Analysis of Airport Data Using Hive & Pig

Case Study

19CSE357 – Big Data Analytics



Date: February 23, 2022

Group Details:

S. No	Name of the Student	Roll No.
1.	DIVESH KOSURI	CB.EN.U4CSE19422
2.	PENUGONDA KOUSHIK	CB.EN.U4CSE19449
3.	RAVELLA ABHINAV	CB.EN.U4CSE19453
4.	SINGADI SHANTHAN REDDY	CB.EN.U4CSE19459

Dataset Description

The main aim of the dataset is to develop a model for the airline data to provide a platform for new analytics based on the following queries as the problem faced is the existing has the ability to analyze limited data from the following databases

In our case study we are dealing with 3 different datasets named airports_mod, Final_airlines, routes

Fields:

Airports_mod:

- Sample: Goroka, Goroka, Papua New Guinea, GKA, AYGA, -6.081689, 145.391881, 5282, 10, U, Pacific/Port_Moresby
- Dataset contains mainly unique Airport ID, Name of the airport, City of the respective airport, Country, 3-letter IATA code, Latitude & Longitude, Altitude, Timezone

Final Airlines:

- Sample: 2,135 Airways, \N,, GNL, GENERAL, United States, N
- In this dataset it contains ID, Name of airline, Shortcut of airline, IATA, ICAO, Callsign, Country

Routes:

- Sample: 2B,410,AER,2965,KZN,2990,,0,CR2
- This dataset contains mainly 3-letter ICAO code, Airline ID, Source airport ID&Code, Destination ID & Code, Halts

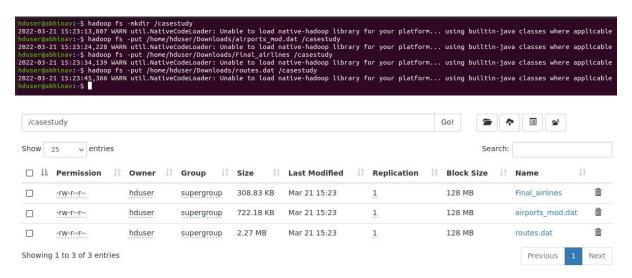
Outcome:

We tried to explore detailed analysis on airline datasets such as listing airports operations, list of airlines having no halts etc., Here we mainly focussed on the processing of big datasets using hive component of Hadoop ecosystem in distributed environment.

At last, it will be useful in accessing and processing their user queries.

Loading the Dataset:

Hive:



Pig:

QUERY 1: Load dataset andusing pig. grunt> airline = LOAD '/user/hduser/proj/airlines.txt' using PigStorage(',') as (airlineID:charArray,airline_name:charArray, airline_alias:charArray, airline_iata:charArray, airline_iata:charArray, airline_icao:charArray,callsign:charArray,territory:charArray, active:charArray);

Query 2:

Filter all the airlines with territory egypt egy = filter airline by territory == 'Egypt'; dump egy;

```
(163,Arab Agricultural Aviation Company,\N,,AGC,AGRICO,Egypt,N)
(317,AMC Airlines,\N,,AMV,,Egypt,Y)
(442,Air Sinai,\N,4D,ASD,AIR SINAI,Egypt,Y)
(610,Air Memphis,\N,,MHS,AIR KEMPHIS,Egypt,N)
(992,Alexandria Airlines,\N,,KHH,,Egypt,N)
(1015,Aleen,\N,,LEM,,Egypt,N)
(1020,Al Ahram Aviation,\N,,LHR,AL AHRAM,Egypt,N)
(1094,Air Cairo,\N,,MSC,,Egypt,N)
(1146,Al Farana Airline,\N,PHR,PHARAOH,Egypt,N)
(1603,Cairo Air Transport Company,\N,,CE,,Egypt,N)
(2143,Egyptair,\N,MS,MSR,EGYPTAIR,Egypt,Y)
(2144,Egyptair Cargo,\N,,MSX,EGYPTAIR,CARGO,Egypt,N)
(2145,Egyptian Air Force,\N,,EGY,Egypt,N)
(2146,Egyptian Aviation,\N,,EJX,,Egypt,N)
(2147,Egyptian Aviation Company,\N,,EMA,,Egypt,N)
(2369,Flash Airlines,\N,,FSH,FLASH,Egypt,N)
(2397,International Air Cargo Corporation,\N,,IAK,AIR CARGO EGYPT,Egypt,N)
```

