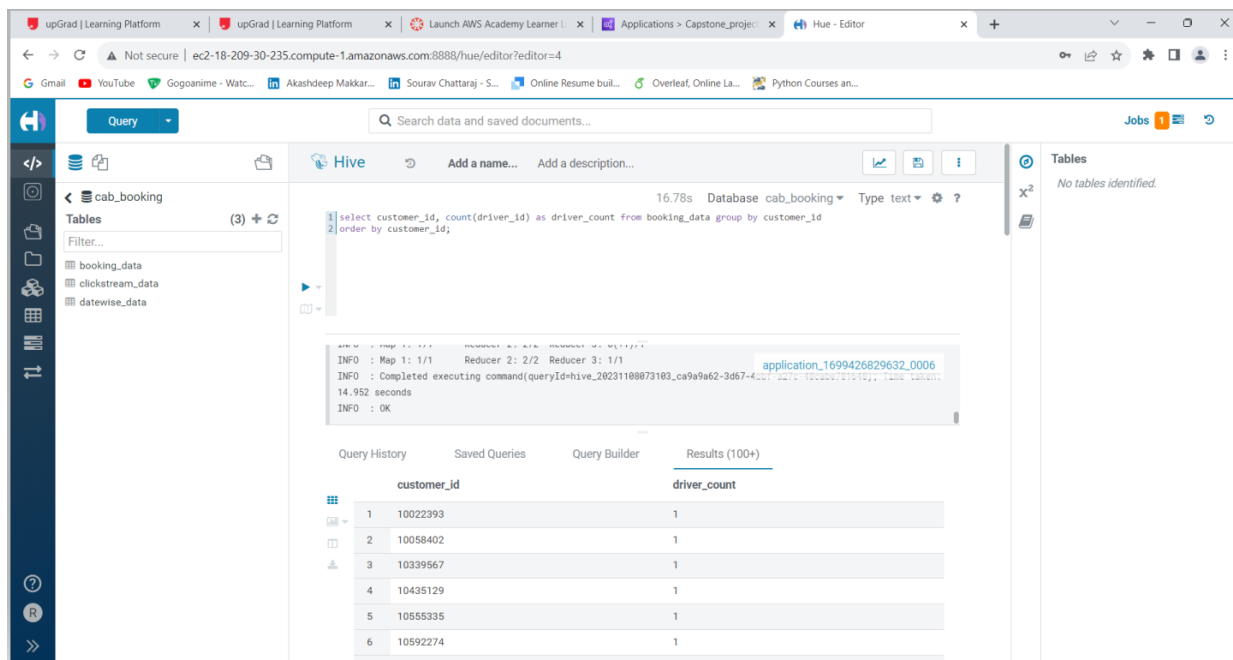# Queries

**1. Hive Query for Task 5:- Calculate the total different drivers for each customer.**
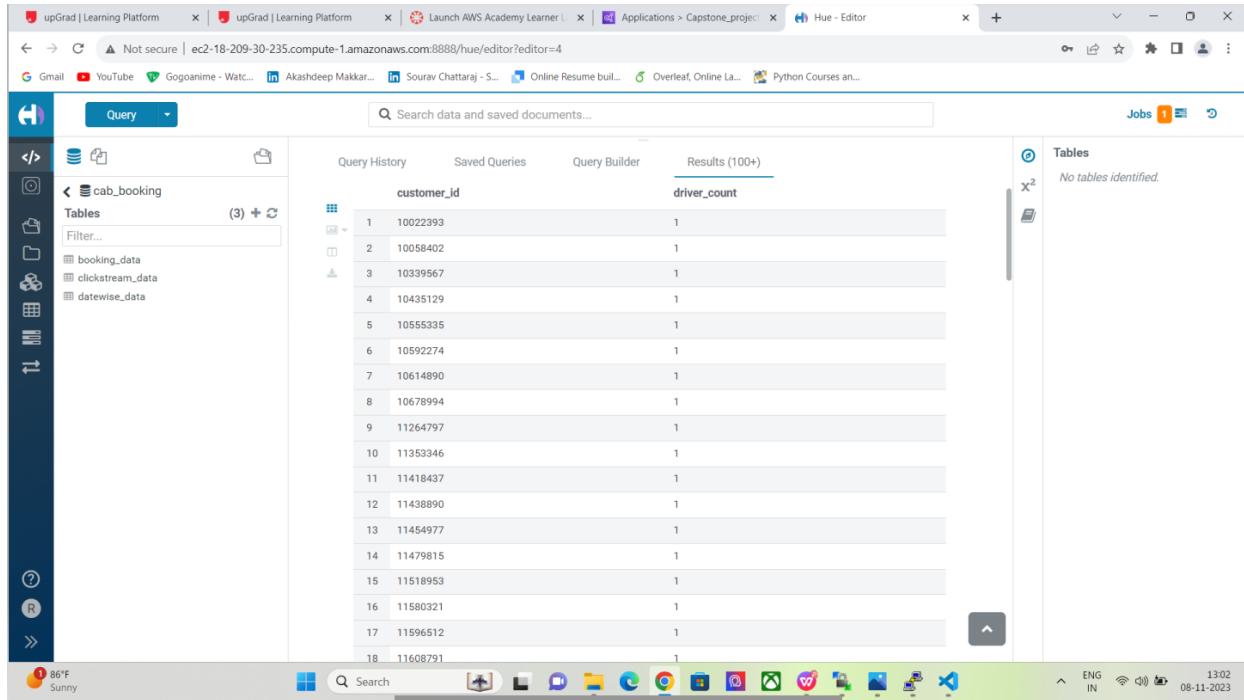
select customer_id, count(driver_id) as driver_count from booking_data group by customer_id order by customer_id;

