

## CSP554—Big Data Technologies

### Assignment #5 (Modules 05)

Worth: 15 points

```
[hadoop@ip-172-31-34-212 ~]$ ls
TestDataGen.class
[hadoop@ip-172-31-34-212 ~]$ java TestDataGen
Magic Number = 155384
[hadoop@ip-172-31-34-212 ~]$ ls
foodplaces155384.txt  foodratings155384.txt  TestDataGen.class
```

**Magic Number: 155384**

Exercise 1) 2 points

Create new versions of the foodratings and foodplaces files by using TestDataGen (as described in assignment #4) and copy them to HDFS (say into /user/hadoop).

Write and execute a sequence of pig latin statements that loads the foodratings file as a relation. Call the relation 'food\_ratings'. The load command should associate a schema with this relation where the first attribute is referred to as 'name' and is of type chararray, the next attributes are referred to as 'f1' through 'f4' and are of type int, and the last field is referred to as 'placeid' and is also of type int.

Execute the describe command on this relation.

Provide the magic number, the load command you wrote and the output of the describe command as the result of this exercise.

**Magic Number: 155384**

```
grunt> food_ratings = LOAD '/user/hadoop/foodratings155384.txt' USING PigStorage(',')
>> AS (
>> name:chararray,
>> f1:int,
>> f2:int,
>> f3:int,
>> f4:int,
>> placeid:int
>> );
21/02/28 23:11:59 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled

grunt> DESCRIBE food_ratings
food_ratings: {name: chararray,f1: int,f2: int,f3: int,f4: int,placeid: int}
```

## Exercise 2) 2 points

Now create another relation with two fields of the initial (food\_ratings) relation: 'name' and 'f4'. Call this relation 'food\_ratings\_subset'.

Store this last relation, food\_ratings\_subset, back to HDFS (perhaps as the file /user/hadoop/fr\_subset)

Also write 6 records of this relation out to the console.

Submit the pig latin statements you used and the six records printed out to the console as the result of this exercise.

**Magic Number: 155384**

```
grunt> food_ratings_subset = FOREACH food_ratings GENERATE name, f4;
```

```
grunt> STORE food_ratings_subset INTO 'fr_subset' USING PigStorage('|');
```

```
21/02/28 23:18:43 INFO tez.TezPigScriptStats: Script Statistics:
```

```
HadoopVersion: 2.10.1-amzn-0
PigVersion: 0.17.0
TezVersion: 0.9.2
UserId: hadoop
FileName:
StartedAt: 2021-02-28 23:18:23
FinishedAt: 2021-02-28 23:18:43
Features: UNKNOWN
```

Success!

```
DAG 0:
      Name: PigLatin:DefaultJobName-0_scope-0
      ApplicationId: job_1614552901913_0001
      TotalLaunchedTasks: 1
      FileBytesRead: 0
      FileBytesWritten: 0
      HdfsBytesRead: 17466
      HdfsBytesWritten: 7027
      SpillableMemoryManager spill count: 0
      Bags proactively spilled: 0
      Records proactively spilled: 0
```

```
DAG Plan:
Tez vertex scope-9
```

```
Vertex Stats:
VertexId Parallelism TotalTasks InputRecords ReduceInputRecords OutputRecords FileBytesRead FileBytesWritten HdfsBytesRead HdfsBytesWritten Alias Feature Outputs
scope-9 1 1 1000 0 1000 0 0 17466 7027 food_ratings,food_ratings_subset hdfs://ip-172-31-34-212.us-east-2.compute.internal:8020/user/hadoop/fr_subset,
```

```
Input(s):
Successfully read 1000 records (17466 bytes) from: "/user/hadoop/foodratings155384.txt"
```

```
Output(s):
Successfully stored 1000 records (7027 bytes) in: "hdfs://ip-172-31-34-212.us-east-2.compute.internal:8020/user/hadoop/fr_subset"
```

```
grunt> food_ratings_subset_6 = LIMIT food_ratings_subset 6;
```

```
grunt> DUMP food_ratings_subset_6;
```

```
(Joy,8)
(Mel,13)
(Joy,15)
(Jill,42)
(Joe,21)
(Mel,29)
```

### Exercise 3) 2 points

Now create another relation using the initial (food\_ratings) relation. Call this relation 'food\_ratings\_profile'. The new relation should only have one record. This record should hold the minimum, maximum and average values for the attributes 'f2' and 'f3'. (So this one record will have 6 fields).

Write the record of this relation out to the console.

Submit the pig latin statements you used and the record printed out to the console as the result of this exercise.