QUERY 3: GET FIRST 3 ACTIVE AIRLINES WITH NAMES IN UPPER CASE. grunt> up_ar = foreach airline generate UPPER(airline_name),active; grunt> a_3 = limit up_ar 3; grunt> dump a_3;

(PRIVATE FLIGHT,Y) (135 AIRWAYS,N) (1TIME <u>A</u>IRLINE,Y)

Query 4:

A = filter filter_airlines by Active == 'Y'; dump A;

```
(19361, Snowbird Airlines, , S8, SBD, , Finland, Y)
(19363, Russkie Krylya, , , KRY, , Russia, Y)
(19367, Kharkiv Airlines, , KH, KHK, , Ukraine, Y)
(19433, XAIR USA, , XA, XAU, XAIR, United States, Y)
(19451, Air Costa, , LB, \N, , India, Y)
(19459, Simrik Airlines, , , RMK, , Nepal, Y)
(19473, XPTO, XPTO , XP, XPT, XPTO, Portugal, Y)
(19474, Royal Flight, , , DME, , Russia, Y)
(19525, BBN-Airways, BlackBurn, , EGH, BBN, United Kingdom, Y)
(19531, Tomsk-Avia, , , TKS, , Russia, Y)
(19541, Malawian Airlines, , , 3M, \N, , Malawi, Y)
(19548, Yeti Airlines , , , NYT, , Nepal, Y)
(19567, Avilu, Avilu' SA, . . . . . . , Switzerland, Y)
(19599, Skyline Ulasim Ticaret A.S., Skyline Ulasim Ticaret A.S., , KCU, Kocoglu, Turk ey, Y)
```

Query 5:

- 1. active_airlines = filter filter_airlines by active == 'Y';
- 2. active_airlines_usa = filter active_airlines by territory == 'United States';
- find_active_airline_names_in_usa = foreach active_airlines_usa generate airline_name;

(Aloha Airlines)
(American Airlines)
(Allegiant Air)
(Alaska Central Express)
(Air Cargo Carriers)
(Airlift International)
(America West Airlines)
(Air Wisconsin)
(Allegheny Commuter Airlines)
(AIr Sunshine)
(ATA Airlines)
(Arrow Air)
(Atlantic Southeast Airlines)
(American Eagle Airlines)

Queries:

Hive:

1. Creating table airport for airports_mod dataset:

create table airports (airport_id int,airport_name string,airport_city string,airport_country string,airport_faa string,airport_icao string,airport_lat double,airport_long double,airport_alt double,airport_timezone double,airport_dst string,airport_tz string) row format delimited fields terminated by ',';

hive> create table airports (airport_id int,airport_name string,airport_city string,airport_country string,airport_faa string,airport_icao string,airport_lat double,airport_long double,airport_ait double,airport_eit string,airport_to string,airpo

2. Creating table final airlines for Final_airlines :

create table final_airlines (airlineID string,airline_name string, airline_alias string, airline_iata string, airline_icao string,callsign string,territory string, active string) row format delimited fields terminated by ',';

hive create table final sirilines (airlineID string, airline_name string, airline_alias string, airline_iata string, airline_icao string, callsign string, territory string, active string) row format delinited fields terminated by ',';

