LEFT JOIN total
ON TRUE
ORDER BY 1,2
--法2:
-- SELECT
--     weather,
--     hour,
--     part_counts / (SELECT all_counts FROM total)::DECIMAL AS probability
-- FROM part
-- ORDER BY 1,2
```