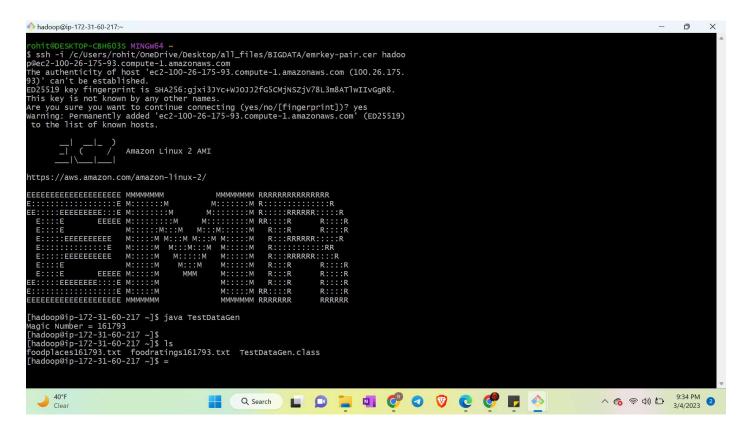
# **CSP-554** Big Data Technologies

## **Bairi Rohith Reddy - Assignment 5**

### Magic Number:



**Magic Number = 161793** 

foodplaces161793.txt

foodratings161793.txt

Exercise 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;

```
hadoop@ip-172-31-60-217:-/pigdemo

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>
```

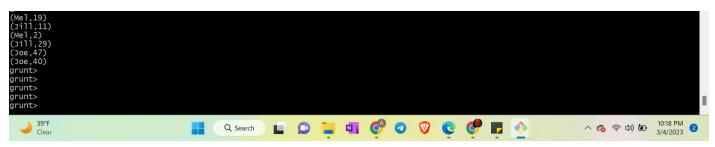
#### Exercise 2)

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;



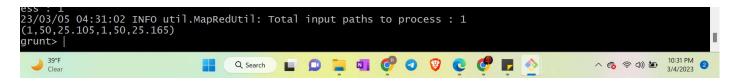
#### Exercise 3)

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;



#### Exercise 4)

food\_ratings\_filtered = FILTER food\_ratings BY (f1<20) AND (f3>5);

fr\_filtered = LIMIT food\_ratings\_filtered 6;

DUMP fr\_filtered;

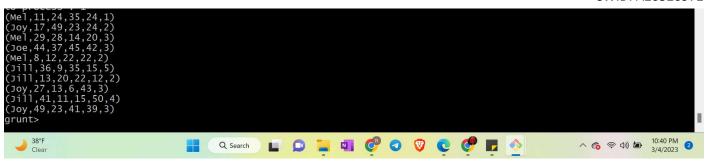


#### Exercise 5)

food\_ratings\_2percent = SAMPLE food\_ratings 0.02;

filtered = LIMIT food\_ratings\_2percent 10;

DUMP filtered;



#### Exercise 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\_ratings BY placeid, food\_places BY placeid;
fr\_result= LIMIT food\_ratings\_w\_place\_names 6;
DUMP fr\_result;



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.

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

Answer: LIMIT

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

Answer: **DISTINCT** 

- 3. 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)

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;