OK

Time taken: 1.869 seconds

The taken: 1.869 seconds

The taken: 1.869 seconds

The taken: 1.869 seconds

The taken: 1.869 seconds

3. Creating table route for routes.dat:

create table routes (route_iata string,route_airid int,route_source_iata string,route_source_airid int,route_des_iata string,route_des_airid int,route_codeshare string,route_stops int,route_equip string) row format delimited fields terminated by ',';

```
hive> show tables;
OK
airports
final_airlines
routes
Time taken: 0.089 seconds, Fetched: 3 row(s)
hive>
```

4. loading data into airport table

load data inpath '/airports_mod.dat' into table airports;

```
hive> load data inpath '/casestudy/airports_mod.dat' into table airports;
Loading data to table default.airports
OK
Time taken: 1.27 seconds
hive>
```

5. loading data into final airlines table

load data inpath '/Final_airlines' into table final_airlines;

6. loading data into route table

load data inpath '/routes.dat' into table routes;

```
hive> load data inpath '/casestudy/airports_mod.dat' into table airports;
Loading data to table default.airports

OK

Time taken: 1.27 seconds
hive> load data inpath '/casestudy/Final_airlines' into table airports;
Loading data to table default.airports

OK

Time taken: 0.363 seconds
hive> load data inpath '/casestudy/routes.dat' into table airports;
Loading data to table default.airports

OK

Time taken: 0.344 seconds
hive>
```

DIVESH KOSURI – CB.EN.U4CSE19422

Hive:

1)Query should return data of all the airlines that are present in United States territory and are active.

select * from final_airlines where territory="United States" and active="Y";

Airline data required is dependent on two attributes territory and active status. So based on these two attributes with "where" clause we can get the expect output. The output data which is generated has active status="Y" and territory="United States".

2)Query should return territories with maximum airlines. select count(airlineID),territory from final_airlines group by territory order by count(airlineID) DESC;

```
MapReduce Jobs Launched:
Stage-Stage-1: HDFS Read: 632486 HDFS Write: 0 SUCCESS Stage-Stage-2: HDFS Read: 632486 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
1080
         United States
439
         Mexico
        United Kingdom
407
318
        Canada
230
        Russia
166
         Spain
131
        Germany
119
        France
93
         Australia
91
         South Africa
90
         Italy
89
        Ukraine
85
         Nigeria
79
         Kazakhstan
70
         China
70
         Sweden
60
         Switzerland
58
         Brazil
52
         Netherlands
         Austria
```

3) Query (Partition)

load data inpath '/user/hduser/proj/airlines.txt' into table final_airlines_t
partition(territory="United States");

```
htve> load data inpath "Juser/Induser/proj/airlines.txt" into table final_airlines t partition(territory="United States");
Loading data to table default.final_airlines_t partition (territory=United States)

OK

Time taken: 4.799 seconds
htve> select * from final_airlines_t limit 10;

OK

I Private flight NULL - N/A United States
2 135 Airways NULL GNL GENERAL United States United States
3 1Time Airline NULL IT ROW NEXTIME South Africa United States
4 2 Sgn No 1 Elementary Flying Training School NULL WYT United Kingdon United States
5 213 Flight Unit NULL TFU RUSSIA United States
6 223 Flight Unit State Airline NULL CHD CHALLOWSK-AVIA Russia United States
7 224th Flight Unit I NULL TFF CARGO UNIT RUSSIA United States
8 247 Jet Ltd NULL TWF CLOUD RUNNER United States
9 30 Aviation NULL SEC SECUREX United States
10 40-Mile Air NULL QS NLA MILE-AIR United States United States
```

RAVELLA ABHINAV – CB.EN.U4CSE19453

Queries:

Hive:

1. Find all the airlines that are active and have an alias names

Query:

create table alias_not_null_airlines as SELECT * FROM final_airlines WHERE airline_alias IS NOT NULL AND active="Y";