**Magic Number: 155384**

```
grunt> food_ratings_group_all = GROUP food_ratings ALL;
grunt> food_ratings_profile = FOREACH food_ratings_group_all GENERATE MIN(food_ratings.f2) AS min_f2, MAX(food_ratings.f2) AS max_f2, AVG(food_ratings.f2) AS avg_f2, MIN(food_ratings.f3) AS min_f3, MAX(food_ratings.f3) AS max_f3, AVG(food_ratings.f3) AS avg_f3;
```

```
grunt> DESCRIBE food_ratings_profile;
food_ratings_profile: {min_f2: int,max_f2: int,avg_f2: double,min_f3: int,max_f3: int,avg_f3: double}
```

```
grunt> DUMP food_ratings_profile;
```

```
Input(s):
Successfully read 1000 records (17466 bytes) from: "/user/hadoop/foodratings155384.txt"

Output(s):
Successfully stored 1 records (28 bytes) in: "hdfs://ip-172-31-34-212.us-east-2.compute.internal:8020/tmp/temp-2142025728/tmp-1889889003"

21/02/28 23:35:49 INFO input.FileInputFormat: Total input files to process : 1
1666157 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
21/02/28 23:35:49 INFO util.MapRedUtil: Total input paths to process : 1
(1,50,24.948,1,50,24.85)
```

### Exercise 4) 2 points

Now create yet another relation from the initial (food\_ratings) relation. This new relation should only include tuples (records) where f1 < 20 and f3 > 5. Call this relation 'food\_ratings\_filtered'.

Write 6 records of this relation out to the console.

Submit the pig latin statements you used and the six records printed out to the console as the result of this exercise.

**Magic Number: 155384**

```
grunt> food_ratings_filtered = FILTER food_ratings BY (f1 < 20) AND (f3 > 5);
grunt> DESCRIBE food_ratings_filtered;
food_ratings_filtered: {name: chararray,f1: int,f2: int,f3: int,f4: int,placeid: int}
grunt> food_ratings_filtered_6 = LIMIT food_ratings_filtered 6;
grunt> DESCRIBE food_ratings_filtered_6;
food_ratings_filtered_6: {name: chararray,f1: int,f2: int,f3: int,f4: int,placeid: int}
grunt> DUMP food_ratings_filtered_6;
```

```

Input(s):
Successfully read 16 records (17466 bytes) from: "/user/hadoop/foodratings155384.txt"

Output(s):
Successfully stored 6 records (117 bytes) in: "hdfs://ip-172-31-34-212.us-east-2.compute.internal:8020/tmp/temp-2142025728/tmp-223768638"

21/02/28 23:45:20 INFO input.FileInputFormat: Total input files to process : 1
2237874 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
21/02/28 23:45:20 INFO util.MapRedUtil: Total input paths to process : 1
(Mel,13,1,47,29,4)
(Joe,19,16,45,5,2)
(Joy,15,33,33,20,1)
(Sam,5,41,49,43,5)
(Mel,12,4,41,5,4)
(Joy,5,50,37,9,1)

```

#### Exercise 5) 2 points

Using the initial (food\_ratings) relation, write and execute a sequence of pig latin statements that creates another relation, call it 'food\_ratings\_2percent', holding a random selection of 2% of the records in the initial relation.

Write 10 of the records out to the console.

Submit the pig latin statements and the records printed out to the console.

**Magic Number: 155384**

```

grunt> food_ratings_2percent = SAMPLE food_ratings 0.02;
grunt> DESCRIBE food_ratings_2percent;
food_ratings_2percent: {name: chararray,f1: int,f2: int,f3: int,f4: int,placeid: int}
grunt> food_ratings_2percent_10 = LIMIT food_ratings_2percent 10;
grunt> DESCRIBE food_ratings_2percent_10;
food_ratings_2percent_10: {name: chararray,f1: int,f2: int,f3: int,f4: int,placeid: int}
grunt> DUMP food_ratings_2percent_10;

```

```

Input(s):
Successfully read 791 records (17466 bytes) from: "/user/hadoop/foodratings155384.txt"

Output(s):
Successfully stored 10 records (200 bytes) in: "hdfs://ip-172-31-34-212.us-east-2.compute.internal:8020/tmp/temp-2142025728/tmp2134899021"

21/02/28 23:52:53 INFO input.FileInputFormat: Total input files to process : 1
2690122 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
21/02/28 23:52:53 INFO util.MapRedUtil: Total input paths to process : 1
(Joe,49,27,39,33,1)
(Sam,13,17,7,45,4)
(Jill,34,43,38,34,2)
(Joe,22,11,19,9,1)
(Jill,44,35,45,24,3)
(Sam,40,33,16,6,2)
(Joy,9,10,44,31,5)
(Joe,32,32,9,18,1)
(Mel,50,10,9,22,5)
(Jill,34,7,22,22,4)

```

#### Exercise 6) 2 points

Write and execute a sequence of pig latin statements that loads the foodplaces file as a relation. Call the relation 'food\_places'. The load command should associate a schema with this relation where the first

attribute is referred to as 'placeid' and is of type int and the second attribute is referred to as 'placename' and is of type chararray.

