



# Logic For Final Submission

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

Printing two columns customer\_id and count of different drivers taken. Group by is done on CustomerId as we need query for each customer. The result is sorted in ascending order. > SELECT customer\_id, count(DISTINCT driver\_id) FROM bookings\_data GROUP BY customer\_id ORDER BY customer\_id ASC;

We can see the output below

```
hive> SELECT customer_id, count(DISTINCT driver_id) FROM bookings_data GROUP BY customer_id ORDER BY customer_id ASC;
Query ID = hadoop_20220123184004_57ca58fa-2d31-47c5-9da6-8fefef140d64
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1642920539798_0041)
        VERTICES
                       MODE
                                   STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
Map 1 ..... container
Reducer 2 ..... container
                                SUCCEEDED
                                                                              0
                                                                                               0
                                SUCCEEDED
                                                2
                                                                                      0
Reducer 3 ..... container
                                                                   0
                                                                           0
                                                                                      0
                                                                                               0
                              SUCCEEDED
                                          ====>>] 100% ELAPSED TIME: 4.63 s
 /ERTICES: 03/03 [==:
OK
NULL
10022393
10058402
10339567
10435129
10555335
10614890
11264797
11353346
11418437
11438890
11454977
11479815
11518953
11580321
11596512
11608791
11655671
11757536
11764909
11869278
11981042
                 1
12106105
12142182
12312603
                 1
12334699
12367832
12856708
12885363
12913608
12914577
12966909
13015449
13229062
```





### Task 6: Calculate the total rides taken by each customer.

Same as above, we print customer\_id and count of total bookings. Group by is done on customer\_id, which is then sorted.

SELECT customer\_id, COUNT(DISTINCT booking\_id) FROM bookings\_data GROUP BY customer\_id ORDER BY customer\_id ASC;

```
hive> SELECT customer_id, COUNT(DISTINCT booking_id) FROM bookings_data GROUP BY customer_id ORDER BY customer_id ASC;
Query ID = hadoop_20220123184656_cfa2832c-5c7f-4d33-9812-052ea44a2afb
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1642920539798_0042)
         VERTICES
                         MODE
                                      STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
Map 1 ...... container
Reducer 2 ..... container
Reducer 3 ..... container
                                   SUCCEEDED
                                                                           0
                                   SUCCEEDED
                                                                                     0
                                                                                                        0
OK
NULL
10022393
10058402
10339567
                  1 1 1
10435129
10555335
10592274
10614890
10678994
11264797
                  1
1
1
11353346
11418437
11438890
11454977
11479815
11518953
11580321
11596512
                  1
11608791
11655671
11757536
11764909
11860278
11981042
12106105
                  1 1 1
12142182
12312693
12334699
                  1
12367832
12856708
12885363
12913608
12914577
                  1
12966909
13015449
13229062
13262795
13356177
13387493
                  1 1 1 1
13389366
13442644
13500355
13590084
13791801
14011511
14143225
14236627
```





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.

#### Find the total visits made by each customer on the booking page

This query will show all the customers who have opened booking page with counts.

This is achieved with use of WHERE clause on page\_id.

select customer\_id, COUNT(page\_id) from click\_stream\_data where

page\_id = 'e7bc5fb2-1231-11eb-adc1-0242ac120002' GROUP BY customer\_id

ORDER BY customer\_id ASC;

# total 'Book Now' button presses

This query will show all the customers who have pressed book\_now button.

This is achieved with use of where clause on button id

select customer\_id, COUNT(button\_id) from click\_stream\_data where
button\_id = 'fcba68aa-1231-11eb-adc1-0242ac120002' GROUP BY
customer\_id ORDER BY customer\_id ASC;

#### Conversion ratio

From the output of the above queries shows that each customers has opened booking page and pressed book\_now max 1 time, we can find number of book button pressed, divided by home page visit. Following query does the same.

select x.result / y.result from (select count(customer\_id) as result
from click\_stream\_data where page\_id = 'e7bc5fb2-1231-11eb-adc10242ac120002') y join (select count(customer\_id) as result from
click\_stream\_data where button\_id = 'fcba68aa-1231-11eb-adc10242ac120002') x on 1=1;