Result:

```
hive> create table alias_not_null_airlines as SELECT * FROM final_airlines WHERE airline_alias IS NOT NULL AND active="Y";

Query ID = hduser_20220322100631_7b70eb96-cde3-4913-a9cd-52c2a86c6138

Total jobs = 3

Launching Job 1 out of 3

Number of reduce tasks is set to 0 since there's no reduce operator

Job running in-process (local Hadoop)

2022-03-22 10:06:35,878 stage-1 map = 0%, reduce = 0%

2022-03-22 10:06:35,878 stage-1 map = 100%, reduce = 0%

Ended Job = job_local216938443_0001

Stage-4 is selected by condition resolver.

Stage-3 is filtered out by condition resolver.

Stage-5 is filtered out by condition resolver.

Moving data to directory hdfs://localhost:54310/user/hive/warehouse/.hive-staging_hive_2022-03-22_10-06-31_052_1734083296723987011-1/-ext-10002

Moving data to directory hdfs://localhost:54310/user/hive/warehouse/alias_not_null_airlines

MapReduce Jobs Launched:

Stage-Stage-1: HDFS Read: 10503019 HDFS Write: 21692 SUCCESS

Total MapReduce CPU Time Spent: 0 msec

OK

Time taken: 5.629 seconds

hive>
```

• Querying the first 10 rows of the resulted table:

Query: SELECT * FROM alias_not_null_airlines LIMIT 10;

```
hive> SELECT * FROM alias_not_null_airlines LIMIT 10;
                                                                                ALL NIPPON
324
         All Nippon Airways
                                   ANA All Nippon Airways NH
                                                                       ANA
                                                                                                  Japan Y
                                            AXM ASIAN EXPRESS
                                                                       Malaysia
576
        AirAsia Air Asia
         Rossiya-Russian Airlines
                                            Pulkovo Aviation Enterprise
                                                                                F۷
641
                                                                                                  PULKOVO Russia
1437
        bmi
                bmi British Midland
                                          BD
                                                     BMA
                                                              MIDLAND United Kingdom
                                                                               BEE-LINE
        Brussels Airlines
1531
                                   SN Brussels Airlines
                                                              SN
                                                                       DAT
                                                                                                  Belgium Y
        Contact Air Contactair C3
Czech Airlines CSA Czech Airlines
Emirates Emirates Airlines
1879
                                                     KIS
                                                              CONTACTAIR
                                                                                Germany Y
                                                              CSA
                                                                       CSA-LINES
                                                                                         Czech Republic Y
United Arab Emirates
1946
2183
                                                     EΚ
                                                              UAE
                                                                       EMIRATES
        easyJet EasyJet Airline U2 EZY
AirAsia X FlyAsianXpress D7
2297
                                                     EASY
                                                              United Kingdom Y
                                                              XANADU Malaysia
2417
                                                     XAX
Time taken: 0.411 seconds, Fetched: 10 row(s)
hive>
```

Explaination:

Job is to find list of airlines with alias names and are still operating (Active). This can be achieved by querying using 'WHERE', 'IS NOT NULL' and 'AND' keywords. In the dataset, all the airlines that has no alias names have 'NULL' as the value in their respective cells. So 'IS' 'NOT' 'NULL' keywords are to be used to fetch all the values rows that have alias names and active status can be directly done using 'WHERE' clause.

2. Find the count of airlines that choose to have routes with 1 stop.

Query: select count(route_airid) from routes where route_stops like "%1%"

Output:

```
hive> select count(route_airid) from routes where route_stops like "%1%";
Query ID = hduser_20220322101423_0447074f-e6cc-4088-9903-196f3f8ede71
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
 set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-03-22 10:14:25,751 Stage-1 map = 100%, reduce = 0%
2022-03-22 10:14:26,790 Stage-1 map = 100%, reduce = 100%
Ended Job = job_local1580320474_0002
MapReduce Jobs Launched:
Stage-Stage-1: HDFS Read: 25765422 HDFS Write: 43384 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
11
Time taken: 3.636 seconds, Fetched: 1 row(s)
hive>
```

Explanation:

Job here is to find the total count of airlines that has one stop in its routes. So we are to query on routes table we already created using one of the aggregate function "count".

