

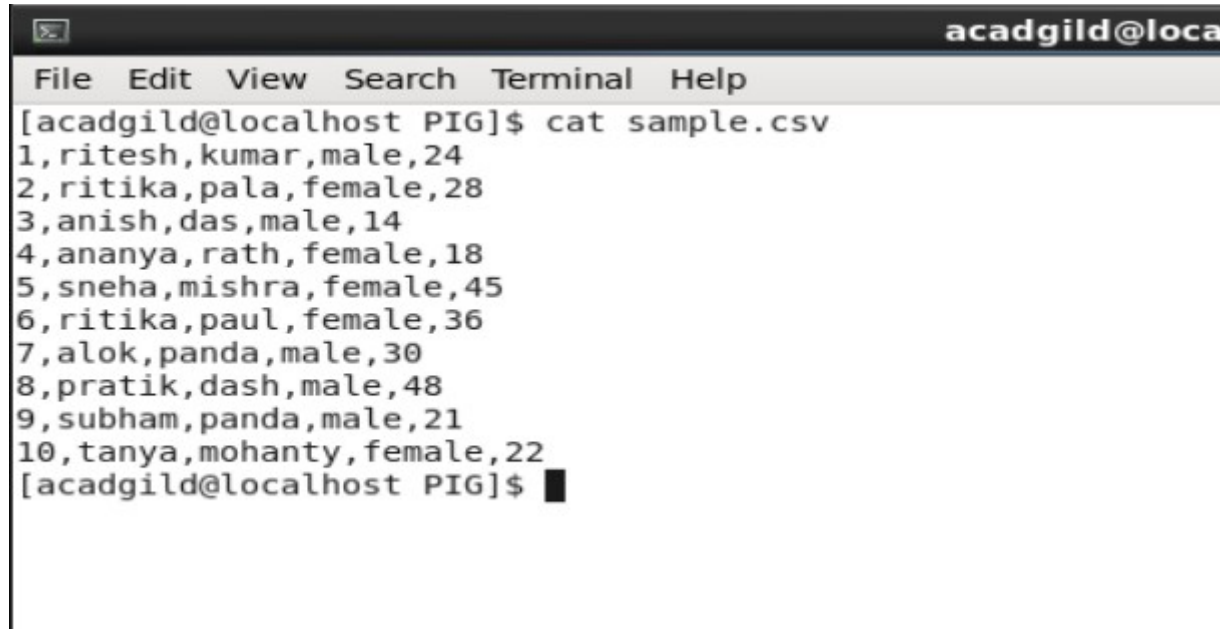
## Bigdata Assignment2.2

Create a sample dataset and implement the below Pig commands on the same dataset.

- 1) Concat
- 2) Tokenize
- 3) Sum
- 4) Min
- 5) Max
- 6) Limit
- 7) Store
- 8) Distinct
- 9) Flatten
- 10) IsEmpty

Created the dataset:-

**cat sample.csv;**



```
acadgild@loca
File Edit View Search Terminal Help
[acadgild@localhost PIG]$ cat sample.csv
1,ritesh,kumar,male,24
2,ritika,pala,female,28
3,anish,das,male,14
4,ananya,rath,female,18
5,sneha,mishra,female,45
6,ritika,paul,female,36
7,alok,panda,male,30
8,pratik,dash,male,48
9,subham,panda,male,21
10,tanya,mohanty,female,22
[acadgild@localhost PIG]$
```

Entered into grunt shell:-

Command - **pig -x local;**

```
acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help
[acadgild@localhost PIG]$ pig -x local
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
18/05/28 20:13:43 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
18/05/28 20:13:43 INFO pig.ExecTypeProvider: Picked LOCAL as the ExecType
2018-05-28 20:13:43,911 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0 (r1746530) compiled Jun 01 2016, 23:10:49
2018-05-28 20:13:43,911 [main] INFO org.apache.pig.Main - Logging error messages to: /home/acadgild/INPUT/PIG/pig_1527518623909.log
2018-05-28 20:13:44,011 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/acadgild/.pigbootstrap not found
2018-05-28 20:13:44,246 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2018-05-28 20:13:44,246 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-28 20:13:44,248 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: file:///
2018-05-28 20:13:44,518 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-28 20:13:44,571 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-default-f74ec77d-bcca-4fa3-ba9a-03eb4a8a6ec7
2018-05-28 20:13:44,571 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
grunt> █
```

