

**SMIT R PATEL**

**19162121031**

**SEM 5**

**PRACTICAL 10**

**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:**

The Load and Store functions in Apache Pig are used to determine how the data goes and comes out of Pig. These functions are used with the load and store operators. Given below is the list of load and store functions available in Pig.

**Function & Description**

1. **PigStorage()**

To load and store structured files.

2. **TextLoader()**

To load unstructured data into Pig.

3. **BinStorage()**

To load and store data into Pig using machine readable format.

4. **Handling Compression**

In Pig Latin, we can load and store compressed data.

# PigStorage()

The **PigStorage()** function loads and stores data as structured text files. It takes a delimiter using which each entity of a tuple is separated as a parameter. By default, it takes '\t' as a parameter.

## Syntax

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

```
grunt> PigStorage(field_delimiter)
```

## Example

Let us suppose we have a file named **student\_data.txt** in the HDFS directory named **/data** with the following content.

```
001,Rajiv,Reddy,9848022337,Hyderabad  
002,siddarth,Battacharya,9848022338,Kolkata  
003,Rajesh,Khanna,9848022339,Delhi  
004,Preethi,Agarwal,9848022330,Pune  
005,Trupthi,Mohanthy,9848022336,Bhuwaneshwar  
006,Archana,Mishra,9848022335,Chennai.
```

We can load the data using the PigStorage function as shown below.

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

```
[cloudera@quickstart Desktop]$ gedit student data1.txt
[cloudera@quickstart Desktop]$ cat student data1.txt
001,Rajiv,Reddy,9848022337,Hyderabad
002,siddarth,Battacharya,9848022338,Kolkata
003,Rajesh,Khanna,9848022339,Delhi
004,Preethi,Agarwal,9848022330,Pune
005,Trupthi,Mohanthy,9848022336,Bhuwaneshwar
006,Archana,Mishra,9848022335,Chennai
[cloudera@quickstart Desktop]$ ls
drwxr-xr-x  2 cloudera cloudera  4096 2021-09-07 00:40 smitrpatel/Practica
lExam
-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 -copyFromLocal student_data1.txt smitrpatel
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrpatel/student_data1.txt
001,Rajiv,Reddy,9848022337,Hyderabad
002,siddarth,Battacharya,9848022338,Kolkata
003,Rajesh,Khanna,9848022339,Delhi
004,Preethi,Agarwal,9848022330,Pune
005,Trupthi,Mohanthy,9848022336,Bhuwaneshwar
006,Archana,Mishra,9848022335,Chennai
[cloudera@quickstart Desktop]$
```

```
student = load'smitrpatel/student_data1.txt' using PigStorage(',') as (id:int, firstname:chararray, lastname:chararray, phone:chararray, city:chararray);
```

```
grunt> student = LOAD 'smitrpatel/student_data1.txt' using PigStorage(',') as
>> (id:int, firstname:chararray, lastname:chararray, phone:chararray, city:chararray);
grunt> |
```

In the above example, we have seen that we have used comma (',') delimiter. Therefore, we have separated the values of a record using (,).

In the same way, we can use the **PigStorage()** function to store the data in to HDFS directory as shown below.

```
grunt> STORE student INTO ' hdfs://localhost:9000/pig_Output/  
' USING PigStorage (',');
```

```
grunt> student = LOAD 'smitrpatel/student_data1.txt' using PigStorage(',') as
>> (id:int, firstname:chararray, lastname:chararray, phone:chararray, city:chararray);
grunt> STORE student into 'smit7output' USING PigStorage(',');
2021-09-07 03:19:429 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN
2021-09-07 03:19:466 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NewPoptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, Partitioner]}
2021-09-07 03:19:492 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.textoutputformat.separator is deprecated. Instead, use output.textoutputformat.separator
```

cloudera@quickstart:~...

mouse pointer outside or press Ctrl+Alt.

```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
0030/jobdetails.jsp?jobid=job_1630982973645_0020
2021-09-07 03:19:24,661 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
2021-09-07 03:19:35,720 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete
2021-09-07 03:19:39,877 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2021-09-07 03:19:39,943 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2021-09-07 03:19:39,944 [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 2021-09-07 03:19:19 2021-09-07 03:19:39 UNKNOWN

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime A
alias Feature Outputs
job_1630982973645_0020 1 0 2 2 2 n/a n/a n/a student MAP_ONLY hdfs://quickstart.cloudera:8028/user/cloudera/smit7output,
```

Input(s):  
Successfully read 7 records (637 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/student\_data1.txt"

Output(s):  
Successfully stored 7 records (236 bytes) in: "hdfs://quickstart.cloudera:8020/user/cloudera/smit7output"

Counters:

Total records written : 7  
Total bytes written : 236  
Spillable Memory Manager spill count : 0  
Total bags proactively spilled: 0  
Total records proactively spilled: 0

Job DAG:  
job\_1630982973645\_0020

```
2021-09-07 03:19:40,033 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 4 time(s).
2021-09-07 03:19:40,033 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt>
```

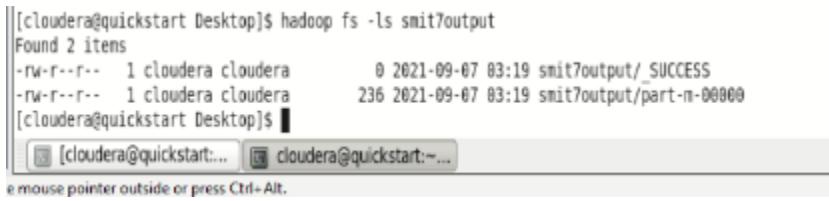
cloudera@quickstart:~...

This will store the data into the given directory. You can verify the data as shown below.

## Verification

You can verify the stored data as shown below. First of all, list out the files in the directory named **pig\_output** using **ls** command as shown below.

```
$ hdfs dfs -ls 'hdfs://localhost:9000/pig_Output/'
```



```
[cloudera@quickstart Desktop]$ hadoop fs -ls smit7output
Found 2 items
-rw-r--r-- 1 cloudera cloudera      0 2021-09-07 03:19 smit7output/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 236 2021-09-07 03:19 smit7output/part-m-00000
[cloudera@quickstart Desktop]$
```

You can observe that two files were created after executing the **Store** statement.

Then, using the **cat** command, list the contents of the file named **part-m-00000** as shown below.

```
$ hdfs dfs -cat 'hdfs://localhost:9000/pig_Output/part-m-00000'
```

```
1,Rajiv,Reddy,9848022337,Hyderabad
2,siddarth,Battacharya,9848022338,Kolkata
3,Rajesh,Khanna,9848022339,Delhi
4,Preethi,Agarwal,9848022330,Pune
5,Trupthi,Mohanthy,9848022336,Bhuwaneshwar
6,Archana,Mishra,9848022335,Chennai
```

```
cloudera@quickstart Desktop]$ hadoop fs -ls smit7output
Found 2 items
-rw-r--r-- 1 cloudera cloudera 0 2021-09-07 03:19 smit7output/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 236 2021-09-07 03:19 smit7output/part-m-00000
[cloudera@quickstart Desktop]$ hadoop fs -cat smit7output/part-m-00000
1,Rajiv,Reddy,9848022337,Hyderabad
2,siddarth,Battacharya,9848022338,Kolkata
3,Rajesh,Khanna,9848022339,Delhi
4,Preethi,Agarwal,9848022338,Pune
5,Trupti,Mohanthy,9848022336,Bhuwaneshwar
6,Archana,Mishra,9848022335,Chennai
....
```

```
[cloudera@quickstart Desktop]$
```

```
[cloudera@quickstart:... [cloudera@quickstart:~...]
```

```
the mouse pointer outside or press Ctrl+Alt.
```

# TextLoader()

The Pig Latin function **TextLoader()** is a Load function which is used to load unstructured data in UTF-8 format.

## Syntax

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

```
grunt> TextLoader()
```

## Example

Let us assume there is a file with named **stu\_data.txt** in the HDFS directory named **/data** as shown below.

```
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
```

```
[cloudera@quickstart Desktop]$ gedit stu_data.txt
[cloudera@quickstart Desktop]$ ^C
[cloudera@quickstart Desktop]$ cat stu_data.txt
cat: stu_data.txt: No such file or directory
[cloudera@quickstart Desktop]$ cat stu_data.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,Konal_Nayak,24,trivendran)
(008,Bharathi_Nambiayar,24,Chennai)

[cloudera@quickstart Desktop]$ hadoop fs -put stu_data.txt smitrpatel
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrpatel/stu_data.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,Konal_Nayak,24,trivendran)
(008,Bharathi_Nambiayar,24,Chennai)

[cloudera@quickstart Desktop]$
```

Now let us load the above file using the **TextLoader()** function.

```
grunt> details = LOAD
'hdfs://localhost:9000/pig_data/stu_data.txt' USING TextLoader();
```

```
grunt> details = load'smitrpatel/stu_data.txt' using TextLoader();
grunt> dump details;
```

You can verify the loaded data using the Dump operator.

```
grunt> dump details;
```

```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
job_1638982973645_0026

2021-09-07 03:52:25,519 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-07 03:52:25,521 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-07 03:52:25,521 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-07 03:52:25,540 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-07 03:52:25,540 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((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_Nambiyar,24,Chennai) )
()
grunt> 
```

# BinStorage()

The **BinStorage()** function is used to load and store the data into Pig using machine readable format. **BinStorage()** in Pig is generally used to store temporary data generated between the MapReduce jobs. It supports multiple locations as input.

## Syntax

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

```
grunt> BinStorage();
```

## Example

Assume that we have a file named **stu\_data.txt** in the HDFS directory **/pig\_data/** as shown below.

### Stu\_data.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
```

Let us load this data into Pig into a relation as shown below.

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

```
grunt> student_details = load'smitrpatel/stu_data.txt' using PigStorage(',') as
>> (id:int, firstname:chararray, age:int, city:chararray);
grunt>
```

Now, we can **store** this relation into the HDFS directory named **/pig\_data** using the **BinStorage()** function.

```
grunt> STORE student_details INTO
  'hdfs://localhost:9000/pig_Output/mydata' USING BinStorage();
```

```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
grunt>
grunt> STORE student_details into 'smitrpatel2' using BinStorage();
2021-09-09 02:50:00,447 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN
2021-09-09 02:50:00,448 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]}
2021-09-09 02:50:00,460 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2021-09-09 02:50:00,463 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-09-09 02:50:00,505 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-09-09 02:50:00,508 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-09-09 02:50:00,521 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.buffer.percent is not set, set to default 0.3
2021-09-09 02:50:01,283 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job5910166312429940354.jar
2021-09-09 02:50:07,447 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job5910166312429940354.jar created
2021-09-09 02:50:07,472 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2021-09-09 02:50:07,473 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2021-09-09 02:50:07,473 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
```

```
File Edit View Search Terminal Help
8,18],student_details[-1,-1] C: R:
2021-09-09 02:50:08,002 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:5030/jobdetails.jsp?jobid=job_1631177379117_0002
2021-09-09 02:50:08,058 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
2021-09-09 02:50:23,932 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete
2021-09-09 02:50:28,155 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2021-09-09 02:50:28,155 [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 2021-09-09 02:50:00 2021-09-09 02:50:28 UNKNOWN
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime A
alias Feature Outputs
job_1631177379117_0002 1 0 4 4 4 n/a n/a n/a n/a student_details MAP_ONLY hdfs://quickstart.c
loudera:8020/user/cloudera/smitrpatel2,
Input(s):
Successfully read 8 records (653 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/stu_data.txt"
Output(s):
Successfully stored 8 records (385 bytes) in: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel2"
Counters:
Total records written : 8
Total bytes written : 385
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1631177379117_0002

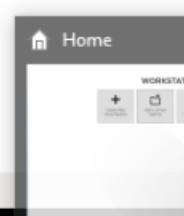
2021-09-09 02:50:28,276 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt>
grunt> S
cloudera@quickstart:... cloudera@quickstart:~...
```

After executing the above statement, the relation is stored in the given HDFS directory. You can see it using the HDFS **ls** command as shown below.

