

SMIT R PATEL

19162121031

SEM = 5

PRACTICAL 9

BIG DATA AND ANALYTICS

AIM- To understand the working of Pig Functions.

Exercise: You need to perform analysis to reach to conclusions about some datasets your manager has provided you, but your company hosts its data over Cloudera, and you do not understand Java concepts clearly, hence need a way to work with these files in Hadoop.

Tasks:

Apache Pig provides various built-in functions namely **eval, load, store, math, string, bag and tuple** functions.

Eval Functions

Given below is the list of **eval** functions provided by Apache Pig.

S.N.	Function & Description
1	<u>AVG()</u> To compute the average of the numerical values within a bag.
2	<u>BagToString()</u> To concatenate the elements of a bag into a string. While concatenating, we can place a delimiter between these values (optional).
3	<u>CONCAT()</u> To concatenate two or more expressions of same type.
4	<u>COUNT()</u> To get the number of elements in a bag, while counting the number of tuples in a bag.

5	<u>COUNT_STAR()</u> It is similar to the COUNT() function. It is used to get the number of elements in a bag.
6	<u>DIFF()</u> To compare two bags (fields) in a tuple.
7	<u>IsEmpty()</u>

	To check if a bag or map is empty.
8	<u>MAX()</u> To calculate the highest value for a column (numeric values or chararrays) in a single-column bag.
9	<u>MIN()</u> To get the minimum (lowest) value (numeric or chararray) for a certain column in a single-column bag.
10	<u>PluckTuple()</u> Using the Pig Latin PluckTuple() function, we can define a string Prefix and filter the columns in a relation that begin with the given prefix.
11	<u>SIZE()</u> To compute the number of elements based on any Pig data type.
12	<u>SUBTRACT()</u> To subtract two bags. It takes two bags as inputs and returns a bag which contains the tuples of the first bag that are not in the second bag.
13	<u>SUM()</u> To get the total of the numeric values of a column in a single-column bag.

14

TOKENIZE()

To split a string (which contains a group of words) in a single tuple and return a bag which contains the output of the split operation.

AVG()

The Pig-Latin **AVG()** function is used to compute the average of the numerical values within a bag. While calculating the average value, the **AVG()** function ignores the NULL values.

Note –

- To get the global average value, we need to perform a **Group All** operation, and calculate the average value using the **AVG()** function.
- To get the average value of a group, we need to group it using the **Group By** operator and proceed with the average function.

Syntax

Given below is the syntax of the **AVG()** function.

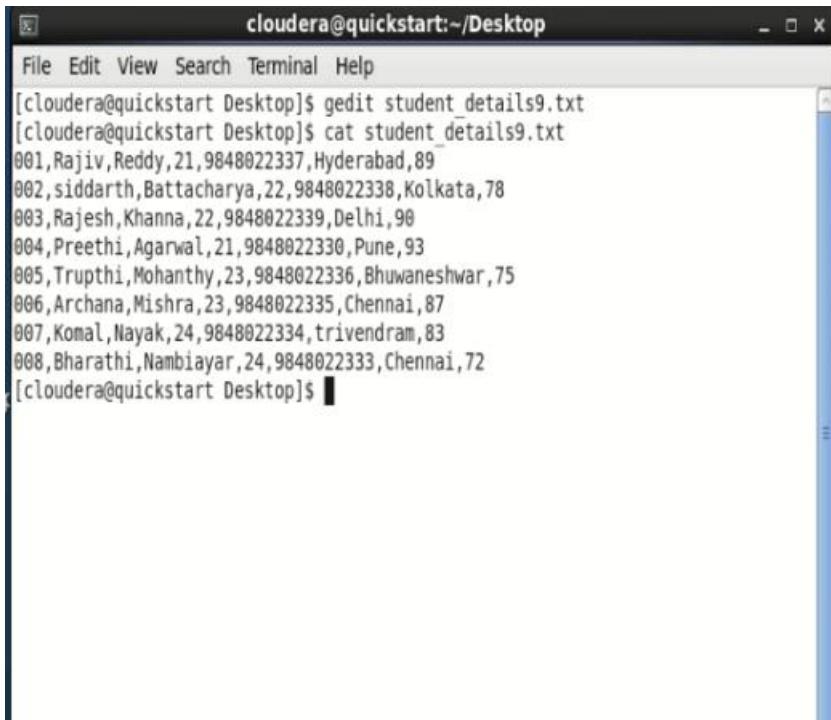
```
grunt> AVG(expression)
```

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below.

student_details.txt

```
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
```



The screenshot shows a terminal window titled "cloudera@quickstart:~/Desktop". The window contains a list of student records. The records are as follows:

```
File Edit View Search Terminal Help
[cloudera@quickstart Desktop]$ gedit student_details9.txt
[cloudera@quickstart Desktop]$ cat student_details9.txt
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
[cloudera@quickstart Desktop]$
```

```

cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
[cloudera@quickstart Desktop]$ gedit student_details9.txt
[cloudera@quickstart Desktop]$ cat student_details9.txt
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiyar,24,9848022333,Chennai,72
[cloudera@quickstart Desktop]$ hadoop fs -ls smitrpatel
Found 8 items
-rw-r--r-- 1 cloudera cloudera      6 2021-08-19 00:43 smitrpatel/ABC.txt
drwxr-xr-x - cloudera cloudera      0 2021-08-19 01:35 smitrpatel/ICT
-rw-r--r-- 1 cloudera cloudera      0 2021-08-17 02:21 smitrpatel/Just_Emp
ty_File.txt
-rw-r--r-- 1 cloudera cloudera      0 2021-08-19 03:13 smitrpatel/Practica
l3
-rw-r--r-- 1 cloudera cloudera    187 2021-08-26 02:05 smitrpatel/customer
s.txt
-rw-r--r-- 1 cloudera cloudera    135 2021-08-24 02:15 smitrpatel/employee
.txt
-rw-r--r-- 1 cloudera cloudera    124 2021-08-26 02:05 smitrpatel/orders.t
xt
-rw-r--r-- 1 cloudera cloudera   236 2021-08-24 00:46 smitrpatel/student_
data.txt
[cloudera@quickstart Desktop]$ hadoop fs -mkdir p9smit
hadoop[cloudera@quickstart Desktop]$ hadoop -put student_details9.txt p9smit
Error: No command named `put` was found. Perhaps you meant `hadoop put`
[cloudera@quickstart Desktop]$ hadoop fs -put student_details9.txt p9smit
[cloudera@quickstart Desktop]$ hadoop fs -ls p9smit
Found 1 items
-rw-r--r-- 1 cloudera cloudera    375 2021-08-28 01:40 p9smit/student_details9.txt
[cloudera@quickstart Desktop]$ 
```

And we have loaded this file into Pig with the relation name **student_details** as shown below.

```

grunt> student_details = LOAD
'hdfs://localhost:9000/pig_data/student_details.txt'
USING PigStorage(',')
  as (id:int, firstname:chararray, lastname:chararray, age:int,
  phone:chararray, city:chararray, gpa:int); 
```

```

grunt>student_details = LOAD 'p9smit/student_details9.txt' using
PigStorage(',') as (id:int, firstname:chararray, lastname:chararray,
age:int, phone:chararray, city:chararray, gpa:int) 
```

```

student@laptop:~/hadoop/cloudera/Desktop/pig_1626474695217.109
grunt> student_details = load 'p9smit/student_details9.txt'
>> using PigStorage(',') as
>> (id:int, firstname:chararray, lastname:chararray, age:int, phone:chararray, city:chararray, gpa:int);
grunt> 
```

Calculating the Average GPA

We can use the built-in function **AVG()** (case-sensitive) to calculate the average of a set of numerical values. Let's group the relation **student_details** using the Group All operator, and store the result in the relation named **student_group_all** as shown below.

```
grunt> student_group_all = Group student_details All;
```

```
grunt> student_details = load 'p9smit/student_details9.txt'
>> using PigStorage(',') as
>> (id:int, firstname:chararray, lastname:chararray, age:int, phone:chararray, city:chararray, gpa:int);
grunt> student_group_all = Group student_details ALL;
grunt> ■
```

This will produce a relation as shown below.

```
grunt> Dump student_group_all;
```

```
File Edit View Search Terminal Help
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.12.0 0.12.0-cdh5.12.0 cloudera 2021-08-28 01:58:18 2021-08-28 01:58:44 GROUP_BY
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime A
alias Feature Outputs
job_1630136589240_0001 1 1 3 3 3 2 2 2 2 student_details,student_group_all GROUP_BY h
dfs://quickstart.cloudera:8020/tmp/temp-750808021/tmp935923112,
Input(s):
Successfully read 8 records (767 bytes) from: "dfs://quickstart.cloudera:8020/user/cloudera/p9smit/student_details9.txt"
Output(s):
Successfully stored 1 records (414 bytes) in: "dfs://quickstart.cloudera:8020/tmp/temp-750808021/tmp935923112"
Counters:
Total records written : 1
Total bytes written : 414
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1630136589240_0001

2021-08-28 01:58:44,934 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-28 01:58:44,936 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-28 01:58:44,936 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.a
ddress
2021-08-28 01:58:44,937 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... Will not generate code.
2021-08-28 01:58:44,952 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-28 01:58:44,952 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(all,{(8,Bharathi,Nambiayar,24,9848022333,Chennai,72),(7,Komal,Nayak,24,9848022334,trivendram,83),(6,Archana,Mishra,23,9848022335,Chennai,87),(5,Trupthi,Moha
nthy,23,9848022336,Bhuwaneshwar,75),(4,Preethi,Agarwal,21,9848022330,Pune,93),(3,Rajesh,Khanna,22,9848022339,Delhi,90),(2,siddarth,Battacharya,22,9848022338,
Kolkata,78),(1,Rajiv,Reddy,21,9848022337,Hyderabad,89)})■
grunt> ■
```

```

File Edit View Search Terminal Help
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.12.0 0.12.0-cdh5.12.0 cloudera 2021-08-28 01:58:18 2021-08-28 01:58:44 GROUP_BY
Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime A
job_1630136589240_0001 1 1 3 3 3 2 2 2 2 student_details,student_group_all GROUP_BY h
dfs://quickstart.cloudera:8020/tmp/temp-750808021/tmp935923112,

Input(s):
Successfully read 8 records (767 bytes) from: "dfs://quickstart.cloudera:8020/user/cloudera/p9smit/student_details9.txt"

Output(s):
Successfully stored 1 records (414 bytes) in: "dfs://quickstart.cloudera:8020/tmp/temp-750808021/tmp935923112"

Counters:
Total records written : 1
Total bytes written : 414
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630136589240_0001

2021-08-28 01:58:44,934 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-28 01:58:44,936 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-28 01:58:44,936 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-28 01:58:44,937 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-28 01:58:44,952 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-28 01:58:44,952 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(all,{(8,Bharathi,Nambiayar,24,9848022333,Chennai,72),(7,Komal,Nayak,24,9848022334,trivendram,83),(6,Archana,Mishra,23,9848022335,Chennai,87),(5,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75),(4,Preethi,Agarwal,21,9848022330,Pune,93),(3,Rajesh,Khanna,22,9848022339,Delhi,90),(2,siddarth,Battacharya,22,9848022338,Kolkata,78),(1,Rajiv,Reddy,21,9848022337,Hyderabad,89)})
grunt> ■

```

```

2021-08-28 01:58:44,934 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-28 01:58:44,936 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-28 01:58:44,936 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-28 01:58:44,937 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-28 01:58:44,952 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-28 01:58:44,952 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(all,{(8,Bharathi,Nambiayar,24,9848022333,Chennai,72),(7,Komal,Nayak,24,9848022334,trivendram,83),(6,Archana,Mishra,23,9848022335,Chennai,87),(5,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75),(4,Preethi,Agarwal,21,9848022330,Pune,93),(3,Rajesh,Khanna,22,9848022339,Delhi,90),(2,siddarth,Battacharya,22,9848022338,Kolkata,78),(1,Rajiv,Reddy,21,9848022337,Hyderabad,89)})
grunt>

```

Let us now calculate the global average GPA of all the students using the **AVG()** function as shown below.

```

grunt> student_gpa_avg = foreach student_group_all
Generate (student_details.firstname,
student_details.gpa), AVG(student_details.gpa);

```

```
File Edit View Search Terminal Help
grunt> student_gpa_avg = foreach student_group_all Generate (student_details.firstname, student_details.gpa), AVG(student_details.gpa);
2021-08-18 23:37:04,820 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:37:04,820 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
grunt> █
```

Command : student_gpa_avg = foreach student_group_all Generate
(student_details.firstname,student_details_gpa), AVG(student_details.gpa);

Verification

Verify the relation **student_gpa_avg** using the **DUMP** operator as shown below.

```
grunt> Dump student_gpa_avg;
```

```
Output(s):
Successfully stored 1 records (127 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-166987612/tmp1306276875"

Counters:
Total records written : 1
Total bytes written : 127
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1629165673222_0020

2021-08-18 23:37:57,411 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-18 23:37:57,412 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:37:57,412 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-18 23:37:57,413 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-18 23:37:57,420 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-18 23:37:57,420 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
({{(Bharathi),(Komal),(Archana),(Trupthi),(Preethi),(Rajesh),(Siddarth),(Rajiv)}},{{(72),(83),(87),(75),(93),(80),(78),(89)}},83.375)
grunt> █
```

Output

It will display the contents of the relation **student_gpa_avg** as follows.

```
({{(Bharathi),(Komal),(Archana),(Trupthi),(Preethi),(Rajesh),(Siddarth),(Rajiv)}},
{{(72),(83),(87),(75),(93),(80),(78),(89)}},83.375)
```

BagToString()

The Pig Latin **BagToString()** function is used to concatenate the elements of a bag into a string. While concatenating, we can place a delimiter between these values (optional).