Loaded the data into alias:-

Command - **data = LOAD 'sample.csv' using PigStorage(',') as (id:int,first:chararray,last:chararray,gender:chararray,age:int);**

```
acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help
grunt> data = LOAD 'sample.csv' using PigStorage(',') as (id:int,first:chararray,last:chararray,gender:chararray,age:int);
2018-05-28 20:17:30,714 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2018-05-28 20:17:30,714 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> dump data;
```

Displayed the data:-

```
acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help
grunt> dump data;
2018-05-28 20:25:06,908 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN

(1,ritesh,kumar,male,24)
(2,ritika,pala,female,28)
(3,anish,das,male,14)
(4,ananya,rath,female,18)
(5,sneha,mishra,female,45)
(6,ritika,paul,female,36)
(7,alok,panda,male,30)
(8,pratik,dash,male,48)
(9,subham,panda,male,21)
(10,tanya,mohanty,female,22)
grunt> █
```

## 1. Concat -

**first\_last = foreach data generate CONCAT(first,' ',last) as name;  
dump first\_last;**

A screenshot of a terminal window titled 'acadgild@localhost:~/INPUT/PIG'. The terminal shows the execution of Pig commands. The first command is 'grunt> first\_last = foreach data generate CONCAT(first,' ',last) as name;', followed by 'grunt> dump first\_last;'. The output displays ten concatenated names: (ritesh kumar), (ritika pala), (anish das), (ananya rath), (sneha mishra), (ritika paul), (alok panda), (pratik dash), (subham panda), and (tanya mohanty). The prompt 'grunt>' is shown at the bottom with a cursor.

```
acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help
grunt> first_last = foreach data generate CONCAT(first,' ',last) as name;
grunt> dump first_last;
(ritesh kumar)
(ritika pala)
(anish das)
(ananya rath)
(sneha mishra)
(ritika paul)
(alok panda)
(pratik dash)
(subham panda)
(tanya mohanty)
grunt> █
```

O/p – Here using CONCAT we concated the two columns first and last. Hence two columns merged data was the output.

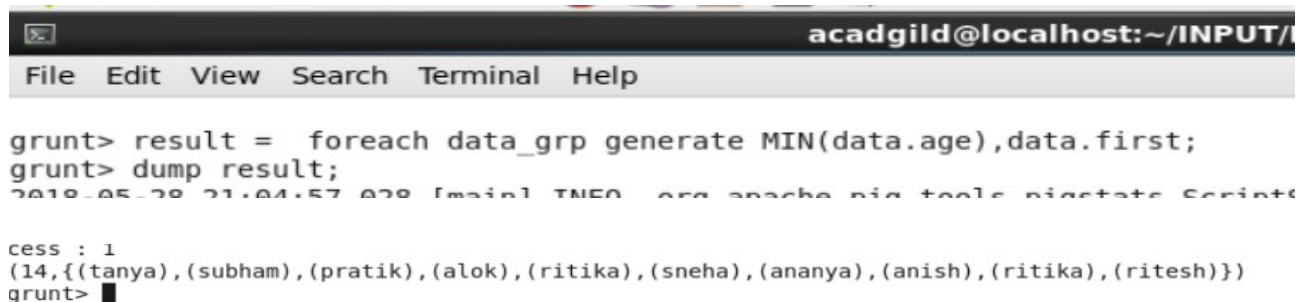
## 2. Tokenize -

**tokenized = foreach first\_last generate TOKENIZE(name);**



#### 4. MIN:-

**result = foreach data\_grp generate MIN(data.age) , data.first;  
dump result;**

A terminal window titled 'acadgild@localhost:~/INPUT/' with a menu bar (File, Edit, View, Search, Terminal, Help). The terminal shows the execution of a Pig script. The first command is 'grunt> result = foreach data\_grp generate MIN(data.age),data.first;'. The second command is 'grunt> dump result;'. The output is '(14,{(tanya),(subham),(pratik),(alok),(ritika),(sneha),(ananya),(anish),(ritika),(ritesh)})'. The prompt 'grunt>' is followed by a cursor.

