

```
rohit@DESKTOP-CBH603S MINGW64 ~  
$ ssh -i /c/Users/rohit/OneDrive/Desktop/all_files/BIGDATA/emrkey-pair.cer hadoop  
p@ec2-100-26-175-93.compute-1.amazonaws.com  
The authenticity of host 'ec2-100-26-175-93.compute-1.amazonaws.com (100.26.175.  
93)' can't be established.  
ED25519 key fingerprint is SHA256:gjxi3JYc+WJOJJ2fG5CmJSzjV7L8m8ATlwiIvGgR8.  
This key is not known by any other names.  
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes  
Warning: Permanently added 'ec2-100-26-175-93.compute-1.amazonaws.com' (ED25519)  
to the list of known hosts.  
  
      _ |   _ |_)  
     _| \___/_|_ Amazon Linux 2 AMI  
  
https://aws.amazon.com/amazon-linux-2/  
  
EEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRRRR  
E:::EEEEEEEEEEEEEE E M:::MM M:::MM R:::R  
EE:::EEEEEEEEEEEEEE E M:::MM M:::MM R:::RRRRRR:::R  
E:::E EEEEE M:::MM M:::MM RR::R R:::R  
E:::E M:::MM:M M:::MM M:::MM R::R R:::R  
E:::EEEEEEEEEEEE M:::M M:::M M:::M R::RRRRRR::R  
E:::EEEEEEEEEEEE M:::M M:::MM M:::MM R::RRRRRR::R  
E:::E M:::MM M:::M M:::MM R::R R:::R  
E:::E EEEEE M:::MM MMM M:::MM R::R R:::R  
EE:::EEEEEEEEEEEE M:::M M:::MM R::R R:::R  
E:::EEEEEEEEEEEE M:::M M:::MM RR::R R:::R  
EEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRR RRRRRR  
  
[hadoop@ip-172-31-60-217 ~]$ java TestDataGen  
Magic Number = 161793  
[hadoop@ip-172-31-60-217 ~]$  
[hadoop@ip-172-31-60-217 ~]$ ls  
foodplaces161793.txt foodratings161793.txt TestDataGen.class  
[hadoop@ip-172-31-60-217 ~]$ =
```

```
hadoop@ip-172-31-60-217:~/pigdemo
grunt>
grunt>
grunt> food_ratings = LOAD '/user/hadoop/foodratings161793.txt' USING P
igStorage(',') AS (name:chararray, f1:int, f2:int, f3:int, f4:int, plac
eid:int);
23/03/05 03:55:39 INFO Configuration.deprecation: yarn.resourcemanager.
system-metrics-publisher.enabled is deprecated. Instead, use yarn.syste
m-metrics-publisher.enabled
grunt>
grunt> DESCRIBE food_ratings;
food_ratings: {name: chararray,f1: int,f2: int,f3: int,f4: int,placeid:
int}
grunt>
```

## Exercise 2)

```

food_ratings_subset = FOREACH food_ratings GENERATE name,f4;
STORE food_ratings_subset INTO '/user/hadoop/fr_subset' USING PigStorage(',');
fr_output = LIMIT food_ratings_subset 6;
dump fr_output;

```

```

(Me1,19)
(Jill,11)
(Me1,2)
(Jill,29)
(Joe,47)
(Joe,40)
grunt>
grunt>
grunt>
grunt>
grunt>

```

## Exercise 3)

```

fr_profile = GROUP food_ratings ALL;
food_ratings_profile = FOREACH fr_profile GENERATE MIN(food_ratings.f2), MAX(food_ratings.f2),
AVG(food_ratings.f2), MIN(food_ratings.f3),MAX(food_ratings.f3), AVG(food_ratings.f3);
DUMP food_ratings_profile;

```

```

23/03/05 04:31:02 INFO util.MapRedUtil: Total input paths to process : 1
(1,50,25.105,1,50,25.165)
grunt>

```

## Exercise 4)

```

food_ratings_filtered = FILTER food_ratings BY (f1<20) AND (f3>5);
fr_filtered = LIMIT food_ratings_filtered 6;
DUMP fr_filtered;