**Screenshot after executing Query:-**

2. **Hive Query for Task 6:- Calculate the total rides taken by each customer.**

select customer_id, count(*) as total_ride from booking_data group by customer_id order by customer_id;

**Screenshot after executing Query:-**
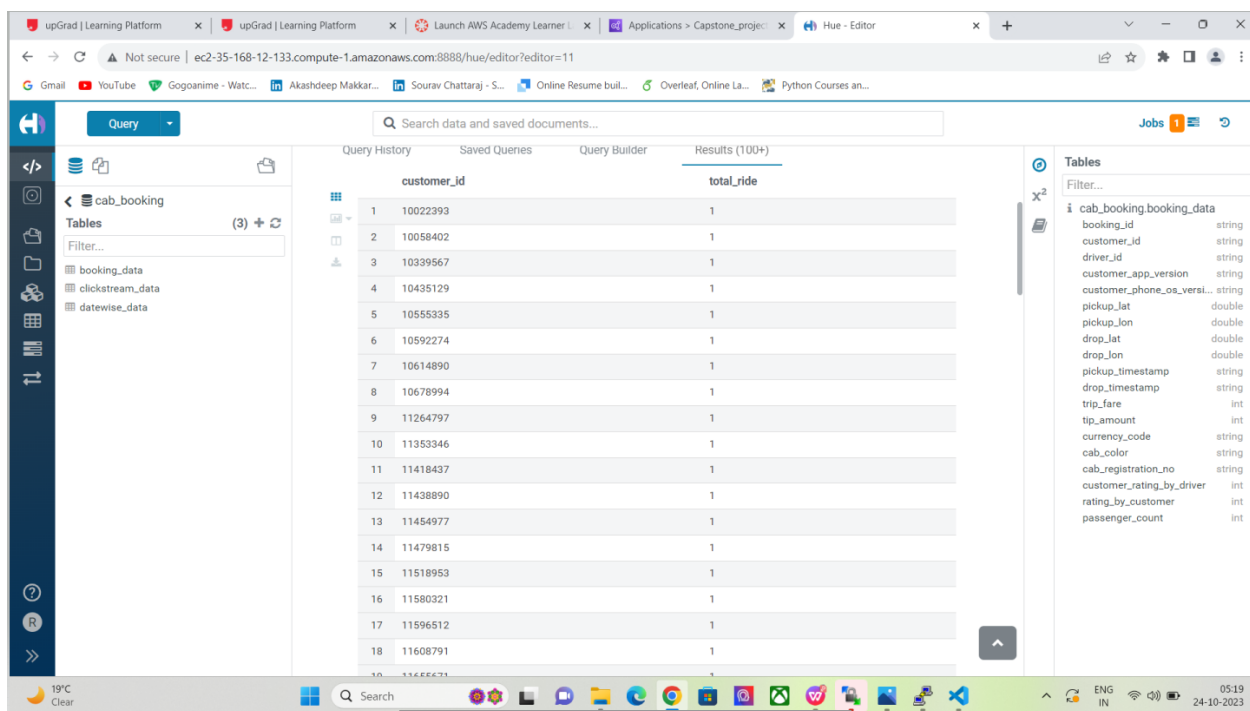
### 3. Hive Query for Task 7:-

**Find total visits by each customer on the booking page and total 'Book Now' button press. This can show the conversion-ratio.**
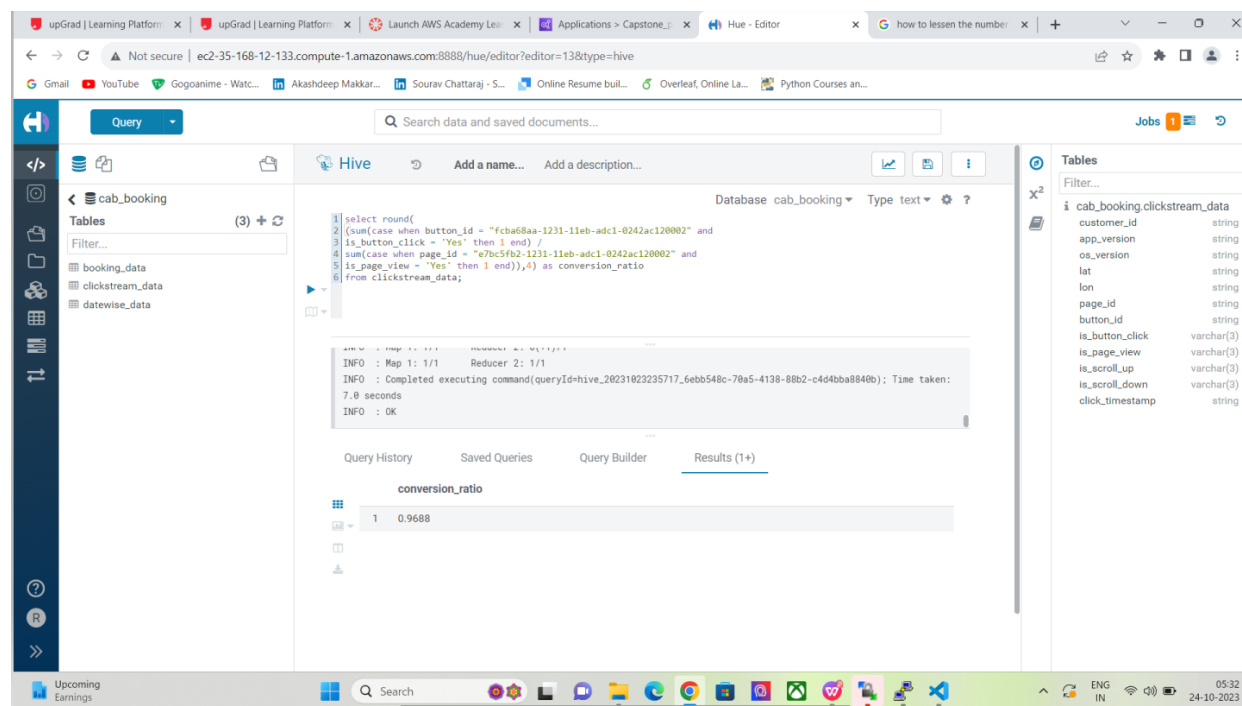**Booking page id is 'e7bc5fb2-1231-11eb-adc1-0242ac120002'**