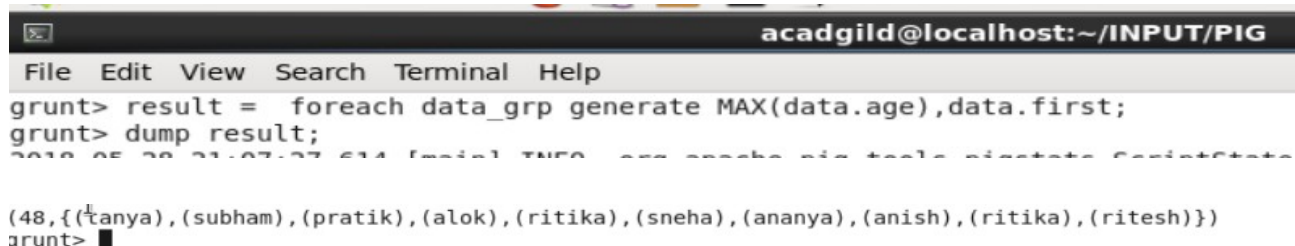
```
acadgild@localhost:~/INPUT/
File Edit View Search Terminal Help

grunt> result =  foreach data_grp generate MIN(data.age),data.first;
grunt> dump result;
2018-05-28 21:04:57 028 [main] INFO org.apache.pig.tools.pigstats.ScriptSt
cess : 1
(14,{(tanya),(subham),(pratik),(alok),(ritika),(sneha),(ananya),(anish),(ritika),(ritesh)})
grunt> █
```

**O/P-** MIN function is used get the minimum value of a certain column in a single-column bag. Here we first made the data into a single column bag then we found the minimum value of the age column.

#### 5. MAX:-

**result = foreach data\_grp generate MAX(data.age) , data.first;  
dump result;**

A terminal window titled 'acadgild@localhost:~/INPUT/PIG' with a menu bar (File, Edit, View, Search, Terminal, Help). The terminal shows the execution of a Pig script. The first command is 'grunt> result = foreach data\_grp generate MAX(data.age),data.first;'. The second command is 'grunt> dump result;'. The output is '(48,{(tanya),(subham),(pratik),(alok),(ritika),(sneha),(ananya),(anish),(ritika),(ritesh)})'. The prompt 'grunt>' is followed by a cursor.

```
acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help

grunt> result =  foreach data_grp generate MAX(data.age),data.first;
grunt> dump result;
2018-05-28 21:07:37 614 [main] INFO org.apache.pig.tools.pigstats.ScriptSt
(48,{(tanya),(subham),(pratik),(alok),(ritika),(sneha),(ananya),(anish),(ritika),(ritesh)})
grunt> █
```

**O/P- MAX** function is used get the minimum value of a certain column in a single-column bag. Here we first made the data into a single column bag then we found the maximum value of the age column.

#### 6. LIMIT:-

**first\_3 = LIMIT data 3;  
dump first\_3;**

```

acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help

grunt> first_3 = LIMIT data 3;
grunt> dump first_3;
2018-05-29 10:07:41,520 [main] WARN org.apache.hadoop.util.NativeCodeLoader - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2018-05-29 10:07:41,559 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: LIMIT
2018-05-29 10:07:41,649 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated.
Instead, use dfs.bytes-per-checksum
2018-05-29 10:07:41,650 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-29 10:07:41,708 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2018-05-29 10:07:41,835 [main] INFO org.apache.pig.impl.util.SpillableMemoryManager - Selected heap (PS Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 489580128, usageThreshold = 489580128
2018-05-29 10:07:42,060 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated.
Instead, use dfs.bytes-per-checksum
2018-05-29 10:07:42,060 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-29 10:07:42,061 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2018-05-29 10:07:42,064 [main] INFO org.apache.pig.builtin.PigStorage - Using PigTextInputFormat
2018-05-29 10:07:42,068 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2018-05-29 10:07:42,073 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
cess : 1
2018-05-29 10:07:42,136 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_0001_m_000001_1' to file:/tmp/temptp1614455752/tmp184462639/temporary/0/task_0001_m_000001
2018-05-29 10:07:42,165 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2018-05-29 10:07:42,189 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2018-05-29 10:07:42,189 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
cess : 1
(1,ritesh,kumar,male,24)
(2,ritika,pala,female,28)
(3,anish,das,male,14)
grunt> █

```

O/p – LIMIT is used to get limited no of tuples from a relation. Here we got 3 no of tuples.

### 7. STORE:-

```
STORE first_3 into '/home/acadgild/INPUT/PIG/first_3';
```

To check the content – **cat part-r-00000;**

```
acadgild@localhost:~/INPUT/PIG
File Edit View Search Terminal Help

grunt> STORE first_3 into '/home/acadgild/INPUT/PIG/first_3' using PigStorage(',');
2016-05-20 10:13:43 750 [main] WARN org.apache.hadoop.conf.Configuration.deprecation

acadgild@localhost:~/INPUT/PIG/first_3
File Edit View Search Terminal Help

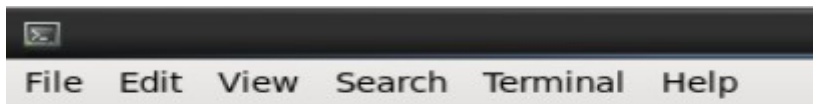
[acadgild@localhost first_3]$ cat part-r-00000
1,ritesh,kumar,male,24
2,ritika,pala,female,28
3,anish,das,male,14
[acadgild@localhost first_3]$
```

O/P – STORE is used to store the loaded data into the file system . Here we , loaded the data first 3 into a folder in the file system and then checked it.

