

Задание 1. Определение пола

В таблице Transactions дана информация о транзакциях клиентов.
В таблице Clients дана информация о клиентах с различной датой актуальности данных.

- Задания:
- Получить таблицу с наиболее актуальными данными по всем клиентам с атрибутивным составом:
 - идентификатор клиента
 - ФИО
 - пол
 - Получить информацию о том, сколько женщин и мужчин совершало покупки в терминалах за январь 2020 в разрезе МСС кодов.
 - (Дополнительно) Выделить фамилию, имя и отчество из ФИО.

Таблица Transactions (пример данных)

| trans_id | client_id | partition_dt | sum_trans | mcc_code |
|----------|-----------|--------------|-----------|----------|
| 1111 | 100000 | 01.01.2020 | 9520 | 5611 |
| 1112 | 200000 | 08.01.2020 | 1623 | 5641 |
| 1113 | 300000 | 09.01.2020 | 1799 | 7832 |
| 1114 | 400000 | 14.01.2020 | 728 | 7538 |
| 1115 | 400000 | 15.01.2020 | 624 | 7832 |
| 1116 | 500000 | 16.01.2020 | 542 | 5641 |
| 1117 | 600000 | 21.01.2020 | 7982 | 7832 |
| 1118 | 700000 | 21.01.2020 | 886 | 7538 |
| 1119 | 700000 | 31.01.2020 | 432 | 5641 |
| 2000 | 800000 | 04.01.2020 | 680 | 5611 |
| 2001 | 100000 | 04.02.2020 | 3472 | 5641 |
| 2002 | 400000 | 05.02.2020 | 572 | 5641 |
| 2003 | 500000 | 06.02.2020 | 783 | 7832 |
| 2004 | 600000 | 07.02.2020 | 9774 | 5611 |
| 2005 | 700000 | 07.02.2020 | 8717 | 5641 |
| 2006 | 800000 | 07.02.2020 | 502 | 7832 |

Таблица Clients (пример данных)

| client_id | partition_dt | client_name |
|-----------|--------------|---------------------------------|
| 100000 | 01.12.2020 | Потапов |
| 200000 | 01.07.2018 | Рыбакова Валерия Львовна |
| 300000 | 01.01.2021 | Трофимов Андрей Ильич |
| 400000 | 01.08.2008 | Балашова Ксения |
| 500000 | 01.06.2007 | Gavrilova Anastasiya Kirillovna |
| 600000 | 01.02.2021 | Быков Семён Семёнович |
| 800000 | 01.01.2020 | Никольская Ульяна Антоновна |
| 800000 | 01.02.2020 | Панфилова Ульяна Антоновна |

Решение:

```
1  -- Получаем информацию о том, сколько женщин и мужчин совершало покупки в терминалах за январь 2020
2  --в разрезе МСС кодов
3  SELECT
4      mcc_code,
5      client_gender,
6      count(client_id) as num_clients
7  FROM
8      (
9          SELECT d.client_id as client_id,
10             c.client_name,
11             c.client_gender,
12             d.partition_dt as partition_dt,
13             d.sum_trans,
14             d.mcc_code
15
16      FROM
17          (
18              -- Получаем таблицу с наиболее актуальными данными по всем клиентам с атрибутивным составом
19              SELECT a.client_id,
20                 b.client_name,
21                 splitByChar(' ',client_name) as fio_array,
22                 arrayElement(fio_array, 1) as client_surname,
23                 -- (Дополнительно) Выделить фамилию, имя и отчество из ФИО
24                 --arrayElement(fio_array, 2) as client_name,
25                 --arrayElement(fio_array, 3) as client_patronymic,
26                 multiIf(client_surname like '%а', 'Ж', client_surname like '%а', 'Ж',
27                    client_surname like '%ая', 'Ж', client_surname like '%aya', 'Ж', 'М')
28                 as client_gender
29
30          FROM
31              (
32                  -- Получаем наиболее актуальные данные
33                  SELECT DISTINCT client_id,
34                     max(partition_dt) as partition_dt
35                  FROM test.Titova_Clients
36                  GROUP BY client_id
37                  ORDER BY client_id
38              ) as a
39
40          LEFT JOIN test.Titova_Clients as b on a.client_id = b.client_id and a.partition_dt = b.partition_dt
41          ) as c
42
43      -- Добавляем данные из таблицы Transactions
44      RIGHT JOIN
45          (
46              SELECT client_id,
47                 partition_dt,
48                 sum_trans,
49                 mcc_code
50              FROM test.Titova_Transactions
51              WHERE partition_dt between '2020-01-01' and '2020-01-31'
52              ) as d on c.client_id = d.client_id
53  )
54
55  GROUP BY mcc_code, client_gender
56  ORDER BY mcc_code
```

Вывод:

| Table | | | + Add Visualization |
|----------|---------------|-------------|---------------------|
| mcc_code | client_gender | num_clients | |
| 5,611 | M | 1 | |
| 5,611 | Ж | 1 | |
| 5,641 | | 1 | |
| 5,641 | Ж | 2 | |

Задание 2. Сводная таблица по терминалам

В таблице Terminal_transactions дана информация о сумме и количестве транзакций терминалов за 8 месяцев 2021 года.

