

## ASSIGNMENT 5.2

### Aviation Data Analysis Using Apache Pig:

There are 2 different datasets, i.e., Delayed\_Flights.csv and Airports.csv.

#### **Delayed\_Flights.csv Dataset**

There are 29 columns in this dataset. Some of them have been mentioned below:

- Year: 1987 – 2008
- Month: 1 – 12
- FlightNum: Flight number
- Canceled: Was the flight canceled?
- CancellationCode: The reason for cancellation.

#### **Airports.csv Dataset**

- iata: the international airport abbreviation code
- name of the airport
- city and country in which airport is located.
- lat and long: the latitude and longitude of the airport

#### **Problem Statement 1:**

Find out the top 5 most visited destinations.

Source Code:

```
REGISTER '/usr/local/pig/lib/piggybank.jar';
```

```
A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.
```

```
storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
```

```
Delayed_flights_limited_records = LIMIT A 10;
```

```
DUMP Delayed_flights_limited_records;
```

```

2017-12-11 04:16:59,126 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-11 04:16:59,126 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(0,2008,1,3,4,2003.0,1955,2211.0,2225,WN,335,N712SW,128.0,150.0,116.0,-14.0,8.0,IAD,TPA,810,4.0,8.0,0,N,0,,,,)
(1,2008,1,3,4,754.0,735,1002.0,1000,WN,3231,N772SW,128.0,145.0,113.0,2.0,19.0,IAD,TPA,810,5.0,10.0,0,N,0,,,,)
(2,2008,1,3,4,628.0,620,804.0,750,WN,448,N428WN,96.0,90.0,76.0,14.0,8.0,IND,BWI,515,3.0,17.0,0,N,0,,,,)
(4,2008,1,3,4,1829.0,1755,1959.0,1925,WN,3920,N464WN,90.0,90.0,77.0,34.0,34.0,IND,BWI,515,3.0,10.0,0,N,0,2.0,0.0,0.0,0.0,32.0)
(5,2008,1,3,4,1940.0,1915,2121.0,2110,WN,378,N726SW,101.0,115.0,87.0,11.0,25.0,IND,JAX,688,4.0,10.0,0,N,0,,,,)
(6,2008,1,3,4,1937.0,1830,2037.0,1940,WN,509,N763SW,240.0,250.0,230.0,57.0,67.0,IND,LAS,1591,3.0,7.0,0,N,0,10.0,0.0,0.0,0.0,47.0)
(10,2008,1,3,4,706.0,700,916.0,915,WN,100,N690SW,130.0,135.0,106.0,1.0,6.0,IND,MCO,828,5.0,19.0,0,N,0,,,,)
(11,2008,1,3,4,1644.0,1510,1845.0,1725,WN,1333,N334SW,121.0,135.0,107.0,80.0,94.0,IND,MCO,828,6.0,8.0,0,N,0,8.0,0.0,0.0,0.0,72.0)
(15,2008,1,3,4,1029.0,1020,1021.0,1010,WN,2272,N263WN,52.0,50.0,37.0,11.0,9.0,IND,MDW,162,6.0,9.0,0,N,0,,,,)
(16,2008,1,3,4,1452.0,1425,1640.0,1625,WN,675,N286WN,228.0,240.0,213.0,15.0,27.0,IND,PHX,1489,7.0,8.0,0,N,0,3.0,0.0,0.0,0.0,12.0)
grunt>

```

B = foreach A generate (int)\$1 as year, (int)\$10 as flight\_num, (chararray)\$17 as origin,(chararray)\$18 as dest;

C = filter B by dest is not null;

D = group C by dest;

E = foreach D generate group, COUNT(C.dest);

F = order E by \$1 DESC;

Result = LIMIT F 5;

