

Задания для стажеров

| Уровень сложности | Задачи по SQL | | |
|-------------------|--|-------|-------------|
| I | Table: daily_weather | | |
| | month | day | temperature |
| | 1 | 1 | -15 |
| | 1 | 2 | -19 |
| | ... | ... | ... |
| | 2 | 1 | -5 |
| | 2 | 2 | 0 |
| | ... | ... | ... |
| | 1. Написать запрос который посчитает среднюю, минимальную и максимальную температуру | | |
| | 2. Написать запрос который посчитает среднюю, минимальную и максимальную температуру в разрезе месяцев | | |
| I | Какое максимальное и минимальное количество строк может быть получено при JOINе двух таблиц по 5 строк каждая (поля not NULL). | | |
| | | MIN | MAX |
| | INNER | | |
| | LEFT | | |
| | RIGHT | | |
| | FULL | | |
| | | | |
| II | Table: mon_salary | | |
| | id | month | salary |
| | 1 | 1 | 100 |
| | 1 | 2 | 150 |
| | ... | ... | ... |
| | 2 | 5 | 200 |
| | 2 | 6 | 180 |
| | ... | ... | ... |
| | Написать sql-запрос отбирающие все id и month у которых salary > avg salary для id | | |

| Уровень сложности | Задачи по SQL | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|-------------------|---|--------------|---------------|-------------|------------|------------|---------|------------|-----------|-------|-------|---------|------------|-----------|--------|------|---------|------------|-----------|--------|------|---------|------------|-----------|-------|-------|---------|------------|-----------|-------|-------|---------|------------|-----------|--------|-------|---------|------------|-----------|--------|-------|---------|------------|-----------|-------|------|-----------|-------------|--------------|---------------|----------|--|--|--|--|--|
| II | Table: VSP_oper_data <table><tr><th>Client_id</th><th>Report_date</th><th>VSP_Number</th><th>Txn_type</th><th>Txn_amount</th></tr><tr><td>1233455</td><td>2017.05.02</td><td>1234/0123</td><td>debit</td><td>10000</td></tr><tr><td>1233455</td><td>2017.05.03</td><td>1236/0123</td><td>credit</td><td>1000</td></tr><tr><td>1233455</td><td>2017.05.04</td><td>1234/0123</td><td>credit</td><td>1000</td></tr><tr><td>1233455</td><td>2017.05.07</td><td>1235/0123</td><td>debit</td><td>15000</td></tr><tr><td>1233456</td><td>2017.05.02</td><td>1234/0123</td><td>debit</td><td>11000</td></tr><tr><td>1233456</td><td>2017.05.03</td><td>1236/0123</td><td>credit</td><td>10000</td></tr><tr><td>1233456</td><td>2017.05.04</td><td>1234/0123</td><td>credit</td><td>10000</td></tr><tr><td>1233456</td><td>2017.06.07</td><td>1237/0123</td><td>debit</td><td>5000</td></tr></table> <p>В таблице VSP_oper_data txn_type принимает значения debit, credit</p> <p>Задание: напишите sql запрос, который для каждого клиента выводит сумму debit, credit операций и последнее посещенное VSP по месяцам. Результат представьте в виде:</p> <table><tr><th>Client_id</th><th>Report_date</th><th>Debit_amount</th><th>Credit_amount</th><th>Last_VSP</th></tr><tr><td></td><td></td><td></td><td></td><td></td></tr></table> | Client_id | Report_date | VSP_Number | Txn_type | Txn_amount | 1233455 | 2017.05.02 | 1234/0123 | debit | 10000 | 1233455 | 2017.05.03 | 1236/0123 | credit | 1000 | 1233455 | 2017.05.04 | 1234/0123 | credit | 1000 | 1233455 | 2017.05.07 | 1235/0123 | debit | 15000 | 1233456 | 2017.05.02 | 1234/0123 | debit | 11000 | 1233456 | 2017.05.03 | 1236/0123 | credit | 10000 | 1233456 | 2017.05.04 | 1234/0123 | credit | 10000 | 1233456 | 2017.06.07 | 1237/0123 | debit | 5000 | Client_id | Report_date | Debit_amount | Credit_amount | Last_VSP | | | | | |
| | Client_id | Report_date | VSP_Number | Txn_type | Txn_amount | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233455 | 2017.05.02 | 1234/0123 | debit | 10000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233455 | 2017.05.03 | 1236/0123 | credit | 1000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233455 | 2017.05.04 | 1234/0123 | credit | 1000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233455 | 2017.05.07 | 1235/0123 | debit | 15000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233456 | 2017.05.02 | 1234/0123 | debit | 11000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233456 | 2017.05.03 | 1236/0123 | credit | 10000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233456 | 2017.05.04 | 1234/0123 | credit | 10000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1233456 | 2017.06.07 | 1237/0123 | debit | 5000 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Client_id | Report_date | Debit_amount | Credit_amount | Last_VSP | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| III | 1. Реализовать задачу из п. 2 тремя различными способами | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Table: daily_weather <table><tr><th>month</th><th>day</th><th>temperature</th></tr><tr><td>1</td><td>1</td><td>-15</td></tr><tr><td>1</td><td>2</td><td>-19</td></tr><tr><td>...</td><td>...</td><td>...</td></tr><tr><td>2</td><td>1</td><td>-5</td></tr><tr><td>2</td><td>2</td><td>0</td></tr><tr><td>...</td><td>...</td><td>...</td></tr></table> <p>2. Написать sql запрос который считает для каждого дня среднюю температуру за 5 предыдущих дней</p> | month | day | temperature | 1 | 1 | -15 | 1 | 2 | -19 | ... | ... | ... | 2 | 1 | -5 | 2 | 2 | 0 | ... | ... | ... | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| month | day | temperature | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 1 | 1 | -15 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 1 | 2 | -19 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ... | ... | ... | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2 | 1 | -5 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2 | 2 | 0 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ... | ... | ... | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| III | <p>Задание: напишите sql запрос, который для каждого клиента из VSP_oper_data выведет долю debit операций клиента к debit операциям всех клиентов по месяцам. Результат в виде таблицы:</p> <table><tr><th>Client_id</th><th>Report_date</th><th>Ratio</th></tr><tr><td></td><td></td><td></td></tr></table> | Client_id | Report_date | Ratio | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Client_id | Report_date | Ratio | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