### 8.DISTINCT:-



**distinct\_data = DISTINCT data;**



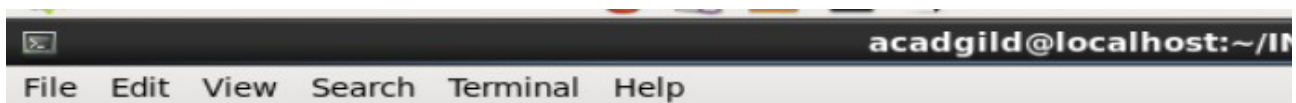
```
grunt> distinct_data = DISTINCT data;
```

```
(1,ritesh,kumar,male,24)
(2,ritika,pala,female,28)
(3,anish,das,male,14)
(4,ananya,rath,female,18)
(5,sneha,mishra,female,45)
(6,ritika,paul,female,36)
(7,alok,panda,male,30)
(8,pratik,dash,male,48)
(9,subham,panda,male,21)
(10,tanya,mohanty,female,22)
grunt> █
```

O/P- DISTINCT is used to remove redundant tuples from a relation. Here , as there were no redundant tuples , the entire data was the output.

## 9. FLATTEN:-

**flatten\_data = foreach data\_grp generate FLATTEN(data);**  
**dump flatten-data;**



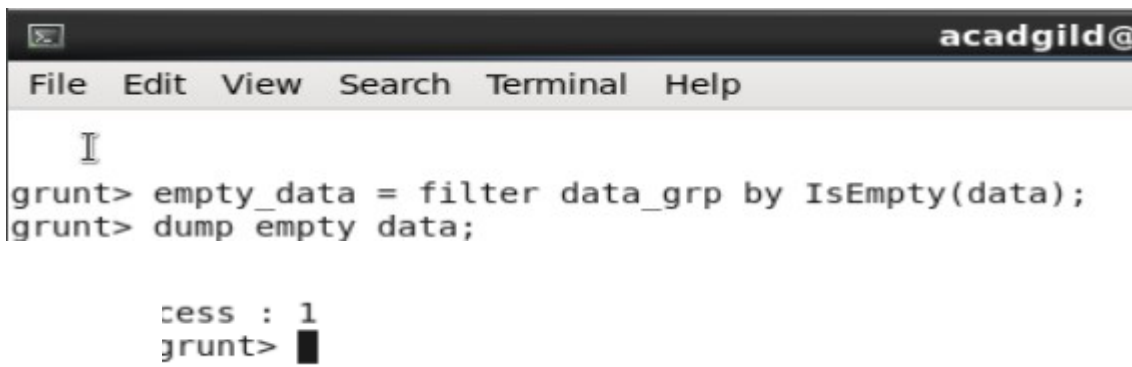
```
grunt> flatten_data = foreach data_grp generate FLATTEN(data);
grunt> dump flatten_data;
```

```
(10,tanya,mohanty,female,22)
(9,subham,panda,male,21)
(8,pratik,dash,male,48)
(7,alok,panda,male,30)
(6,ritika,paul,female,36)
(5,sneha,mishra,female,45)
(4,ananya,rath,female,18)
(3,anish,das,male,14)
(2,ritika,pala,female,28)
(1,ritesh,kumar,male,24)
grunt> █
```

O/P – This flattens the nested schema into unnested schema. Here , we converted the data\_grp bag into flatten-data.

#### 10. IsEmpty:-

**empty\_data = filter data\_grp by IsEmpty(data);**  
**dump empty\_data;**

A screenshot of a terminal window with a dark background. The title bar shows a window icon and the text 'acadgild@'. The menu bar includes 'File', 'Edit', 'View', 'Search', 'Terminal', and 'Help'. The terminal content shows a cursor icon followed by the command 'grunt> empty\_data = filter data\_grp by IsEmpty(data);' on one line, and 'grunt> dump empty\_data;' on the next line. Below these, the output 'cess : 1' is displayed, followed by 'grunt>' and a black cursor block.

```
acadgild@  
File Edit View Search Terminal Help  
I  
grunt> empty_data = filter data_grp by IsEmpty(data);  
grunt> dump empty_data;  
  
cess : 1  
grunt> █
```

O/P – ISEmpty is used to check if empty map or bags are present. Here , the output is nothing since there are no empty bags.