```
$ hdfs dfs -ls
hdfs://localhost:9000/pig_Output/mydata/
```

```
[cloudera@quickstart Desktop]$ hdfs dfs -ls smitrpatel2
Found 2 items
-rw-r--r-- 1 cloudera cloudera 0 2021-09-09 02:50 smitrpatel2
atel2/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 385 2021-09-09 02:50 smitrpatel2/part-m-00000
[cloudera@quickstart Desktop]$
```

```
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrpatel2/part-m-00000
Rajiv_Reddy
Hyderabad
Siddarth_Battacharya
Kolkata
Rajesh_Khanna
Delhi
Preethi_Agarwal
Bhubaneshwar
Archana_Mishra
Chennai
Komal_Nayak
trivendram
Sarathini_Nambiayar
Chennai[cloudera@quickstart Desktop]$
```



Now, load the data from the file **part-m-00000**.

```
grunt> result = LOAD 'hdfs://localhost:9000/pig_Output/b/part-m_00000' USING BinStorage();
```

```
grunt> result = load 'smitrpatel2/part-m-00000' using BinStorage();
2021-09-09 03:01:34,935 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:01:34,935 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:01:35,072 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:01:35,072 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> dump result;
```

```
[cloudera@quickstart:...]
```

Verify the contents of the relation as shown below

```
grunt> Dump result;
```

```
(1,Rajiv_Reddy,21,Hyderabad)
(2,siddarth_Battacharya,22,Kolkata)
(3,Rajesh_Khanna,22,Delhi)
(4,Preethi_Agarwal,21,Pune)
(5,Trupthi_Mohanthy,23,Bhuwaneshwar)
(6,Archana_Mishra,23,Chennai)
(7,Komal_Nayak,24,trivendram)
(8,Bharathi_Nambiayar,24,Chennai)
```

```
Applications Places System 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  A
alias  Feature Outputs
job 1631177379117 0004 1     4     4     4     n/a     n/a     n/a     result  MAP_ONLY      hdfs://quickstart.cloudera:80
20/tmp/temp937174768/tmp1920425366,0

Input(s):
Successfully read 8 records (774 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel2/part-m-00000"

Output(s):
Successfully stored 8 records (304 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp1920425366"

Counters:
Total records written : 8
Total bytes written : 304
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0004

2021-09-09 03:02:33,198 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:02:33,199 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:02:33,199 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.a
ddress
2021-09-09 03:02:33,200 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:02:33,208 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:02:33,208 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(1,Rajiv_Reddy,21,Hyderabad )
(2,siddarth_Battacharya,22,Kolkata )
(3,Rajesh_Khanna,22,Delhi )
(4,Preeti_Agarwal,21,Pune )
(5,Trupthi_Mohanthy,23,Bhuvaneshwar )
(6,Archana_Mishra,23,Chennai )
(7,Komal_Nayak,24,trivendram )
(8,Bharathi_Nambiar,24,Chennai)
grun> [cloudera@quickstart:~] cloudera@quickstart:~]
```

# Handling Compression

We can load and store compressed data in Apache Pig using the functions **BinStorage()** and **TextLoader()**.

## Example

Assume we have a file named **employee.txt.zip** in the HDFS directory **/pigdata/**. Then, we can load the compressed file into pig as shown below.

Using **PigStorage**:

```
grunt> data = LOAD  
'hdfs://localhost:9000/pig_data/employee.txt.zip'  
USING PigStorage('');
```

Using **TextLoader**:

```
grunt> data = LOAD  
'hdfs://localhost:9000/pig_data/employee.txt.zip'  
USING TextLoader;
```



```
grunt> data = load 'pig data/employee.txt.zip' USING TextLoader();  
2021-09-09 03:06:56,795 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:06:56,796 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:06:56,943 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:06:56,943 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
grunt> data = load 'pig data/employee.txt.zip' USING PigStorage('');  
2021-09-09 03:07:00,775 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:07:00,775 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:07:00,929 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:07:00,929 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
grunt>
```

In the same way, we can store the compressed files into pig as shown below.

Using **PigStorage**:

```
grunt> store data INTO  
'hdfs://localhost:9000/pig_Output/data.bz' USING  
PigStorage('');
```

# Bag & Tuple Functions

Given below is the list of Bag and Tuple functions.

## Function & Description

### 1. TOBAG()

To convert two or more expressions into a bag.

### 2. TOP()

To get the top N tuples of a relation.

### 3. TOTUPLE()

To convert one or more expressions into a tuple.

### 4. TOMAP()

To convert the key-value pairs into a Map.

## TOBAG()

The **TOBAG()** function of Pig Latin converts one or more expressions to individual tuples. And these tuples are placed in a bag.

## Syntax

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

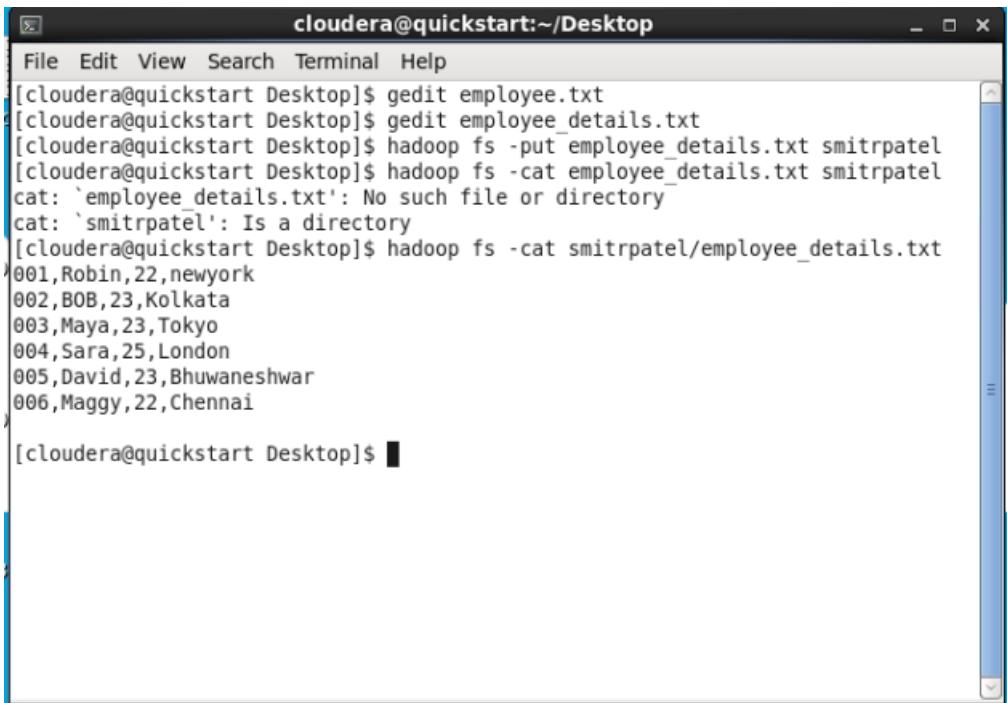
```
TOBAG (expression [, expression ...])
```

## Example

Assume we have a file named **employee\_details.txt** in the HDFS directory **/pig\_data/**, with the following content.

### **employee\_details.txt**

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
```



```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
[cloudera@quickstart Desktop]$ gedit employee.txt
[cloudera@quickstart Desktop]$ gedit employee_details.txt
[cloudera@quickstart Desktop]$ hadoop fs -put employee_details.txt smitrpatel
[cloudera@quickstart Desktop]$ hadoop fs -cat employee_details.txt smitrpatel
cat: 'employee_details.txt': No such file or directory
cat: 'smitrpatel': Is a directory
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrpatel/employee_details.txt
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai

[cloudera@quickstart Desktop]$
```

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

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

```
Details at logfile: /home/cloudera/Desktop/pig_1631179507531.log
grunt> emp_data = load 'smitrpatel/employee_details.txt' USING PigStorage(',') as (id:int, name:chararray, age:int, city:chararray);
2021-09-09 03:23:01,901 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:23:01,901 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:23:02,111 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:23:02,111 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
```

Let us now convert the id, name, age and city, of each employee (record) into a tuple as shown below.

```
tobag = FOREACH emp_data GENERATE TOBAG (id, name, age, city);
```

```
grunt> emp_data = load 'smitrpatel/employee_details.txt' USING PigStorage(',') as (id:int, name:chararray, age:int, city:chararray);
2021-09-09 03:23:01,901 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:23:01,901 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:23:02,111 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:23:02,111 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:23:06,841 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:23:06,841 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:23:07,056 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:23:07,056 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> [cloudera@quickstart:...]
```

## Verification

You can verify the contents of the **tobag** relation using the **Dump** operator as shown below.

```
grunt> DUMP tobag;
```

```
Applications Places System Thu Sep 9, 3:24 AM cloudera
File Edit View Search Terminal Help
cloudera@quickstart:~/Desktop
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime A
alias Feature Outputs
job_1631177379117_0007 1 0 4 4 4 n/a n/a n/a n/a emp_data,tobag MAP_ONLY hdfs://quickstart.clo
uder@8020/tmp/temp937174768/tmp-2039231129,,

Input(s):
Successfully read 7 records (530 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details.txt"

Output(s):
Successfully stored 7 records (212 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-2039231129"

Counters:
Total records written : 7
Total bytes written : 212
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0007

2021-09-09 03:24:51,362 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 3 time(s).
2021-09-09 03:24:51,362 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:24:51,363 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:24:51,363 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:24:51,363 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:24:51,374 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:24:51,374 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
{{(1),(Robin),(22),(newyork )}}
{{(2),(BOB),(23),(Kolkata )}}
{{(3),(Maya),(23),(Tokyo )}}
{{(4),(Sara),(25),(London )}}
{{(5),(David),(23),(Bhuvaneshwar )}}
{{(6),(Maggy),(22),(Chennai )}}
{{(),(),(),()}}
grunt> [cloudera@quickstart:~/Desktop]
```

# TOP()

The **TOP()** function of Pig Latin is used to get the top **N** tuples of a bag. To this function, as inputs, we have to pass a relation, the number of tuples we want, and the column name whose values are being compared. This function will return a bag containing the required columns.

## Syntax

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

```
grunt> TOP (topN, column, relation)
```

## Example

Assume we have a file named **employee\_details.txt** in the HDFS directory **/pig\_data/**, with the following content.