A1 = load '/home/acadgild/airline\_usecase/airports.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',','NO\_MULTILINE','UNIX','SKIP\_INPUT\_HEADER');

```

2017-12-11 04:13:11,090 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-11 04:13:11,090 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(00M,Thigpen ,Bay Springs,MS,USA,31.95376472,-89.23450472)
(00R,Livingston Municipal,Livingston,TX,USA,30.68586111,-95.01792778)
(00V,Meadow Lake,Colorado Springs,CO,USA,38.94574889,-104.5698933)
(01G,Perry-Warsaw,Perry,NY,USA,42.74134667,-78.05208056)
(01J,Hilliard Airpark,Hilliard,FL,USA,30.6880125,-81.90594389)
(01M,Tishomingo County,Belmont,MS,USA,34.49166667,-88.20111111)
(02A,Gragg-Wade ,Clanton,AL,USA,32.85048667,-86.61145333)
(02C,Capitol,Brookfield,WI,USA,43.08751,-88.17786917)
(02G,Columbiana County,East Liverpool,OH,USA,40.67331278,-80.64140639)
(03D,Memphis Memorial,Memphis,MO,USA,40.44725889,-92.22696056)
grunt>

```

A2 = foreach A1 generate (chararray)\$0 as dest, (chararray)\$2 as city, (chararray)\$4 as country;

joined\_table = join Result by \$0, A2 by dest;

```
grunt> A = load '/home/acadgild/airline usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage('','NO_MULTILINE','UNIX','SKIP_INPUT_HEADER');
2017-12-14 16:49:56,532 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-14 16:49:56,532 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-14 16:49:56,532 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> B = foreach A generate (int)$1 as year, (int)$10 as flight_num, (chararray)$17 as origin,(chararray) $18 as dest;
grunt> C = filter B by dest is not null;
grunt> D = group C by dest;
grunt> E = foreach D generate group, COUNT(C.dest);
grunt> F = order E by $1 DESC;
grunt> Result = LIMIT F 5;
grunt> A1 = load '/home/acadgild/airline usecase/airports.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage('','NO_MULTILINE','UNIX','SKIP_INPUT_HEADER');
2017-12-14 17:52:48,174 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-14 17:52:48,175 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-14 17:52:48,175 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> A2 = foreach A1 generate (chararray)$0 as dest, (chararray)$2 as city, (chararray)$4 as country;
grunt>
grunt> joined_table = join Result by $0, A2 by dest;
grunt>
```

**In Line 1:** We are registering the *piggybank* jar in order to use the CSVExcelStorage class.

In relation **A**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and headers.

In relation **B**, we are generating the columns that are required for processing and explicitly typecasting each of them.

In relation **C**, we are filtering the null values from the “dest” column.

In relation **D**, we are grouping relation C by “dest.”

In relation **E**, we are generating the grouped column and the count of each.

Relation **F** and **Result** is used to order and limit the result to top 5.

These are the steps to find the top 5 most visited destinations. However, adding few more steps in this process, we will be using another table to find the city name and country as well.

In relation **A1**, we are loading another table to which we will look-up and find the city as well as the country.

In relation **A2**, we are generating dest, city, and country from the previous relation.

In relation **joined\_table**, we are joining Result and A2 based on a common column, i.e., “dest”

Finally, using dump, we are printing the result.

DUMP joined\_table;

```
2017-12-14 18:57:18,041 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-14 18:57:18,041 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(ATL,106898,ATL,Atlanta,USA)
(DEN,63003,DEN,Denver,USA)
(DFW,70657,DFW,Dallas-Fort Worth,USA)
(LAX,59969,LAX,Los Angeles,USA)
(ORD,108984,ORD,Chicago,USA)
grunt> DESCRIBE joined_table;
joined table: {Result::group: chararray,long,A2::dest: chararray,A2::city: chararray,A2::country: chararray}
grunt> █
```

## Problem Statement 2:

Which month has seen the most number of cancellations due to bad weather?