**Book Now button id is 'fcba68aa-1231-11eb-adc1-0242ac120002'.**
**Also, calculate the conversion ratio as described in the Tasks segment.**

```
select round(
(sum(case when button_id = "fcba68aa-1231-11eb-adc1-0242ac120002" and
is_button_click = 'Yes' then 1 end) /
sum(case when page_id = "e7bc5fb2-1231-11eb-adc1-0242ac120002" and
is_page_view = 'Yes' then 1 end)),4) as conversion_ratio
from clickstream_data;
```
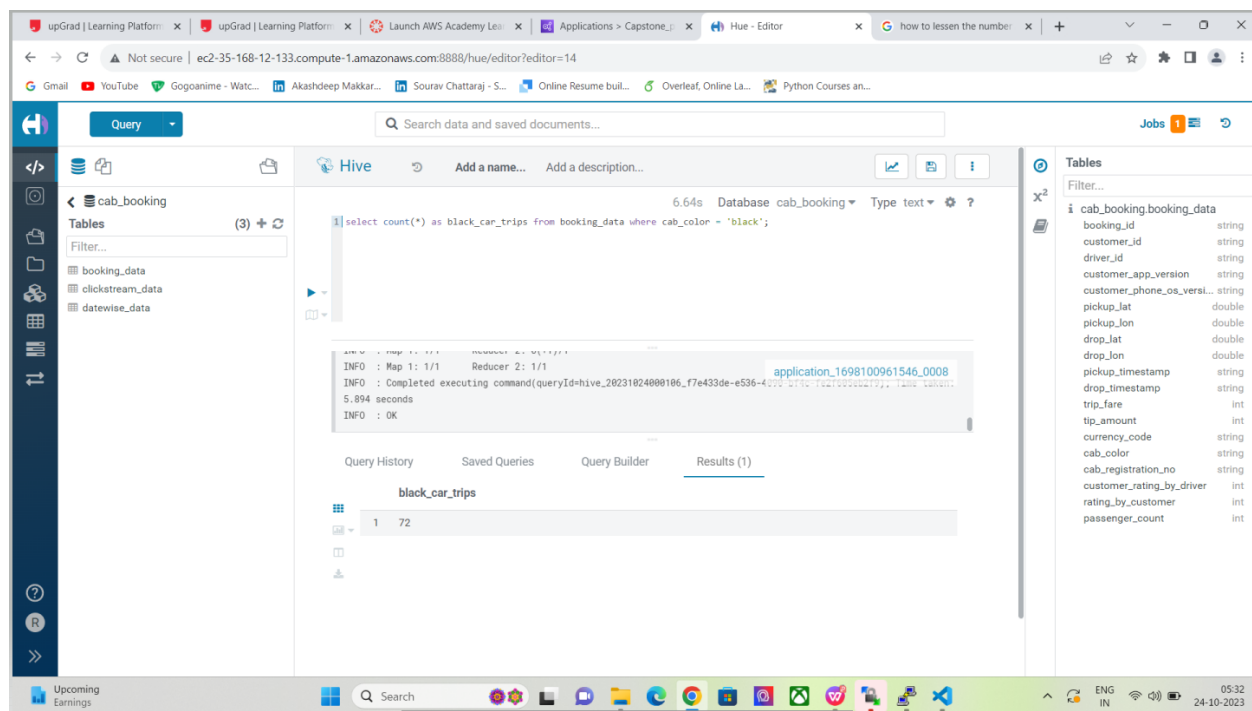
**Screenshot after executing Query:-**