### **employee\_details.txt**

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

```
[cloudera@quickstart Desktop]$ gedit employee_details2.txt
[cloudera@quickstart Desktop]$ hadoop fs -put employee_details2.txt smitrpatel
[cloudera@quickstart Desktop]$ hadoop fs -cat smitrpatel/employee_details2.txt
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai

[cloudera@quickstart Desktop]$
```

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

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

```
grunt> emp_data = load 'smitrpatel/employee_details2.txt' USING PigStorage(',') as (id:int, name:chararray, age:int, city:chararray);
2021-09-09 03:27:40,857 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:27:40,857 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:27:41,087 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:27:41,087 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:27:41,087 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt>
```

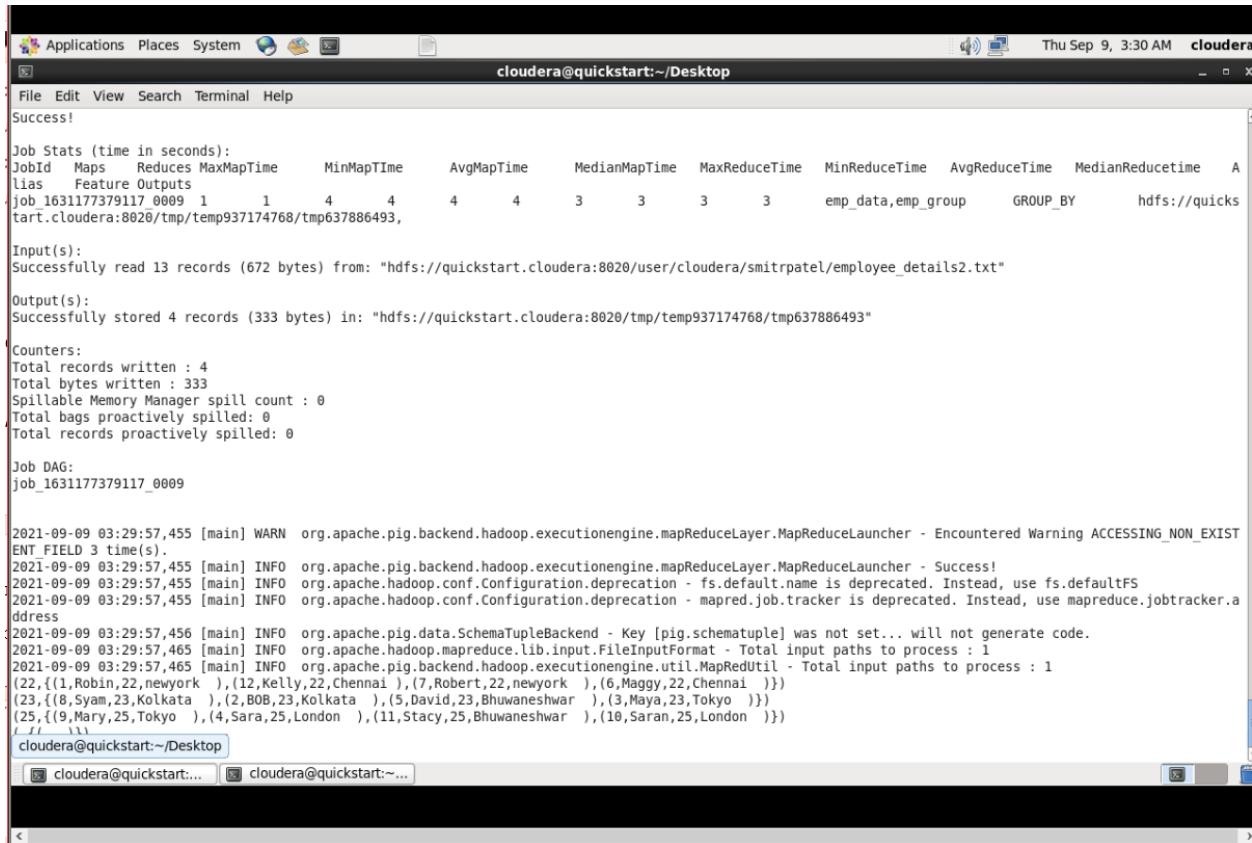
Group the relation **emp\_data** by age, and store it in the relation **emp\_group**

```
grunt> emp_group = Group emp_data BY age;
```

```
grunt> emp_group = Group emp_data BY age;
2021-09-09 03:28:16,550 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:28:16,550 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:28:16,762 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:28:16,762 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:28:16,762 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt>
```

Verify the relation **emp\_group** using the Dump operator as shown below.

```
grunt> Dump emp_group;
```



```
Success!
Job Stats (time in seconds):
JobId  Maps Reduces MaxMapTime     MinMapTime     AvgMapTime     MedianMapTime   MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReducetime   A
alias  Feature Outputs
job 1631177379117 0009  1      1      4      4      4      3      3      3      3      emp_data,emp_group      GROUP_BY      hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp637886493,
Input(s):
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"
Output(s):
Successfully stored 4 records (333 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp637886493"
Counters:
Total records written : 4
Total bytes written : 333
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1631177379117_0009

2021-09-09 03:29:57,455 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 3 time(s).
2021-09-09 03:29:57,455 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:29:57,455 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:29:57,455 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:29:57,456 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:29:57,465 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:29:57,465 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(22,{(1,Robin,22,newyork ),(12,Kelly,22,Chennai ),(7,Robert,22,newyork ),(6,Maggy,22,Chennai )})
(23,{(8,Syam,23,Kolkata ),(2,BOB,23,Kolkata ),(5,David,23,Bhuwaneshwar ),(3,Maya,23,Tokyo )})
(25,{(9,Mary,25,Tokyo ),(4,Sara,25,London ),(11,Stacy,25,Bhuwaneshwar ),(10,Saran,25,London )})
cloudera@quickstart:~/Desktop
```

Now, you can get the top two records of each group arranged in ascending order (**based on id**) as shown below.

```
grunt> data_top = FOREACH emp_group {
  top = TOP(2, 0, emp_data);
  GENERATE top;
}
```

```
grunt> data_top = FOREACH emp_group{top = TOP(2, 0, emp_data); GENERATE top; }
2021-09-09 03:31:40,343 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:31:40,343 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:31:40,555 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:31:40,555 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> [REDACTED]
```

In this example we are retrieving the top 2 tuples of a group having greater id. Since we are retrieving top 2 tuples basing on the **id**, we are passing the index of the column name id as second parameter of TOP() function.

### Verification

You can verify the contents of the **data\_top** relation using the **Dump** operator as shown below.

```
grunt> Dump data_top;
```

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

Job Stats (time in seconds):
JobId  Maps  Reduces MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  A
alias  Feature Outputs
job_1631177379117_0010  1  1  3  3  3  3  3  3  3  data_top,emp_data,emp_group  GROUP_BY,COMBINER  h
dfs://quickstart.cloudera:8020/tmp/temp937174768/tmp1656750643,

Input(s):
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"

Output(s):
Successfully stored 4 records (187 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp1656750643"

Counters:
Total records written : 4
Total bytes written : 187
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0010

2021-09-09 03:32:41,725 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 3 time(s).
2021-09-09 03:32:41,726 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:32:41,726 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:32:41,726 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:32:41,726 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:32:41,734 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:32:41,734 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
{{(7,Robert,22,newyork ),(12,Kelly,22,Chennai )}}
{{(5,David,23,Bhuwaneshwar ),(8,Syam,23,Kolkata )}}
{{(10,Saran,25,London ),(11,Stacy,25,Bhuwaneshwar )}}
{{(,,,)}}
grunt> █ cloudera@quickstart:~/Desktop
cloudera@quickstart:~/Desktop
```

# TOTUPLE()

The **TOTUPLE()** function is used convert one or more expressions to the data type **tuple**.

## Syntax

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

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

## Example

Assume we have a file named **employee\_details.txt** in the HDFS directory **/pig\_data/**, with the following content.

### employee\_details.txt

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
```

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

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

```
grunt> emp_data = load 'smitrpatel/employee_details.txt' USING PigStorage(',') as (id:int, name:chararray, age:int, city:chararray);
2021-09-09 03:33:38,819 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:33:38,819 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:33:39,042 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:33:39,042 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
cloudera@quickstart:~/Desktop'
```



Let us now convert the id, name and age of each student (record) into a tuple.

```
grunt> totuple = FOREACH emp_data GENERATE TOTUPLE (id, name, age);
```

```
grunt> totuple = FOREACH emp_data GENERATE TOTUPLE (id, name, age);
2021-09-09 03:35:17,380 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:35:17,380 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:35:17,618 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:35:17,618 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> ■
```

## Verification

You can verify the contents of the **totuple** schema using the **Dump** operator as shown below.

```
grunt> DUMP totuple;
```

```
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_1631177379117_0011  1      0      3      3      3      n/a      n/a      n/a      n/a      emp_data,totuple      MAP_ONLY      hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-142375159, 

Input(s):
Successfully read 7 records (530 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details.txt"

Output(s):
Successfully stored 7 records (105 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-142375159"

Counters:
Total records written : 7
Total bytes written : 105
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0011
0

2021-09-09 03:36:32,030 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:36:32,031 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:36:32,031 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:36:32,032 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:36:32,040 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:36:32,040 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((1,Robin,22))
((2,BOB,23))
((3,Maya,23))
((4,Sara,25))
((5,David,23))
((6,Maggy,22))
((,,))
grunt> ■
```

# TOMAP()

The **TOMAP()** function of Pig Latin is used to convert the key-value pairs into a Map.

## Syntax

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

```
grunt> TOMAP(key-expression, value-expression [, key-expression, value-expression ...])
```

## Example

Assume we have a file named **employee\_details.txt** in the HDFS directory **/pig\_data/**, with the following content.

### **employee\_details.txt**

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
```

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

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

Let us now take the name and age of each record as key-value pairs and convert them into map as shown below.

```
grunt> tomap = FOREACH emp_data GENERATE TOMAP(name, age);
```

```
grunt> tomap = FOREACH emp_data GENERATE TOMAP(name, age);
2021-09-09 03:37:57,220 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:37:57,221 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:37:57,443 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:37:57,443 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> ■
```

cloudera@quickstart:... cloudera@quickstart:~... cloudera@quickstart:~/Desktop

## Verification

You can verify the contents of the **tomap** relation using the **Dump** operator as shown below.

```
grunt> DUMP tomap;
```

```
Applications Places System Thu Sep 9, 3:39 AM 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
job_1631177379117_0012 1 0 3 3 3 n/a n/a n/a emp_data,tomap MAP_ONLY hdfs://quickstart.clo
uder:8020/tmp/temp937174768/tmp1974067546,
Input(s):
Successfully read 7 records (530 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details.txt"
Output(s):
Successfully stored 7 records (100 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp1974067546"
Counters:
Total records written : 7
Total bytes written : 100
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1631177379117_0012

2021-09-09 03:38:44,240 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:38:44,241 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:38:44,241 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.a
ddress
2021-09-09 03:38:44,241 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:38:44,249 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:38:44,249 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
([Robin#22])
([Bob#23])
([Maya#23])
([Sara#25])
([David#23])
([Maggy#22])
([null#])
grunt> ■
```

cloudera@quickstart:... cloudera@quickstart:~... cloudera@quickstart:~/Desktop

# String Functions

We have the following String functions in Apache Pig.

## Functions & Description

### 1. ENDSWITH(string, testAgainst)

To verify whether a given string ends with a particular substring.

### 2. STARTSWITH(string, substring)

Accepts two string parameters and verifies whether the first string starts with the second.

### 3. SUBSTRING(string, startIndex, stopIndex)

Returns a substring from a given string.

### 4. EqualsIgnoreCase(string1, string2)

To compare two strings ignoring the case.

### 5. INDEXOF(string, 'character', startIndex)

Returns the first occurrence of a character in a string, searching forward from a start index.

### 6. LAST\_INDEX\_OF(expression)

Returns the index of the last occurrence of a character in a string, searching backward from a start index.

### 7. LCFIRST(expression)

Converts the first character in a string to lower case.

### 8. UCFIRST(expression)

Returns a string with the first character converted to upper case.

#### 9. UPPER(expression)

UPPER(expression) Returns a string converted to upper case.

#### 10. LOWER(expression)

Converts all characters in a string to lower case.

#### 11.REPLACE(string, 'oldChar', 'newChar');

To replace existing characters in a string with new characters.

#### 12.STRSPLIT(string, regex, limit)

To split a string around matches of a given regular expression.

#### 13.STRSPLITTOBAG(string, regex, limit)

Similar to the STRSPLIT() function, it splits the string by given delimiter and returns the result in a bag.

#### 14. TRIM(expression)

Returns a copy of a string with leading and trailing whitespaces removed.

#### 15. LTRIM(expression)

Returns a copy of a string with leading whitespaces removed.

#### 16. RTRIM(expression)

Returns a copy of a string with trailing whitespaces removed.

# ENDSWITH()

This function accepts two String parameters, it is used to verify whether the first string ends with the second string.

## Syntax

```
grunt> ENDSWITH(string1, string2)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name age and city.