```

```

(Me1,1,18,48,19,2)
(Me1,19,3,24,20,4)
(Joe,13,33,30,6,3)
(Joe,4,19,39,1,5)
(Jill,11,46,21,36,3)
(Jill,17,21,47,15,3)
grunt>

```

## Exercise 5)

```

food_ratings_2percent = SAMPLE food_ratings 0.02;
filtered = LIMIT food_ratings_2percent 10;
DUMP filtered;

```

```

(1,China Bistro,Sam,18,5,4,38,1)
(1,China Bistro,Joe,8,43,29,25,1)
(1,China Bistro,Sam,12,7,16,48,1)
(1,China Bistro,Mel,20,11,6,24,1)
(1,China Bistro,Joe,40,47,3,8,1)
(1,China Bistro,Sam,39,11,2,6,1)
grunt>

```

### Exercise 6)

```
food_places = LOAD '/user/hadoop/foodplaces161793.txt' USING PigStorage(',') AS (placeid: int,
placeid: int,
placename: chararray);
```

```
DESCRIBE food_places;
```

```

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;
```

```
fr_result= LIMIT food_ratings_w_place_names 6;
```

```
DUMP fr_result;
```

```

(1,China Bistro,Sam,18,5,4,38,1)
(1,China Bistro,Joe,8,43,29,25,1)
(1,China Bistro,Sam,12,7,16,48,1)
(1,China Bistro,Mel,20,11,6,24,1)
(1,China Bistro,Joe,40,47,3,8,1)
(1,China Bistro,Sam,39,11,2,6,1)
grunt>

```

**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.

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

Answer: LIMIT

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

Answer: DISTINCT

- 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?

Answer: **(f1, f2, f3, f4)** is incorrect because a pig file cannot load data properly without explicitly mentioning of the data types for the columns.

At the same time (f1:byte array, f2:integer, f3:byte array, f4:integer) also doesn't seem correct because the f3 column in the file contains an integer in decimal format but the data type used is byte array. This can cause errors while performing arithmetic operations as it pig has to convert this data using the ASCII values for storing and there is a possibility of errors while this is done.

But option A and B looks same and pig accepts Char array and String interchangeably.

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

5. Pig Latin is a \_\_\_\_\_ language. Select the best choice to fill in the blank.

Answer: **data flow**

6. 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: **relB = FILTER relA by \$0 < 20**

### **Documentation of commands executed**

#### **Excercise 1:-**

```
food_ratings = LOAD '/user/hadoop/foodratings161793.txt' USING PigStorage(',') AS
(name:chararray, f1:int, f2:int, f3:int, f4:int, placeid:int);
```

```
DESCRIBE food_ratings;
```

#### **Excercise 2:-**

```
food_ratings_subset = FOREACH food_ratings name, f4;
STORE food_ratings_subset INTO '/user/hadoop/fr_subset' USING PigStorage(',');
fr_output = LIMIT food_ratings_subset 6;
dump fr_output;
```

#### **Excercise 3:-**

```
fr_profile = GROUP food_ratings ALL;  
food_ratings_profile = FOREACH fr_profile GENERATE MIN(food_ratings.f2),  
MAX(food_ratings.f2), AVG(food_ratings.f2), MIN(food_ratings.f3), MAX(food_ratings.f3),  
AVG(food_ratings.f3);  
DUMP food_ratings_profile;
```

#### Excercise 4:-

```
food_ratings_filtered = FILTER food_ratings BY (f1<20) AND (f3>5);  
fr_filtered = LIMIT food_ratings_filtered 6;  
DUMP fr_filtered;
```

#### Excercise 5:-

```
food_ratings_2percent = SAMPLE food_ratings 0.02;  
filtered = LIMIT food_ratings_2percent 10;  
DUMP filtered;
```

#### Excercise 6:-

```
food_places = LOAD '/user/hadoop/foodplaces161793.txt' USING PigStorage(',') AS  
(placeid:int, placename:chararray);  
DESCRIBE food_places;  
food_ratings_w_place_names = JOIN food_places BY placeid, food_ratings BY placeid;  
fr_result = LIMIT food_ratings_w_place_names 6;  
DUMP fr_result;
```