# Logic For Final Submission

## Active EMR cluster:



## S3 files:

Validation already completed for task 1 to task 4 in LogicFirstSubmission.pdf

# Capstone Project - Validation Document

**Note**: Please note that records retrieved from the Kafka topic and related metrics given below can vary a bit.

## Data Ingestion with Sqoop

Please check the number of records that are imported after the Sqoop Job

```
Number of records retrieved - 1000
```

## Bookings Table Count

Please check the number of records in the bookings table

```
Number of records - 1000
```

## Clickstream Table Count

Please check the number of records in the clickstream table

```
Number of records - 2984
```

## Bookings Aggregates Table Count

Please check the number of records in the bookings aggregates table

```
Number of records - 289
```

Enable HUE UI for Hive execution: add port 8888 in both master and slave security group

Open HUE link:



HUE: query check in UI, select cabrides database and execute table query

**Task 5**: Calculate the total number of different drivers for each customer

**QUERY**:

<mark>SELECT CUSTOMER_ID, COUNT(DISTINCT driver_id) AS TOTAL_NUMBER_OF_DRIVERS
FROM BOOKINGS_DETAIL
GROUP BY CUSTOMER_ID
ORDER BY CUSTOMER_ID;</mark>

**EXPLANATION:**

Check multiple drivers for each customer by grouping customer and count total individual drivers for each and sort by customer id

**OUTPUT**: each customer is having only one driver , so total number of drivers is "1"

**VALIDATION:** Exact Match

1. When you run the query to calculate the total number of different drivers for each customer, you would get an output as shown below:

```
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2020-11-17 12:23:06,034 Stage-1 map = 0%,  reduce = 0%
2020-11-17 12:23:12,394 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 3.27 sec
2020-11-17 12:23:20,727 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 7.69 sec
MapReduce Total cumulative CPU time: 7 seconds 690 msec
Ended Job = job_1605615116654_0005
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1   Cumulative CPU: 7.69 sec   HDFS Read: 43007 HDFS Write: 11000 SUCCESS
Total MapReduce CPU Time Spent: 7 seconds 690 msec
OK
10022393        1
10058402        1
10339567        1
10435129        1
10555335        1
10592274        1
10614890        1
10678994        1
11264797        1
11353346        1
11418437        1
11438890        1
11454977        1
11479815        1
11518953        1
11580321        1
11596512        1
11608791        1
11655671        1
11757536        1
11764909        1
11860278        1
11981042        1
12106105        1
12142182        1
12312603        1
12334699        1
12367832        1
12856708        1
12885363        1
12913608        1
12914577        1
12966909        1
13015449        1
13229062        1
```

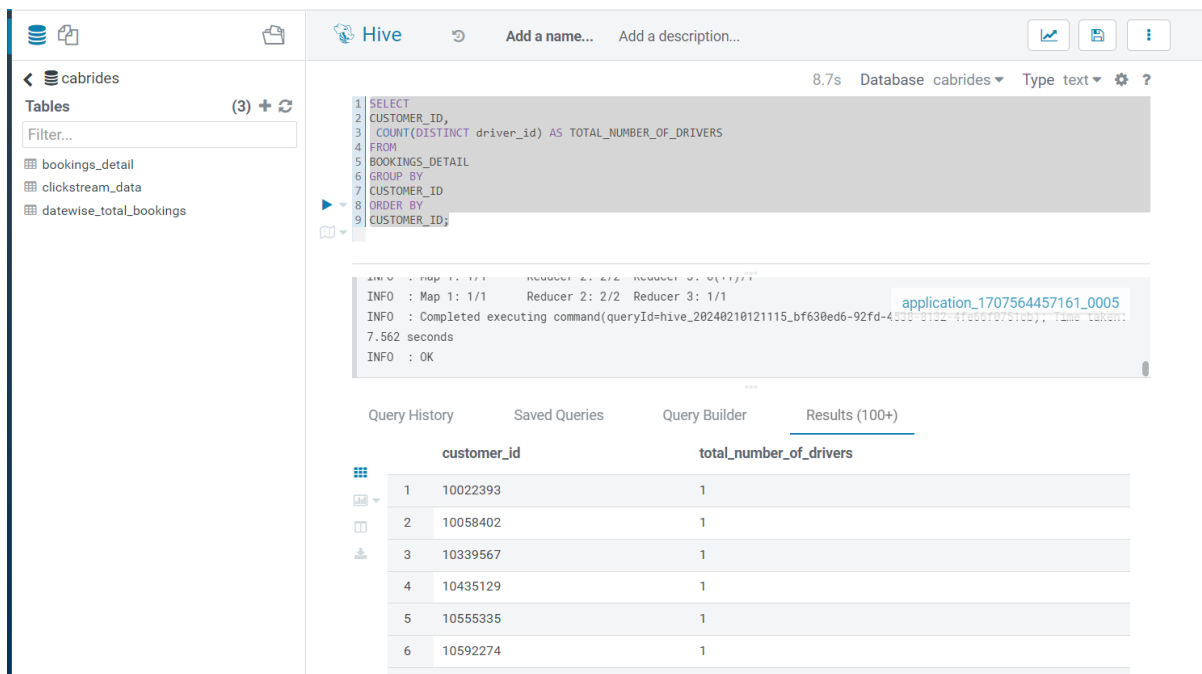- **Task 6**:

  Calculate the total rides taken by each customer.
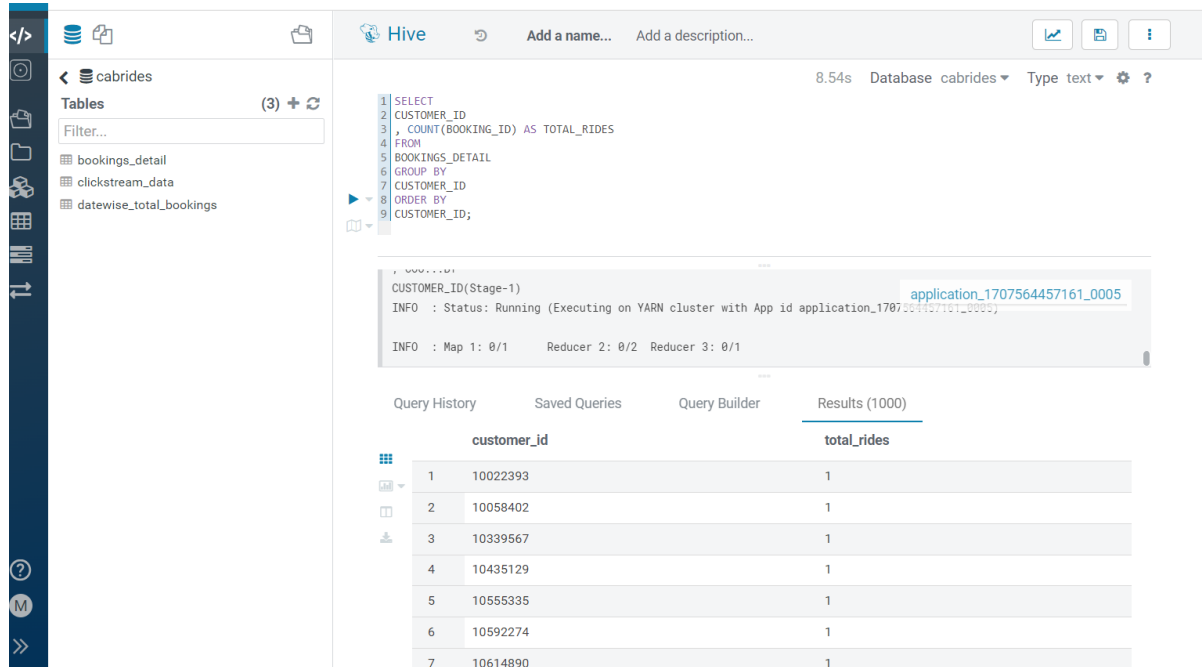  **QUERY:**
  SELECT CUSTOMER_ID, COUNT(BOOKING_ID) AS TOTAL_RIDES
  FROM BOOKINGS_DETAIL
  GROUP BY CUSTOMER_ID
  ORDER BY CUSTOMER_ID;

  **EXPLANATION:**
  Check total rides for each customer by grouping customer, count total rides   and sort by customer id

**OUTPUT**: each customer is having only one ride , so total ride is "1" for each customer



**VALIDATION:** Exact Match

- **Task 7**:

Find the total visits made by each customer on the booking page and the total 'Book Now' button presses. This can show the conversion ratio.
The booking page id is 'e7bc5fb2-1231-11eb-adc1-0242ac120002'.
The Book Now button id is 'fcba68aa-1231-11eb-adc1-0242ac120002'. You also need to calculate the conversion ratio as part of this task. Conversion ratio can be calculated as Total 'Book Now' Button Press/Total Visits made by customer on the booking page.
**QUERY:**
```
SELECT
SUM(CASE WHEN PAGE_ID = 'e7bc5fb2-1231-11eb-adc1-0242ac120002' THEN 1 ELSE 0
END) AS
TOTAL_PAGE_VISITS,
SUM(CASE WHEN BUTTON_ID = 'fcba68aa-1231-11eb-adc1-0242ac120002' THEN 1 ELSE 0
END) AS
TOTAL_BUTTON_PRESSED,
ROUND(CAST(SUM(CASE WHEN BUTTON_ID = 'fcba68aa-1231-11eb-adc1-0242ac120002'
THEN 1 ELSE 0 END) AS FLOAT) /
CAST(SUM(CASE WHEN PAGE_ID = 'e7bc5fb2-1231-11eb-adc1-0242ac120002' THEN 1
ELSE
0 END) AS FLOAT), 4) AS CONVERSION_RATIO
FROM CLICKSTREAM_DATA;
```

**EXPLANATION:**
- This analysis is useful for understanding customer behavior, such as how often they book a ride after visiting the booking page.
- Tracking the number of times the booking page is visited (PAGE_ID = 'e7bc5fb2-1231-11eb-adc1-0242ac120002') provides the total count of page visits.
- Recording the number of times the booking button is clicked (BUTTON_ID = 'fc6ba68a-1231-11eb-adc1-0242ac120002') indicates the total number of rides that have been booked.
- The conversion ratio is simply the number of total bookings divided by the number of total page visits.
- The conversion ratio is an essential Key Performance Indicator (KPI) for the company, showing a high likelihood of booking when a customer visits the booking page, which is approximately 98% in this instance.

**OUTPUT:**



**VALIDATION:** close to Match (0.9852), since Kafka had extra 16 records compare to validation it should get the conversion ratio as 0.9688.

- **Task 8**: Calculate the count of all trips done on black cabs.

    **QUERY:**
    ```
    SELECT COUNT(BOOKING_ID) AS TOTAL_TRIPS_BLACK_CABS
    FROM BOOKINGS_DETAIL
    WHERE UPPER(CAB_COLOR) ='BLACK';
    ```

    **EXPLANATION:**
    - This assessment aids in determining the overall number of journeys made using Black taxis.
    - Categorizing by CAB_COLOR can reveal the total count or proportion of trips completed by taxis of a particular color.
    - This might reveal any obscure trends in booking related to the color of the taxis, such as a customer's preference for booking taxis of a certain color.

**OUTPUT:** since case sensitive, handles upper from cab colour



**VALIDATION:** Exact Match

Count of all trips done on black cabs -72.

- **Task 9**: Calculate the total amount of tips given date wise to all drivers by customers.

**QUERY:**

SELECT DATE(PICKUP_TIMESTAMP) TRIP_DATE
, ROUND(SUM(TIP_AMOUNT),0) AS TOTAL_TIP_AMOUNT
FROM BOOKINGS_DETAIL
GROUP BY DATE(PICKUP_TIMESTAMP)
ORDER BY TRIP_DATE;