### 4. Hive Query for Task 8:- Calculate the count of all trips done on Black cabs.

```
select count(*) as black_car_trips from booking_data where cab_color = 'black';
```

**Screenshot after executing Query:-**



5. **Hive Query for Task 9:- Calculate the total tip amount on a given date to all drivers by customers.**

select date_format(pickup_timestamp,'YYYY-MM-dd') as datewise, sum(tip_amount) as total_tip
from booking_data
group by date_format(pickup_timestamp,'YYYY-MM-dd')
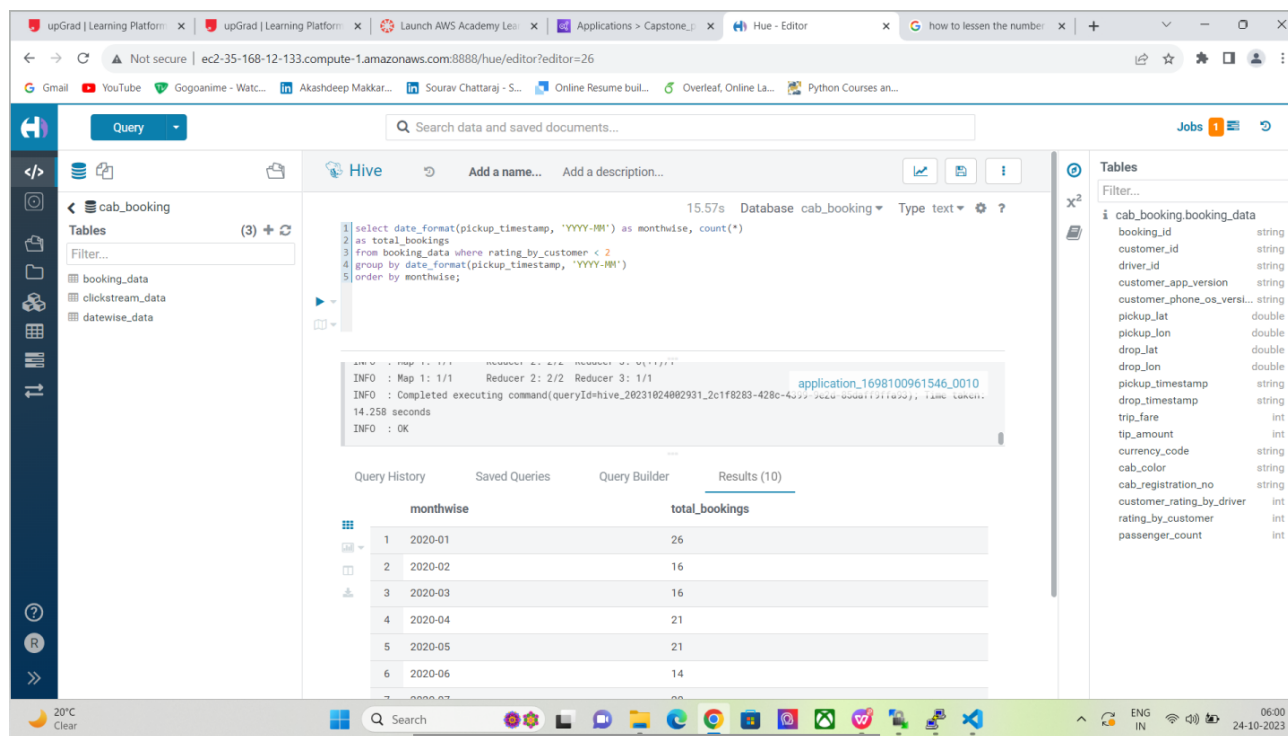order by datewise;