**emp.txt**

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

```
grunt> emp_data = load 'smitrpatel/employee_details2.txt' USING PigStorage(',') as (id:int, name:chararray, age:int, city:chararray);
2021-09-09 03:40:07,740 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:40:07,740 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:40:07,992 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:40:07,992 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> [cloudera@quickstart:...]
```

Following is an example of **ENDSWITH()** function, in this example we are verifying, weather the name of every employee ends with the character **n**.

```
grunt> emp_endswith = FOREACH emp_data GENERATE  
  (id, name), ENDSWITH ( name, 'n' );
```

```
grunt> emp_endswith = FOREACH emp_data GENERATE (id, name), ENDSWITH(name, 'n');  
2021-09-09 03:41:38,263 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:41:38,263 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:41:38,523 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:41:38,523 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
grunt> █  
cloudera@quickstart:~$ [cloudera@quickstart:~$]  
cloudera@quickstart:~$ █
```

The above statement verifies weather the name of the employee ends with the letter n. Since the names of the employees **Saran** and **Robin** ends with the letter n for these two tuples **ENDSWITH()** function returns the Boolean value **'true'** and for remaining tuples the value will be **'false'**.

The result of the statement will be stored in the relation named **emp\_endswith**.

Verify the content of the relation **emp\_endswith**, using the Dump operator as shown below.

```
grunt> Dump emp_endswith;
```

```
((1,Robin),true)  
((2,BOB),false)  
((3,Maya),false)  
((4,Sara),false)  
((5,David),false)  
((6,Maggy),false)  
((7,Robert),false)  
((8,Syam),false)  
((9,Mary),false)  
((10,Saran),true)  
((11,Stacy),false)  
((12,Kelly),false)
```

```
Applications Places System cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"
Output(s):
Successfully stored 13 records (194 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-1312275600"

Counters:
Total records written : 13
Total bytes written : 194
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0013

2021-09-09 03:42:42,392 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning UDF_WARNING_2 1 time(s).
2021-09-09 03:42:42,392 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:42:42,393 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:42:42,393 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:42:42,394 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:42:42,402 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:42:42,402 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((1,Robin),true)
((2,BOB),false)
((3,Maya),false)
((4,Sara),false)
((5,David),false)
((6,Maggy),false)
((7,Robert),false)
((8,Syam),false)
((9,Mary),false)
((10,Saran),true)
((11,Stacy),false)
((12,Kelly),false)
((,),)
grunt> █ cloudera@quickstart:~/Desktop
cloudera@quickstart:... [cloudera@quickstart:...]
```

# STARTSWITH()

This function accepts two string parameters. It verifies whether the first string starts with the second.

## Syntax

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

```
grunt> STARTSWITH(string, substring)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### emp.txt

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

## Example

Following is an example of the **STARTSWITH()** function. In this example, we have verified whether the names of all the employees start with the substring “Ro”.

```
grunt> startswith_data = FOREACH emp_data GENERATE  
(id, name), STARTSWITH (name, 'Ro');
```

```
grunt> startswith emp = FOREACH emp_data GENERATE (id, name), STARTSWITH(name, 'Ro');  
2021-09-09 03:47:21,576 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:47:21,576 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:47:21,829 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:47:21,829 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
grunt> dump startswith_data;
```

The above statement parses the names of all the employees if any of these names starts with the substring ‘Ro’. Since the names of the employees ‘Robin’ and ‘Robert’ starts with the substring ‘Ro’ for these two tuples the **STARTSWITH()** function returns the Boolean value ‘true’ and for remaining tuples the value will be ‘false’.

The result of the statement will be stored in the relation named **startswith\_data**. Verify the content of the relation **startswith\_data**, using the Dump operator as shown below.

```
grunt> Dump startswith_data;
```

```
grunt> dump startswith_data;  
2021-09-09 03:47:38,730 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:47:38,730 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:47:38,901 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 1003: Unable to find an operator for alias startswith_data  
Details at logfile: /home/cloudera/Desktop/pig_1631179507531.log  
grunt> dump startswith_emp;
```

```

Applications Places System Terminal Help
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
grunt> dump startswith_emp;
2021-09-09 03:47:56,978 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:47:56,978 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:47:57,137 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN
2021-09-09 03:47:57,139 [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]
2021-09-09 03:47:57,140 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for emp.data: $2, $3
2021-09-09 03:47:57,148 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2021-09-09 03:47:57,150 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-09-09 03:47:57,150 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-09-09 03:47:57,179 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-09-09 03:47:57,181 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-09-09 03:47:57,198 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-09-09 03:47:57,434 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job2109952258476592765.jar
2021-09-09 03:48:01,172 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job2109952258476592765.jar created
2021-09-09 03:48:01,180 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2021-09-09 03:48:01,180 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2021-09-09 03:48:01,180 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2021-09-09 03:48:01,180 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2021-09-09 03:48:01,191 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-09-09 03:48:01,192 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:48:01,194 [JobControl] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-09-09 03:48:01,201 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:48:01,309 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:48:01,309 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:48:01,312 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
2021-09-09 03:48:01,349 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:1
2021-09-09 03:48:01,382 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1631177379117_0015
2021-09-09 03:48:01,442 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1631177379117_0015
2021-09-09 03:48:01,447 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8088/proxy/application_1631177379117_0015/
cloudera@quickstart:~/Desktop

```

```

((1,Robin),true)
((2,BOB),false)
((3,Maya),false)
((4,Sara),false)
((5,David),false)
((6,maggy),false)
((7,Robert),true)
((8,Syam),false)
((9,Mary),false)
((10,Saran),false)
((11,Stacy),false)
((12,Kelly),false)

```

```
cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"
Output(s):
Successfully stored 13 records (194 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-499131371"
Counters:
Total records written : 13
Total bytes written : 194
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0015

2021-09-09 03:48:16,893 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning UDF_WARNING_2 1 time(s).
2021-09-09 03:48:16,893 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:48:16,894 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:48:16,894 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:48:16,894 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:48:16,902 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:48:16,902 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((1,Robin),true)
((2,BOB),false)
((3,Maya),false)
((4,Sara),false)
((5,David),false)
((6,Maggy),false)
((7,Robert),true)
((8,Syam),false)
((9,Mary),false)
((10,Saran),false)
((11,Stacy),false)
((12,Kelly),false)
((,,),)
grunt> [cloudera@quickstart:~/Desktop]
cloudera@quickstart:~/Desktop
cloudera@quickstart:~/Desktop [cloudera@quickstart:~/Desktop]
```

# SUBSTRING()

This function returns a substring from the given string.

## Syntax

Given below is the syntax of the **SUBSTRING()** function. This function accepts three parameters one is the column name of the string we want. And the other two are the start and stop indexes of the required substring.

```
grunt> SUBSTRING(string, startIndex, stopIndex)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name age and city.

**emp.txt**

```
001,Robin,22,newyork
002,Stacy,25,Bhuwaneshwar
003,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Following is an example of the **SUBSTRING()** function. This example fetches the sub strings that starts with 0<sup>th</sup> letter and ends with 2<sup>nd</sup> letter from the employee names.

```
grunt> substring_data = FOREACH emp_data GENERATE
(id, name), SUBSTRING(name, 0, 2);
```

```
grunt> substring_emp = FOREACH emp_data GENERATE (id, name), SUBSTRING(name,0,2);
2021-09-09 03:51:24,046 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:51:24,046 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:51:24,316 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:51:24,316 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
grunt> dump substring_data;
```

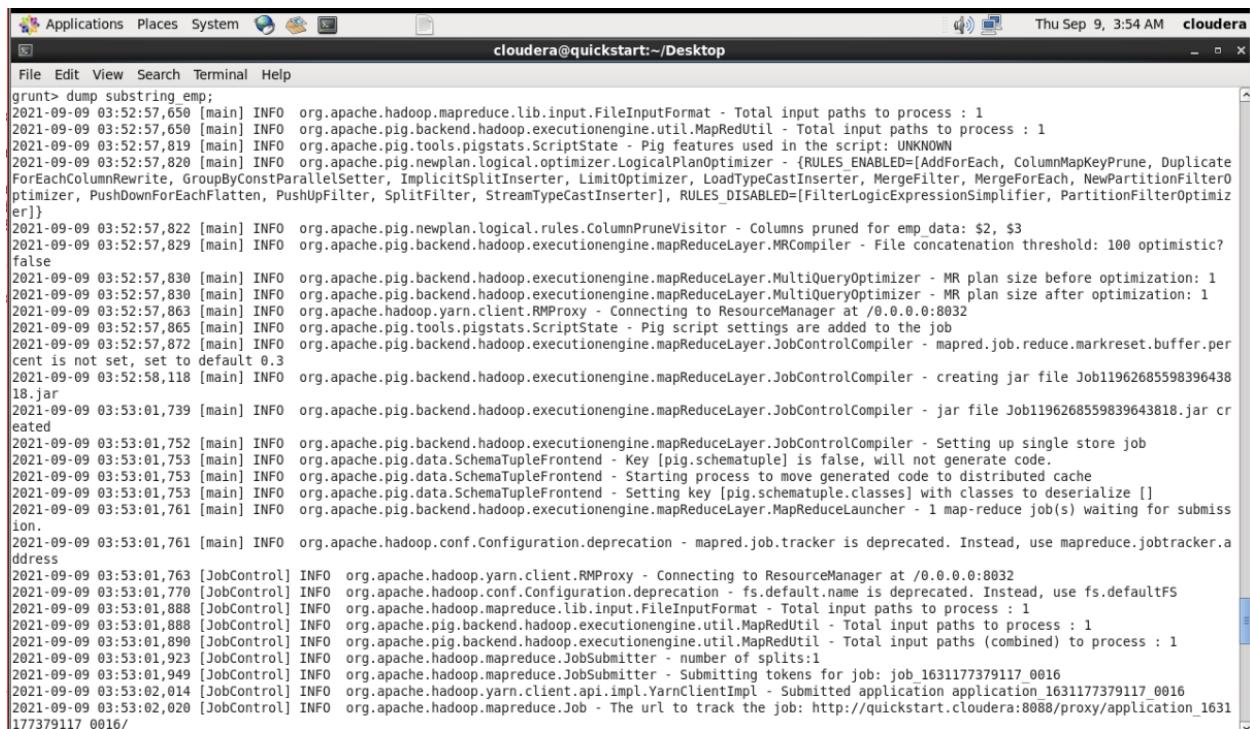


The above statement fetches the required substrings from the names of the employees. The result of the statement will be stored in the relation named **substring\_data**.

Verify the content of the relation **substring\_data**, using the Dump operator as shown below.