Generally bags are disordered and can be arranged by using **ORDER BY**

operator. **Syntax**

Given below is the syntax of the **BagToString()** function.

```
grunt> BagToString(vals:bag [,
```

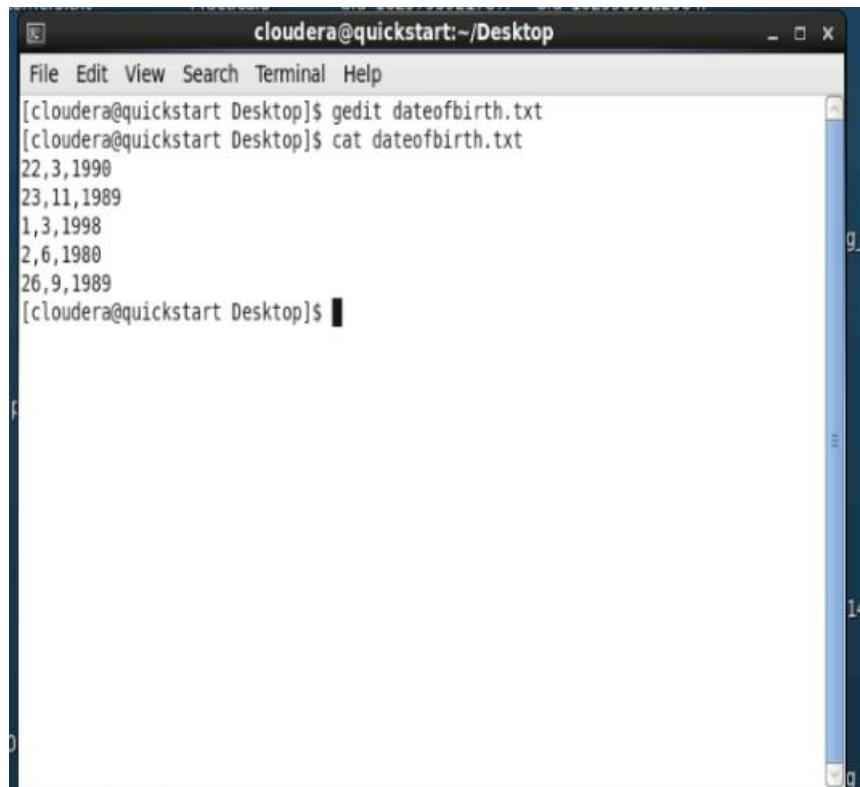
```
delimiter:chararray])
```

Example

Assume that we have a file named **dateofbirth.txt** in the HDFS directory **/pig_data/** as shown below. This file contains the date-of-births.

dateofbirth.txt

```
22,3,1990
23,11,1989
1,3,1998
2,6,1980
26,9,1989
```



The screenshot shows a terminal window titled "cloudera@quickstart:~/Desktop". The window contains the following text:

```
File Edit View Search Terminal Help
[cloudera@quickstart Desktop]$ gedit dateofbirth.txt
[cloudera@quickstart Desktop]$ cat dateofbirth.txt
22,3,1990
23,11,1989
1,3,1998
2,6,1980
26,9,1989
[cloudera@quickstart Desktop]$
```

The terminal window has a dark background with light-colored text. The scroll bar on the right side shows the number "14" at the bottom.

And we have loaded this file into Pig with the relation name **dob** as shown below.

```
grunt> dob = LOAD  
'hdfs://localhost:9000/pig_data/dateofbirth.txt'  
USING PigStorage(',')  
as (day:int, month:int, year:int);
```

```
>> (day:int month:int, year:int);  
2021-08-28 03:06:53,227 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 12  
00: <line 3, column 9> mismatched input 'month' expecting RIGHT_PAREN  
Details at logfile: /home/cloudera/Desktop/pig_1630144798039.log  
nt:int, year:int);  
grunt> dob = load 'p9smi/dateofbirth.txt' using PigStorage(',') as (day:int, mont:int, year:int);  
grunt> ■
```

Converting Bag to String

Using the **bagtostring()** function, we can convert the data in the bag to string. Let us group the **dob** relation. The group operation will produce a bag containing all the tuples of the relation.

Group the relation **dob** using the **Group All** operator, and store the result in the relation named **group_dob** as shown below.

```
grunt> group_dob = Group dob All;
```

```
grunt> dob = load 'p9smi/dateofbirth.txt' using PigStorage(',') as (day:int, mont:int, year:int);  
grunt> group_dob = Group dob ALL;  
grunt> ■
```

It will produce a relation as shown below.

```
grunt> Dump group_dob;
```

```
Output(s):
Successfully stored 1 records (51 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-166
6987612/tmp283669265"

Counters:
Total records written : 1
Total bytes written : 51
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1629165673222_0021
```

```
2021-08-18 23:43:09,294 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapRedu
ceLayer.MapReduceLauncher - Success!
2021-08-18 23:43:09,294 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs
.default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:43:09,295 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - ma
pred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-18 23:43:09,295 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.sch
ematuple] was not set... will not generate code.
2021-08-18 23:43:09,306 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat
- Total input paths to process : 1
2021-08-18 23:43:09,306 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.Ma
pRedUtil - Total input paths to process : 1
(all,{{(26,9,1989),(2,6,1980),(1,3,1998),(23,11,1989),(22,3,1990)}})
grunt> █
```

```
(all,{{(26,9,1989),(2,6,1980),(1,3,1998),(23,11,1989),(22,3,1990)}}
} )
```

Here, we can observe a bag having all the date-of-births as tuples of it. Now, let's convert the bag to string using the function **BagToString()**.

```
grunt> dob_string = foreach group_dob Generate BagToString(dob);
```

```
grunt>
grunt>
grunt> dob_string = foreach group_dob Generate BagToString(dob);
grunt> █
```

Verification

Verify the relation **dob_string** using the **DUMP** operator as shown below.

```
grunt> Dump dob_string;
```

```
2021-08-28 03:11:40,573 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input
paths to process : 1
2021-08-28 03:11:40,573 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total
input paths to process : 1
(26_9_1989_2_6_1980_1_3_1998_23_11_1989_22_3_1990)
grunt>
```

```

cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
2021-08-28 03:11:40,459 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:
adoopVersion PigVersion UserId StartedAt FinishedAt Features
.6.0-cdh5.12.0 0.12.0-cdh5.12.0 cloudera 2021-08-28 03:11:15 2021-08-28 03:11:40 GROUP_BY

uccess!

ob Stats (time in seconds):
obId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime
ias Feature Outputs
ob 1630136589240 0003 1 3 3 3 2 2 2 2 dob,dob_string,group_dob GROUP_BY hc
/quickstart.cloudera:8020/tmp/temp60001356/tmp93024786,

nput(s):
uccessfully read 5 records (441 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/p9smit/dateofbirth.txt"

utput(s):
uccessfully stored 1 records (55 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp60001356/tmp93024786"

ounters:
otal records written : 1
otal bytes written : 55
pillable Memory Manager spill count : 0
otal bags proactively spilled: 0
otal records proactively spilled: 0

ob DAG:
ob_1630136589240_0003

2021-08-28 03:11:40,555 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-28 03:11:40,557 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-28 03:11:40,557 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracke
dress
2021-08-28 03:11:40,558 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-28 03:11:40,573 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-28 03:11:40,573 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
26_9_1989_2_6_1980_1_3_1998_23_11_1989_22_3_1990)
rants

```

Output

It will produce the following output, displaying the contents of the relation

dob_string. (26_9_1989_2_6_1980_1_3_1998_23_11_1989_22_3_1990)

CONCAT()

The **CONCAT()** function of Pig Latin is used to concatenate two or more expressions of the same type.

Syntax

```
grunt> CONCAT (expression, expression,  
[...expression])
```

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below.

student_details.txt

```
001,Rajiv,Reddy,21,9848022337,Hyderabad,89  
002,siddarth,Battacharya,22,9848022338,Kolkata,78  
003,Rajesh,Khanna,22,9848022339,Delhi,90  
004,Preethi,Agarwal,21,9848022330,Pune,93  
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75  
006,Archana,Mishra,23,9848022335,Chennai,87  
007,Komal,Nayak,24,9848022334,trivendram,83  
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
```

And we have loaded this file into Pig with the relation name **student_details** as shown below.

```
grunt> student_details = LOAD  
'hdfs://localhost:9000/pig_data/student_details.txt'  
USING PigStorage(',')  
as (id:int, firstname:chararray, lastname:chararray, age:int,  
phone:chararray, city:chararray, gpa:int);
```

Concatenating Two Strings

We can use the **CONCAT()** function to concatenate two or more expressions. First of all, verify the contents of the **student_details** relation using the **Dump** operator as shown below.

```
grunt> Dump student_details;  
( 1,Rajiv,Reddy,21,9848022337,Hyderabad,89 )  
( 2,siddarth,Battacharya,22,9848022338,Kolkata,78 )  
( 3,Rajesh,Khanna,22,9848022339,Delhi,90 )  
( 4,Preethi,Agarwal,21,9848022330,Pune,93 )  
( 5,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75 )  
( 6,Archana,Mishra,23,9848022335,Chennai,87 )  
( 7,Komal,Nayak,24,9848022334,trivendram,83 )  
( 8,Bharathi,Nambiayar,24,9848022333,Chennai,72 )
```

And, verify the schema using **describe** operator as shown below.

```
grunt> Describe student_details;
```

```
student_details: {id: int, firstname: chararray, lastname:  
chararray, age: int,
```

```
phone: chararray, city: chararray, gpa: int}
```

In the above schema, you can observe that the name of the student is represented using two chararray values namely **firstname** and **lastname**. Let us concatenate these two values using the **CONCAT()** function.

```
grunt> student_name_concat = foreach student_details
Generate CONCAT (firstname, lastname);
```

```
File Edit View Search Terminal Help
grunt> student_name_concat = foreach student_details Generate CONCAT(firstname, lastname); █
grunt> █
```

Verification

Verify the relation **student_name_concat** using the **DUMP** operator as shown below.

```
grunt> Dump student_name_concat;
```

```
File Edit View Search Terminal Help
Counters:
Total records written : 8
Total bytes written : 166
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1629165673222_0022

2021-08-18 23:45:58,663 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-18 23:45:58,663 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:45:58,664 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-18 23:45:58,664 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-18 23:45:58,671 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-18 23:45:58,671 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(RajivReddy)
(SiddarthBattacharya)
(RajeshKhanna)
(PreethiAgarwal)
(TrupthiMohanthy)
(ArchanaMishra)
(KomalNayak)
(BharathiNambiayar)
grunt> █
```

Output

It will produce the following output, displaying the contents of the relation **student_name_concat**.

```
(RajivReddy)
(siddarthBattacharya)
```

```
(RajeshKhanna)
(PreethiAgarwal)
(TrupthiMohanty)
(ArchanaMishra)
(KomalNayak)
(BharathiNambiar)
```

We can also use an optional delimiter between the two expressions as shown below.

```
grunt> CONCAT(firstname, '_', lastname);
```

Now, let us concatenate the first name and last name of the student records in the **student_details** relation by placing ‘_’ between them as shown below.

```
grunt> student_name_concat = foreach student_details
  GENERATE CONCAT(firstname, '_', lastname);
```

```
File Edit View Search Terminal Help
grunt> student_name_concat = foreach student_details Generate CONCAT(firstname, '_', lastname);
grunt> dump student_name_concat;
```

Verification

Verify the relation **student_name_concat** using the **DUMP** operator as shown below.

```
grunt> Dump student_name_concat;
```

Counters:

```
Total records written : 8
Total bytes written : 174
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
```

Job DAG:

```
job_1629165673222_0023
```

```
2021-08-18 23:47:25,395 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-18 23:47:25,395 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:47:25,395 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-18 23:47:25,396 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-18 23:47:25,401 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-18 23:47:25,401 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Rajiv_Reddy)
(Siddarth_Battacharya)
(Rajesh_Khanna)
(Preethi_Agarwal)
(Trupthi_Mohanthy)
(Archana_Mishra)
(Komal_Nayak)
(Bharathi_Nambiayar)
grunt> ■
```

Output

It will produce the following output, displaying the contents of the relation **student_name_concat** as follows.

```
(Rajiv_Reddy)
(siddarth_Battacharya)
(Rajesh_Khanna)
(Preethi_Agarwal)
(Trupthi_Mohanthy)
(Archana_Mishra)
(Komal_Nayak)
(Bharathi_Nambiayar)
```

COUNT()

The **COUNT()** function of Pig Latin is used to get the number of elements in a bag. While counting the number of tuples in a bag, the **COUNT()** function ignores (will not count) the tuples having a NULL value in the FIRST FIELD.

Note –

- To get the global count value (total number of tuples in a bag), we need to perform a **Group All** operation, and calculate the count value using the **COUNT()** function.
- To get the count value of a group (Number of tuples in a group), we need to group it using the **Group By** operator and proceed with the count function.