Source code:

```
REGISTER '/home/acadgild/airline_usecase/piggybank.jar';
```

```
A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
```

```
B = foreach A generate (int)$2 as month, (int)$10 as flight_num, (int)$22 as cancelled, (chararray)$23 as cancel_code;
```

```
C = filter B by cancelled == 1 AND cancel_code == 'B';
```

```
D = group C by month;
```

```
E = foreach D generate group, COUNT(C.cancelled);
```

```
F = order E by $1 DESC;
```

```
Result = limit F 1;
```

```

2017-12-14 18:59:10,753 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.userlog.retain.hours is deprecated. Instead, use mapreduce.job.userlog.retain.hours
2017-12-14 18:59:10,754 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.local.dir.minspacekill is deprecated. Instead, use mapreduce.tasktracker.local.dir.minspacekill
grunt> B = foreach A generate (int)$2 as month,(int)$10 as flight_num,(int)$22 as cancelled,(chararray)$23 as cancel_code;
grunt> C = filter B by cancelled == 1 AND cancel_code == 'B';
grunt> D = group C by month;
grunt> E = foreach D generate group, COUNT(C.cancelled);
grunt> F= order E by $1 DESC;
grunt> Result = limit F 1;
grunt>

```

**In Line 1:** We are registering *piggybank* jar in order to use the CSVExcelStorage class.

In relation **A**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and header.

In relation **B**, we are generating the columns which are required for processing and explicitly typecasting each of them.

In relation **C**, we are filtering the data based on cancellation and cancellation code, i.e., canceled = 1 means flight have been canceled and cancel\_code = 'B' means the reason for cancellation is "weather." So relation C will point to the data which consists of canceled flights due to bad weather.

In relation **D**, we are grouping the relation C based on every month.

In relation **E**, we are finding the count of canceled flights every month.

Relation **F** and **Result** is for ordering and finding the top month based on cancellation.

Finally, using dump, we are printing the result.

**DUMP Result;**

```

2017-12-14 19:01:25,285 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-14 19:01:25,285 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(12,250)
grunt> DESCRIBE Result;
Result: {group: int,long}
grunt>

```

### Problem Statement 3:

Top ten origins with the highest AVG departure delay

Source code:

```
REGISTER '/home/acadgild/airline_usecase/piggybank.jar';
```

```
A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
```

```
B1 = foreach A generate (int)$16 as dep_delay, (chararray)$17 as origin;
```

```
C1 = filter B1 by (dep_delay is not null) AND (origin is not null);
```

```
D1 = group C1 by origin;
```

```
E1 = foreach D1 generate group, AVG(C1.dep_delay);
```

```
Result = order E1 by $1 DESC;
```

```
Top_ten = limit Result 10;
```

```
Lookup = load '/home/acadgild/airline_usecase/airports.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
```

```
Lookup1 = foreach Lookup generate (chararray)$0 as origin, (chararray)$2 as city, (chararray)$4 as country;
```

```
Joined = join Lookup1 by origin, Top_ten by $0;
```

```
Final = foreach Joined generate $0,$1,$2,$4;
```

```
Final_Result = ORDER Final by $3 DESC;
```

```

grunt> A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(' ', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
2017-12-14 19:02:41,285 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-14 19:02:41,285 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-14 19:02:41,285 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> B1 = foreach A generate (int)$16 as dep_delay, (chararray)$17 as origin;
grunt> C1 = filter B1 by (dep_delay is not null) AND (origin is not null);
grunt> D1 = group C1 by origin;
grunt> E1 = foreach D1 generate group, AVG(C1.dep_delay);
grunt> Result = order E1 by $1 DESC;
grunt> Top_ten = limit Result 10;
grunt> Lookup = load '/home/acadgild/airline_usecase/airports.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(' ', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
2017-12-14 19:03:36,681 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-14 19:03:36,681 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-14 19:03:36,681 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> Lookup1 = foreach Lookup generate (chararray)$0 as origin, (chararray)$2 as city, (chararray)$4 as country;
grunt> Joined = join Lookup1 by origin, Top_ten by $0;
grunt> Final = foreach Joined generate $0,$1,$2,$4;
grunt> Final_Result = ORDER Final by $3 DESC;
grunt>

```

Explanation of first 3 lines are the same as explained in the previous 2 problem statements.

In relation **C1**, we are removing the null values fields present if any.

In relation **D1**, we are grouping the data based on column “origin.”

In relation **E1**, we are finding average delay from each unique origin.