**EXPLANATION:**

- The date function isolates the date from the datetime_stamp value and labels it as Pickup_Date.
- The SUM function, in conjunction with GROUP BY Pickup Date, calculates the aggregate tip amount for each Pickup date.
- The ORDER BY clause, applied to the Pickup Date alias, arranges the results in chronological order based on the date of pickup.
- Analyzing this data can provide insights into whether customers tend to tip more on certain occasions or specific days.
- Based on these insights, management can introduce promotions for customers or drivers, inferring that a higher tip amount correlates with greater customer satisfaction on particular days.
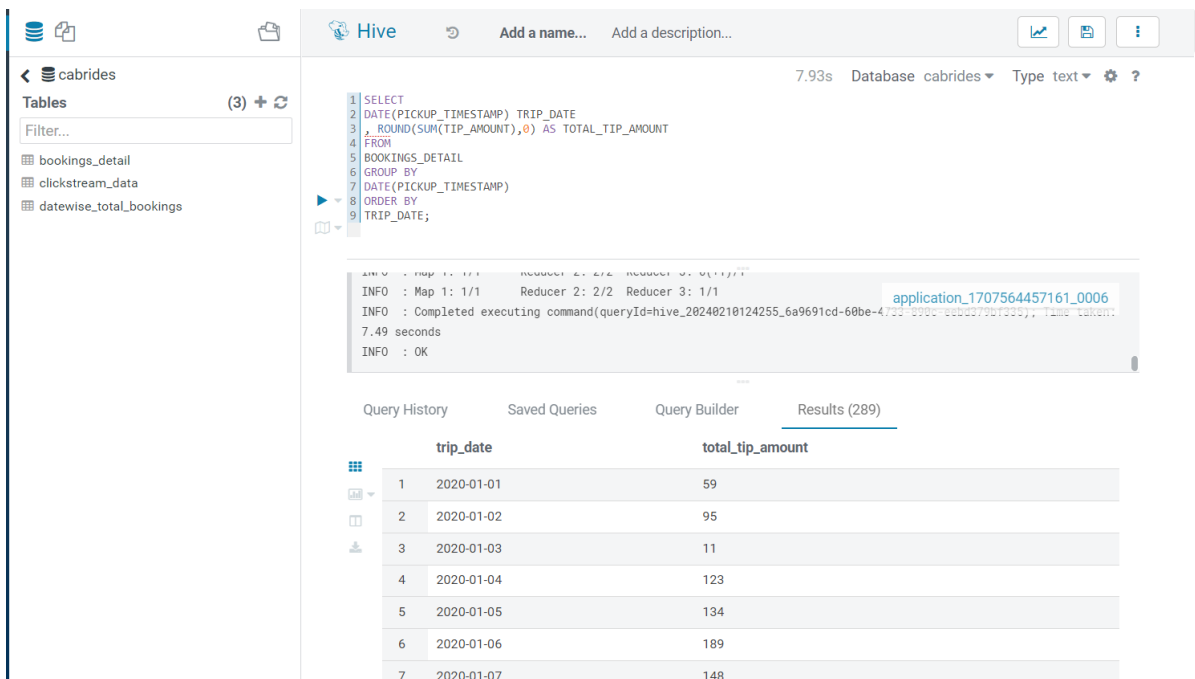
**OUTPUT:**



```
1  SELECT
2  DATE(PICKUP_TIMESTAMP) TRIP_DATE
3  , ROUND(SUM(TIP_AMOUNT),0) AS TOTAL_TIP_AMOUNT
4  FROM
5  BOOKINGS_DETAIL
6  GROUP BY
7  DATE(PICKUP_TIMESTAMP)
8  ORDER BY
9  TRIP_DATE;
```

INFO  : Map 1: 1/1      Reducer 2: 2/2  Reducer 3: 0(+1)/1
INFO  : Map 1: 1/1      Reducer 2: 2/2  Reducer 3: 1/1
INFO  : Completed executing command(queryId=hive_20240210124255_6a9691cd-60be-4... application_1707564457161_0006
7.49 seconds
INFO  : OK

| | trip_date | total_tip_amount |
|---|---|---|
| 1 | 2020-01-01 | 59 |
| 2 | 2020-01-02 | 95 |
| 3 | 2020-01-03 | 11 |
| 4 | 2020-01-04 | 123 |
| 5 | 2020-01-05 | 134 |
| 6 | 2020-01-06 | 189 |
| 7 | 2020-01-07 | 148 |

VALIDATION: Exact Match

| | |
|---|---|
| 2020-01-01 | 59 |
| 2020-01-02 | 95 |
| 2020-01-03 | 11 |
| 2020-01-04 | 123 |
| 2020-01-05 | 134 |
| 2020-01-06 | 189 |
| 2020-01-07 | 148 |
| 2020-01-08 | 111 |
| 2020-01-09 | 48 |
| 2020-01-10 | 77 |
| 2020-01-11 | 81 |
| 2020-01-12 | 109 |
| 2020-01-14 | 142 |
| 2020-01-15 | 338 |
| 2020-01-16 | 155 |
| 2020-01-17 | 296 |
| 2020-01-18 | 240 |
| 2020-01-20 | 210 |
| 2020-01-21 | 5 |
| 2020-01-23 | 148 |
| 2020-01-24 | 472 |
| 2020-01-25 | 98 |
| 2020-01-26 | 209 |
| 2020-01-27 | 231 |
| 2020-01-28 | 567 |

- **Task 10**:

Calculate the total count of all the bookings with ratings lower than 2 as given by customers in a particular month.