Syntax

Given below is the syntax of the **COUNT()** function.

```
grunt> COUNT(expression)
```

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below.

student_details.txt

```
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
```

And we have loaded this file into Pig with the relation named **student_details** as shown below.

```
grunt> student_details = LOAD
  'hdfs://localhost:9000/pig_data/student_details.txt'
  USING PigStorage(',')
  as (id:int, firstname:chararray, lastname:chararray, age:int,
  phone:chararray, city:chararray, gpa:int);
```

Calculating the Number of Tuples

We can use the built-in function **COUNT()** (case sensitive) to calculate the number of tuples in a relation. Let us group the relation **student_details** using the **Group All** operator, and store the result in the relation named **student_group_all** as shown below.

```
grunt> student_group_all = Group student_details All;
```

It will produce a relation as shown below.

```
grunt> Dump student_group_all;

(all, {(8,Bharathi,Nambiayar,24,9848022333,Chennai,72),
, (7,Komal,Nayak,24,9848022 334,trivendram,83),
(6,Archana,Mishra,23,9848022335,Chennai,87),
(5,Trupthi,Mohan thy,23,9848022336,Bhuwaneshwar,75),
(4,Preethi,Agarwal,21,9848022330,Pune,93),
(3 ,Rajesh,Khanna,22,9848022339,Delhi,90),
(2,siddarth,Battacharya,22,9848022338,Ko lkata,78),
(1,Rajiv,Reddy,21,9848022337,Hyderabad,89) })
```

Let us now calculate number of tuples/records in the relation.

```
grunt> student_count = foreach student_group_all
Generate COUNT(student_details.gpa);
```

```
File Edit View Search Terminal Help
grunt> student_count = foreach student_group_all Generate COUNT(student_details.gpa);
grunt> █
```

Verification

Verify the relation **student_count** using the **DUMP** operator as shown below.

```
grunt> Dump student_count;
```

```
Output(s):
Successfully stored 1 records (6 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-1666
987612/tmp-493324534"
```

```
Counters:
Total records written : 1
Total bytes written : 6
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
```

```
Job DAG:
job_1629165673222_0024
```

```
2021-08-18 23:48:56,748 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapRedu
ceLayer.MapReduceLauncher - Success!
2021-08-18 23:48:56,749 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs
.default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:48:56,749 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - ma
pred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-18 23:48:56,749 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.sch
ematuple] was not set... will not generate code.
2021-08-18 23:48:56,757 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat
- Total input paths to process : 1
2021-08-18 23:48:56,757 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.Ma
pRedUtil - Total input paths to process : 1
(8)
grunt> █
```

Output

It will produce the following output, displaying the contents of the relation **student_count**.

COUNT_STAR()

The **COUNT_STAR()** function of Pig Latin is similar to the **COUNT()** function. It is used to get the number of elements in a bag. While counting the elements, the **COUNT_STAR()** function includes the NULL values.

Note –

- To get the global count value (total number of tuples in a bag), we need to perform a **Group All** operation, and calculate the count_star value using the **COUNT_STAR()** function.
- To get the count value of a group (Number of tuples in a group), we need to group it using the **Group By** operator and proceed with the count_star function.

Syntax

Given below is the syntax of the **COUNT_STAR()** function.

```
grunt> COUNT_STAR(expression)
```

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below. This file contains an empty record.

student_details.txt

```
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
```

And we have loaded this file into Pig with the relation name **student_details** as shown below.

```
grunt> student_details = LOAD
  'hdfs://localhost:9000/pig_data/student_details.txt'
  USING PigStorage(',')
  as (id:int, firstname:chararray, lastname:chararray, age:int,
  phone:chararray, city:chararray, gpa:int);
```

Calculating the Number of Tuples

We can use the built-in function **COUNT_STAR()** to calculate the number of tuples in a relation. Let us group the relation **student_details** using the **Group All** operator, and store the result in the relation named **student_group_all** as shown below.

```
grunt> student_group_all = Group student_details All;
```

It will produce a relation as shown below.

```
grunt> Dump student_group_all;

(all,{{8,Bharathi,Nambiayar,24,9848022333,Chennai,72),
 , (7,Komal,Nayak,24,9848022 334,trivendram,83),
 (6,Archana,Mishra,23,9848022335,Chennai,87),
 (5,Trupthi,Mohan thy,23,9848022336,Bhuwaneshwar,75),
 (4,Preethi,Agarwal,21,9848022330,Pune,93),
 (3 ,Rajesh,Khanna,22,9848022339,Delhi,90),
 (2,siddarth,Battacharya,22,9848022338,Ko lkata,78),
 (1,Rajiv,Reddy,21,9848022337,Hyderabad,89),
 ( , , , , , ))}
```

Let us now calculate the number of tuples/records in the relation.

```
grunt> student_count = foreach student_group_all
Generate COUNT_STAR(student_details.gpa);
```

```
File Edit View Search Terminal Help
grunt> student_count = foreach student_group_all Generate COUNT_STAR(student_details.gpa);
grunt> ■
```

Verification

Verify the relation **student_count** using the **DUMP** operator as shown below.

```
grunt> Dump student_count;
```

```
Output(s):
Successfully stored 1 records (6 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-166987612/tmp152862587"
```

```
Counters:
Total records written : 1
Total bytes written : 6
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
```

```
Job DAG:
job_1629165673222_0025
```

```
2021-08-18 23:50:47,381 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-18 23:50:47,382 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - f .default.name is deprecated. Instead, use fs.defaultFS
2021-08-18 23:50:47,382 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - m pred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-18 23:50:47,383 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schema] was not set... will not generate code.
2021-08-18 23:50:47,389 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-18 23:50:47,389 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.M pReduceUtil - Total input paths to process : 1
(8)
grunt> ■
```

Output

It will produce the following output, displaying the contents of the relation **student_count**.

9

Since we have used the function **COUNT_STAR()**, it included the null tuple and returned 9.

DIFF()

The **DIFF()** function of Pig Latin is used to compare two bags (fields) in a tuple. It takes two fields of a tuple as input and matches them. If they match, it returns an empty bag. If they do not match, it finds the elements that exist in one field (bag) and not found in the other, and returns these elements by wrapping them within a bag.

Syntax

Given below is the syntax of the **DIFF()** function.

```
grunt> DIFF (expression, expression)
```

Example

Generally the **DIFF()** function compares two bags in a tuple. Given below is its example, here we create two relations, cogroup them, and calculate the difference between them.

Assume that we have two files namely **emp_sales.txt** and **emp_bonus.txt** in the HDFS directory **/pig_data/** as shown below. The **emp_sales.txt** contains the details of the employees of the sales department and the **emp_bonus.txt** contains the employee details who got bonus.

emp_sales.txt

```
1,Robin,22,25000,sales
2,BOB,23,30000,sales
3,Maya,23,25000,sales
4,Sara,25,40000,sales
5,David,23,45000,sales
6,Maggy,22,35000,sales
```

emp_bonus.txt

```
1,Robin,22,25000,sales
2,Jaya,23,20000,admin
3,Maya,23,25000,sales
4,Alia,25,50000,admin
5,David,23,45000,sales
6,Omar,30,30000,admin
```

```
[cloudera@quickstart Desktop]$ gedit emp_sales.txt
[cloudera@quickstart Desktop]$ cat emp_sales.txt
1,Robin,22,25000,sales
2,BOB,23,30000,sales
3,Maya,23,25000,sales
4,Sara,25,40000,sales
5,David,23,45000,sales
6,Maggy,22,35000,sales
[cloudera@quickstart Desktop]$ cat emp_bonus.txt
1,Robin,22,25000,sales
2,Jaya,23,20000,admin
3,Maya,23,25000,sales
4,Alia,25,50000,admin
5,David,23,45000,sales
6,Omar,30,30000,admin
[cloudera@quickstart Desktop]$ hadoop fs -put emp_bonus.txt smitrapatel
[cloudera@quickstart Desktop]$ hadoop fs -put emp_sales.txt smitrapatel
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrapatel/emp_bonus.txt
```

And we have loaded these files into Pig, with the relation names **emp_sales** and **emp_bonus** respectively.

```
grunt> emp_sales = LOAD  
'hdfs://localhost:9000/pig_data/emp_sales.txt' USING  
PigStorage(',')  
as (sno:int, name:chararray, age:int, salary:int,  
dept:chararray);  
  
grunt> emp_bonus = LOAD  
'hdfs://localhost:9000/pig_data/emp_bonus.txt' USING  
PigStorage(',')  
as (sno:int, name:chararray, age:int, salary:int,  
dept:chararray);
```

Group the records/tuples of the relations **emp_sales** and **emp_bonus** with the key **sno**, using the COGROUP operator as shown below.

```
grunt> cogroup_data = COGROUP emp_sales by sno, emp_bonus by sno;
```

```
grunt> emp_sales = load 'smitrpatel/emp_sales.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);  
grunt> emp_bonus = load 'smitrpatel/emp_bonus.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);  
grunt> cogroup_data = COGROUP emp_sales BY sno, emp_bonus BY sno;  
grunt> dump cogroup_data;S
```

Verify the relation **cogroup_data** using the **DUMP** operator as shown below.

```
grunt> Dump cogroup_data;
```

```

Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReducetime  A
alias  Feature Outputs
job_1630371021450_0002  2      1      25      25      25      25      10      10      10      10      cogroup_data,emp_bonus,emp_sales      COGROUP hdfs:
//quickstart.cloudera:8020/tmp/temp-864208229/tmp1446958649,                                    

[Input(s):
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"

[Output(s):
Successfully stored 6 records (367 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp1446958649"

Counters:
Total records written : 6
Total bytes written : 367
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0002

2021-08-30 18:31:20,619 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 18:31:20,620 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 18:31:20,620 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 18:31:20,621 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 18:31:20,629 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 18:31:20,629 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
[1,{(1,Robin,22,25000,sales )}, {(1,Robin,22,25000,sales )}]
[2,{(2,BOB,23,30000,sales )}, {(2,Jaya,23,20000,admin )}]
[3,{(3,Maya,23,25000,sales )}, {(3,Maya,23,25000,sales )}]
[4,{(4,Sara,25,40000,sales )}, {(4,Alia,25,50000,admin )}]
[5,{(5,David,23,45000,sales )}, {(5,David,23,45000,sales )}]
[6,{(6,Maggy,22,35000,sales )}, {(6,Omar,30,30000,admin )}]
grunt> █ cloudera@quickstart:~/Desktop

```

```

(1, { (1,Robin,22,25000,sales )}, { (1,Robin,22,25000,sales )})
(2, { (2,BOB,23,30000,sales )}, { (2,Jaya,23,20000,admin )})
(3, { (3,Maya,23,25000,sales )}, { (3,Maya,23,25000,sales )})
(4, { (4,Sara,25,40000,sales )}, { (4,Alia,25,50000,admin )})
(5, { (5,David,23,45000,sales )}, { (5,David,23,45000,sales )})
(6, { (6,Maggy,22,35000,sales )}, { (6,Omar,30,30000,admin )})

```

Calculating the Difference between Two Relations

Let us now calculate the difference between the two relations using **DIFF()** function and store it in the relation **diff_data** as shown below.

```

grunt> diff_data = FOREACH cogroup_data GENERATE
  DIFF(emp_sales,emp_bonus);

```

```

2021-08-30 18:31:20,621 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key
2021-08-30 18:31:20,629 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileIn
2021-08-30 18:31:20,629 [main] INFO org.apache.pig.backend.hadoop.executionengine
[1,{(1,Robin,22,25000,sales )}, {(1,Robin,22,25000,sales )}]
[2,{(2,BOB,23,30000,sales )}, {(2,Jaya,23,20000,admin )}]
[3,{(3,Maya,23,25000,sales )}, {(3,Maya,23,25000,sales )}]
[4,{(4,Sara,25,40000,sales )}, {(4,Alia,25,50000,admin )}]
[5,{(5,David,23,45000,sales )}, {(5,David,23,45000,sales )}]
[6,{(6,Maggy,22,35000,sales )}, {(6,Omar,30,30000,admin )}]
grunt> diff_data = FOREACH cogroup_data Generate DIFF(emp_sales, emp_bonus);
grunt> █

```

Verification

Verify the relation **diff_data** using the DUMP operator as shown below.

```

grunt> Dump diff_data;

```

```

File Edit View Search Terminal Help

Job Stats (time in seconds):
JobId  Maps Reduces MaxMapTime      MinMapTime      AvgMapTime      MedianMapTime      MaxReduceTime      MinReduceTime      AvgReduceTime      MedianReducetime      A
alias  Feature Outputs
job_1630371021450 0003 2      1      26      25      25      25      12      12      12      cogroup_data,diff_data,emp_bonus,emp_sales      COGRO
UP      hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp423696791,      COGRO

Input(s):
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"

Output(s):
Successfully stored 6 records (188 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp423696791"

Counters:
Total records written : 6
Total bytes written : 188
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0003

2021-08-30 18:36:32,819 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 18:36:32,820 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 18:36:32,820 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.a
ddress
2021-08-30 18:36:32,821 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 18:36:32,830 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 18:36:32,830 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
({})
({{(2,BOB,23,30000,sales),(2,Jaya,23,20000,admin)})}
({})
({{(4,Sara,25,40000,sales),(4,Alia,25,50000,admin)})}
({})
({{(6,Maggy,22,35000,sales),(6,Omar,30,30000,admin)}})
grunt> 

```

```

({ })
({{2,BOB,23,30000,sales),(2,Jaya,23,20000,admin)})}
({})
({{(4,Sara,25,40000,sales),(4,Alia,25,50000,admin)})}
({})
({{(6,Maggy,22,35000,sales),(6,Omar,30,30000,admin)}})

```

The **diff_data** relation will have an empty tuple if the records in **emp_bonus** and **emp_sales** match. In other cases, it will hold tuples from both the relations (tuples that differ).

For example, if you consider the records having **sno** as **1**, then you will find them same in both the relations **((1,Robin,22,25000,sales), (1,Robin,22,25000,sales))**. Therefore, in the **diff_data** relation, which is the result of **DIFF()** function, you will get an empty tuple for **sno 1**.