Relations named **Result** and **Top\_ten** are ordering the results in descending order and printing the top ten values.

These steps are good enough to find the top ten origins with the highest average departure delay.

However, rather than generating just the code of origin, we will be following a few more steps to find some more details like country and city.

In the relation **Lookup**, we are loading another table to which we will look up and find the city as well as the country.

In the relation **Lookup1**, we are generating the destination, city, and country from the previous relation.

In the relation **Joined**, we are joining relation **Top\_ten** and **Lookup1** based on common a column, i.e., “origin.”

In the relation **Final**, we are generating required columns from the **Joined** table.

Finally, using **dump**, we are printing the result.

DUMP Final\_Result;

```

2017-12-14 19:05:28,692 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-14 19:05:28,692 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(CMX,Hancock,USA,116.1470588235294)
(PLN,Pellston,USA,93.76190476190476)
(SPI,Springfield,USA,83.84873949579831)
(ALO,Waterloo,USA,82.2258064516129)
(MQT,NA,USA,79.55665024630542)
(ACY,Atlantic City,USA,79.3103448275862)
(MOT,Minot,USA,78.66165413533835)
(HHH,NA,USA,76.53005464480874)
(EGE,Eagle,USA,74.12891986062718)
(BGM,Binghamton,USA,73.15533980582525)
grunt> DESCRIBE Final Result;
Final Result: {Lookup1::origin: chararray,Lookup1::city: chararray,Lookup1::country: chararray,double}
grunt> █

```

## Problem Statement 4:

Which route (origin & destination) has seen the maximum diversion?

Source code:

```
REGISTER '/home/acadgild/airline_usecase/piggybank.jar';
```

```
A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
```

```
B = FOREACH A GENERATE (chararray)$17 as origin, (chararray)$18 as dest, (int)$24 as diversion;
```

```
C = FILTER B BY (origin is not null) AND (dest is not null) AND (diversion == 1);
```

```
D = GROUP C by (origin,dest);
```

```
E = FOREACH D generate group, COUNT(C.diversion);
```

```
F = ORDER E BY $1 DESC;
```

```
Result = limit F 10;
```



```

grunt> A = load '/home/acadqild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage('
', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
2017-12-14 19:06:34,174 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is depre
cated. Instead, use mapreduce.job.counters.max
2017-12-14 19:06:34,175 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated.
Instead, use dfs.bytes-per-checksum
2017-12-14 19:06:34,175 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instea
d, use fs.defaultFS
grunt> B = FOREACH A GENERATE (chararray)$17 as origin, (chararray)$18 as dest, (int)$24 as diversion;
grunt> C = FILTER B BY (origin is not null) AND (dest is not null) AND (diversion == 1);
grunt> D = GROUP C by (origin,dest);
grunt> E = FOREACH D generate group, COUNT(C.diversion);
grunt> F = ORDER E BY $1 DESC;
grunt> Result = limit F 10;
grunt>

```

**In Line 1:** We are registering *piggybank* jar in order to use CSVExcelStorage class.

In relation **A**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and headers.

In relation **B**, we are generating the columns which are required for processing and explicitly type-casting each of them.

In relation **C**, we are filtering the data based on “not null” and diversion =1. This will remove the null records, if any, and give the data corresponding to the diversion taken.

In relation **D**, we are grouping the data based on origin and destination.

Relation **D** finds the count of diversion taken per unique origin and destination.

Relations **F** and **Result** orders the result and produces top 10 results.

Finally, using dump, we are printing the result.

**DUMP Result;**

```

2017-12-14 19:08:20,758 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-14 19:08:20,758 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to pro
cess : 1
((ORD,LGA),39)
((DAL,HOU),35)
((DFW,LGA),33)
((ATL,LGA),32)
((ORD,SNA),31)
((SLC,SUN),31)
((MIA,LGA),31)
((BUR,JFK),29)
((HRL,HOU),28)
((BUR,DFW),25)
grunt> DESCRIBE Result;
Result: {group: (origin: chararray,dest: chararray),long}
grunt>

```