```
grunt> Dump substring_data;
```

```
grunt> dump substring_data;
2021-09-09 03:51:58,948 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:51:58,948 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:51:59,108 [main] ERROR org.apache.pig.tools.grunt.Grunt - ERROR 1003: Unable to find an operator for alias substring_data
Details at logfile: /home/cloudera/Desktop/pig_1631179507531.log
```



```
File Edit View Search Terminal Help
cloudera@quickstart:~/Desktop
grunt> dump substring emp;
2021-09-09 03:52:57,650 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:52:57,650 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:52:57,819 [main] INFO  org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN
2021-09-09 03:52:57,820 [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])
2021-09-09 03:52:57,822 [main] INFO  org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for emp_data: $2, $3
2021-09-09 03:52:57,829 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2021-09-09 03:52:57,830 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-09-09 03:52:57,830 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-09-09 03:52:57,863 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-09-09 03:52:57,865 [main] INFO  org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-09-09 03:52:57,872 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-09-09 03:52:58,118 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job1196268559839643818.jar
2021-09-09 03:53:01,739 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job1196268559839643818.jar created
2021-09-09 03:53:01,752 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2021-09-09 03:53:01,753 [main] INFO  org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2021-09-09 03:53:01,753 [main] INFO  org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2021-09-09 03:53:01,753 [main] INFO  org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2021-09-09 03:53:01,761 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-09-09 03:53:01,761 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:53:01,763 [JobControl] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2021-09-09 03:53:01,770 [JobControl] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:53:01,888 [JobControl] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:53:01,888 [JobControl] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-09-09 03:53:01,890 [JobControl] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
2021-09-09 03:53:01,923 [JobControl] INFO  org.apache.hadoop.mapreduce.JobSubmitter - number of splits:1
2021-09-09 03:53:01,949 [JobControl] INFO  org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1631177379117_0016
2021-09-09 03:53:02,014 [JobControl] INFO  org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1631177379117_0016
2021-09-09 03:53:02,020 [JobControl] INFO  org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8088/proxy/application_1631177379117_0016/
```

```
Applications Places System cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"
Output(s):
Successfully stored 13 records (242 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-1204248295"

Counters:
Total records written : 13
Total bytes written : 242
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0016

2021-09-09 03:53:17,561 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning UDF_WARNING_2 1 time(s).
2021-09-09 03:53:17,561 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 03:53:17,562 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 03:53:17,562 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 03:53:17,563 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 03:53:17,571 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 03:53:17,571 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((1,Robin),Ro)
((2,BOB),BO)
((3,Maya),Ma)
((4,Sara),Sa)
((5,David),Da)
((6,Maggy),Ma)
((7,Robert),Ro)
((8,Syam),Sy)
((9,Mary),Ma)
((10,Saran),Sa)
((11,Stacy),St)
((12,Kelly),Ke)
(...)

cloudera@quickstart:~/Desktop
```

# EqualsIgnoreCase()

The **EqualsIgnoreCase()** function is used to compare two strings and verify whether they are equal. If both are equal this function returns the Boolean value **true** else it returns the value **false**.

## Syntax

Given below is the syntax of the function **EqualsIgnoreCase()**

```
grunt> EqualsIgnoreCase(string1, string2)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name age and city.

### emp.txt

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Given below is an example of the **EqualsIgnoreCase()** function. In this example we are comparing the names of every employees with the string value **'Robin'**.

```
grunt> equals_data = FOREACH emp_data GENERATE  
(id, name), EqualsIgnoreCase(name, 'Robin');
```

```
grunt> equals_data = FOREACH emp_data GENERATE (id, name), EqualsIgnoreCase(name, 'Robin');  
2021-09-09 03:56:55,886 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:56:55,886 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:56:56,161 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:56:56,162 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
grunt> dump equals_data;  
2021-09-09 03:57:02,179 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 03:57:02,179 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 03:57:02,362 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN  
2021-09-09 03:57:02,363 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]}  
2021-09-09 03:57:02,364 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for emp_data: $2, $3  
2021-09-09 03:57:02,372 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false  
2021-09-09 03:57:02,374 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1  
2021-09-09 03:57:02,374 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1  
2021-09-09 03:57:02,374 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032  
cloudera@quickstart:~/Desktop$ [n] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job  
[ ] cloudera@quickstart:... [ ] cloudera@quickstart:... [ ]
```

The above statement compares the string “**Robin**” (case insensitive) with the names of the employees, if the value matches it returns **true** else it returns **false**. In short, this statement searches the employee record whose name is ‘**Robin**’

The result of the statement will be stored in the relation named **equals\_data**. Verify the content of the relation **equals\_data**, using the Dump operator as shown below.

```
grunt> Dump equals_data;
```

```
((1,Robin),true)  
((2,BOB),false)  
((3,Maya),false)  
((4,Sara),false)  
((5,David),false)  
((6,Maggy),false)  
((7,Robert),false)  
((8,Syam),false)  
((9,Mary),false)  
((10,Saran),false)  
((11,Stacy),false)  
((12,Kelly),false)
```

```
pRedUtil - Total input paths to process : 1
((1,Robin),true)
((2,BOB),false)
((3,Maya),false)
((4,Sara),false)
((5,David),false)
((6,Maggy),false)
((7,Robert),false)
((8,Syam),false)
((9,Mary),false)
((10,Saran),false)
((11,Stacy),false)
((12,Kelly),false)
grunt> █
```

# INDEXOF()

The **INDEXOF()** function accepts a string value, a character and an index (integer). It returns the first occurrence of the given character in the string, searching forward from the given index.

## Syntax

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

```
grunt> INDEXOF(string, 'character', startIndex)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### emp.txt

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Given below is an example of the **INDEXOF()** function. In this example, we are finding the occurrence of the letter 'r' in the names of every employee using this function.

```
grunt> indexof_data = FOREACH emp_data GENERATE  
(id, name), INDEXOF(name, 'r', 0);
```

```
grunt> indexof data = FOREACH emp data GENERATE (id, name),INDEXOF(name, 'r', 0);  
2021-09-09 04:02:12,711 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS  
2021-09-09 04:02:12,711 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracke  
ddress  
2021-09-09 04:02:12,816 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 04:02:12,816 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-09-09 04:02:13,097 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-09-09 04:02:13,097 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
grunt> dump indexof_data;
```

The above statement parses the name of each employee and returns the index value at which the letter 'r' occurred for the first time. If the name doesn't contain the letter 'r' it returns the value **-1**

The result of the statement will be stored in the relation named **indexof\_data**. Verify the content of the relation **indexof\_data**, using the Dump operator as shown below.

```
grunt> Dump indexof_data;
```

```
((1,Robin),-1)  
((2,BOB),-1)  
((3,Maya),-1)  
((4,Sara),2)  
((5,David),-1)  
((6,Maggy),-1)  
((7,Robert),4)  
((8,Syam),-1)  
((9,Mary),2)  
((10,Saran),2)  
((11,Stacy),-1)  
((12,Kelly),-1)
```

```
Applications Places System cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
cloudera@quickstart:~/Desktop
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"

Output(s):
Successfully stored 13 records (206 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-935082370"

Counters:
Total records written : 13
Total bytes written : 206
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0019

2021-09-09 04:03:02,638 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning UDF_WARNING_1 1 times.
2021-09-09 04:03:02,638 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 04:03:02,639 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 04:03:02,639 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 04:03:02,640 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 04:03:02,647 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 04:03:02,647 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((1,Robin),-1)
((2,BOB),-1)
((3,Maya),-1)
((4,Sara),2)
((5,David),-1)
((6,Maggy),-1)
((7,Robert),4)
((8,Syan),-1)
((9,Mary),2)
((10,Saran),2)
((11,Stacy),-1)
((12,Kelly),-1)
((,),)
grunt> [cloudera@quickstart:~/Desktop]
cloudera@quickstart:~/Desktop [cloudera@quickstart:...]
```

# LAST\_INDEX\_OF()

The **LAST\_INDEX\_OF()** function accepts a string value and a character. It returns the last occurrence of the given character in the string, searching backward from the end of the string.

## Syntax

Given below is the syntax of the **LAST\_INDEX\_OF()** function

```
grunt> LAST_INDEX_OF(string, 'character')
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### emp.txt

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Given below is an example of the **LAST\_INDEX\_OF()** function. In this example, we are going to find the occurrence of the letter 'g' from the end, in the names of every employee.

```
grunt> last_index_data = FOREACH emp_data GENERATE  
(id, name), LAST_INDEX_OF(name, 'g');
```

```
1|grunt> last_index_data = FOREACH emp_data GENERATE (id, name), LAST_INDEX_OF(name, 'g');  
2|2021-09-09 04:05:10,933 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
3|2021-09-09 04:05:10,933 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
4|2021-09-09 04:05:11,207 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
5|2021-09-09 04:05:11,208 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
6|grunt> dump last_index_data;
```

The above statement parses the name of each employee from the end and returns the index value at which the letter 'g' occurred for the first time. If the name doesn't contain the letter 'g' it returns the value -1

The result of the statement will be stored in the relation named **last\_index\_data**. Verify the content of the relation **last\_index\_data** using the Dump operator as shown below.

```
grunt> Dump last_index_data;
```

```
((1,Robin),-1)  
((2,BOB),-1)  
((3,Maya),-1)  
((4,Sara),-1)  
((5,David),-1)  
((6,Maggy),3)  
((7,Robert),-1)  
((8,Syam),-1)  
((9,Mary),-1)  
((10,Saran),-1)  
((11,Stacy),-1)  
((12,Kelly),-1)
```

```
Applications Places System cloudera@quickstart:~/Desktop
File Edit View Search Terminal Help
cloudera@quickstart:~/Desktop
Successfully read 13 records (672 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/smitrpatel/employee_details2.txt"
Output(s):
Successfully stored 13 records (206 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp937174768/tmp-1663720490"

Counters:
Total records written : 13
Total bytes written : 206
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1631177379117_0020

2021-09-09 04:06:03,822 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning UDF_WARNING_1 1 times.
2021-09-09 04:06:03,822 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-09-09 04:06:03,822 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-09-09 04:06:03,823 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-09-09 04:06:03,823 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-09-09 04:06:03,831 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-09-09 04:06:03,831 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((1,Robin),-1)
((2,BOB),-1)
((3,Maya),-1)
((4,Sara),-1)
((5,David),-1)
((6,Maggy),3)
((7,Robert),-1)
((8,Syam),-1)
((9,Mary),-1)
((10,Saran),-1)
((11,Stacy),-1)
((12,Kelly),-1)
((.,),)
grunt> [cloudera@quickstart:...]
[cloudera@quickstart:...]
```

# LCFIRST()

This function is used to convert the first character of the given string into lowercase.

