



Queries

<Hive Query for Task 5>

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

SELECT customer_id, count(DISTINCT driver_id) FROM bookings_data GROUP BY customer id ORDER BY customer id ASC;

```
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
                                                              2
                                                                        0
                                                                                  0
                                                                                                    0
                                  SUCCEEDED
                                                  2
                                                                                           0
Reducer 3 ..... container
                                                                        0
                                                                                  0
                                                                                           0
                                                                                                    0
                                 SUCCEEDED
 /ERTICES: 03/03 [==:
                                             ===>>] 100% ELAPSED TIME: 4.63 s
OK
NULL
10022393
                 1
10058402
10339567
                 1
10435129
10555335
10592274
10614890
10678994
11264797
11353346
11418437
11438890
11454977
11479815
11518953
11580321
11596512
11608791
11655671
                 1
11757536
11764909
11869278
11981042
                 1
12106105
                 1
12142182
12312603
12334699
12367832
12856708
12885363
12913608
12914577
12966909
13015449
13229062
```





<Hive Query for Task 6>

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

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
                                 SUCCEEDED
                                                                                                   A
Reducer 2 ..... container Reducer 3 ..... container
                                                  2
                                  SUCCEEDED
                                                                       0
                                                                                                   0
                                  SUCCEEDED
                                                                                                   0
OK
NULL
10022393
                 1
10058402
10339567
10435129
10555335
10592274
18614898
10678994
11264797
11353346
11418437
11438890
11454977
                 1
11479815
11518953
                 1
11580321
11596512
11608791
11655671
11757536
                 1
11764909
11860278
11981042
12106105
12142182
12312603
12334699
                 1111
12367832
12856708
12885363
12913608
12914577
12966909
13015449
13229062
                 1 1 1
13262795
13356177
13387493
13389366
13442644
13500355
13590084
                 1
13791801
13798100
                 1
14011511
                 1
14143225
```





<Hive Query for Task 7>

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>

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

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

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

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;

<Screenshot after executing Query>

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)
                                    STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
         VERTICES
                       MODE
                                 SUCCEEDED
Map 1 ..... container
                                                                                         0
                                                                                                 0
                                                 1
                                                             1
                                                                       0
                                                                                 0
Reducer 2 ..... container
                                 SUCCEEDED
                                                 2
                                                                       0
                                                                                 0
                                                                                         0
                                                                                                 0
                                 SUCCEEDED
                                                                                         0
                                                                                                 0
Reducer 3 ..... container
                                                 1
                                                             1
VERTICES: 03/03 [================>>] 100% ELAPSED TIME: 4.64 s
OK
10168879
                 1
10276292
10405598
                 1
                 1
10463231
10707209
10917583
                 1
10985972
11234701
                 1
11372759
11439057
11459135
                 1
11617260
                 1
11702141
                 1
11970941
                 1
12252116
12275339
12388855
12609914
                 1
12635200
12648576
                 1
12731678
13014916
                 1
13042136
                 1
13066424
13125118
                 1
13172005
                 1
13219572
                 1
13222167
                 1
13288349
                 1
13593893
                 1
13785948
                 1
13867614
                 1
13948107
14004235
                 1
14111800
                 1
14147392
14171711
                 1
14197474
14281485
                 1
14329925
                 1
```

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)
         VERTICES
                        MODE
                                     STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
                                  SUCCEEDED
Map 1 ..... container
                                                  1
                                                              1
                                                                        a
                                                                                  A
                                                                                                    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
11988474
12089943
12269901
                 1
12452446
                  1
12635200
                  1
12636650
```

Conversion ratio

```
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 = 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
                                           SUCCEEDED
                                                                1
                                                                               1
                                                                                           0
                                                                                                        0
                                                                                                                   0
                                                                                                                              0
Map 1 ..... container
Map 3 ..... container
                                           SUCCEEDED
                                                                1
                                                                                           0
                                                                                                        0
                                                                                                                   0
                                                                                                                              0
Reducer 4 ..... container
                                           SUCCEEDED
                                                                                                        0
Reducer 2 ..... container
                                                                                           0
                                                                                                        0
                                                                                                                              0
                                           SUCCEEDED
                                                                               1
                                                                                                                   0
 ERTICES: 04/04 [=
                                                            =>>] 100% ELAPSED TIME: 5.47 s
0.985207100591716
```

<Hive Query for Task 8>

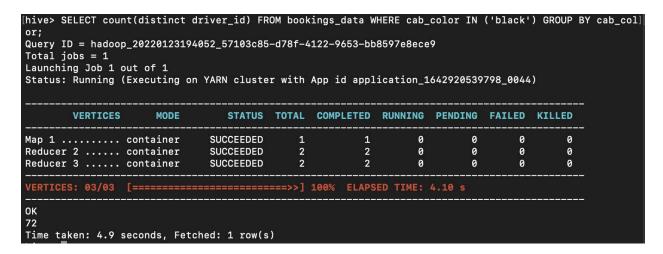
Task 8: Calculate the count of all trips done on black cabs.





SELECT count(distinct driver_id) FROM bookings_data WHERE cab_color IN ('black') GROUP BY cab_color;

<Screenshot after executing Query>



<Hive Query for Task 9>

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

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;





```
| State | Select | Se
```

<Hive Query for Task 10>

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

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;

<Screenshot after executing Query>



<Hive Query for Task 11>





Task 11: Calculate the count of total iOS users.

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
                                   STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
                      MODE
Map 1 ..... container
Reducer 2 ..... container
                                SUCCEEDED
                                                                                              0
                                               2
                                                           2
                                SUCCEEDED
                                                                    0
                                                                             0
                                                                                     0
                                                                                              0
                                                           2
                                                                    0
                                                                             0
Reducer 3 ..... container
                                SUCCEEDED
                                               2
                                                                                     0
                                                                                              0
VERTICES: 03/03 [==
OK
1515
Time taken: 5.006 seconds, Fetched: 1 row(s)
hive>
```