**QUERY:**

SELECT DATE_FORMAT(PICKUP_TIMESTAMP, 'yyyy-MM') TRIP_MONTH
, COUNT(BOOKING_ID) AS NO_OF_BOOKINGS
FROM BOOKINGS_DETAIL
WHERE RATING_BY_CUSTOMER < 2
GROUP BY DATE_FORMAT(PICKUP_TIMESTAMP, 'yyyy-MM')
ORDER BY TRIP_MONTH;

**EXPLANATION**:
- DATE_FORMAT function formats datetimestamp value in the specified format like yyyy-MM in this case which results like 2023-06.
- WHERE clause is used to filter bookings where rating given by customers is less than 2 which indicates customers dissatisfaction.
- ORDER BY clause with Trip month alias is used to show output in ascending order of pickup month.
- This analysis could help to understand number of trips by month where customers were not happy.
- Also could give insight or a hidden pattern in dissatisfactory rides in a specific month or period which could be n number of factors like low rating because of AC was not on during summery time, cab reached late on pickup point due to traffic on a rainy day/season etc.
- Based on this analysis, instructions can be given to driver to make customers happy and take care of things which could lead to low customer rating.

**OUTPUT:**



VALIDATION: Exact Match

```
Total MapReduce CPU Time Spent: 7 seconds 970 msec
OK
2020-01 26
2020-02 16
2020-03 16
2020-04 21
2020-05 21
2020-06 14
2020-07 20
2020-08 32
2020-09 21
2020-10 15
```

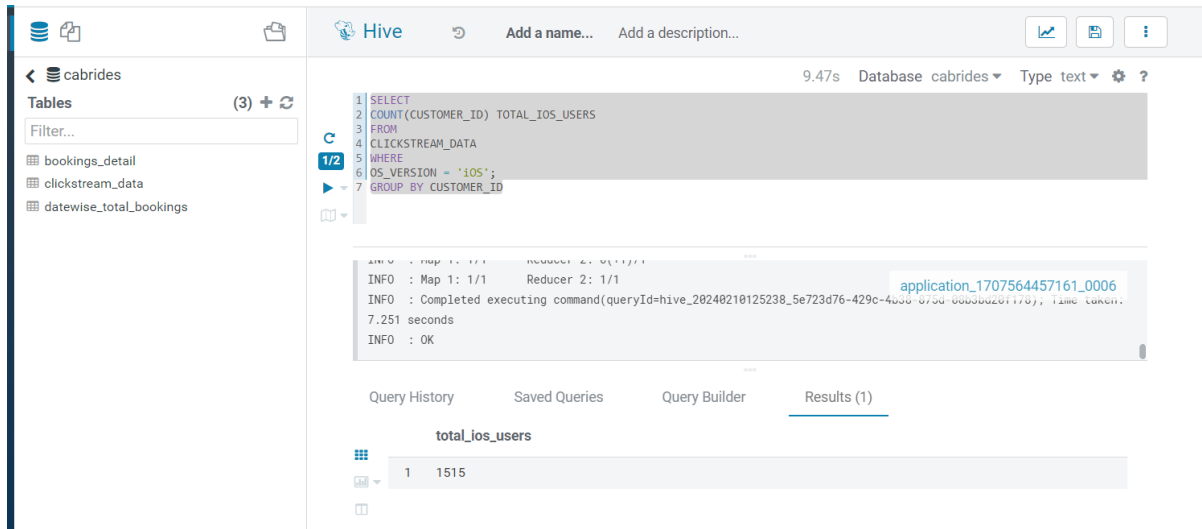- **Task 11**: Calculate the count of total iOS users.

**QUERY:**

```
SELECT COUNT(CUSTOMER_ID) TOTAL_IOS_USERS
FROM CLICKSTREAM_DATA
WHERE  OS_VERSION = 'iOS';
GROUP BY CUSTOMER_ID;
```

**EXPLANATION**:
- The DISTINCT option within the COUNT function is utilized to tally the distinct number of customers using iOS devices.
- The WHERE clause filters for events originating from iOS devices.
- The ORDER BY clause organizes the data in chronological order by the pickup month, using Trip month as an alias.
- This analysis offers a broad overview of the quantity or proportion of customers who use a particular type of device or operating system.
- For instance, should the company roll out updates to the iOS and Android mobile applications, this data could estimate the customer base that would be affected.

**OUTPUT:**



**VALIDATION:** close to Match(1515) , since Kafka had extra 16 records compare to validation You should get the count of all iOS users as 1503.