## Syntax

Following is the syntax of the **LCFIRST()** function.

```
grunt> LCFIRST(expression)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### **emp.txt**

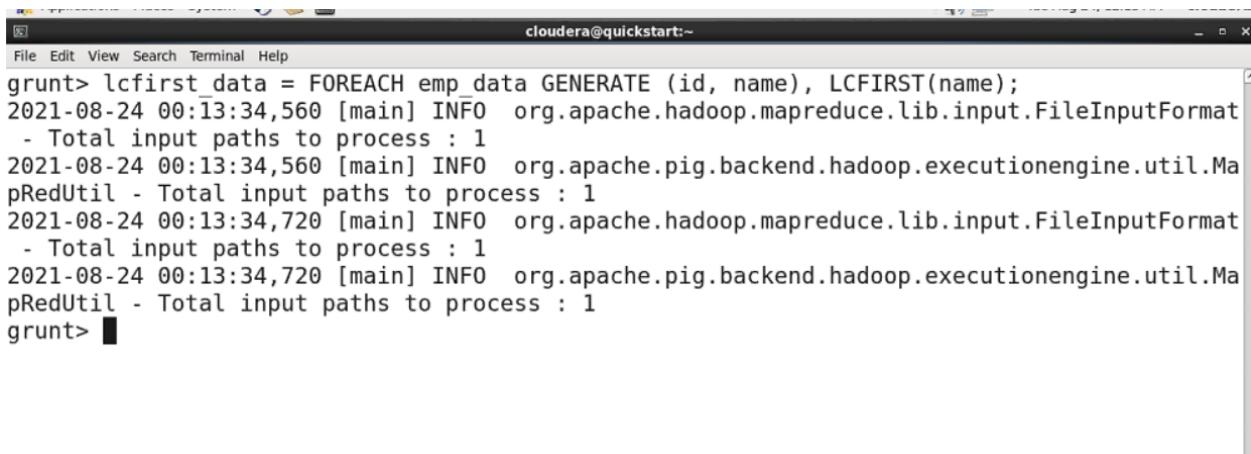
```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Given below is an example of the **LCFIRST()** function. In this example, we have converted all the first letters of the names of the employees to lowercase.

```
grunt> Lcfirst_data = FOREACH emp_data GENERATE  
(id, name), LCFIRST(name);
```



```
grunt> Lcfirst_data = FOREACH emp_data GENERATE (id, name), LCFIRST(name);  
2021-08-24 00:13:34,560 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:13:34,560 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
2021-08-24 00:13:34,720 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:13:34,720 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
grunt> .
```

The result of the statement will be stored in the relation named **Lcfirst\_data**. Verify the content of the relation **Lcfirst\_data**, using the **Dump** operator as shown below.

```
grunt> Dump Lcfirst_data;
```

```
((1,Robin),robin)  
((2,BOB),bob)  
((3,Maya),maya)  
((4,Sara),sara)  
((5,David),david)  
((6,Maggy),maggy)  
((7,Robert),robert)  
((8,Syam),syam)  
((9,Mary),mary)  
((10,Saran),saran)  
((11,Stacy),stacy)  
((12,Kelly),kelly)
```

```
pRedUtil - Total input paths to process : 1
((1,Robin),robin)
((2,BOB),b0B)
((3,Maya),maya)
((4,Sara),sara)
((5,David),david)
((6,Maggy),maggy)
((7,Robert),robert)
((8,Syam),syam)
((9,Mary),mary)
((10,Saran),saran)
((11,Stacy),stacy)
((12,Kelly),kelly)
grunt> █
```

# UCFIRST()

This function accepts a string, converts the first letter of it into uppercase, and returns the result.

## Syntax

Here is the syntax of the function **UCFIRST()** function.

```
grunt> UCFIRST(expression)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### emp.txt

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Following is an example of the **UCFIRST()** function. In this example, we are trying to convert the first letters of the names of the cities, to which the employees belong to, to uppercase.

```
grunt> ucfirst_data = FOREACH emp_data GENERATE  
(id,city), UCFIRST(city);
```

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
grunt> ucfirst_data = FOREACH emp_data GENERATE (id, name), UCFIRST(name);  
2021-08-24 00:14:28,591 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:14:28,591 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.M  
apRedUtil - Total input paths to process : 1  
2021-08-24 00:14:28,741 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:14:28,741 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.M  
apRedUtil - Total input paths to process : 1  
grunt> ■
```

The result of the statement will be stored in the relation named **ucfirst\_data**. Verify the content of the relation **ucfirst\_data**, using the Dump operator as shown below.

In our example, the first letter of the name of the city “**newyork**” is in lowercase. After applying UCFIRST() function, it turns into “**NEWYORK**”

```
grunt>Dump ucfirst_data;
```

```
((1,newyork),Newyork)  
((2,Kolkata),Kolkata)  
((3,Tokyo),Tokyo)  
((4,London),London)  
((5,Bhuwaneshwar),Bhuwaneshwar)  
((6,Chennai),Chennai)  
((7,newyork),Newyork)  
((8,Kolkata),Kolkata)  
((9,Tokyo),Tokyo)  
((10,London),London)  
((11,Bhuwaneshwar),Bhuwaneshwar)  
((12,Chennai),Chennai)
```

```
pRedUtil - Total input paths to process : 1
((1,newyork),Newyork)
((2,Kolkata),Kolkata)
((3,Tokyo),Tokyo)
((4,London),London)
((5,Bhuwaneshwar),Bhuwaneshwar)
((6,Chennai),Chennai)
((7,newyork),Newyork)
((8,Kolkata),Kolkata)
((9,Tokyo),Tokyo)
((10,London),London)
((11,Bhuwaneshwar),Bhuwaneshwar)
((12,Chennai),Chennai)
grunt> █
```

# UPPER()

This function is used to convert all the characters in a string to uppercase.

## Syntax

The syntax of the **UPPER()** function is as follows –

```
grunt> UPPER(expression)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/**. This file contains the employee details such as id, name, age, and city.

### emp.txt

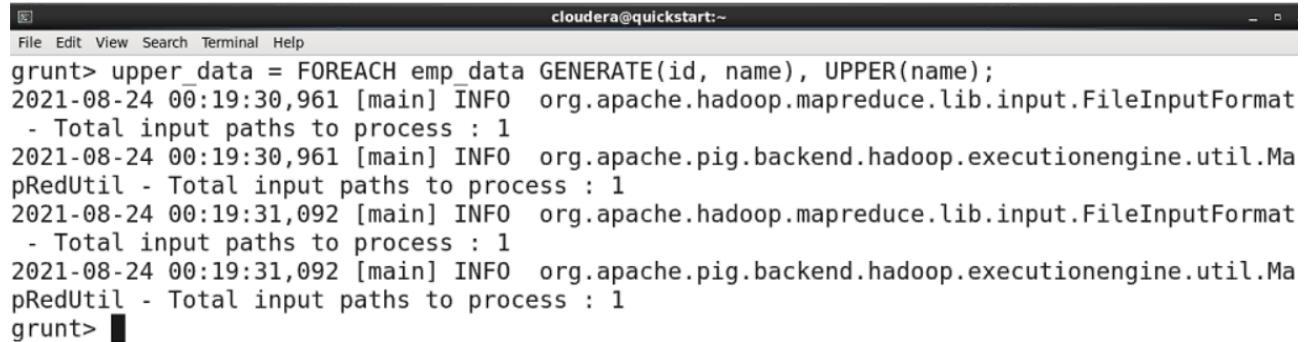
```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Given below is an example of the **UPPER()** function. In this example, we have converted the names of all the employees to upper case.

```
grunt> upper_data = FOREACH emp_data GENERATE  
(id, name), UPPER(name);
```



A terminal window titled "cloudera@quickstart:~" showing the execution of a Pig Latin script. The script defines a relation `upper_data` using the `FOREACH` operator to convert employee names to uppercase. The terminal output shows the command entered and the corresponding log messages from the Hadoop ecosystem, including INFO logs for MapReduce input format and Pig execution engine.

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
grunt> upper_data = FOREACH emp_data GENERATE(id, name), UPPER(name);  
2021-08-24 00:19:30,961 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:19:30,961 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
2021-08-24 00:19:31,092 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:19:31,092 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
grunt> █
```

The above statement converts the names of all the employees to uppercase and returns the result.

The result of the statement will be stored in a relation named **upper\_data**. Verify the content of the relation **upper\_data**, using the **Dump** operator as shown below.

```
grunt> Dump upper_data;
```

```
((1,Robin),ROBIN)  
((2,BOB),BOB)  
((3,Maya),MAYA)  
((4,Sara),SARA)  
((5,David),DAVID)  
((6,Maggy),MAGGY)  
((7,Robert),ROBERT)  
((8,Syam),SYAM)  
((9,Mary),MARY)  
((10,Saran),SARAN)  
((11,Stacy),STACY)  
((12,Kelly),KELLY)
```

```
pRedUtil - Total input paths to process : 1
((1,Robin),ROBIN)
((2,BOB),BOB)
((3,Maya),MAYA)
((4,Sara),SARA)
((5,David),DAVID)
((6,Maggy),MAGGY)
((7,Robert),ROBERT)
((8,Syam),SYAM)
((9,Mary),MARY)
((10,Saran),SARAN)
((11,Stacy),STACY)
((12,Kelly),KELLY)
grunt> █
```

# LOWER()

This function is used to convert all the characters in a string to lowercase.

## Syntax

Following is the syntax of the **LOWER()** function.

```
grunt> LOWER(expression)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### emp.txt

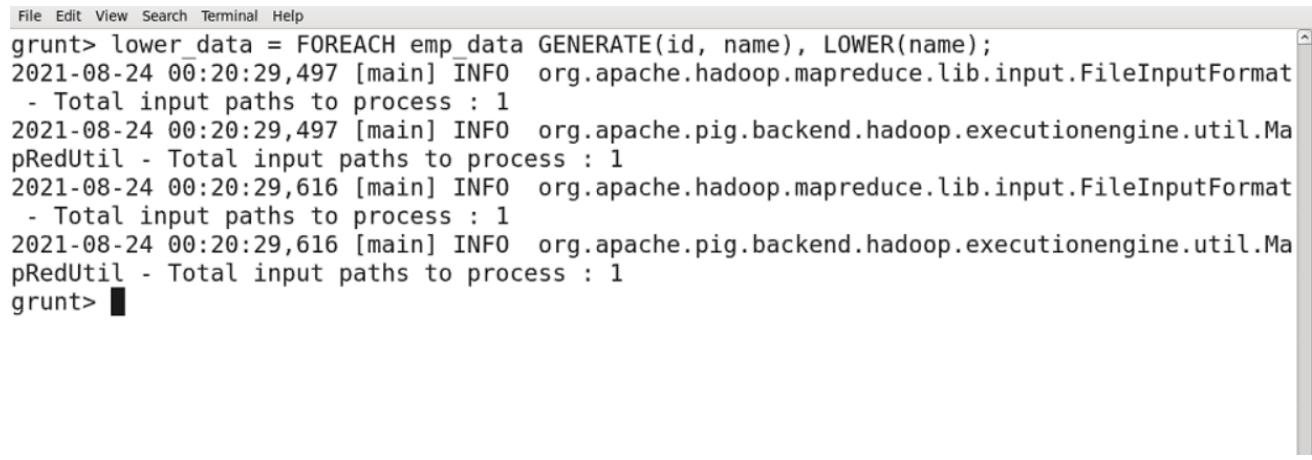
```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Given below is an example of the **LOWER()** function. In this example, we have converted the names of all the employees to lowercase.

```
grunt> lower_data = FOREACH emp_data GENERATE  
(id, name), LOWER(name);
```



```
File Edit View Search Terminal Help  
grunt> lower_data = FOREACH emp_data GENERATE(id, name), LOWER(name);  
2021-08-24 00:20:29,497 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:20:29,497 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
2021-08-24 00:20:29,616 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:20:29,616 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
grunt> █
```

The above statement converts the names of all the employees to lowercase and returns the result.

The result of the statement will be stored in the relation named **lower\_data**. Verify the content of the relation **lower\_data**, using the Dump operator.

```
grunt> Dump lower_data;
```

```
((1,Robin),robin)  
((2,BOB),bob)  
((3,Maya),maya)  
((4,Sara),sara)  
((5,David),david)  
((6,Maggy),maggy)  
((7,Robert),robert)  
((8,Syam),syam)  
((9,Mary),mary)  
((10,Saran),saran)  
((11,Stacy),stacy)  
((12,Kelly),kelly)
```

```
pRedUtil - Total input paths to process : 1
((1,Robin),robin)
((2,BOB),bob)
((3,Maya),maya)
((4,Sara),sara)
((5,David),david)
((6,Maggy),maggy)
((7,Robert),robert)
((8,Syam),syam)
((9,Mary),mary)
((10,Saran),saran)
((11,Stacy),stacy)
((12,Kelly),kelly)
grunt> █
```

# REPLACE()

This function is used to replace all the characters in a given string with the new characters.

## Syntax

Given below is the syntax of the **REPLACE()** function. This function accepts three parameters, namely,

- **string** – The string that is to be replaced. If we want to replace the string within a relation, we have to pass the column name the string belongs to.
- **regEXP** – Here we have to pass the string/regular expression we want to replace.
- **newChar** – Here we have to pass the new value of the string.