IsEmpty()

The **IsEmpty()** function of Pig Latin is used to check if a bag or map is empty. **Syntax**

Given below is the syntax of the **IsEmpty()** function.

```
grunt> IsEmpty(expression)
```

Example

Assume that we have two files namely **emp_sales.txt** and **emp_bonus.txt** in the HDFS directory **/pig_data/** as shown below. The **emp_sales.txt** contains the details of the employees of the sales department and the **emp_bonus.txt** contains the employee details who got bonus.

emp_sales.txt

```
1,Robin,22,25000,sales
2,BOB,23,30000,sales
3,Maya,23,25000,sales
4,Sara,25,40000,sales
5,David,23,45000,sales
6,Maggy,22,35000,sales
```

emp_bonus.txt

```
1,Robin,22,25000,sales
2,Jaya,23,20000,admin
3,Maya,23,25000,sales
4,Alia,25,50000,admin
5,David,23,45000,sales
6,Omar,30,30000,admin
```

```
[cloudera@quickstart Desktop]$ gedit emp_sales.txt
[cloudera@quickstart Desktop]$ cat emp_sales.txt
1,Robin,22,25000,sales
2,BOB,23,30000,sales
3,Maya,23,25000,sales
4,Sara,25,40000,sales
5,David,23,45000,sales
6,Maggy,22,35000,sales
[cloudera@quickstart Desktop]$ cat emp_bonus.txt
1,Robin,22,25000,sales
2,Jaya,23,20000,admin
3,Maya,23,25000,sales
4,Alia,25,50000,admin
5,David,23,45000,sales
6,Omar,30,30000,admin
[cloudera@quickstart Desktop]$ hadoop fs -put emp_bonus.txt smitrapatel
[cloudera@quickstart Desktop]$ hadoop fs -put emp_sales.txt smitrapatel
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrapatel/emp_bonus.txt
```

And we have loaded these files into Pig, with the relation names **emp_sales** and **emp_bonus** respectively, as shown below.

```

grunt> emp_sales = LOAD
  'hdfs://localhost:9000/pig_data/emp_sales.txt'
  USING PigStorage(',')
  as (sno:int, name:chararray, age:int, salary:int,
  dept:chararray);

grunt> emp_bonus = LOAD
  'hdfs://localhost:9000/pig_data/emp_bonus.txt'
  USING PigStorage(',')
  as (sno:int, name:chararray, age:int, salary:int,
  dept:chararray);

```

Let us now group the records/tuples of the relations **emp_sales** and **emp_bonus** with the key **age**, using the **cogroup** operator as shown below.

```
grunt> cogroup_data = COGROUP emp_sales by age, emp_bonus by age;
```

```

File Edit View Search Terminal Help
grunt> cogroup_data = COGROUP emp_sales by age, emp_bonus by age;
grunt> █

```

Verify the relation **cogroup_data** using the **DUMP** operator as shown below.

```
grunt> Dump cogroup_data;
```

```

File Edit View Search Terminal Help
Job Stats (time in seconds):
JobId  Maps  Reduces MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReducetime  A
alias  Feature Outputs
job_1630371021450_0002  2      1      25      25      25      10      10      10      10      cogroup_data,emp_bonus,emp_sales      COGROUP hdfs:
//quickstart.cloudera:8020/tmp/temp-864208229/tmp1446958649,                                    

[Input(s):
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"

[Output(s):
Successfully stored 6 records (367 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp1446958649"

Counters:
Total records written : 6
Total bytes written : 367
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0002

2021-08-30 18:31:20,619 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 18:31:20,620 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 18:31:20,620 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 18:31:20,621 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 18:31:20,629 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 18:31:20,629 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
[1,{(1,Robin,22,25000,sales )}, {(1,Robin,22,25000,sales )}]
[2,{(2,BOB,23,30000,sales )}, {(2,Jaya,23,20000,admin )}]
[3,{(3,Maya,23,25000,sales )}, {(3,Maya,23,25000,sales )}]
[4,{(4,Sara,25,40000,sales )}, {(4,Alia,25,50000,admin )}]
[5,{(5,David,23,45000,sales )}, {(5,David,23,45000,sales )}]
[6,{(6,Maggy,22,35000,sales )}, {(6,Omar,30,30000,admin )}]
grunt> █

```

cloudera@quickstart:~/Desktop

[Cloudera Live : Welco...] [cloudera@quickstart:...] [cloudera@quickstart:~...]

The COGROUP operator groups the tuples from each relation according to age. Each group depicts a particular age value.

For example, if we consider the 1st tuple of the result, it is grouped by age 22. And it contains two bags, the first bag holds all the tuples from the first relation (student_details in this case) having age 22, and the second bag contains all the tuples from the second relation (employee_details in this case) having age 22. In case a relation doesn't have tuples having the age value 22, it returns an empty bag.

Getting the Groups having Empty Bags

Let's list such empty bags from the **emp_sales** relation in the group using the **IsEmpty()** function.

```
grunt> isempty_data = filter cogroup_data by IsEmpty(emp_sales);
```

```
2021-08-30 18:36:32,819 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduc
2021-08-30 18:36:32,820 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.
2021-08-30 18:36:32,820 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - map
ddress
2021-08-30 18:36:32,821 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.sche
2021-08-30 18:36:32,830 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat
2021-08-30 18:36:32,830 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.Map
({})
({{(2,BOB,23,30000,sales),(2,Jaya,23,20000,admin )}})
({})
({{(4,Sara,25,40000,sales),(4,Alia,25,50000,admin )}})
({})
({{(6,Maggy,22,35000,sales),(6,Omar,30,30000,admin )}})
grunt> isempty_data = filter cogroup_data by IsEmpty(emp_sales);
grunt> dump isempty_data;
```

Verification

Verify the relation **isempty_data** using the DUMP operator as shown below. The **emp_sales** relation holds the tuples that are not there in the relation **emp_bonus**.

```
grunt> Dump isempty_data;
```

File Edit View Search Terminal Help

2021-08-30 18:49:51,032 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:

HadoopVersion	PigVersion	UserId	StartedAt	FinishedAt	Features
2.6.0-cdh5.12.0	0.12.0-cdh5.12.0	cloudera	2021-08-30 18:48:43	2021-08-30 18:49:51	COGROUP,FILTER

Success!

Job Stats (time in seconds):

JobId	Maps	Reduces	MaxMapTime	MinMapTime	AvgMapTime	MedianMapTime	MaxReduceTime	MinReduceTime	AvgReduceTime	MedianReducetime	A
job_1630371021450_0005	2	1	25	25	25	11	11	11	11	cogroup_data,emp_bonus,emp_sales,isempty_data	COGROUP
UP											

Input(s):

Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"

Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"

Output(s):

Successfully stored 0 records in: "hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp-1229756602"

Counters:

Total records written : 0

Total bytes written : 0

Spillable Memory Manager spill count : 0

Total bags proactively spilled: 0

Total records proactively spilled: 0

Job DAG:

job_1630371021450_0005

2021-08-30 18:49:51,203 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!

2021-08-30 18:49:51,203 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS

2021-08-30 18:49:51,203 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address

2021-08-30 18:49:51,204 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.

2021-08-30 18:49:51,212 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1

2021-08-30 18:49:51,212 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1

grunt> █

Cloudera Live : Welco... cloudera@quickstart... cloudera@quickstart:~...

MAX()

The Pig Latin **MAX()** function is used to calculate the highest value for a column (numeric values or chararrays) in a single-column bag. While calculating the maximum value, the **Max()** function ignores the NULL values.

Note –

- To get the global maximum value, we need to perform a **Group All** operation, and calculate the maximum value using the **MAX()** function.
- To get the maximum value of a group, we need to group it using the **Group By** operator and proceed with the maximum function.

Syntax

Given below is the syntax of the **Max()** function.

```
grunt> Max (expression)
```

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below.

student_details.txt

```
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
```

File Edit View Search Terminal Help

008,Bharathi,Nambiayar,24,9848022333,Chennai,72

```
[cloudera@quickstart Desktop]$ hadoop fs -put student_details.txt smitrpatel
put: `smitrpatel/student_details.txt': File exists
[cloudera@quickstart Desktop]$ hadoop fs -put student_details.txt smitrpatel
put: `smitrpatel/student_details.txt': File exists
[cloudera@quickstart Desktop]$ cat student_details.txt
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72

[cloudera@quickstart Desktop]$ hadoop fs -put student_details.txt smitrpatel
put: `smitrpatel/student_details.txt': File exists
[cloudera@quickstart Desktop]$ hadoop fs -ls smitrpatel
Found 9 items
-rw-r--r-- 1 cloudera cloudera 194 2021-08-25 17:32 smitrpatel/customers.txt
-rw-r--r-- 1 cloudera cloudera 144 2021-08-30 18:24 smitrpatel/emp_bonus.txt
-rw-r--r-- 1 cloudera cloudera 144 2021-08-30 18:25 smitrpatel/emp_sales.txt
-rw-r--r-- 1 cloudera cloudera 304 2021-08-25 18:08 smitrpatel/employee.txt
-rw-r--r-- 1 cloudera cloudera 275 2021-08-25 18:08 smitrpatel/employee_contact.txt
-rw-r--r-- 1 cloudera cloudera 124 2021-08-25 17:32 smitrpatel/orders.txt
-rw-r--r-- 1 cloudera cloudera 240 2021-08-25 18:33 smitrpatel/student_data1.txt
-rw-r--r-- 1 cloudera cloudera 76 2021-08-25 18:33 smitrpatel/student_data2.txt
-rw-r--r-- 1 cloudera cloudera 351 2021-08-25 18:44 smitrpatel/student_details.txt
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrpatel/student_details.txt
001,Rajiv,Reddy,21,9848022337,Hyderabad
002,siddarth,Battacharya,22,9848022338,Kolkata
003,Rajesh,Khanna,22,9848022339,Delhi
004,Preethi,Agarwal,21,9848022330,Pune
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar
006,Archana,Mishra,23,9848022335,Chennai
007,Komal,Nayak,24,9848022334,trivendram
008,Bharathi,Nambiayar,24,9848022333,Chennai
[cloudera@quickstart Desktop]$
```

cloudera@quickstart:~/Desktop

And we have loaded this file into Pig with the relation name **student_details** as shown below.

```
grunt> student_details = LOAD
  'hdfs://localhost:9000/pig_data/student_details.txt'
  USING PigStorage(',')
  as (id:int, firstname:chararray, lastname:chararray, age:int,
  phone:chararray, city:chararray, gpa:int);
```

```
2021-08-30 19:34:36,468 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
  rray, city:chararray, gpa:int);smitrpatel/student_details.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, age:int, phone:char)
grunt> student_group_all = Group student_details All;
grunt> dump student_group_all;
2021-08-30 19:37:18,574 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2021-08-30 19:37:18,578 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, Duplicate
ForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NewPartitionFilterOptim
izer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]}
2021-08-30 19:37:18,621 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic?
false
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-08-30 19:37:18,658 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-08-30 19:37:18,659 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.per
cent is not set, set to default 0.3
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of
```

Calculating the Maximum GPA

We can use the built-in function **MAX()** (case-sensitive) to calculate the maximum

value from a set of given numerical values. Let us group the relation **student_details** using the **Group All** operator, and store the result in the relation named **student_group_all** as shown below.

```
grunt> student_group_all = Group student_details All;
```

```
2021-08-30 19:34:36,468 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
array, city:chararray, gpa:int);mitrpatel/student_details.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, age:int, phone:char)
grunt> student_group_all = Group student_details All;
grunt> dump student_group_all;
2021-08-30 19:37:18,574 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2021-08-30 19:37:18,578 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]}
2021-08-30 19:37:18,621 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic?
false
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-08-30 19:37:18,658 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-08-30 19:37:18,659 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of
```

This will produce a relation as shown below.

```
grunt> Dump student_group_all;
```

```
File Edit View Search Terminal Help

Success!

Job Stats (time in seconds):
JobId  Maps  Reduces MaxMapTime      MinMapTime      AvgMapTime      MedianMapTime      MaxReduceTime      MinReduceTime      AvgReduceTime      MedianReduceTime
alias  Feature Outputs
job_1630371021450_0007  1      1      13      13      13      13      11      11      11      11      student_details,student_group_all      GROUP_BY
dfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp566066035

Input(s):
Successfully read 8 records (746 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/student_details.txt"

Output(s):
Successfully stored 1 records (418 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp566066035"

Counters:
Total records written : 1
Total bytes written : 418
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0007

2021-08-30 19:38:18,127 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 8 time(s).
2021-08-30 19:38:18,127 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 19:38:18,127 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 19:38:18,127 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 19:38:18,128 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 19:38:18,134 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 19:38:18,134 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(all,{(8,Bharathi,Nambiayar,24,9848022333,Chennai),(7,Komal,Nayak,24,9848022334,trivendram),(6,Archana,Mishra,23,9848022335,Chennai),(5,Trupthi,Mohany,23,9848022336,Bhuwaneshwar),(4,Preethi,Agarwal,21,9848022330,Pune),(3,Rajesh,Khanna,22,9848022339,Delhi),(2,siddarth,Battacharya,22,9848022338,Calcutta),(1,Rajiv,Reddy,21,9848022337,Hyderabad)})}

grunt> cloudera@quickstart:~/Desktop
[Cloudera Live : Welco... cloudera@quickstart:... cloudera@quickstart:~...]
```

Let us now calculate the global maximum of GPA, i.e., maximum among the GPA values of all the students using the **MAX()** function as shown below.

```
grunt> student_gpa_max = foreach student_group_all
Generate  (student_details.firstname,
student_details.gpa), MAX(student_details.gpa);
```

```
2021-08-30 19:30:10,154 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
(all,{(8,Bharathi,Nambyayar,24,9848022333,Chennai),(7,Komal,Nayak,24,9848022334,trivendram),(6,Archana,Mishra,23,9848022335,Chennai),
hy,23,9848022336,Bhuwaneshwar,),(4,Preethi,Agarwal,21,9848022330,Pune),(3,Rajesh,Khanna,22,9848022339,Delhi),(2,siddarth,Battacharya
lkata,),(1,Rajiv,Reddy,21,9848022337,Hyderabad,)})}
grunt> student_gpa_max = foreach student_group_all Generate (student_details.firstname, student_details.gpa), MAX(student_details.gpa);
grunt> ■
```

Verification

Verify the relation **student_gpa_max** using the **DUMP** operator as shown below.

```
grunt> Dump student_gpa_max;
```

```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.12.0 0.12.0-cdh5.12.0 cloudera 2021-08-30 19:43:09 2021-08-30 19:44:03 GROUP_BY
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime
lias Feature Outputs
job 1630371021450_0008 1 1 11 11 11 11 11 11 11 11 student_details,student_gpa_max,student_group_all
ROUP_BY hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp1937260669,
Input(s):
Successfully read 8 records (746 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/student_details.txt"
Output(s):
Successfully stored 1 records (111 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp1937260669"
Counters:
Total records written : 1
Total bytes written : 111
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1630371021450_0008

2021-08-30 19:44:03,275 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 8 time(s).
2021-08-30 19:44:03,275 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 19:44:03,276 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 19:44:03,276 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker
address
2021-08-30 19:44:03,276 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 19:44:03,282 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 19:44:03,282 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
({{(Bharathi),(Komal),(Archana),(Trupthi),(Preethi),(Rajesh),(siddarth),(Rajiv)},{{(),(),(),(),(),(),(),()}}}),)
grunt> ■
```

Output

It will produce the following output, displaying the contents of the relation **student_gpa_max**.

```
2021-08-30 19:44:03,282 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
({{(Bharathi),(Komal),(Archana),(Trupthi),(Preethi),(Rajesh),(siddarth),(Rajiv)},{{(),(),(),(),(),(),(),()}}}),)
```

MIN()

The **MIN()** function of Pig Latin is used to get the minimum (lowest) value (numeric or chararray) for a certain column in a single-column bag. While calculating the minimum value, the **MIN()** function ignores the NULL values.