**Screenshot after executing Query:-**

**6. Hive Query for Task 10:- Calculate the count of all the bookings with a rating below 2 in a particular month.**

select date_format(pickup_timestamp, 'YYYY-MM') as monthwise, count(*)
as total_bookings
from booking_data where rating_by_customer < 2
group by date_format(pickup_timestamp, 'YYYY-MM')
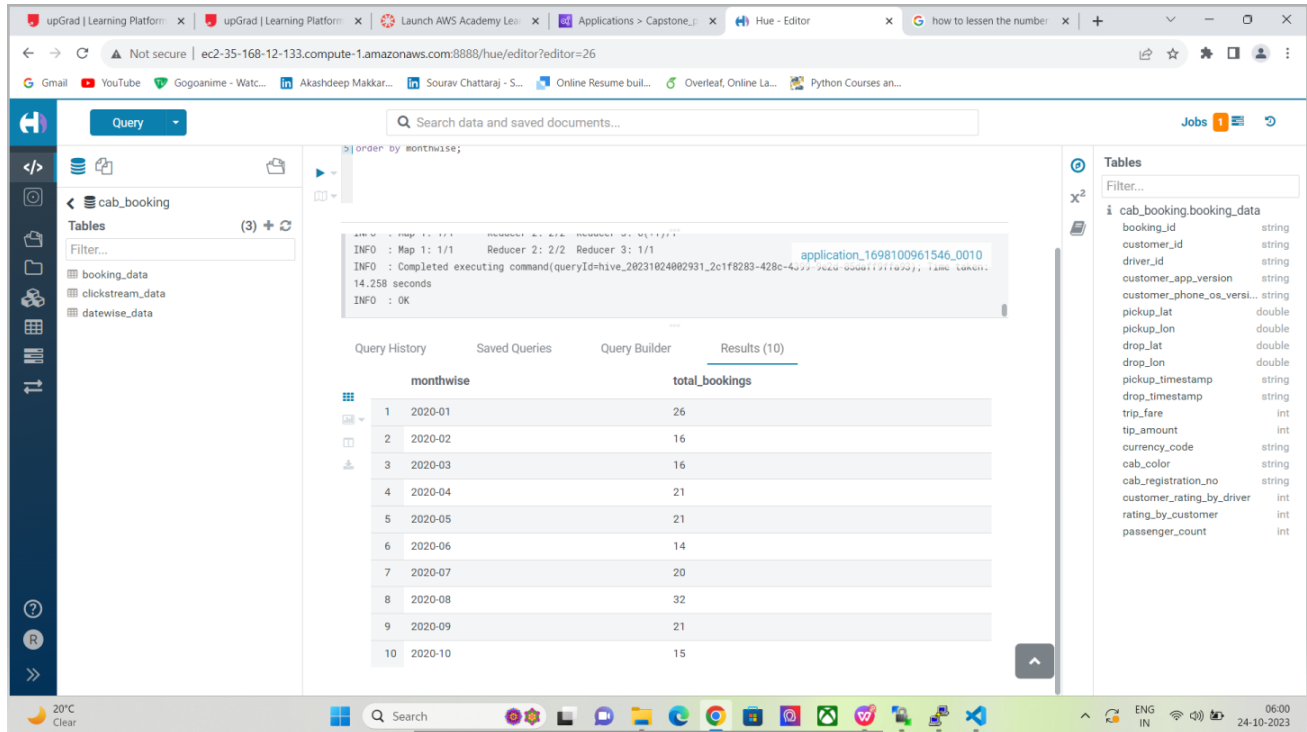order by monthwise;

**Screenshot after executing Query:-**

**7. Hive Query for Task 11:- Calculate the count of total iOS users.**

select count(*) AS TOTAL_USERS from clickstream_data where os_version = 'iOS';

**Screenshot after executing Query:-**