```
grunt> REPLACE(string, 'regExp', 'newChar');
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### **emp.txt**

```
001,Robin,22,newyork
002,BOB,23,Kolkata
003,Maya,23,Tokyo
004,Sara,25,London
005,David,23,Bhuwaneshwar
006,Maggy,22,Chennai
007,Robert,22,newyork
008,Syam,23,Kolkata
009,Mary,25,Tokyo
010,Saran,25,London
011,Stacy,25,Bhuwaneshwar
012,Kelly,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

```
grunt> emp_data = LOAD  
'hdfs://localhost:9000/pig_data/emp1.txt' USING  
PigStorage(',')  
as (id:int, name:chararray, age:int, city:chararray);
```

Following is an example of the **REPLACE()** function. In this example, we have replaced the name of the city **Bhuwaneshwar** with a shorter form **Bhuw**.

```
grunt> replace_data = FOREACH emp_data GENERATE  
(id,city), REPLACE(city, 'Bhuwaneshwar', 'Bhuw');
```

```
grunt> replace_data = FOREACH emp_data GENERATE(id, city), REPLACE(city, 'Bhuwaneshwar',  
'Bhuw');  
2021-08-24 00:22:11,300 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:22:11,300 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
2021-08-24 00:22:11,428 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:22:11,428 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma  
pRedUtil - Total input paths to process : 1  
grunt> █
```

The above statement replaces the string '**Bhuwaneshwar**' with '**Bhuw**' in the column named **city** in the **emp\_data** relation and returns the result. This result is stored in the relation named **replace\_data**. Verify the content of the relation **replace\_data** using the Dump operator as shown below.

```
grunt> Dump replace_data;
```

```
((1,newyork),newyork)  
((2,Kolkata),Kolkata)  
((3,Tokyo),Tokyo)  
((4,London),London)  
((5,Bhuwaneshwar),Bhuw)  
((6,Chennai),Chennai)  
((7,newyork),newyork)  
((8,Kolkata),Kolkata)  
((9,Tokyo),Tokyo)  
((10,London),London)  
((11,Bhuwaneshwar),Bhuw)  
((12,Chennai),Chennai)
```

```
pRedUtil - Total input paths to process : 1
((1,newyork),newyork)
((2,Kolkata),Kolkata)
((3,Tokyo),Tokyo)
((4,London),London)
((5,Bhuwaneshwar),Bhuw)
((6,Chennai),Chennai)
((7,newyork),newyork)
((8,Kolkata),Kolkata)
((9,Tokyo),Tokyo)
((10,London),London)
((11,Bhuwaneshwar),Bhuw)
((12,Chennai),Chennai)
grunt> █
```

# STRSPLIT()

This function is used to split a given string by a given delimiter.

## Syntax

The syntax of **STRSPLIT()** is given below. This function accepts a string that is needed to be split, a regular expression, and an integer value specifying the limit (the number of substrings the string should be split). This function parses the string and when it encounters the given regular expression, it splits the string into **n** number of substrings where **n** will be the value passed to **limit**.

```
grunt> STRSPLIT(string, regex, limit)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### **emp.txt**

```
001,Robin_Smith,22,newyork
002,BOB_Wilson,23,Kolkata
003,Maya_Reddy,23,Tokyo
004,Sara_Jain,25,London
005,David_Miller,23,Bhuwaneshwar
006,Maggy_Moore,22,Chennai
007,Robert_Scott,22,newyork
008,Syam_Ketavarapu,23,Kolkata
009,Mary_Carter,25,Tokyo
010,Saran_Naidu,25,London
011,Stacy_Green,25,Bhuwaneshwar
012,Kelly_Moore,22,Chennai
```

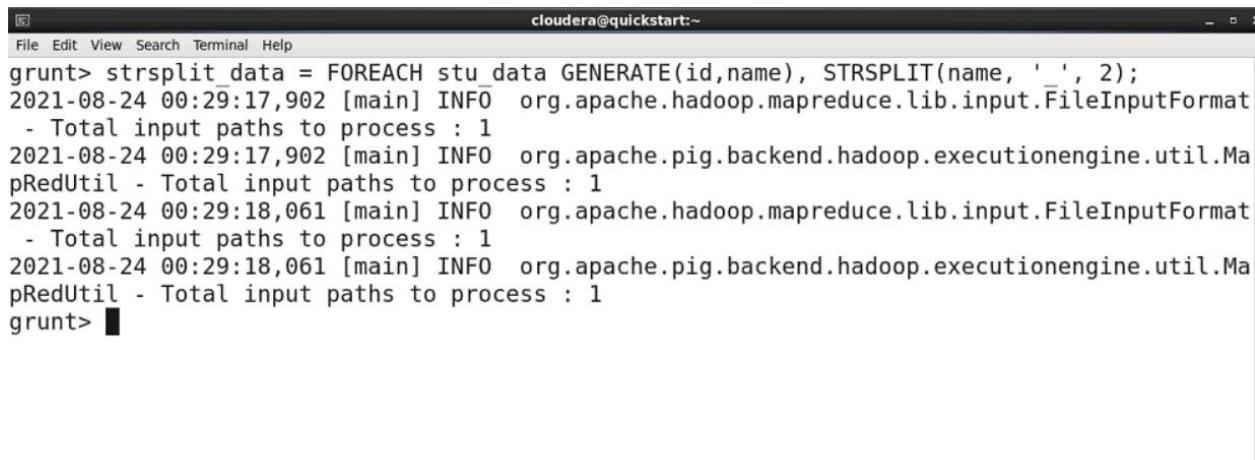
And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Following is an example of the **STRSPLIT()** function. If you observe the emp.txt file, you can find that, in the **name** column, we have the names and surnames of the employees separated by the delimiter '\_'.

In this example, we are trying to split the name and surname of the employees using **STRSPLIT()** function.

```
grunt> strsplit_data = FOREACH emp_data GENERATE  
(id, name), STRSPLIT (name, '_', 2);
```



```
File Edit View Search Terminal Help  
cloudera@quickstart:~  
grunt> strsplit_data = FOREACH stu_data GENERATE(id,name), STRSPLIT(name, ' ', 2);  
2021-08-24 00:29:17,902 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:29:17,902 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil  
- Total input paths to process : 1  
2021-08-24 00:29:18,061 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat  
- Total input paths to process : 1  
2021-08-24 00:29:18,061 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil  
- Total input paths to process : 1  
grunt> █
```

The result of the statement will be stored in the relation named **strsplit\_data**. Verify the content of the relation **strsplit\_data**, using the **Dump** operator as shown below.

```
grunt> Dump strsplit_data;  
  
((1,Robin_Smith),(Robin,Smith))  
((2,BOB_Wilson),(BOB,Wilson))  
((3,Maya_Reddy),(Maya,Reddy))  
((4,Sara_Jain),(Sara,Jain))  
((5,David_Miller),(David,Miller))  
((6,Maggy_Moore),(Maggy,Moore))  
((7,Robert_Scott),(Robert,Scott))  
((8,Syam_Ketavarapu),(Syam,Ketavarapu))  
((9,Mary_Carter),(Mary,Carter))  
((10,Saran_Naidu),(Saran,Naidu))  
((11,Stacy_Green),(Stacy,Green))  
((12,Kelly_Moore),(Kelly,Moore))
```

```
pRedUtil - Total input paths to process : 1
((1,Rajiv_Reddy),(Rajiv,Reddy))
((2,siddarth_Battacharya),(siddarth,Battacharya))
((3,Rajesh_Khanna),(Rajesh,Khanna))
((4,Preethi_Agarwal),(Preethi,Agarwal))
((5,Trupthi_Mohanthy),(Trupthi,Mohanthy))
((6,Archana_Mishra),(Archana,Mishra))
((7,Komal_Nayak),(Komal,Nayak))
((8,Bharathi_Nambiayar),(Bharathi,Nambiayar))
grunt> █
```

# STRSPLITTOBAG()

This function is similar to the **STRSPLIT()** function. It splits the string by a given delimiter and returns the result in a bag.

## Syntax

The syntax of **STRSPLITTOBAG()** is given below. This function accepts a string that is needed to be split, a regular expression, and an integer value specifying the limit (the number of substrings the string should be split). This function parses the string and when it encounters the given regular expression, it splits the string into **n** number of substrings where **n** will be the value passed to **limit**.

```
grunt> STRSPLITTOBAG(string, regex, limit)
```

## Example

Assume that there is a file named **emp.txt** in the **HDFS** directory **/pig\_data/** as shown below. This file contains the employee details such as id, name, age, and city.

### emp.txt

```
001,Robin_Smith,22,newyork
002,BOB_Wilson,23,Kolkata
003,Maya_Reddy,23,Tokyo
004,Sara_Jain,25,London
005,David_Miller,23,Bhuwaneshwar
006,Maggy_Moore,22,Chennai
007,Robert_Scott,22,newyork
008,Syam_Ketavarapu,23,Kolkata
009,Mary_Carter,25,Tokyo
010,Saran_Naidu,25,London
011,Stacy_Green,25,Bhuwaneshwar
012,Kelly_Moore,22,Chennai
```

And, we have loaded this file into Pig with a relation named **emp\_data** as shown below.

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

Following is an example of the **STRSPLITTOBAG()** function. If you observe the emp.txt file, you can find that, in the **name** column, we have name and surname of the employees separated by the delimiter “\_”.

In this example we are trying to split the name and surname of the employee, and get the result in a bag using **STRSPLITTOBAG()** function.

```
grunt> strsplittobag_data = FOREACH emp_data GENERATE  
(id, name), STRSPLITTOBAG (name, '_', 2);
```

The result of the statement will be stored in the relation named **strsplittobag\_data**. Verify the content of the relation **strsplittobag\_data**, using the Dump operator as shown below.

```
grunt> Dump strsplittobag_data;
```

# TRIM()

The **TRIM()** function accepts a string and returns its copy after removing the unwanted spaces before and after it.

## Syntax

Here is the syntax of the **TRIM()** function.

```
grunt> TRIM(expression)
```

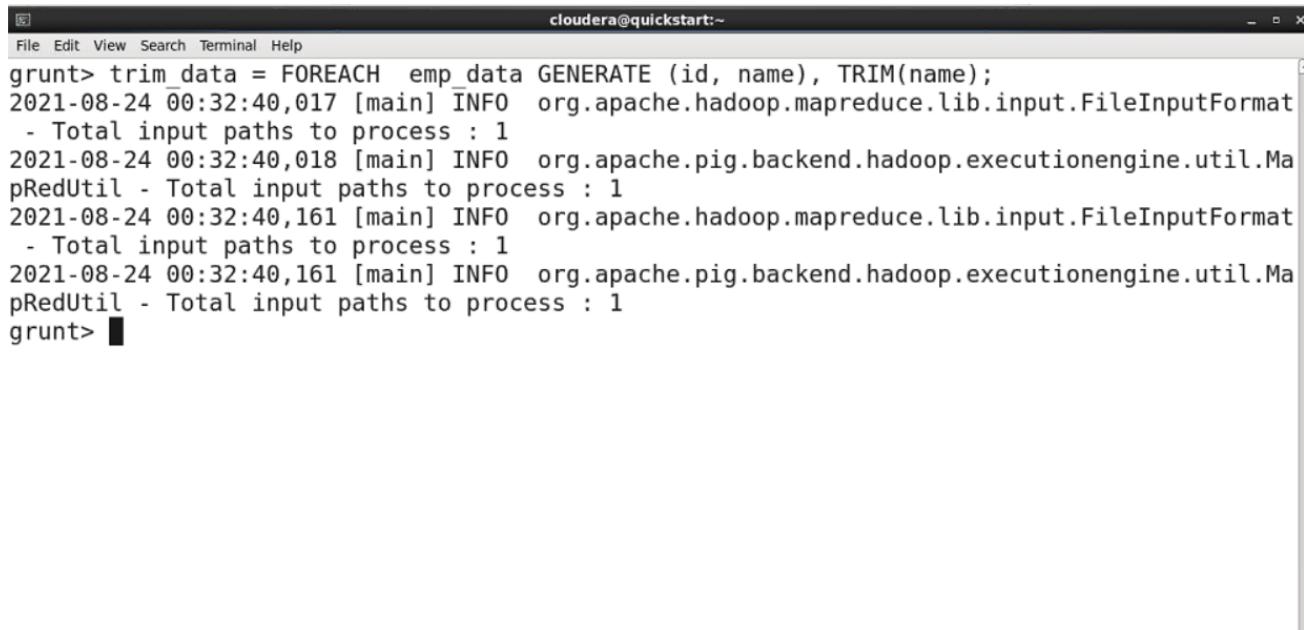
## Example

Assume we have some unwanted spaces before and after the names of the employees in the records of the **emp\_data** relation.

```
grunt> Dump emp_data;  
  
(1, Robin ,22,newyork)  
(2,BOB,23,Kolkata)  
(3, Maya ,23,Tokyo)  
(4,Sara,25,London)  
(5, David ,23,Bhuwaneshwar)  
(6,maggy,22,Chennai)  
(7,Robert,22,newyork)  
(8, Syam ,23,Kolkata)  
(9,Mary,25,Tokyo)  
(10, Saran ,25,London)  
(11, Stacy,25,Bhuwaneshwar)  
(12, Kelly ,22,Chennai)
```