Note –

- To get the global minimum value, we need to perform a **Group All** operation, and calculate the minimum value using the **MIN()** function.
- To get the minimum value of a group, we need to group it using the **Group By** operator and proceed with the minimum function.

Syntax

Given below is the syntax of the **MIN()** function.

```
grunt> MIN(expression)
```

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below.

student_details.txt

```
001,Rajiv,Reddy,21,9848022337,Hyderabad,89
002,siddarth,Battacharya,22,9848022338,Kolkata,78
003,Rajesh,Khanna,22,9848022339,Delhi,90
004,Preethi,Agarwal,21,9848022330,Pune,93
005,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar,75
006,Archana,Mishra,23,9848022335,Chennai,87
007,Komal,Nayak,24,9848022334,trivendram,83
008,Bharathi,Nambiayar,24,9848022333,Chennai,72
```

And we have loaded this file into Pig with the relation named **student_details** as shown below.

```
grunt> student_details = LOAD
  'hdfs://localhost:9000/pig_data/student_details.txt'
  USING PigStorage(',')
  as (id:int, firstname:chararray, lastname:chararray, age:int,
  phone:chararray, city:chararray, gpa:int);
```

```
2021-08-30 19:34:36,468 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-08-30 19:34:36,468 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Input paths: rray, city:chararray, gpa:int);mitrpatel/student_details.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, age:int, phone:char)
grunt> student_group_all = Group student_details All;
grunt> dump student_group_all;
2021-08-30 19:37:18,574 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2021-08-30 19:37:18,578 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, Duplicate
ForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeForEach, NewPartitionFilterOptimizer,
PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]}
2021-08-30 19:37:18,621 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic?
false
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-08-30 19:37:18,658 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-08-30 19:37:18,659 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.per
cent is not set, set to default 0.3
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of
```

Calculating the Minimum GPA

We can use the built-in function **MIN()** (case sensitive) to calculate the minimum value from a set of given numerical values. Let us group the relation **student_details** using the **Group All** operator, and store the result in the relation named **student_group_all** as shown below

```
grunt> student_group_all = Group student_details All;
```

```
2021-08-30 19:34:36,468 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
array, city:chararray, gpa:int);mitrpatel/student_details.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, age:int, phone:char)
grunt> student_group_all = Group student_details All;
grunt> dump student_group_all;
2021-08-30 19:37:18,574 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2021-08-30 19:37:18,578 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, Duplicate
ForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeForEach, NewPartitionFilterOpti
mizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimiz
er]}
2021-08-30 19:37:18,621 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic?
false
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-08-30 19:37:18,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-08-30 19:37:18,658 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-08-30 19:37:18,659 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.per
cent is not set, set to default 0.3
2021-08-30 19:37:18,664 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of
```

It will produce a relation as shown below.

```
grunt> Dump student_group_all;
```

```
File Edit View Search Terminal Help

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTIme AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime
alias Feature Outputs
job_1630371021450_0007 1 1 13 13 13 13 11 11 11 11 student_details,student_group_all GROUP_BY
dfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp566066035,

Input(s):
Successfully read 8 records (746 bytes) from: "dfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/student_details.txt"

Output(s):
Successfully stored 1 records (418 bytes) in: "dfs://quickstart.cloudera:8020/tmp/temp-864208229/tmp566066035"

Counters:
Total records written : 1
Total bytes written : 418
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0007

2021-08-30 19:38:18,127 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 8 time(s).
2021-08-30 19:38:18,127 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 19:38:18,127 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 19:38:18,127 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.
ddress
2021-08-30 19:38:18,128 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 19:38:18,134 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 19:38:18,134 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(all,{(8,Bharathi,Nambiar,24,9848022333,Chennai,), (7,Komal,Nayak,24,9848022334,trivendram,), (6,Archana,Mishra,23,9848022335,Chennai,), (5,Trupthi,Mohan
hy,23,9848022336,Bhuwaneshwar,), (4,Preethi,Agarwal,21,9848022330,Pune,), (3,Rajesh,Khanna,22,9848022339,Delhi,), (2,siddarth,Battacharya,22,9848022338,Al
lkata,), (1,Rajiv,Reddy,21,9848022337,Hyderabad,)},)
cloudera@quickstart:~/Desktop
grunt> [Cloudera Live : Welco... cloudera@quickstart:... cloudera@quickstart:...]
```

Let us now calculate the global minimum of GPA, i.e., minimum among the GPA values of all the students using the **MIN()** function as shown below.

```
grunt> student_gpa_min = foreach student_group_all
Generate  (student_details.firstname,
student_details.gpa), MIN(student_details.gpa);
```

```
\l{kata },(1,Rajiv,Reddy,21,9848022337,Hyderabad,))  
grunt> student_gpa_min = foreach student_group_all Generate (student_details.firstname, student_details.gpa), MIN(student_details.gpa);  
2021-08-30 19:55:20,054 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS  
2021-08-30 19:55:20,054 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.a  
ddress  
grunt> student_gpa_min = foreach student_group_all Generate (student_details.firstname, student_details.gpa), MIN(student_details.gpa);  
grunt> dump student_gpa_min;|
```

Verification

Verify the relation **student_gpa_min** using the **DUMP** operator as shown below.

```
grunt> Dump student_gpa_min;
```

Output

It will produce the following output, displaying the contents of the relation **student_gpa_min**.

```
2021-08-30 19:57:33,727 [main] INFO
org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil -
Total input paths to process : 1
({{ (Bharathi), (Komal), (Archana), (Trupthi), (Preethi), (Rajesh), (sid
darth), (Rajiv) }}, { (), (), (), (), (), (), (), () }), )
```

PluckTuple()

After performing operations like join to differentiate the columns of the two schemas, we use the function **PluckTuple()**. To use this function, first of all, we have to define a string Prefix and we have to filter for the columns in a relation that begin with that prefix.

Syntax

Given below is the syntax of the **PluckTuple()** function.

```
DEFINE pluck PluckTuple(expression1)
DEFINE pluck PluckTuple(expression1,expression3)
pluck(expression2)
```

Example

Assume that we have two files namely **emp_sales.txt** and **emp_bonus.txt** in the HDFS directory **/pig_data/**. The **emp_sales.txt** contains the details of the employees of the sales department and the **emp_bonus.txt** contains the employee details who got bonus.

emp_sales.txt

```
1,Robin,22,25000,sales
2,BOB,23,30000,sales
3,Maya,23,25000,sales
4,Sara,25,40000,sales
5,David,23,45000,sales
6,Maggy,22,35000,sales
```

emp_bonus.txt

```
1,Robin,22,25000,sales
2,Jaya,23,20000,admin
3,Maya,23,25000,sales
4,Alia,25,50000,admin
5,David,23,45000,sales
6,Omar,30,30000,admin
```

And we have loaded these files into Pig, with the relation names **emp_sales** and **emp_bonus** respectively.

```
grunt> emp_sales = LOAD
'hdfs://localhost:9000/pig_data/emp_sales.txt'
USING PigStorage(',')
as (sno:int, name:chararray, age:int, salary:int,
dept:chararray);

grunt> emp_bonus = LOAD
'hdfs://localhost:9000/pig_data/emp_bonus.txt'
USING PigStorage(',')
as (sno:int, name:chararray, age:int, salary:int,
dept:chararray);
```

```
grunt>emp_sales = load 'smitrpatel/emp_sales.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);
```

```
grunt> emp_sales = load 'smitrpatel/emp_sales.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);
grunt> emp_bonus = load 'smitrpatel/emp_bonus.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);
grunt> join_data = join emp_sales by sno, emp_bonus by sno;
grunt> dump join_data;
```

Join these two relations using the **join** operator as shown below.

```
grunt> join_data = join emp_sales by sno, emp_bonus by sno;
```

```
grunt> emp_sales = load 'smitrpatel/emp_sales.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);
grunt> emp_bonus = load 'smitrpatel/emp_bonus.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salary:int, dept:chararray);
grunt> join_data = join emp_sales by sno, emp_bonus by sno;
grunt> dump join_data;
```

Verify the relation **join_data** using the **Dump** operator.

```
grunt> Dump join_data;
```

```
File Edit View Search Terminal Help

Job Stats (time in seconds):
JobId  Maps   Reduces  MaxMapTime   MinMapTime   AvgMapTime   MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReducetime
alias  Feature Outputs
job_1630371021450_0011  2      1      25      24      24      24      11      11      11      11      emp_bonus,emp_sales,join_data  HASH_JOIN      hdfs
//quickstart.cloudera:8020/tmp/temp-340327799/tmp-1824817083

Input(s):
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"

Output(s):
Successfully stored 6 records (326 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-340327799/tmp-1824817083"

Counters:
Total records written : 6
Total bytes written : 326
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0011

2021-08-30 20:04:28,154 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 20:04:28,154 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 20:04:28,155 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 20:04:28,155 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 20:04:28,164 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 20:04:28,164 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(1,Robin,22,25000,sales ,1,Robin,22,25000,sales )
(2,BOB,23,30000,sales ,2,Jaya,23,20000,admin )
(3,Maya,23,25000,sales ,3,Maya,23,25000,sales )
(4,Sara,25,40000,sales ,4,Alia,25,50000,admin )
(5,David,23,45000,sales ,5,David,23,45000,sales )
(6,Maggy,22,35000,sales ,6,Omar,30,30000,admin )
grunt>
```

Using **PluckTuple()** Function

Now, define the required expression by which you want to differentiate the columns using **PluckTuple()** function.

```
grunt> DEFINE pluck PluckTuple('a::');
```

```
";" ...  
Details at logfile: /home/cloudera/Desktop/pig_1630378251018.log  
grunt> DEFINE pluck PluckTuple('a::');  
grunt> data = foreach join_data generate Flatten(pluck(*));  
grunt> █
```

Filter the columns in the `join_data` relation as shown below.

```
grunt> data = foreach join_data generate FLATTEN(pluck(*));
```

```
";" ...
Details at logfile: /home/cloudera/Desktop/pig_1630378251018.log
grunt> DEFINE pluck PluckTuple('a:');
grunt> data = foreach join_data generate Flatten(pluck(*));
grunt> 
```

Describe the relation named **data** as shown below.

```
grunt> Describe data;
```

Since we have defined the expression as “`a::`”, the columns of the `emp_sales` schema are plucked as `emp_sales::column name` and the columns of the `emp_bonus` schema are plucked as `emp_bonus::column name`

SIZE()

The **SIZE()** function of Pig Latin is used to compute the number of elements based on any Pig data type.

Syntax

Given below is the syntax of the **SIZE()** function.

```
grunt> SIZE (expression)
```

The return values vary according to the data types in Apache Pig.

Data type Value

int, long, float, double :- For all these types, the size function returns 1.

Char array :- For a char array the size() function returns the number of characters in the array.

Byte array :- For a bytearray the size() function returns the number of bytes in the array.

Tuple :- For a tuple the size() function returns number of fields in the tuple.

Bag :- For a bag the size() function returns number of tuples in the bag.

Map :- For a map the size() function returns the number of key/value pairs in the map.

Example

Assume that we have a file named **employee.txt** in the HDFS directory **/pig_data/** as shown below.

employee.txt

```
1,John,2007-01-24,250
2,Ram,2007-05-27,220
3,Jack,2007-05-06,170
3,Jack,2007-04-06,100
4,Jill,2007-04-06,220
5,Zara,2007-06-06,300
5,Zara,2007-02-06,350
```

```
[cloudera@quickstart Desktop]$ gedit employee.txt
[cloudera@quickstart Desktop]$ hadoop -put employee.txt smitrptael
Error: No command named '-put' was found. Perhaps you meant 'hadoop put'
[cloudera@quickstart Desktop]$ hadoop fs -put employee.txt smitrptael
[cloudera@quickstart Desktop]$ hadoop fs -put employee.txt smitrptael
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrptael/employee.txt
1,John,2007-01-24,250
2,Ram,2007-05-27,220
3,Jack,2007-05-06,170
3,Jack,2007-04-06,100
4,Jill,2007-04-06,220
5,Zara,2007-06-06,300
5,Zara,2007-02-06,350
[cloudera@quickstart Desktop]$
```

And we have loaded this file into Pig with the relation name **employee_data** as shown below.

```
grunt> employee_data = LOAD
'hdfs://localhost:9000/pig_data/ employee.txt' USING
PigStorage(',',)
as (id:int, name:chararray, workdate:chararray,
daily_typing_pages:int);
```

```
grunt> employee_data = load 'smitrptael/employee.txt' using PigStorage(',') as (id:int, name:chararray, workdate:chararray, daily_typing_pages:int);
grunt> size = FOREACH employee_data Generate SIZE(name);
grunt> dump size;
```

Calculating the Size of the Type

To calculate the size of the type of a particular column, we can use the **SIZE()** function. Let's calculate the size of the name type as shown below.

