

# Problem Statement 1

Find out the top 5 most visited destinations.

```
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)$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');
```

```
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;
```

```
dump joined_table;
```

```
grunt> REGISTER '/home/acadgild/airline_usecase/piggybank.jar';  
2018-05-15 16:40:50,885 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum  
2018-05-15 16:40:50,892 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS  
2018-05-15 16:40:50,906 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 101: file '/home/acadgild/airline_usecase/piggybank.jar' does not exist.  
Details at logfile: /home/acadgild/Desktop/pig_1526382604355.log
```

```

2018-05-15 16:44:44,512 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 16:44:44,518 [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');
2018-05-15 16:51:22,437 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 16:51:22,437 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS

```

```

grunt> A1 = load '/home/acadgild/airline_usecase/airports.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage('','NO_MULTILINE','UNIX','SKIP_INPUT_HEADER');
2018-05-15 16:51:22,437 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 16:51:22,437 [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> joined_table = join Result by $0, A2 by dest;
grunt> dump joined_table;
2018-05-15 16:51:59,114 [main] WARN org.apache.hadoop.util.NativeCodeLoader - Unable to load native-hadoop library for your platform... using builtin-java c

```

```

2018-05-15 16:52:46,252 [main] INFO org.apache.pig.piggybank.storage.CSVExcelStorage - (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> █

```

## Problem Statement 2

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

**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;**

**dump Result;**

```
grunt> REGISTER '/home/acadgild/airline_usecase/piggybank.jar';
2018-05-15 16:58:33,965 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 16:58:33,965 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-15 16:58:33,970 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 101: file '/home/acadgild/airline_usecase/piggybank.jar' does not exist.
Details at logfile: /home/acadgild/Desktop/pig_1526382604355.log
grunt> A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
2018-05-15 16:58:48,805 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 16:58:48,812 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
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>
grunt> D = group C by month;
grunt>
grunt> E = foreach D generate group, COUNT(C.cancelled);
grunt>
grunt> F = order E by $1 DESC;
grunt> Result = limit F 1;
grunt> dump Result;
```

```
2018-05-15 17:00:19,
(12,250)
grunt> █
```

---

## Problem Statement 3

Top ten origins with the highest AVG departure delay

```
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\_INPU  
T\_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;

dump Final\_Result;

```
grunt> REGISTER '/home/acadgild/airline_usecase/piggybank.jar';
2018-05-15 17:02:18,333 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 17:02:18,333 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-15 17:02:18,342 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 101: file '/home/acadgild/airline_usecase/piggybank.jar' does not exist.
Details at logfile: /home/acadgild/Desktop/pig_1526382604355.log
grunt> A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_IN
PUT_HEADER');
2018-05-15 17:02:31,391 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 17:02:31,392 [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>
grunt> C1 = filter B1 by (dep_delay is not null) AND (origin is not null);
grunt>
grunt> D1 = group C1 by origin;
grunt>
grunt> E1 = foreach D1 generate group, AVG(C1.dep_delay);
grunt>
grunt> Result = order E1 by $1 DESC;
grunt>
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_IN
PUT_HEADER');
2018-05-15 17:03:05,689 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 17:03:05,689 [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>
grunt> Joined = join Lookup1 by origin, Top_ten by $0;
grunt>
grunt> Final = foreach Joined generate $0,$1,$2,$4;
grunt>
grunt> Final_Result = ORDER Final by $3 DESC;
grunt> dump Final_Result;
```

```
2018-05-15 17:03:07,007 [main] INFO org
(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>
```

## Problem Statement 4

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

```
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;
```

```
dump Result;
```

```
grunt> REGISTER '/home/acadgild/airline_usecase/piggybank.jar';
2018-05-15 17:10:15,739 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 17:10:15,740 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-15 17:10:15,746 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 101: file '/home/acadgild/airline_usecase/piggybank.jar' does not exist.
Details at logfile: /home/acadgild/Desktop/pig_1526382604355.log
grunt> A = load '/home/acadgild/airline_usecase/DelayedFlights.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');
2018-05-15 17:10:32,109 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-15 17:10:32,109 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt>
grunt> B = FOREACH A GENERATE (chararray)$17 as origin, (chararray)$18 as dest, (int)$24 as diversion;
grunt>
grunt> C = FILTER B BY (origin is not null) AND (dest is not null) AND (diversion == 1);
grunt>
grunt> D = GROUP C by (origin,dest);
grunt>
grunt> E = FOREACH D generate group, COUNT(C.diversion);
grunt>
grunt> F = ORDER E BY $1 DESC;
grunt> Result = limit F 10;
grunt> dump Result;
```

```
2018-05-15 17:11:
((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)
```