

Module 3 Assignment | Lab 4: Hadoop with Hive

Trang Tran

CPS, Northeastern University

ALY6110 | Data Management and Big Data

Professor Andrew Kinley

Aug 05, 2023

1. `SELECT c_id, c_lname FROM customers WHERE c_state = 'New York';`

Select `c_id` and `c_lname` columns from the `customers` table (as picture below) where the value in the `c_state` column is 'New York'.

```
hive> describe customers;
OK
c_id          int
c_fname       string
c_lname       string
c_street      string
c_city        string
c_state       string
c_zip         string
c_yob         int
c_gender      string
credit_card   string
internet      string
mobile        string
```

```
> -- 1: Display customers (ID and Name) from the state of New York
> SELECT c_id, c_lname
> FROM customers
> WHERE c_state = 'New York';
OK
171      Hensley
177      Wheeler
180      Kramer
181      Tate
```

2. `SELECT count(c_id) AS total, c_zip FROM customers GROUP BY c_zip;`

Calculate the total count of customers (`c_id`), store as `total` for each unique `c_zip` (zip code) from the `customers` table, and group the results by `c_zip`.

3. `SELECT c.c_id, c.c_lname, p.transaction_id, p.transaction_date FROM customers c`
`JOIN payments p ON c.c_id = p.c_id`
`WHERE p.late = 'TRUE';`

Retrieve `c_id`, `c_lname`, `transaction_id`, and `transaction_date` from the `customers` and `payments` tables, joining them on matching ID values, where the `p.late` equals `'TRUE'`.

'c' and 'p' are used as an alias for the 'customers' and 'payment' tables.

```
OK
172 Foster 401129380-4 23-8-2015
174 McCormick 851112155-2 4-12-2019
175 King 108659198-9 13-9-2016
178 Noble 694690715-8 31-5-2014
179 Matthews 318241713-5 15-7-2016
181 Tate 146268743-8 23-8-2020
183 Silva 270161074-X 19-7-2016
183 Silva 559786593-4 13-12-2018
```

4. `SELECT c_city, COUNT(c_id) FROM customers WHERE credit_card = 'TRUE' GROUP BY c_city;`

Count the number of customers (c_id) with a valid credit card (credit_card = 'TRUE') in each unique c_city from the 'customers' table and group the results by c_city.

```
OK
Chicago 2
Dallas 1
New York 3
San Diego 1
Austin 1
Chattanooga 1
Detroit 1
```

5. `SELECT c.c_state, COUNT(c.c_id) as total FROM customers c JOIN payments p ON c.c_id = p.c_id WHERE p.late = 'TRUE' GROUP BY c.c_state;`

Count the number of customers (c_id) with late payments (p.late = 'TRUE') in each unique c_state from the 'customers' and 'payments' tables, grouping the results by c_state values.

```
OK
Illinois 2
New York 1
Michigan 2
Tennessee 1
Texas 2
```

6. `SELECT * FROM customers WHERE mobile = 'FALSE';`

9. `SELECT c.c_lname, c.c_city, c.c_state FROM customers c JOIN payments p ON c.c_id = p.c_id WHERE p.late = 'TRUE' and c.c_yob > 1985;`

Select c_lname, c_city, and c_state from the 'customers' table, joining it with the 'payments' table on matching c_id values, where the payment is marked as late and the customer's year of birth (c_yob) is greater than 1985.

```
OK
McCormick      Chicago Illinois
King    Chicago Illinois
Noble    Hixson  Tennessee
Tate     New York      New York
```

10. `SELECT c.c_id, c.c_lname, c.c_city, c.c_state`
`FROM customers c JOIN payments p ON c.c_id = p.c_id`
`WHERE p.late = 'FALSE' AND c.internet = 'TRUE'`
`GROUP BY c.c_id, c.c_lname, c.c_city, c.c_state;`

Select columns c_id, c_lname, c_city, c_state from the 'customers' table, joining it with the 'payments' table on matching c_id values, where the payment is not late (p.late = 'FALSE') and the internet = 'TRUE'. Group the results by c_id, c_lname, c_city, c_state.

```
OK
173    Livingston    Chattanooga    Tennessee
182    Estrada Dallas    Texas
183    Silva    Detroit Michigan
Time taken: 1.659 seconds, Fetched: 3 row(s)
hive> -- exit hive
hive> exit
```