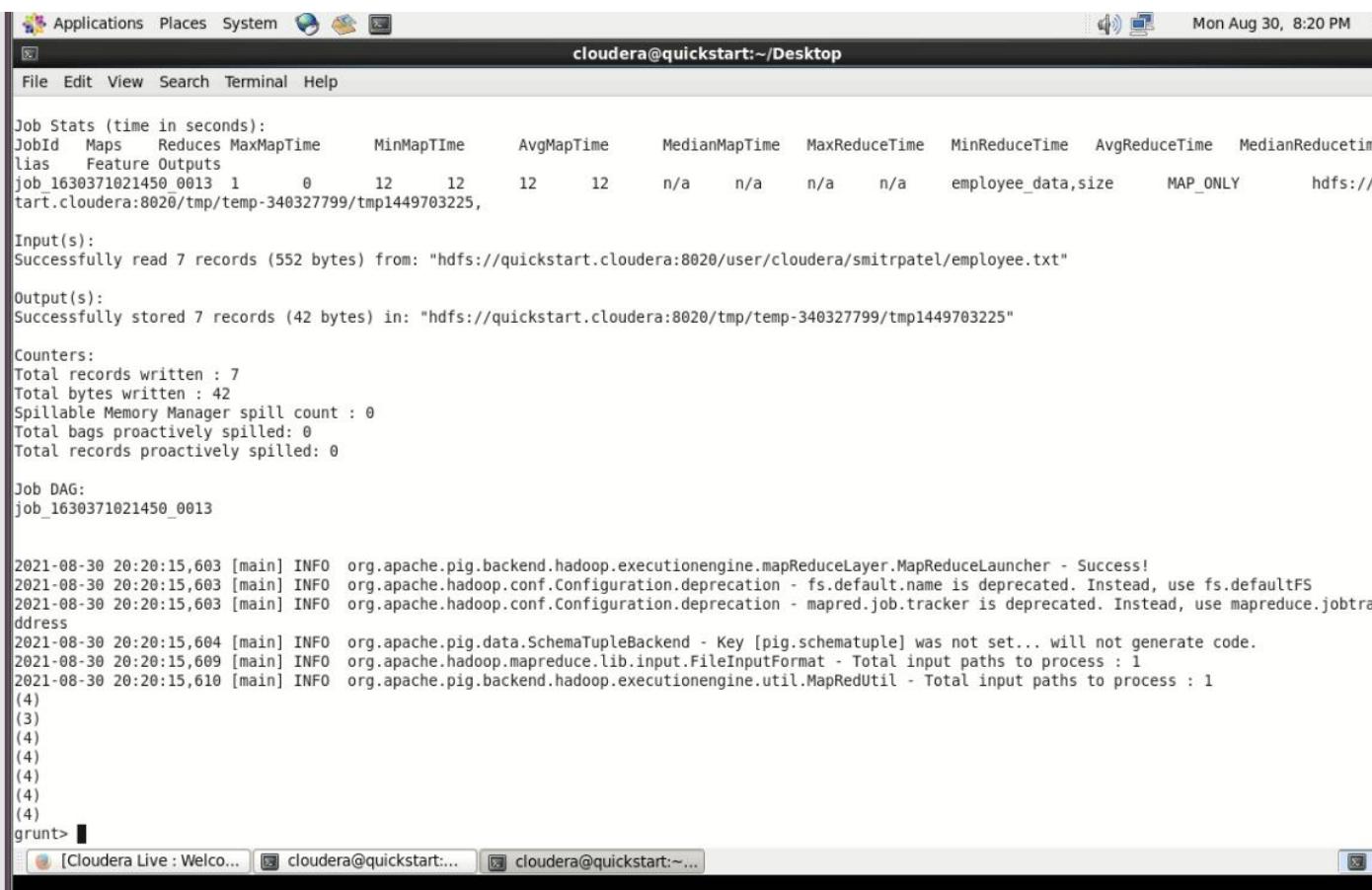
```
grunt> size = FOREACH employee_data GENERATE SIZE(name);
```

```
grunt> employee_data = load 'smitrptael/employee.txt' using PigStorage(',') as (id:int, name:chararray, workdate:chararray, daily_typing_pages:int);
grunt> size = FOREACH employee_data Generate SIZE(name);
grunt> dump size;
```

Verification

Verify the relation **size** using the **DUMP** operator as shown below.

```
grunt> Dump size;
```



```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime     MinMapTime     AvgMapTime     MedianMapTime   MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime
job_1630371021450_0013 1          0           12           12           n/a          n/a          n/a          employee_data.size      MAP_ONLY      hdfs://
tarf.cloudera:8020/tmp/temp-340327799/tmp1449703225

Input(s):
Successfully read 7 records (552 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee.txt"

Output(s):
Successfully stored 7 records (42 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-340327799/tmp1449703225"

Counters:
Total records written : 7
Total bytes written : 42
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0013

2021-08-30 20:20:15,603 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 20:20:15,603 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 20:20:15,603 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 20:20:15,604 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 20:20:15,609 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 20:20:15,610 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(4)
(3)
(4)
(4)
(4)
(4)
grunt> [Cloudera Live : Welco... cloudera@quickstart:... cloudera@quickstart:~...]
```

Output

It will produce the following output, displaying the contents of the relation **size** as follows. In the example, we have calculated the size of the **name** column. Since it is of varchar type, the **SIZE()** function gives you the number of characters in the name of each employee.

2021-08-30 20:20:15,609 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1

2021-08-30 20:20:15,610 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1

(4)

(3)

(4)

(4)

(4)

(4)

(4)

SUBTRACT()

The **SUBTRACT()** function of Pig Latin is used to subtract two bags. It takes two bags as inputs and returns a bag which contains the tuples of the first bag that are not in the second bag.

Syntax

Given below is the syntax of the **SUBTRACT()** function.

```
grunt> SUBTRACT(expression, expression)
```

Example

Assume that we have two files namely **emp_sales.txt** and **emp_bonus.txt** in the HDFS directory **/pig_data/** as shown below. The **emp_sales.txt** contains the details of the employees of the sales department and the **emp_bonus.txt** contains the employee details who got bonus.

emp_sales.txt

```
1,Robin,22,25000,sales
2,BOB,23,30000,sales
3,Maya,23,25000,sales
4,Sara,25,40000,sales
5,David,23,45000,sales
6,Maggy,22,35000,sales
```

emp_bonus.txt

```
1,Robin,22,25000,sales
2,Jaya,23,20000,admin
3,Maya,23,25000,sales
4,Alia,25,50000,admin
5,David,23,45000,sales
6,Omar,30,30000,admin
```

And we have loaded these files into Pig, with the relation names **emp_sales** and **emp_bonus** respectively.

```
grunt> emp_sales = LOAD
'hdfs://localhost:9000/pig_data/emp_sales.txt'
USING PigStorage(',')
  as (sno:int, name:chararray, age:int, salary:int,
dept:chararray);

grunt> emp_bonus = LOAD
'hdfs://localhost:9000/pig_data/emp_bonus.txt'
USING PigStorage(',')
  as (sno:int, name:chararray, age:int, salary:int,
dept:chararray);
```

```
grunt> emp_sales = load'smitrpatel/emp_sales.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salaray:int, dept:chararray);
grunt> emp_bonus = load'smitrpatel/emp_bonus.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salaray:int, dept:chararray);
grunt> cogroup_data = COGROUP emp_sales by sno, emp_bonus by sno;
grunt> dump cogroup_data;
```

Let us now group the records/tuples of the relations **emp_sales** and **emp_bonus** with the key **sno**, using the **COGROUP** operator as shown below.

```
grunt> cogroup_data = COGROUP emp_sales by sno, emp_bonus by sno;
```

```
grunt> emp_sales = load'smitrpatel/emp_sales.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salaray:int, dept:chararray);
grunt> emp_bonus = load'smitrpatel/emp_bonus.txt' using PigStorage(',') as (sno:int, name:chararray, age:int, salaray:int, dept:chararray);
grunt> cogroup_data = COGROUP emp_sales by sno, emp_bonus by sno;
grunt> dump cogroup_data;
```

Verify the relation **cogroup_data** using the **DUMP** operator as shown below. **grunt> Dump cogroup_data;**

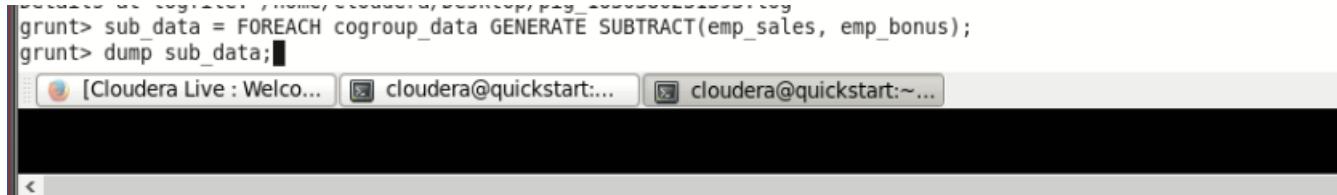
```
File Edit View Search Terminal Help
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
alias  Feature Outputs
job_1630371021450_0014 2 1 22 22 22 22 9 9 9 9 cogroup_data,emp_bonus,emp_sales COGRO
//quickstart.cloudera:8020/tmp/temp-918163753/tmp963766796,
Input(s):
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"
Output(s):
Successfully stored 6 records (367 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp963766796"
Counters:
Total records written : 6
Total bytes written : 367
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1630371021450_0014

2021-08-30 20:28:49,938 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED
ecting to job history server
2021-08-30 20:28:50,611 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 20:28:50,614 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 20:28:50,614 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.job
address
2021-08-30 20:28:50,614 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 20:28:50,631 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 20:28:50,631 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(1,{{1,Robin,22,25000,sales }},{{1,Robin,22,25000,sales }})
(2,{{2,BOB,23,30000,sales }},{{2,Jaya,23,20000,admin }})
(3,{{3,Maya,23,25000,sales }},{{3,Maya,23,25000,sales }})
(4,{{4,Sara,25,40000,sales }},{{4,Alia,25,50000,admin }})
(5,{{5,David,23,45000,sales }},{{5,David,23,45000,sales }})
(6,{{6,Maggy,22,35000,sales }},{{6,Omar,30,30000,admin }})
grunt>
```

Subtracting One Relation from the Other

Let us now subtract the tuples of **emp_bonus** relation from **emp_sales** relation. The resulting relation holds the tuples of **emp_sales** that are not there in **emp_bonus**.

```
grunt> sub_data = FOREACH cogroup_data GENERATE
SUBTRACT (emp_sales, emp_bonus);
```

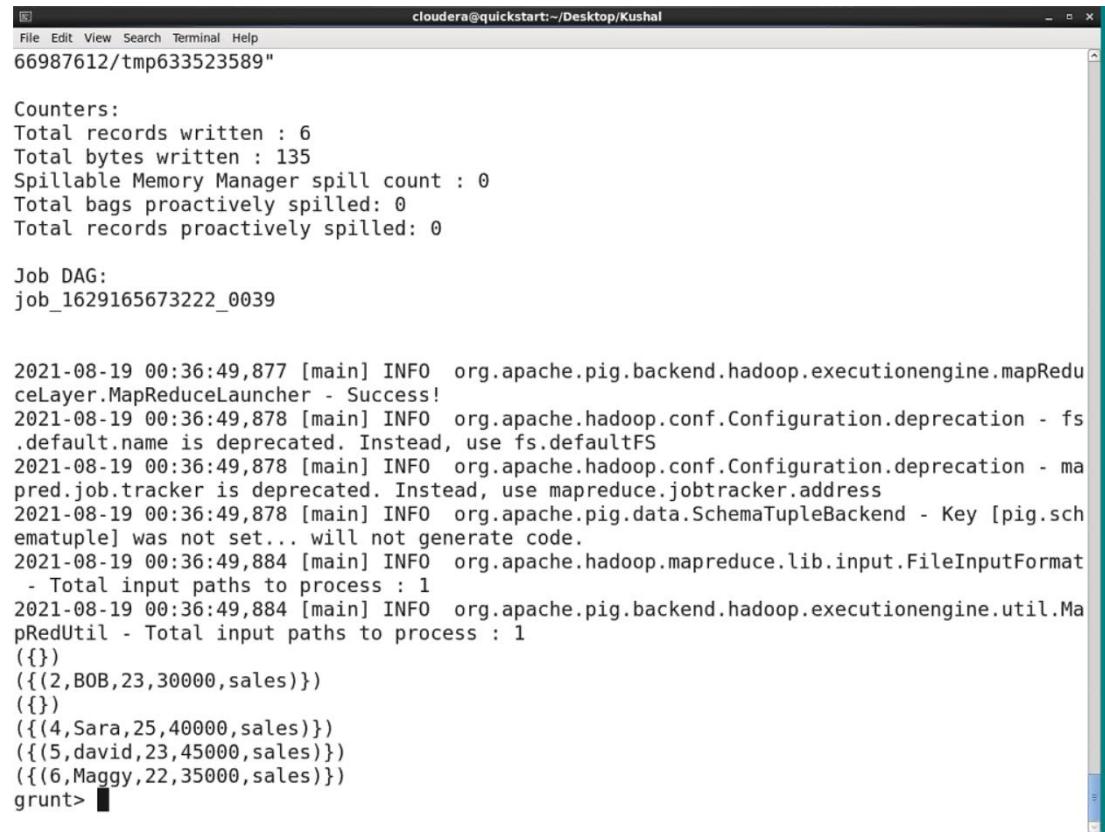


```
grunt> sub_data = FOREACH cogroup_data GENERATE SUBTRACT(emp_sales, emp_bonus);
grunt> dump sub_data;
```