Задание:

Необходимо написать SQL-запрос, который переводит эту таблицу в pivot и для каждого терминала посчитать долю суммы транзакций каждого месяца от общей суммы в формате:

| | | | | | | | |
|-------------|------------|------------|------------|-----|------------|------------|-----|
| Terminal_id | Sum_1_2020 | Sum_2_2020 | Sum_3_2020 | ... | Prc_1_2020 | Prc_2_2020 | ... |
|-------------|------------|------------|------------|-----|------------|------------|-----|

К таблице Terminal_transactions можно обращаться только один раз

Таблица Terminal_transactions (пример данных)

| terminal_id | partition_dt | sum_trans |
|-------------|--------------|-----------|
| 555555 | 01.01.2021 | 3537690 |
| 555555 | 01.02.2021 | 4677980 |
| 555555 | 01.03.2021 | 5678989 |
| 555555 | 01.04.2021 | 356657 |
| 555555 | 01.05.2021 | 456757 |
| 555555 | 01.06.2021 | 5676899 |
| 555555 | 01.07.2021 | 6789020 |
| 555555 | 01.08.2021 | 5678789 |

Решение:

```
1 SELECT terminal_id,
2     SUM(CASE WHEN partition_mnth = '1' THEN sum_trans END) as Sum_1_2020,
3     SUM(CASE WHEN partition_mnth = '2' THEN sum_trans END) as Sum_2_2020,
4     SUM(CASE WHEN partition_mnth = '3' THEN sum_trans END) as Sum_3_2020,
5     SUM(CASE WHEN partition_mnth = '4' THEN sum_trans END) as Sum_4_2020,
6     SUM(CASE WHEN partition_mnth = '5' THEN sum_trans END) as Sum_5_2020,
7     SUM(CASE WHEN partition_mnth = '6' THEN sum_trans END) as Sum_6_2020,
8     SUM(CASE WHEN partition_mnth = '7' THEN sum_trans END) as Sum_7_2020,
9     SUM(CASE WHEN partition_mnth = '8' THEN sum_trans END) as Sum_8_2020,
10    SUM(CASE WHEN partition_mnth = '9' THEN sum_trans END) as Sum_9_2020,
11    SUM(CASE WHEN partition_mnth = '10' THEN sum_trans END) as Sum_10_2020,
12    SUM(CASE WHEN partition_mnth = '11' THEN sum_trans END) as Sum_11_2020,
13    SUM(CASE WHEN partition_mnth = '12' THEN sum_trans END) as Sum_12_2020,
14    (Sum_1_2020 / total_sum)*100 as Prc_1_2020,
15    (Sum_2_2020 / total_sum)*100 as Prc_2_2020,
16    (Sum_3_2020 / total_sum)*100 as Prc_3_2020,
17    (Sum_4_2020 / total_sum)*100 as Prc_4_2020,
18    (Sum_5_2020 / total_sum)*100 as Prc_5_2020,
19    (Sum_6_2020 / total_sum)*100 as Prc_6_2020,
20    (Sum_7_2020 / total_sum)*100 as Prc_7_2020,
21    (Sum_8_2020 / total_sum)*100 as Prc_8_2020,
22    (Sum_9_2020 / total_sum)*100 as Prc_9_2020,
23    (Sum_10_2020 / total_sum)*100 as Prc_10_2020,
24    (Sum_11_2020 / total_sum)*100 as Prc_11_2020,
25    (Sum_12_2020 / total_sum)*100 as Prc_12_2020,
26    sum(sum_trans) as total_sum
27 FROM
28 (
29     SELECT terminal_id,
30            EXTRACT(MONTH FROM partition_dt) as partition_mnth,
31            sum_trans
32     FROM test.Terminal_Transactions
33 )
34 GROUP BY terminal_id
```

Вывод:

Table

+ Add Visualization

| terminal_id | Sum_1_2020 | Sum_2_2020 | Sum_3_2020 | Sum_4_2020 | Sum_5_2020 | Sum_6_2020 | Sum_7_2020 | Sum_8_2020 | Sum_9_2020 | Sum_10_2020 | Sum_11_2020 | Sum_12_2020 |
|-------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|-------------|-------------|-------------|
| 555,555 | 3,537,690 | 4,677,980 | 5,678,989 | 356,657 | 456,757 | 5,676,899 | 6,789,020 | 5,678,789 | | | | |

Table

+ Add Visualization

| Sum_12_202 | Prc_1_2020 | Prc_2_2020 | Prc_3_2020 | Prc_4_2020 | Prc_5_2020 | Prc_6_2020 | Prc_7_2020 | Prc_8_2020 | Prc_9_2020 | Prc_10_2020 | Prc_11_2020 | Prc_12_2020 | total_sum |
|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|-------------|-------------|-------------|------------|
| | 10.77 | 14.24 | 17.29 | 1.09 | 1.39 | 17.28 | 20.66 | 17.29 | | | | | 32,852,781 |

Задание 3. Определение брендов магазинов

По названиям терминалов можно определить, к какому бренду относится магазин. Но часто эти названия написаны по-разному для одного и того же бренда и есть вероятность что названия терминалов у разных брендов будут похожи.