#### **SCREENSHOTS OF THE ABOVE QUERIES**

Find the total visits made by each customer on the booking page

```
hive> select customer_id, COUNT(page_id) from click_stream_data where page_id = 'e7bc5fb2-1231-11eb-adc
1-0242ac120002' GROUP BY customer_id ORDER BY customer_id ASC;
Query ID = hadoop_20220123193554_9c91bd39-7778-4ec0-b2ce-bb0d4ea0fe51
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1642920539798_0044)
          VERTICES
                                        STATUS TOTAL COMPLETED RUNNING PENDING FAILED
                          MODE
                                                                                                    KILLED
Map 1 ..... container
                                    SUCCEEDED
                                                                                                          0
Reducer 2 ..... container
                                                      2
                                                                   2
                                                                                        0
                                                                                                 0
                                                                                                          0
                                    SUCCEEDED
                                                                             0
Reducer 3 ..... container
                                    SUCCEEDED
                                                      1
                                                                   1
                                                                             0
                                                                                        0
                                                                                                 0
                                                                                                          0
OK
10168879
                   1
10276292
                   1
10405598
                   1
10463231
                   1
10707209
                   1
10917583
                   1
10985972
                   1
                   1
11234701
11372759
                   1
11439057
11459135
                   1
11617260
11702141
11970941
                   1
12252116
                   1
12275339
                   1
12388855
                   1
12609914
12635200
                   1
12648576
12731678
                   1
13014916
13042136
                   1
                   1
13066424
13125118
13172005
                   1
13219572
13222167
                   1
13288349
13593893
13785948
                   1
13867614
13948107
                   1
14004235
                   1
14111800
14147392
                   1
14171711
                   1
14197474
                   1
14281485
                   1
14329925
```

total 'Book Now' button presses





```
hive> select customer_id, COUNT(button_id) from click_stream_data where button_id = 'fcba68aa-1231-11eb
-adc1-0242ac120002' GROUP BY customer_id ORDER BY customer_id ASC;
Query ID = hadoop_20220123193657_509a4d4b-c37d-4529-b855-4ac50828806c
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1642920539798_0044)
                                     STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
         VERTICES
                        MODE
Map 1 ..... container
                                  SUCCEEDED
                                                  1
                                                              1
                                                                        a
                                                                                  A
                                                                                           0
                                                                                                   A
Reducer 2 ..... container
                                  SUCCEEDED
                                                  2
                                                              2
                                                                        0
                                                                                  0
                                                                                           0
                                                                                                   0
Reducer 3 ..... container
                                  SUCCEEDED
                                                                                  0
                                                                                           0
                                                                                                   0
                                                  1
                                                              1
                                                                        0
VERTICES: 03/03 [=======
                                     =======>>] 100% ELAPSED TIME: 4.14 s
OK
10097931
                 1
10276292
                 1
10303507
                 1
10318382
10405598
10697432
                 1
10800309
11037726
                 1
11235483
11439057
11651952
                 1
11970941
11980742
                 1
                 1
11988474
12089943
12269901
                 1
12452446
                  1
12635200
                  1
12636650
```

#### Conversion ratio = 0.9852

hive> select x.result / y.result from (select count(customer\_id) as result from click\_stream\_data where page\_id = 'e7bc5fb2-1231-11eb-adc1-0242ac120002') y join (select count(customer\_id) as result from cli ck\_stream\_data where button\_id = 'fcba68aa-1231-11eb-adc1-0242ac120002') x on 1=1; Warning: Map Join MAPJOIN[21][bigTable=?] in task 'Reducer 2' is a cross product Query ID = hadoop\_20220123193742\_90956543-08ab-4a6d-919f-2e35c0c2af63 Total jobs = 1Launching Job 1 out of 1 Status: Running (Executing on YARN cluster with App id application\_1642920539798\_0044) VERTICES MODE STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED SUCCEEDED Map 1 ..... container SUCCEEDED Map 3 ..... container Reducer 4 ..... container SUCCEEDED Reducer 2 ..... container SUCCEEDED /ERTICES: 04/04 [=: =>>] 100% ELAPSED TIME: 5.47 s 0.985207100591716





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

We find black cabs by using WHERE clause on cab\_colour. We also group by cab\_color to count different driver\_id.