Verification

Verify the relation **sub_data** using the DUMP operator as shown below. The **emp_sales** relation holds the tuples that are not there in the relation **emp_bonus**.

```
grunt> Dump sub_data;
```



```
File Edit View Search Terminal Help
66987612/tmp633523589"

Counters:
Total records written : 6
Total bytes written : 135
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1629165673222_0039

2021-08-19 00:36:49,877 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-19 00:36:49,878 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-19 00:36:49,878 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-19 00:36:49,878 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-19 00:36:49,884 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-19 00:36:49,884 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
({})
({{2,BOB,23,30000,sales}})
({})
({{4,Sara,25,40000,sales}})
({{5,david,23,45000,sales}})
({{6,Maggy,22,35000,sales}})
grunt>
```

```
({ })
({ (2,BOB,23,30000,sales) })
({ })
({ (4,Sara,25,40000,sales) })
({ })
({ (6,Maggy,22,35000,sales) })
```

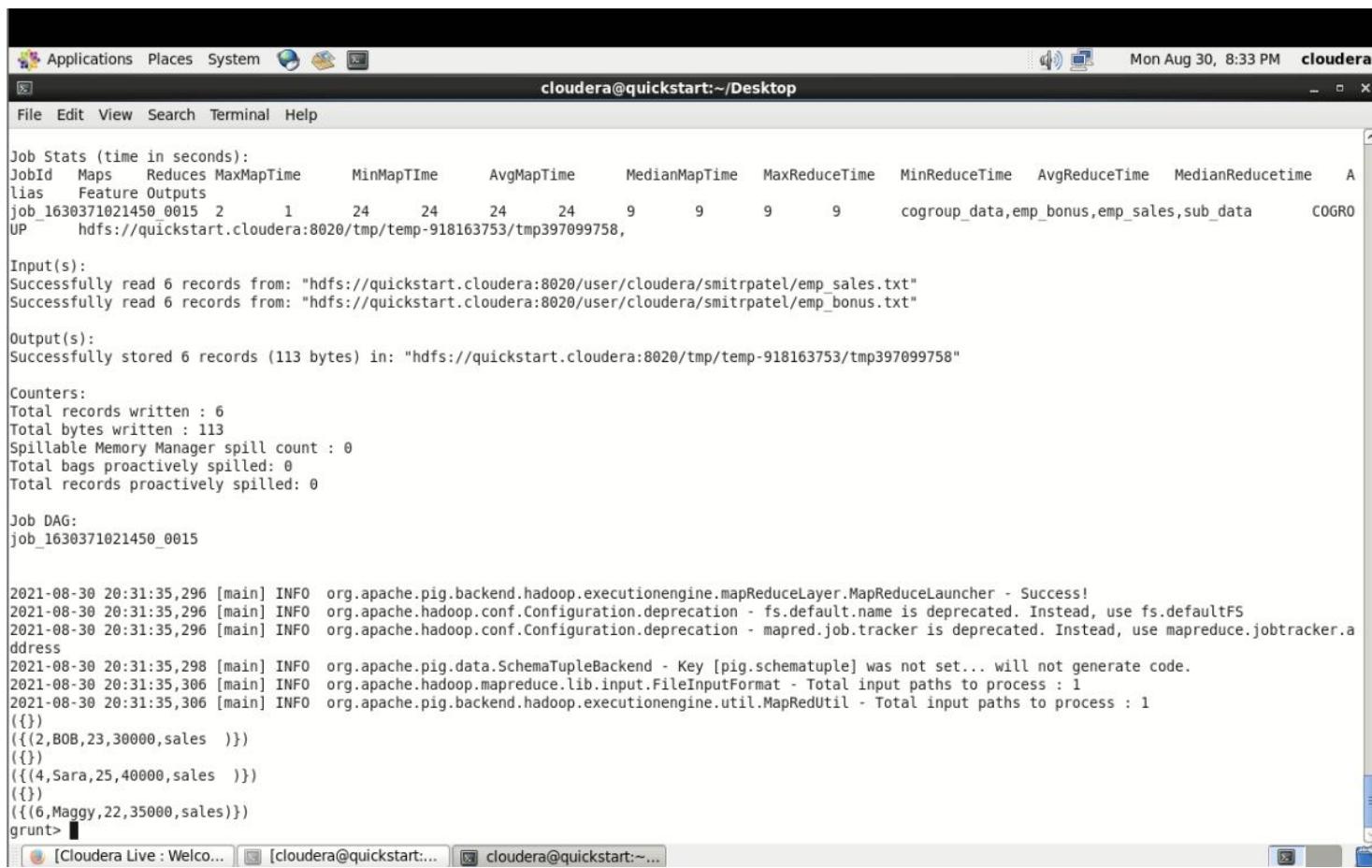
In the same way, let us subtract the **emp_sales** relation from **emp_bonus** relation

as shown below.

```
grunt> sub_data = FOREACH cogroup_data GENERATE  
SUBTRACT (emp_bonus, emp_sales);
```

Verify the contents of the **sub_data** relation using the Dump operator as shown below.

```
grunt> Dump sub_data;
```



```
Applications Places System  Mon Aug 30, 8:33 PM cloudera  
File Edit View Search Terminal Help  
cloudera@quickstart:~/Desktop  
Job Stats (time in seconds):  
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime A  
alias Feature Outputs  
job_1630371021450_0015 2 1 24 24 24 24 9 9 9 9 cogroup_data,emp_bonus,emp_sales,sub_data COGRO  
UP hdfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp397099758,  
Input(s):  
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_sales.txt"  
Successfully read 6 records from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/emp_bonus.txt"  
Output(s):  
Successfully stored 6 records (113 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp397099758"  
Counters:  
Total records written : 6  
Total bytes written : 113  
Spillable Memory Manager spill count : 0  
Total bags proactively spilled: 0  
Total records proactively spilled: 0  
Job DAG:  
job_1630371021450_0015  
  
2021-08-30 20:31:35,296 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!  
2021-08-30 20:31:35,296 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS  
2021-08-30 20:31:35,296 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.a  
ddress  
2021-08-30 20:31:35,298 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.  
2021-08-30 20:31:35,306 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-08-30 20:31:35,306 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
({})  
({{2,BOB,23,30000,sales }})  
({})  
({{4,Sara,25,40000,sales }})  
({})  
({{6,Maggy,22,35000,sales}})  
grunt>
```

SUM()

You can use the **SUM()** function of Pig Latin to get the total of the numeric values of a column in a single-column bag. While computing the total, the **SUM()** function ignores the NULL values.

Note –

- To get the global sum value, we need to perform a **Group All** operation, and calculate the sum value using the **SUM()** function.
- To get the sum value of a group, we need to group it using the **Group By** operator and proceed with the sum function.

Syntax

Given below is the syntax of the **SUM()** function.

```
grunt> SUM(expression)
```

Example

Assume that we have a file named **employee.txt** in the HDFS directory **/pig_data/** as shown below.

employee.txt

```
1, John, 2007-01-24, 250
2, Ram, 2007-05-27, 220
3, Jack, 2007-05-06, 170
3, Jack, 2007-04-06, 100
4, Jill, 2007-04-06, 220
5, Zara, 2007-06-06, 300
5, Zara, 2007-02-06, 350
```

And we have loaded this file into Pig with the relation name **employee_data** as shown below.

```
grunt> employee_data = LOAD
'hdfs://localhost:9000/pig_data/ employee.txt' USING
PigStorage(',')
as (id:int, name:chararray, workdate:chararray,
daily_typing_pages:int);
```

```
2021-08-30 20:31:35,306 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
({})
({{2,BOB,23,30000,sales }})
({})
({{4,Sara,25,40000,sales }})
({})
({{6,Maggy,22,35000,sales}})
grunt> employee_data = load 'smitrpatel/employee.txt' using PigStorage(',') as (id:int, name:chararray, workdate:chararray, daily_typing_pages:int);
grunt> employee_group = Group employee_data all;
```

Calculating the Sum of All GPA

To demonstrate the **SUM()** function, let's try to calculate the total number of pages typed daily of all the employees. We can use the Apache Pig's built-in function **SUM()** (case sensitive) to calculate the sum of the numerical values. Let us group the relation **employee_data** using the **Group All** operator, and store the result in the relation

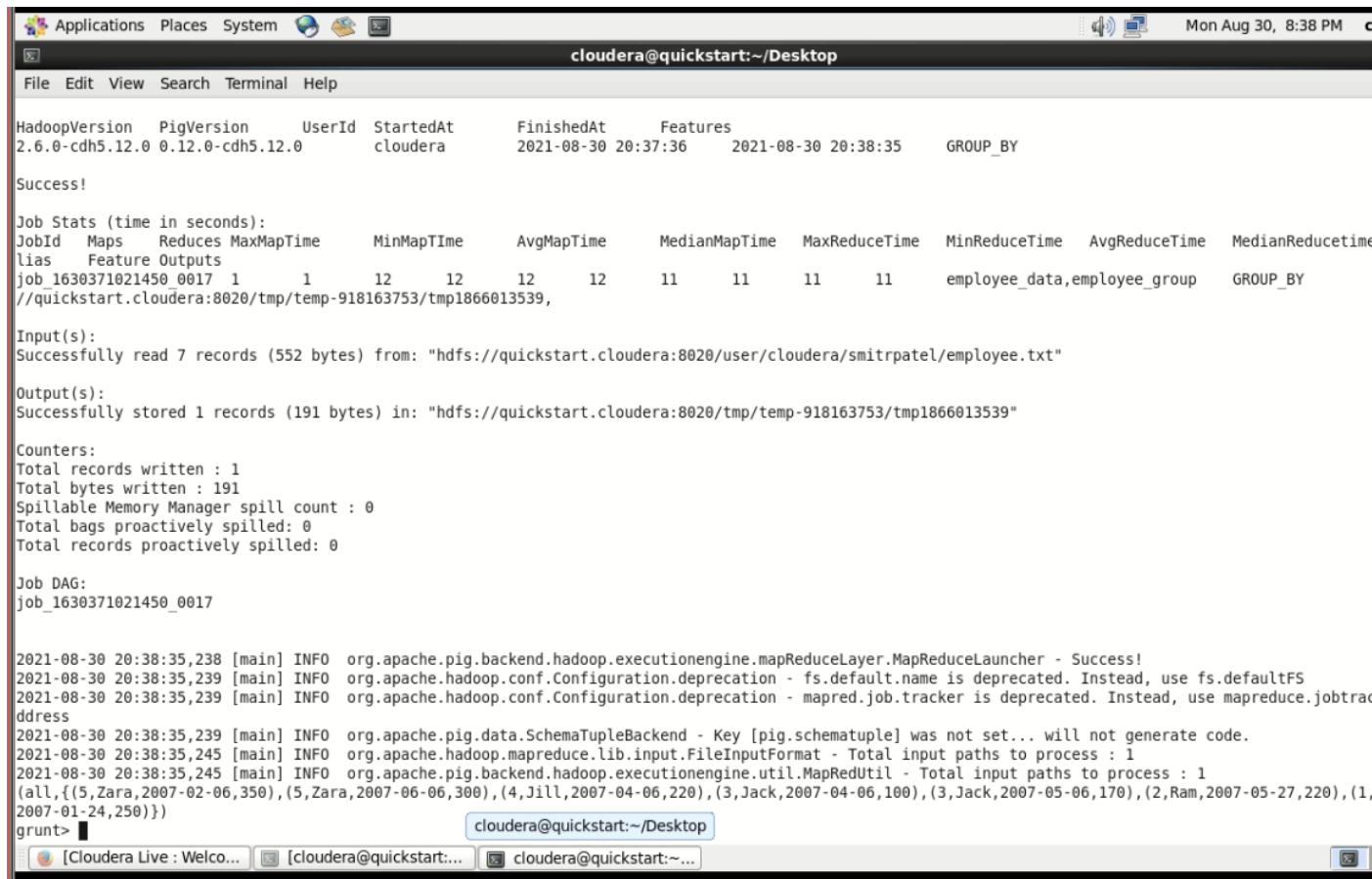
named `employee_group` as shown below.

```
grunt> employee_group = Group employee_data all;
```

```
2021-08-30 20:31:35,306 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
({})
({{2,BOB,23,30000,sales }})
({})
({{4,Sara,25,40000,sales }})
({})
({{6,Maggy,22,35000,sales}})
grunt> employee_data = load 'smitrpatel/employee.txt' using PigStorage(',') as (id:int, name:chararray, workdate:chararray, daily_typing_pages:int)
grunt> employee_group = Group employee_data all;
```