Задание:

Необходимо написать регулярные выражения, которые будут правильно определять бренд магазинов "Азбука Вкуса" и "Белорусский трикотаж".

Примеры некоторых наименований терминалов для «Азбуки вкуса» (для белорусского трикотажа необходимо самостоятельно придумать возможные варианты):

AZBUKA VKUSA 82
AZBUKA VKUSA 131 PROSP
AZBUKA VKUSOV
AB DAILY.
OOO AZBUKA VKUSA
AV MARKET
SP_AV AZBUKAVKUSA
"AZBUKA VKUSA" SHOP
AV AZBUKAVKUSA..
AV DAILY..
AZBUKA VKUSA PROSPEKT
AV .AZBUKAVKUSA
MAGAZIN AZBUKA VKUSA
MINIMARKET AZBUKA VKUSA
AV AZBUKAVKUSA.
AZBUKAVKUSA
"AZBUKAVKUSA"
AZBUKA_VKUSA 80
AV.MARKET
"AZBUKA VKUSA" SHOP GO

Примеры наименований других магазинов, которые не должны выбираться с помощью регулярных выражений для азбуки вкуса:

AZBUKA
AVE
VKUSVILL
MAGAZIN
DALY MARKT
ABC

Решение:

```
1 SELECT terminal_name,  
2     case when multiMatchAny(replace(replace(terminal_name, ' ', ' '), ' ', ' '),  
3     ['AZBUKA\s\S', 'AZBUKA_\s', 'AB\s\S', 'AV\s\S', 'AV.\s', 'AV\s', 'AZBUKAVKUSA']) = '1'  
4     then 'AZBUKA VKUSA'  
5     when multiMatchAny(replace(replace(terminal_name, ' ', ' '), ' ', ' '),  
6     ['BELORUSSKIY\s\S', 'BELORUSSKIY_\s', 'BEL\s\S', 'BEL_\s', 'BEL.'])= '1'  
7     then 'BELORUSSKIY TRIKOTAZH' end as brand_name  
8 FROM test.Terminal_name  
9 ORDER BY brand_name DESC
```

Вывод:

Table + Add Visualization

| terminal_name | brand_name |
|------------------------|--------------|
| "AZBUKA VKUSA" SHOP | AZBUKA VKUSA |
| "AZBUKA VKUSA" SHOP GO | AZBUKA VKUSA |

< 1 2 >

Задание 4. Определение брендов магазинов

Таблица Merchants является ненормализованной. Она содержит атрибут listofterminal – список терминалов этого магазина с информацией о времени работы.

Задание:

Необходимо привести таблицу к виду магазин – терминал – время работы

Таблица Merchants (пример данных)

| merchant | listofterminal |
|----------|--|
| 11111 | 11222222: 12:00-20:00; 11222226: 14:00-17:00 |
| 12344 | 333333; 444444: 9:00-23:00; 555555: 9:00-23:00 |
| 22222 | 22222221: 9:00-22:00 |

Решение:

```
1 SELECT merchant,
2     terminal_array2[1] AS terminal,
3     replace(concat(terminal_array2[2],':', terminal_array2[3],':',terminal_array2[4]), '::~', '')
4     as terminal_working_hours
5
6 FROM
7     (
8         -- Разбиваем терминал и время работы на отдельные элементы
9         SELECT merchant,
10             splitByString(':', terminal_array) as terminal_array2
11         FROM
12             (
13                 -- Раскладываем каждое отдельное значение массива в строку
14                 SELECT merchant, terminal_array
15                 FROM
16                     (
17                         -- Разбиваем listofterminal на массив с элементами
18                         SELECT merchant,
19                             splitByString(';', listofterminal) as terminal_array
20                         FROM test.Merchants
21                     )
22                 ARRAY JOIN terminal_array
23             )
24     )
25
```

Вывод:

Table + Add Visualization

| merchant | terminal | terminal_working_hours |
|----------|----------|------------------------|
| 11,111 | 11222222 | 12:00-20:00 |
| 11,111 | 11222226 | 14:00-17:00 |
| 12,344 | 333333 | |
| 12,344 | 444444 | 9:00-23:00 |
| 12,344 | 555555 | 9:00-23:00 |
| 22,222 | 22222221 | 9:00-22:00 |