- 1. select all the route id's which have their no of stops equal to 1
- 2. Add the aggregate function "count" to count the no of ids that are resulted as a result of first query.
- **3.** Find all airports in the world which lie at an altitude greater than 5000 ft.

Query: create table high_alt_airports as select * from airports where airport_alt > 5000;

Result:

```
hive> create table high_alt_airports as select * from airports where airport_alt > 5000;

Query ID = hduser_20220322101742_05oc5cee-fd82-4137-9922-1ebec2490a38

Total jobs = 3

Launching Job 1 out of 3

Number of reduce tasks is set to 0 since there's no reduce operator
Job running in-process (local Hadoop)

202-03-22 10:17:45,255 Stage-1 map = 100%, reduce = 0%

Ended Job = job_local2075807615_0003

Stage-4 is selected by condition resolver.

Stage-3 is filtered out by condition resolver.

Stage-5 is filtered out by condition resolver.

Moving data to directory hdfs://localhost:54310/user/hive/warehouse/.hive-staging_hive_2022-03-22_10-17-42_629_8769889742244920012-1/-ext-10002

Moving data to directory hdfs://localhost:54310/user/hive/warehouse/high_alt_airports

MapReduce Jobs Launched:

Stage-Stage-1: HDFS Read: 16313974 HDFS Write: 52226 SUCCESS

Total MapReduce CPU Time Spent: 0 msec

OK

Time taken: 3.178 seconds

htve>
```

Subquery:

select * from high_alt_airports limit 10;

Explaination:

Job is to find out the list of all the airports at higher altitudes (alt > 5000 ft) we use a binary operator ">" to select all the airports that have their airport alt > 5000.

- 1. We first use the select clause to find all the airports above air_alt > 5000, create a new table high_alt_airports and store the result of above query in that new table.
- 2. Now we query the table for 10 airports with altitude above 5000 using 'LIMIT' keyword.

P. Sai Koushik

1. Find list of Airports operating in the Country India;

create table india_opert_airport as select * from airports where airport_country LIKE '% India%';

```
hive> create table india_opert_airport as select * from airports where airport_country LIKE '%India%';
Query ID = hduser_20220322103609_79bca6ed-adca-446e-b81f-3ada52ca92cf
Total jobs = 3
Launching Job 1 out of 3
Number of reduce tasks is set to 0 since there's no reduce operator
Job running in-process (local Hadoop)
2022-03-22 10:36:13,652 Stage-1 map = 100%, reduce = 0%
Ended Job = job_local294702227_0001
Stage-4 is selected by condition resolver.
Stage-3 is filtered out by condition resolver.
Stage-5 is filtered out by condition resolver.
Moving data to directory hdfs://localhost:54310/user/hive/warehouse/.hive-staging_hive_2022-03-22_10-36-10_002_6773573062564491529-1/-ext-10002
Moving data to directory hdfs://localhost:54310/user/hive/warehouse/india_opert_airport
MapReduce Jobs Launched:
Stage-Stage-I: HDFS Read: 5649832 HDFS Write: 23793 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 3.937 seconds
hive>
```

2. select * from india_opert_airport limit 10; (to show only first 10 values due to large nos of data)

Here we display the first 10 airports operating in India.

3. Find the list of Airlines having zero stops

create table stop as select * from routes where route_stops LIKE '%0';

4. select * from stop limit 10;

SHANTHAN REDDY – CB.EN.U4CSE19459 HIVE QUERIES:

1. select * from final airlines where territory="United States" and active="Y";

Explain: Query should return data of all the airlines that are present in United States territory and are active.

Airline data required is dependent on two attributes, territory, and active status. So based on these two attributes with "where" clause we can get the expect output.