It will produce a relation as shown below.

```
grunt> Dump employee_group;
```



The screenshot shows a terminal window on a Linux desktop. The title bar says "cloudera@quickstart:~/Desktop". The terminal content displays the execution of a Pig Latin script to group employee data by name. The output shows the job statistics, including the number of maps and reduce tasks, and the resulting grouped data.

```
Applications Places System  Mon Aug 30, 8:38 PM
File Edit View Search Terminal Help
cloudera@quickstart:~/Desktop
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.12.0 0.12.0-cdh5.12.0 cloudera 2021-08-30 20:37:36 2021-08-30 20:38:35 GROUP_BY
Success!
Job Stats (time in seconds):
JobId  Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime
alias  Feature Outputs
job_1630371021450_0017 1 1 12 12 12 11 11 11 11 employee_data,employee_group GROUP_BY
//quickstart.cloudera:8020/tmp/temp-918163753/tmp1866013539,
Input(s):
Successfully read 7 records (552 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee.txt"
Output(s):
Successfully stored 1 records (191 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp1866013539"
Counters:
Total records written : 1
Total bytes written : 191
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1630371021450_0017

2021-08-30 20:38:35,238 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 20:38:35,239 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 20:38:35,239 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 20:38:35,239 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 20:38:35,245 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 20:38:35,245 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(all,{{5,Zara,2007-02-06,350),(5,Zara,2007-06-06,300),(4,Jill,2007-04-06,220),(3,Jack,2007-04-06,100),(3,Jack,2007-05-06,170),(2,Ram,2007-05-27,220),(1,John,2007-01-24,250)})
grunt> cloudera@quickstart:~/Desktop
[Cloudera Live : Welco... [cloudera@quickstart:... [cloudera@quickstart:~...
```

```
2021-08-30 20:38:35,245 [main] INFO
org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total
input paths to process : 1

(all,{{5,Zara,2007-02-06,350),(5,Zara,2007-06-06,300),(4,Jill,2007-04-06,220),(3,Jack,2007-04-06,100),(3,Jack,2007-05-06,170),(2,Ram,2007-05-27,220),(1,John,2007-01-24,250)})
```

Let us now calculate the global sum of the pages typed daily.

```
grunt> student_workpages_sum = foreach employee_group
  Generate
    (employee_data.name,employee_data.daily_typing_pages),SUM(employee_data.daily_typing_pages);
```

```
File Edit View Search Terminal Help
grunt> student_workpages_sum = foreach employee_group Generate (employee_data.name, employee_data.daily_typeing_pages), SUM(employee_data.daily_typeing_pages);
grunt> █
```

Verification

Verify the relation **student_workpages_sum** using the **DUMP** operator as shown below.

```
grunt> Dump student_workpages_sum;
```

```
pRedUtil - Total input paths to process : 1
(({(Zara),(Zara),(Jill),(Jack),(Jack),(Ram),(John)},{(250),(300),(220),(100),(170),(220),(250)}),1510)
grunt> █
```

Output

It will produce the following output, displaying the contents of the relation **student_workpages_sum** as follows.

```
(({ (Zara), (Zara), (Jill) , (Jack) , (Jack) , (Ram) , (John) } ,
{ (350) , (300) , (220) , (100) , (170) , (220) , (250) } ),1610)
```

TOKENIZE()

The **TOKENIZE()** function of Pig Latin is used to split a string (which contains a group of words) in a single tuple and returns a bag which contains the output of the split operation.

Syntax

Given below is the syntax of the **TOKENIZE()** function.

```
grunt> TOKENIZE(expression [, 'field_delimiter'])
```

As a delimiter to the **TOKENIZE()** function, we can pass space [], double quote ["], coma [,], parenthesis [()], star [*].

Example

Assume that we have a file named **student_details.txt** in the HDFS directory **/pig_data/** as shown below. This file contains the details of a student like id, name, age and city. If we closely observe, the name of the student includes first and last names separated by space [].

student_details.txt

```
001,Rajiv Reddy,21,Hyderabad
002,siddarth Battacharya,22,Kolkata
003,Rajesh Khanna,22,Delhi
004,Preethi Agarwal,21,Pune
005,Trupthi Mohanthy,23,Bhuwaneshwar
006,Archana Mishra,23 ,Chennai
007,Komal Nayak,24,trivendram
008,Bharathi Nambiayar,24,Chennai
```

We have loaded this file into Pig with the relation name **student_details** as shown below.

```
grunt> student_details = LOAD
'hdfs://localhost:9000/pig_data/student_details.txt'
USING PigStorage(',')
as (id:int, name:chararray, age:int, city:chararray);
```

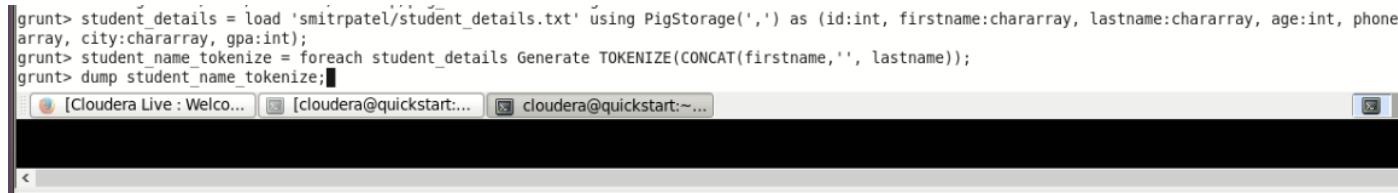
```
grunt> student_details = load 'smitrpatel/student_details.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, age:int, phone:chararray, city:chararray, gpa:int);
grunt> student_name_tokenize = foreach student_details Generate TOKENIZE(CONCAT(firstname,'', lastname));
grunt> dump student_name_tokenize;
```



Tokenizing a String

We can use the **TOKENIZE()** function to split a string. As an example let us split the name using this function as shown below.

```
grunt> student_name_tokenize = foreach student_details
Generate TOKENIZE(name);
```



```

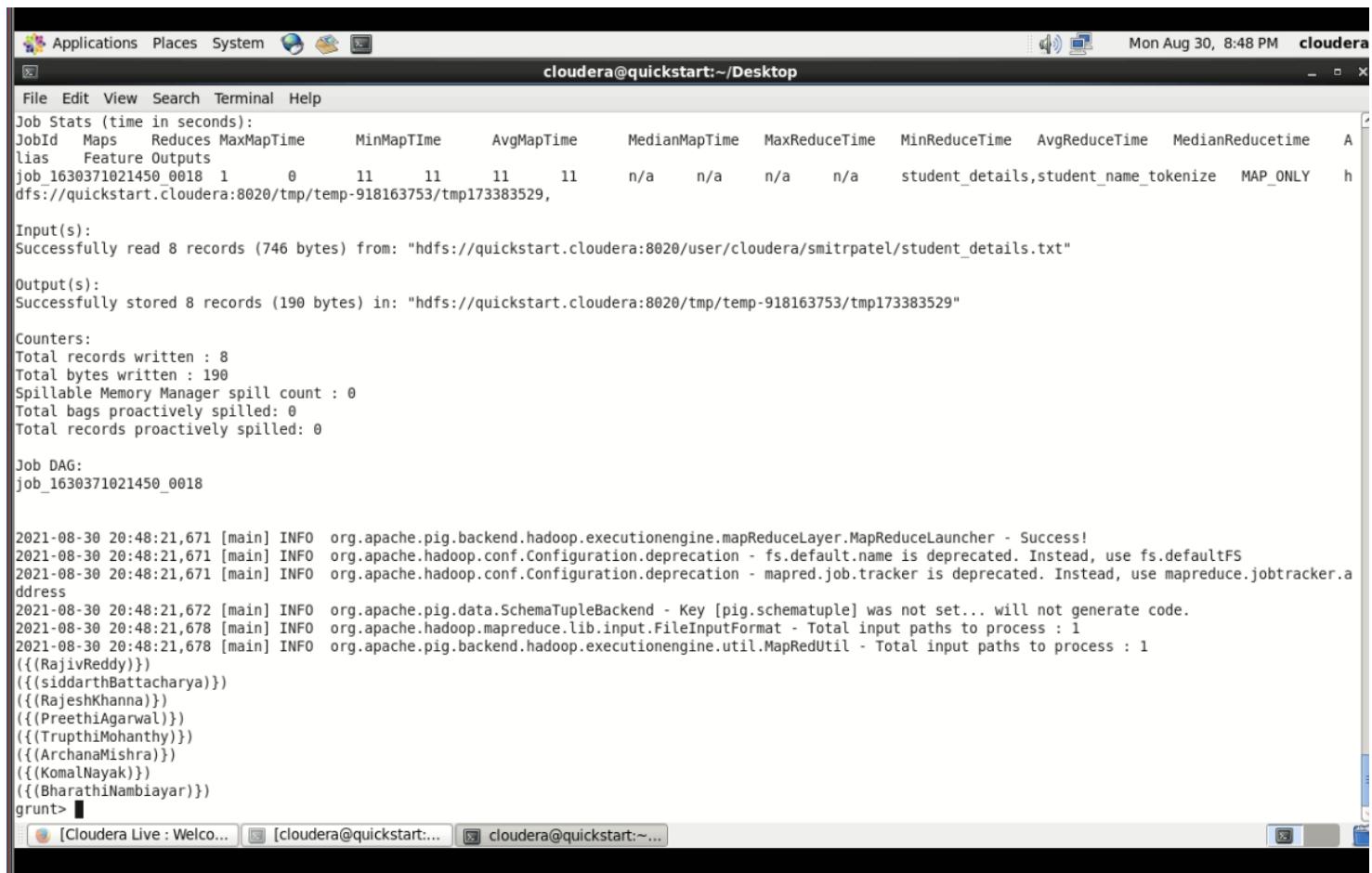
grunt> student_details = load 'smitrpatel/student_details.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, age:int, phone
array, city:chararray, gpa:int);
grunt> student_name_tokenize = foreach student_details Generate TOKENIZE(CONCAT(firstname,'', lastname));
grunt> dump student_name_tokenize;

```

Verification

Verify the relation **student_name_tokenize** using the **DUMP** operator as shown below.

```
grunt> Dump student_name_tokenize;
```



```

File Edit View Search Terminal Help
Job Stats (time in seconds):
JobID Maps Reduces MaxMapTime      MinMapTime      AvgMapTime      MedianMapTime      MaxReduceTime      MinReduceTime      AvgReduceTime      MedianReducetime
Job 1630371021450 0018 1          0          11          11          11          n/a          n/a          n/a          n/a          student_details,student_name_tokenize      MAP_ONLY
dfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp173383529,
Input(s):
Successfully read 8 records (746 bytes) from: "dfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/student_details.txt"
Output(s):
Successfully stored 8 records (190 bytes) in: "dfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp173383529"
Counters:
Total records written : 8
Total bytes written : 190
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1630371021450_0018

2021-08-30 20:48:21,671 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 20:48:21,671 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 20:48:21,671 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-30 20:48:21,672 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 20:48:21,678 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 20:48:21,678 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
{{(RajivReddy)}}
{{(siddarthBattacharya)}}
{{(RajeshKhanna)}}
{{(PreethiAgarwal)}}
{{(TrupthiMohanty)}}
{{(ArchanaMishra)}}
{{(KomalNayak)}}
{{(BharathiNambiar)}}
grunt> 

```

Output

It will produce the following output, displaying the contents of the relation **student_name_tokenize** as follows.

```

File Edit View Search Terminal Help
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime     MinMapTime     AvgMapTime     MedianMapTime   MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime
job_1630371021450_0018 1      0      11      11      11      11      n/a      n/a      n/a      n/a      student_details,student_name_tokenize  MAP_ON
dfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp173383529,

Input(s):
Successfully read 8 records (746 bytes) from: "dfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/student_details.txt"

Output(s):
Successfully stored 8 records (190 bytes) in: "dfs://quickstart.cloudera:8020/tmp/temp-918163753/tmp173383529"

Counters:
Total records written : 8
Total bytes written : 190
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1630371021450_0018

2021-08-30 20:48:21,671 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-08-30 20:48:21,671 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-30 20:48:21,671 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.job.tracker
2021-08-30 20:48:21,672 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-08-30 20:48:21,678 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-08-30 20:48:21,678 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
{{{{RajivReddy}}}}
{{{{siddarthBattacharya}}}}
{{{{RajeshKhanna}}}}
{{{{PreethiAgarwal}}}}
{{{{TrupthiMohanty}}}}
{{{{ArchanaMishra}}}}
{{{{KomalNayak}}}}
{{{{BharathiNambiar}}}}
grunt> 

```

Other Delimiters

In the same way, including space [], the TOKENIZE() function accepts double quote [" "], coma [,], parenthesis [()], star [*] as delimiters.

Example

Suppose there is a file named **details.txt** with students details like id, name, age, and city. Under the name column this file contains the first name and the last name of the students separated by various delimiters as shown below.

details.txt

```

001,"siddarth""Battacharya",22,Kolkata
002,Rajesh*Khanna,22,Delhi
003,(Preethi) (Agarwal),21,Pune

```

We have loaded this file into Pig with the relation name **details** as shown below.

```

grunt> details = LOAD
'hdfs://localhost:9000/pig_data/details.txt'
USING PigStorage(',')
as (id:int, name:chararray, age:int, city:chararray);

```

Now, try to separate the first name and the last name of the students using TOKENIZE() as follows.

```
grunt> tokenize_data = foreach details Generate TOKENIZE(name);
```

```
File Edit View Search Terminal Help
grunt> tokenize_data = foreach details generate TOKENIZE(name);
grunt> █
```

On verifying the **tokenize_data** relation using dump operator you will get the following result.

```
grunt> Dump tokenize_data;
```

```
Successfully stored 4 records (92 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-166
5987612/tmp1149022552"
```

```
Counters:
Total records written : 4
Total bytes written : 92
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
```

```
Job DAG:
job_1629165673222_0042
```

```
2021-08-19 00:47:30,353 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapRedu
ceLayer.MapReduceLauncher - Success!
2021-08-19 00:47:30,354 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs
.default.name is deprecated. Instead, use fs.defaultFS
2021-08-19 00:47:30,354 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - ma
pred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-19 00:47:30,355 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.sch
ematuple] was not set... will not generate code.
2021-08-19 00:47:30,366 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
- Total input paths to process : 1
2021-08-19 00:47:30,366 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma
pRedUtil - Total input paths to process : 1
{{(siddarth),(Battacharya)}}
{{(Rajesh),(Khanna)}}
{{(Preethi),(Agarwal)}}
()
grunt> █
```