SELECT count(distinct driver\_id) FROM bookings\_data WHERE cab\_color IN ('black') GROUP BY cab\_color;

#### Screenshot:

```
[hive> SELECT count(distinct driver_id) FROM bookings_data WHERE cab_color IN ('black') GROUP BY cab_col]
or;
Query ID = hadoop_20220123194052_57103c85-d78f-4122-9653-bb8597e8ece9
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1642920539798_0044)
        VERTICES
                       MODE
                                   STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
Map 1 ..... container Reducer 2 ..... container
                                SUCCEEDED
                                                1
                                                           1
                                                                              0
                                                                                              0
                                                2
                                                           2
                                SUCCEEDED
                                                                    0
                                                                              0
                                                                                      0
                                                                                              0
                                                2
                                                           2
Reducer 3 ..... container
                                SUCCEEDED
                                                                    0
                                                                              0
                                                                                              0
VERTICES: 03/03
                                            ==>>] 100% ELAPSED TIME: 4.10 s
OK
72
Time taken: 4.9 seconds, Fetched: 1 row(s)
```



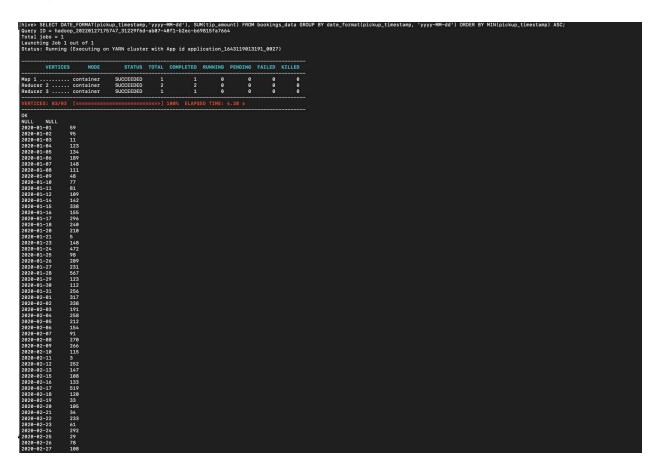


# Task 9: Calculate the total amount of tips given date wise to all drivers by customers.

We use SUM method over GROUP BY to find SUM of given INT field (here: tip\_amount). Group by is done on pickup\_timestamp.

#### QUERY:

SELECT DATE\_FORMAT(pickup\_timestamp,'yyyy-MM-dd'), SUM(tip\_amount) FROM bookings\_data GROUP BY date\_format(pickup\_timestamp, 'yyyy-MM-dd') ORDER BY MIN(pickup\_timestamp) ASC;



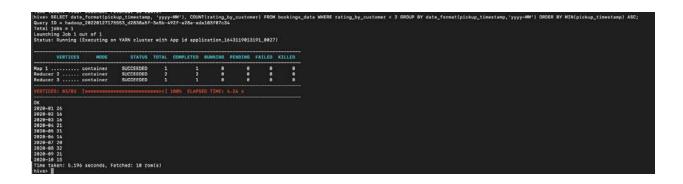




# Task 10: Calculate the total count of all the bookings with ratings lower than 2 as given by customers in a particular month.

To find rating lesser than 2, we use WHERE clause. Which is grouped for each month-year as each month uniquely identifies with month and year values.

SELECT date\_format(pickup\_timestamp, 'yyyy-MM'), COUNT(rating\_by\_customer) FROM bookings\_data WHERE rating\_by\_customer < 2 GROUP BY date\_format(pickup\_timestamp,'yyyy-MM') ORDER BY MIN(pickup\_timestamp) ASC;







# <Hive Query for Task 11>

#### Task 11: Calculate the count of total iOS users.

We use data from click\_stream\_data as there can be customers who haven't booked any cabs. We use WHERE clause on os\_version to filter the 'iOS' lines, which is grouped by to find the total count.

SELECT count(distinct customer\_id) FROM click\_stream\_data WHERE os\_version in ('iOS') GROUP BY os\_version;

```
hive> SELECT count(distinct customer_id) FROM click_stream_data WHERE os_version in ('iOS')
    > GROUP BY os_version;
Query ID = hadoop_20220123195611_f4c7058b-6726-483d-ad49-9e1f3dab84a4
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1642920539798_0044)
        VERTICES
                     MODE
                                 STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
                                                        1
2
Map 1 ..... container
                              SUCCEEDED
                                                                                  0
                                                                                          0
                                             2
                              SUCCEEDED
                                                                 0
                                                                          0
                                                                                  0
                                                                                          0
Reducer 2 ..... container
                                                        2
Reducer 3 ..... container
                              SUCCEEDED
                                             2
                                                                 0
                                                                          0
                                                                                  0
                                                                                          0
OK
1515
Time taken: 5.006 seconds, Fetched: 1 row(s)
hive>
```