Execute the describe command on this relation.

Now perform a join between the initial place\_ratings relation and the food\_places relation on the placeid attributes to create a new relation called 'food\_ratings\_w\_place\_names'. This new relation should have all the attributes (columns) of both relations. The new relation will allow us to work with place ratings and place names together.

Write 6 records of this relation out to the console.

Submit the pig latin statements you used and the six records printed out to the console as the result of this exercise.

**Magic Number: 155384**

```
grunt> food_places = LOAD '/user/hadoop/foodplaces155384.txt' USING PigStorage(',')
>> AS (
>> placeid:int,
>> placename:chararray
>> );
21/02/28 23:58:36 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
grunt> DESCRIBE food_places;
food_places: {placeid: int,placename: chararray}
```

```
grunt> food_ratings_w_place_names = JOIN food_ratings BY placeid, food_places BY placeid;
grunt> DESCRIBE food_ratings_w_place_names;
food_ratings_w_place_names: {food_ratings::name: chararray,food_ratings::f1: int,food_ratings::f2: int,food_ratings::f3: int,food_ratings::f4: int,food_ratings::placeid: int,food_places::placeid: int,food_places::placename: chararray}
```

```
grunt> food_ratings_w_place_names_6 = LIMIT food_ratings_w_place_names 6;
grunt> DESCRIBE food_ratings_w_place_names_6;
food_ratings_w_place_names_6: {food_ratings::name: chararray,food_ratings::f1: int,food_ratings::f2: int,food_ratings::f3: int,food_ratings::f4: int,food_ratings::placeid: int,food_places::placeid: int,food_places::placename: chararray}
grunt> DUMP food_ratings_w_place_names_6;
```

```
Input(s):
Successfully read 1000 records (17466 bytes) from: "/user/hadoop/foodratings155384.txt"
Successfully read 5 records (59 bytes) from: "/user/hadoop/foodplaces155384.txt"

Output(s):
Successfully stored 6 records (210 bytes) in: "hdfs://ip-172-31-34-212.us-east-2.compute.internal:8020/tmp/temp-2142025728/tmp-556639745"

21/03/01 00:07:50 INFO input.FileInputFormat: Total input files to process : 1
3587404 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
21/03/01 00:07:50 INFO util.MapRedUtil: Total input paths to process : 1
(Joy,48,18,49,25,1,1,China Bistro)
(Joe,44,31,15,31,1,1,China Bistro)
(Sam,43,12,9,24,1,1,China Bistro)
(Sam,21,45,4,8,1,1,China Bistro)
(Joe,38,21,47,3,1,1,China Bistro)
(Sam,24,25,27,27,1,1,China Bistro)
```

Exercise 7) (3 points) Identify the one correct answer for each the following questions. These questions are similar to the ones you might find on the mid-term covering Pig. Each is worth ½ point.

- I. Which keyword is used to select a certain number of rows from a relation when forming a new relation?

Answer:   A  

Choices:

**A. LIMIT**

B. DISTINCT

C. UNIQUE

D. SAMPLE

II. Which keyword returns only unique rows for a relation when forming a new relation?

Choices:

Answer:   C  

A. SAMPLE

B. FILTER

**C. DISTINCT**

D. SPLIT

III. Assume you have an HDFS file with a large number of records similar to the examples below

- Mel, 1, 2, 3
- Jill, 3, 4, 5

Which of the following would NOT be a correct pig schema for such a file?

Choices:

Answer:   B  

A. (f1: CHARARRAY, f2: INT, f3: INT, f4: INT)

**B. (f1: STRING, f2: INT, f3: INT, f4: INT)**

C. (f1, f2, f3, f4)

D. (f1: BYTEARRAY, f2: INT, f3: BYTEARRAY, f4: INT)

IV. Which one of the following statements would create a relation (relB) with two columns from a relation (relA) with 4 columns? Assume the pig schema for relA is as follows:

(f1: INT, f2, f3, f4: FLOAT)

Answer:   B  

Choices:

A. relB = GROUP relA GENERATE f1, f3;

**B. relB = FOREACH relA GENERATE \$0, f3;**

C. relB = FOREACH relA GENERATE f1, f5;

D. relB = FOREACH relA SELECT f1, f3;

V. Pig Latin is a   B   language. Select the best choice to fill in the blank.

Choices:

- A. functional
- B. data flow**
- C. procedural
- D. declarative

VI. Given a relation (relA) with 4 columns and pig schema as follows: (f1: INT, f2, f3, f4: FLOAT) which one statement will create a relation (relB) having records all of whose first field is less than 20

Answer: **A**

Choices:

- A. relB = FILTER relA by \$0 < 20**
- B. relB = GROUP relA by f1 < 20
- C. relB = FILTER relA by \$1 < 20
- D. relB = FOREACH relA GENERATE f1 < 20