| Уровень сложности | Scala или python (используя spark, pandas или стандартные коллекции) |
|-------------------|---|
| I | <p>Необходимо на любом знакомом языке программирования написать алгоритм:</p> <p>Вывести построчно целые числа от 1 до N не более C чисел в строке</p> |
| II | <p>## Задача Написать приложение на scala или python, которое в локальном режиме выполняет следующее: По имеющимся данным о рейтингах книг посчитать агрегированную статистику по ним.</p> <ol style="list-style-type: none"> 1. Прочитать csv файл: book.csv 2. Вывести схему для dataframe полученного из п.1 3. Вывести количество записей 4. Вывести информацию по книгам у которых рейтинг выше 4.50 5. Вывести средний рейтинг для всех книг. 6. Вывести агрегированную информацию по количеству книг в диапазонах: <p>0 - 1 1 - 2 2 - 3 3 - 4 4 - 5</p> |
| III | <p>## Задача Написать приложение (опционально spark приложение), которое в локальном режиме выполняет следующее: По имеющимся данным о рейтингах фильмов (MovieLens: 100 000 рейтингов) посчитать агрегированную статистику по ним. Использовать архив data.zip</p> <p>## Описание данных Имеются следующие входные данные: Архив с рейтингами фильмов. Файл README содержит описания файлов. Файл u.data содержит все оценки, а файл u.item — список всех фильмов. (используются только эти два файла) id_film=32</p> <ol style="list-style-type: none"> 1. Прочитать данные файлы. 2. создать выходной файл в формате json, где <p>Поле "Toy Story" нужно заменить на название фильма, соответствующего id_film и указать для заданного фильма количество поставленных оценок в следующем порядке: "1", "2", "3", "4", "5". То есть сколько было единиц, двоек, троек и т.д.</p> <p>В поле "hist_all" нужно указать то же самое только для всех фильмов общее количество поставленных оценок в том же порядке: "1", "2", "3", "4", "5".</p> <p>Пример решения:</p> <pre>{ "Toy Story": [134, 123, 782, 356, 148], "hist_all": [134, 123, 782, 356, 148] }</pre> |

| Уровень сложности | Scala или python (используя spark, pandas или стандартные коллекции) |
|-------------------|--|
| III | <p>Входные данные:</p> <p>Customer.csv – информация о клиентах</p> <p>Имя поля: формат</p> <p>id: Int, name: String, email: String, joinDate: Date, status: String</p> <p>Product.csv – информация о товарах</p> <p>id: Int name: String price: Double numberOfProducts: Int</p> <p>Order.csv – информация о заказах</p> <p>customerID: Int orderID: Int productID: Int numberOfProduct: Int – кол-во товара в заказе orderDate: Date status: String</p> <p>Необходимо определить самый популярный продукт у клиента</p> <p>Итоговое множество должно содержать поля: customer.name, product.name</p> <p>Результат записать в csv-файл</p> |