Using the **TRIM()** function, we can remove these heading and tailing spaces from the names, as shown below.

```
grunt> trim_data = FOREACH emp_data GENERATE  
(id,TRIM(name));
```



```
cloudera@quickstart:~
```

```
File Edit View Search Terminal Help
```

```
grunt> trim_data = FOREACH emp_data GENERATE (id, name), TRIM(name);
```

```
2021-08-24 00:32:40,017 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
```

```
- Total input paths to process : 1
```

```
2021-08-24 00:32:40,018 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma
```

```
pRedUtil - Total input paths to process : 1
```

```
2021-08-24 00:32:40,161 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
```

```
- Total input paths to process : 1
```

```
2021-08-24 00:32:40,161 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma
```

```
pRedUtil - Total input paths to process : 1
```

```
grunt> █
```

The above statement returns the copy of the names by removing the heading and tailing spaces from the names of the employees. The result is stored in the relation named **trim\_data**. Verify the result of the relation **trim\_data** using the Dump operator as shown below.

```
grunt> Dump trim_data;
```

```
((1, Robin ),Robin)
```

```
((2,BOB),BOB)
```

```
((3, Maya ),Maya)
```

```
((4,Sara),Sara)
```

```
((5, David ),David)
```

```
((6,maggy),maggy)
```

```
((7,Robert),Robert)
```

```
((8, Syam ),Syam)
```

```
((9,Mary),Mary)
```

```
((10, Saran ),Saran)
```

```
((11, Stacy),Stacy)
```

```
((12, Kelly ),Kelly)
```

```
pRedUtil - Total input paths to process : 1
((1,Robin),Robin)
((2,BOB),BOB)
((3,Maya),Maya)
((4,Sara),Sara)
((5,David),David)
((6,Maggy),Maggy)
((7,Robert),Robert)
((8,Syam),Syam)
((9,Mary),Mary)
((10,Saran),Saran)
((11,Stacy),Stacy)
((12,Kelly),Kelly)
grunts ■
```

# LTRIM()

The function **LTRIM()** is same as the function **TRIM()**. It removes the unwanted spaces from the left side of the given string (heading spaces).

## Syntax

Here is the syntax of the LTRIM() function.

```
grunt> LTRIM(expression)
```

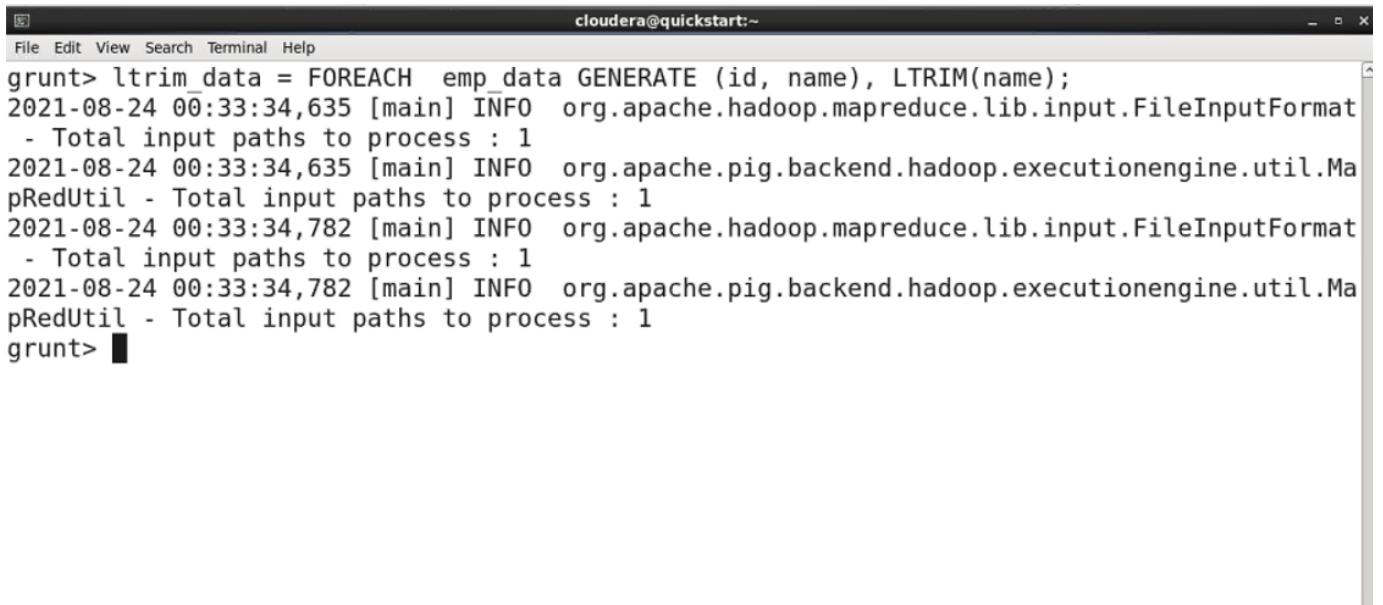
## Example

Assume we have some unwanted spaces before and after the names of the employees in the records of the **emp\_data** relation.

```
grunt> Dump emp_data;  
  
(1, Robin ,22,newyork)  
(2, BOB,23,Kolkata)  
(3, Maya ,23,Tokyo)  
(4, Sara,25,London)  
(5, David ,23,Bhuwaneshwar)  
(6, maggy,22,Chennai)  
(7, Robert,22,newyork)  
(8, Syam ,23,Kolkata)  
(9, Mary,25,Tokyo)  
(10, Saran ,25,London)  
(11, Stacy,25,Bhuwaneshwar)  
(12, Kelly ,22,Chennai)
```

Using the **LTRIM()** function, we can remove the heading spaces from the names as shown below.

```
grunt> ltrim_data = FOREACH emp_data GENERATE  
(id,name), LTRIM(name);
```



```
cloudera@quickstart:~
```

```
File Edit View Search Terminal Help
```

```
grunt> ltrim_data = FOREACH emp_data GENERATE (id, name), LTRIM(name);
```

```
2021-08-24 00:33:34,635 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
```

```
- Total input paths to process : 1
```

```
2021-08-24 00:33:34,635 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil
```

```
- Total input paths to process : 1
```

```
2021-08-24 00:33:34,782 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
```

```
- Total input paths to process : 1
```

```
2021-08-24 00:33:34,782 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil
```

```
- Total input paths to process : 1
```

```
grunt> █
```

The above statement returns the copy of the names by removing the leading spaces from the names of the employees. The result is stored in the relation named **ltrim\_data**. Verify the result of the relation **ltrim\_data** using the Dump operator as shown below.

```
grunt> Dump ltrim_data;
```

```
((1, Robin ),Robin )
```

```
((2,BOB),BOB)
```

```
((3, Maya ),Maya )
```

```
((4,Sara),Sara)
```

```
((5, David ),David )
```

```
((6,maggy),maggy)
```

```
((7,Robert),Robert)
```

```
((8, Syam ),Syam )
```

```
((9,Mary),Mary)
```

```
((10, Saran),Saran)
```

```
((11, Stacy),Stacy)
```

```
((12, Kelly ),Kelly )
```

```
pRedUtil - Total input paths to process : 1
```

```
((1,Robin),Robin)
```

```
((2,BOB),BOB)
```

```
((3,Maya),Maya)
```

```
((4,Sara),Sara)
```

```
((5,David),David)
```

```
((6,Maggy),Maggy)
```

```
((7,Robert),Robert)
```

```
((8,Syam),Syam)
```

```
((9,Mary),Mary)
```

```
((10,Saran),Saran)
```

```
((11,Stacy),Stacy)
```

```
((12,Kelly),Kelly)
```

```
grunt> █
```

# RTRIM()

The function **RTRIM()** is same as the function **TRIM()**. It removes the unwanted spaces from the right side of a given string (tailing spaces).

## Syntax

The syntax of the **RTRIM()** function is as follows –

```
grunt> RTRIM(expression)
```

## Example

Assume we have some unwanted spaces before and after the names of the employees in the records of the **emp\_data** relation as shown below.

```
grunt> Dump emp_data;
```

```
(1, Robin ,22,newyork)
(2, BOB,23,Kolkata)
(3, Maya ,23,Tokyo)
(4, Sara,25,London)
(5, David ,23,Bhuwaneshwar)
(6, maggy,22,Chennai)
(7, Robert,22,newyork)
(8, Syam ,23,Kolkata)
(9, Mary,25,Tokyo)
(10, Saran ,25,London)
(11, Stacy,25,Bhuwaneshwar)
(12, Kelly ,22,Chennai)
```

Using the **RTRIM()** function, we can remove the heading spaces from the names as shown below

```
grunt> rtrim_data = FOREACH emp_data GENERATE
(id, name), RTRIM(name);
```

```

grunt> rtrim_data = FOREACH emp_data GENERATE (id, name), RTRIM(name);
2021-08-24 00:34:26,440 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
- Total input paths to process : 1
2021-08-24 00:34:26,440 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma
pRedUtil - Total input paths to process : 1
2021-08-24 00:34:26,590 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat
- Total input paths to process : 1
2021-08-24 00:34:26,590 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.Ma
pRedUtil - Total input paths to process : 1
grunt> █

```

The above statement returns the copy of the names by removing the **tailing** spaces from the names of the employees. The result is stored in the relation named **rtrim\_data**. Verify the result of the relation **rtrim\_data** using the **Dump** operator as shown below.

```
grunt> Dump rtrim_data;
```

```

((1, Robin ), Robin)
((2,BOB),BOB)
((3, Maya ), Maya)
((4,Sara),Sara)
((5, David ), David)
((6,maggy),maggy)
((7,Robert),Robert)
((8, Syam ), Syam)
((9,Mary),Mary)
((10, Saran ), Saran)
((11, Stacy), Stacy)
((12, Kelly ), Kelly)

```

```

pRedUtil - Total input paths to process : 1
((1,Robin),Robin)
((2,BOB),BOB)
((3,Maya),Maya)
((4,Sara),Sara)
((5,David),David)
((6,Maggy),Maggy)
((7,Robert),Robert)
((8,Syam),Syam)
((9,Mary),Mary)
((10,Saran),Saran)
((11,Stacy),Stacy)
((12,Kelly),Kelly)
grunt> █

```