The output data which is generated has active status="Y" and territory="United States."

```
hduser@shanthan-VirtualBox: ~
19845
          FTI Fluggesellschaft
                                                                                  Germany N
Time taken: 1.672 seconds, Fetched: 6048 row(s)
hive> select * from final_airlines where territory="United States" and active="Y";
                                         Q5
AQ
NULL
10
22
          40-Mile Air NULL
Aloha Airlines NULL
                                                   MLA
                                                             MILE-AIR
                                                                                 United States
                                                   AAH
                                                             ALOHA United States
AAL AMERICAN
          American Airlines
                                                   AA
                                                                                            United States
35
          Allegiant Air
                            NULL
                                                   AAY
                                                             ALLEGIANT
                                                                                 United States
          Alaska Central Express
                                                                       ACE AIR United States
                                         NULL
                                                   KO
109
                                                             AER
          Air Cargo Carriers
                                                                        NIGHT CARGO
                                                                                            United States
                                                   2Q
          Airlift International
                                                                        AIRLIFT United States
                                                   ΗР
          America West Airlines
Air Wisconsin NULL
                                                                       CACTUS United States
281
                                         NULL
                                                             AWE
                                                   AWI
                                                             AIR WISCONSIN
                                                                                 United States
282
          Allegheny Commuter Airlines
Air Sunshine NULL
287
                                                   NULL
                                                                                  ALLEGHENY
                                                                                                      United States
                                                             AIR SUNSHINE
295
                                                   RSI
                                                                                 United States
                                                             AMTRAN United States
BIG A United States
315
          ATA Airlines
                              NULL
                                                   AMT
397
          Arrow Air
                              NULL
          Atlantic Southeast Airlines
American Eagle Airlines NULL
                                                             EV
EGF
                                                                       ASQ ACEY
EAGLE FLIGHT
452
                                                   NULL
                                                                                            United States
                                                   MQ
CYD
                                                                                            United States
659
                                                             CYCLONE United States
792
          Access Air
                              NULL
                                         ZA
          Air Florida
882
                                                             AIR FLORIDA
                                                             GIANT United States
CITRUS United States
928
          Atlas Air
                              NULL
          AirTran Airways NULL
Bemidji Airlines
1316
                                                   TRS
1442
                                         NULL
                                                             вмј
                                                                       BEMIDJI United States
          Bering Air NULL
Cape Air NULL
Chautauqua Airlines
1472
                                                             BERING AIR
                                                                                 United States
                                                                       United States
1629
                              NULL
                                                   KAP
                                                             CAIR
                                                   RP
                                                                                            United States
United States
                                         NULL
                                                             CHQ
U.S.
                                                                       CHAUTAUOUA
1739
                              NULL
1814
          Coastal Air
                                         DQ
                                                                   Virgin Islands
                                                  CJC COLGAN United States
COMAIR United States Y
1821
          Colgan Air
                              NULL
1828
          Comair NULL
                              OH
                                         COM
          CommutAir
                              NULL
                                                   UCA
                                                             COMMUTAIR
                                                                                  United States
1843
1860
          Compass Airlines
                                                   СР
                                                                        Compass Rose
                                                                                            United States
          Continental Airlines
Continental Express
Continental Micronesia
1881
                                         NULL
                                                             COA
                                                                        CONTINENTAL
                                                                                            United States
                                                                        JETLINK United States
                                                   CO
1883
                                         NULL
                                                                                            United States
1884
                                                                        AIR MIKE
          Crown Airways NULL (
Delta Air Lines NULL DL (
Evergreen International Airlines
1931
                                                   CRO
                                                             CROWN AIRWAYS United States
                                                                       United States
                                                   DAL
2009
                                                             DELTA
                                                                                            EVERGREEN
                                                                                                                 United States
2261
                                                                        ΕZ
                                                             NULL
                                                                                 EIA
2293
                                                                                  LONGHORN
                                                                                                     United States
          Express One International
                                                             ΕO
          ExpressJet
                                                             JET LINK
2295
                             NULL
                                        XF
                                                                                  United States
          Florida West International Airways
Freedom Air NULL FP FRE
Freedom Airlines NULL
                                                                                                                 United States
                                                                       RF
                                                                                            FLO WEST
2404
                                                             NULL
                                                                                 FWL
2454
                                                   FRE
                                                             FREEDOM United States
                                                                       RELEUUM AIR United States
FRONTIER FLIGHT United States
FRONTIER-AIR United Character
United States
2456
                                                             FRL
          Frontier Airlines NULL
Frontier Flying Service NULL
GoJet Airlines NULL G7
2468
                                                             FFT
                                                             FTA
2470
                                                   GJS
2577
                                                             GATEWAY United States
          Great Lakes Airlines
                                                                                            United States
2607
                                         NULL
                                                   ZK
                                                             GLA
                                                                        LAKES AIR
                                                                                 GFT
          Gulfstream International Airlines
Hageland Aviation Services NU
2645
2657
                                                             NULL
                                                                                            GULF FLIGHT
                                                                                                                 United States
                                                                                 HAGELAND
                                                                                                      United States
                                                             Нб
2688
          Hawaiian Airlines
                                                                        HAWAIIAN
                                                                                            United States
                            Horizon Airlines
                                                                                 HORIZON AIR
2778
          Horizon Air
                                                                        QXE
                                                                                                      United States
```

2. select count(airlineID),territory from final_airlines group by territory order by count(airlineID) DESC;

Explain: To find where most people depend on airways for travel ordering territories based on total number of airlines. The query orders the airlines from highest to least.

```
hduser@shanthan-VirtualBox: -
19676 Rainbow Air Polynesia Rainbow Air POL RX RPO Rainbow Air United States Y
19678 Rainbow Air US Rainbow Air US RM RNY Rainbow Air United States Y
19774 Spike Airlines Aero Spike SO SAL Spike Air United States Y
19804 All America All America A2 AL2 United States Y
Time taken: 0.453 seconds, Fetched: 141 row(s)
hive> select count(airlineID),territory from final_airlines group by territory order by count(airlineID) DESC;
Query ID = hduser_20220322101538_9cc17915-c935-47e4-8db3-e88b05999414
Total jobs = 2
Launching Job 1 out of 2
 Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
The order to set a contract number of reducers.
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
    set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-03-22 10:15:41,105 Stage-1 map = 100%, reduce = 0%
2022-03-22 10:15:42,133 Stage-1 map = 100%, reduce = 100%
Ended Job = job_local2008815317_0001
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
    set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
    set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-03-22 10:15:43,549 Stage-2 map = 100%, reduce = 100%
Ended Job = job_local748215050_0002
MapReduce Jobs Launched:
 Stage-Stage-1: HDFS Read: 1897506 HDFS Write: 0 SUCCESS Stage-Stage-2: HDFS Read: 1897506 HDFS Write: 0 SUCCESS Total MapReduce CPU Time Spent: 0 msec
  1080
                                  United States
   439
                                  United Kingdom
  407
  318
230
                                  Canada
                                  Russia
 166
                                   Spain
  131
119
                                  Germany
                                   France
 93
91
90
                                  Australia
                                   South Africa
                                   Italy
 89
85
79
                                  Ukraine
                                  Nigeria
Kazakhstan
                                  China
```

4. select airport_country,count(*) as cnt from airports group by airport_country ORDER BY cnt DESC;

```
Aduser@shanthan-VirtualBox:-

FAILED: ParseException line 4:0 missing EOF at 'select' near 'cnt'
Nives select airport_country.count(*) as cnt from airports group by airport_country ORDER BY cnt DESC;
(Query ID = Induit___202023221011__2026107-44c4-40e0_9080-74380470239

Launching Job 1 out of 2

Launching Job 2 out of 3

Launching Job 2 out of 4

Launching Job 2 out of 4

Launching Job 2 out of 5

Launching Job 2 out of 5

Launching Job 2 out of 6

Launching Job 2 out of 7

Launching Job 2 out of 8

Launching Job 2 out of 8

Launching Job 2 out of 9

Launchi
```