# Bigdata Assignment2.2

Created the dataset:-
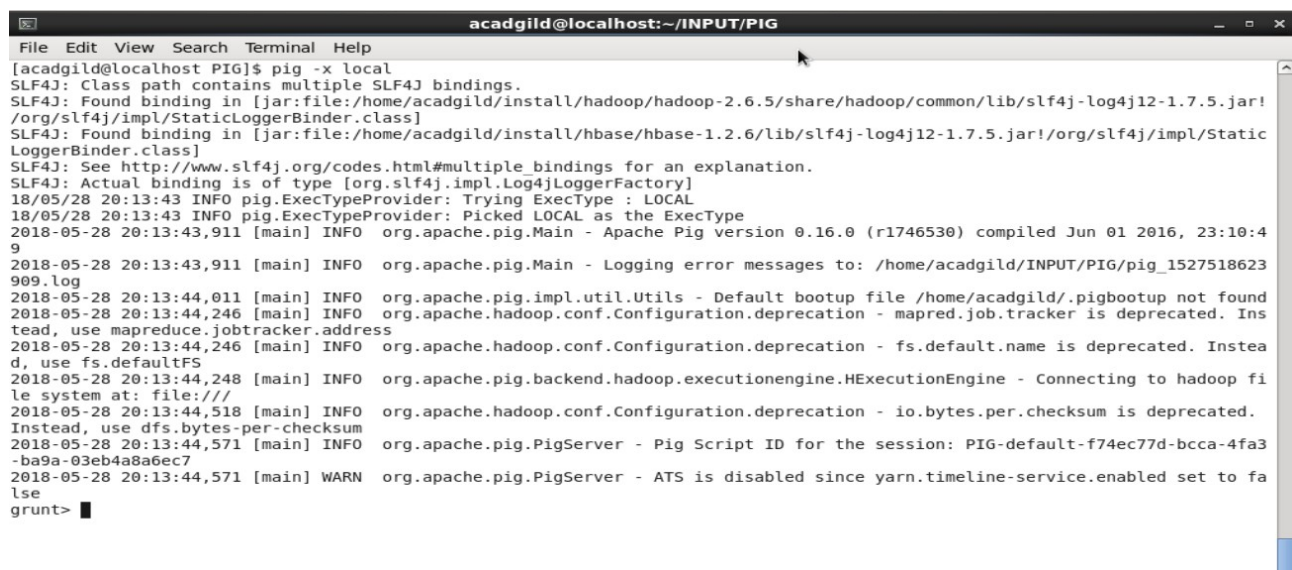
**cat sample.csv;**



Entered into grunt shell:-

Command - **pig -x local;**

## Loaded the data into alias:-

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



## Displayed the 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>
```

## 1. Concat -

**first_last = foreach data generate CONCAT(first,' ',last) as name;
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);**

```
acadgild@localhost:~/INPUT/
File  Edit  View  Search  Terminal  Help
grunt> tokenized = foreach first_last generate TOKENIZE(name);
grunt> dump tokenized;

---- - -
({(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 TOKENIZE when a chararray is given as an argument, this method will split the chararray and return a bag with a tuple for each chararray that results from the split.
We slitted the first_last charaaray.

3.SUM:-

**data_grp = GROUP data all;**
**result = foreach data_grp generate SUM(data.age);**

**dump result;**



**O/P- SUM** function is used get the total of the numeric values of a column in a single-column bag. Here we first made the data into a single column bag then we found the sum of the age column.

4. <u>MIN:-</u>

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



**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.<u>MAX:-</u>

**result = foreach data_grp generate MAX(data.age) , data.first;**
<u>**dump result;**</u>

File   Edit   View   Search   Terminal   Help

```
grunt> result =  foreach data_grp generate MAX(data.age),data.first;
grunt> dump result;
```

```
(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;**

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. Instea
d, use fs.defaultFS
2018-05-29 10:07:41,708 [main] INFO  org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEa
ch, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, Mer
geForEach, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTyp
eCastInserter]}
2018-05-29 10:07:41,835 [main] INFO  org.apache.pig.impl.util.SpillableMemoryManager - Selected heap (PS Old Gen) of size 699
400192 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. Instea
d, 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 pro
cess : 1
2018-05-29 10:07:42,136 [main] INFO  org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attem
pt__0001_m_000001_1' to file:/tmp/temp1614455752/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 pro
cess : 1
(1,ritesh,kumar,male,24)
(2,ritika,pala,female,28)
(3,anish,das,male,14)
grunt> ▌
```
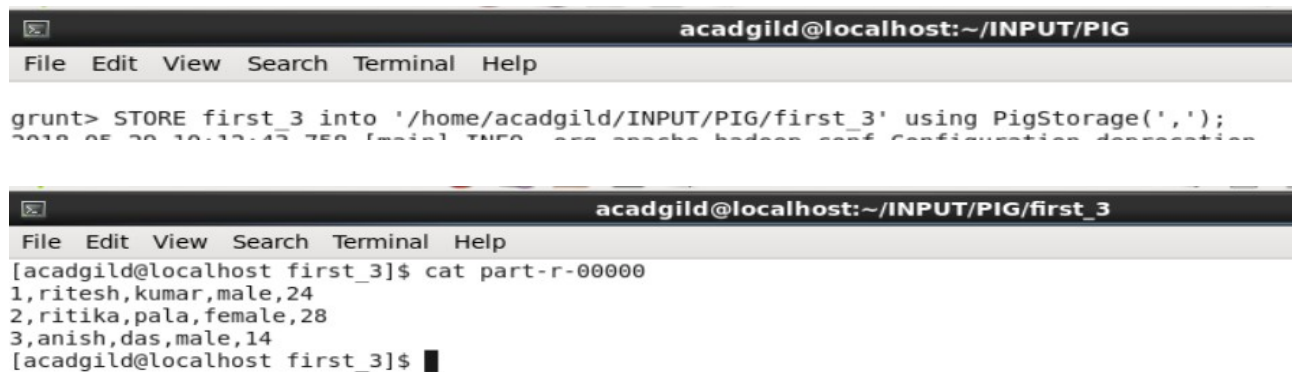
O/p – LIMIT is used to get limited no of tuples from a relation.Here we got 3 no of tuples.

7. STORE:-

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

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



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;**



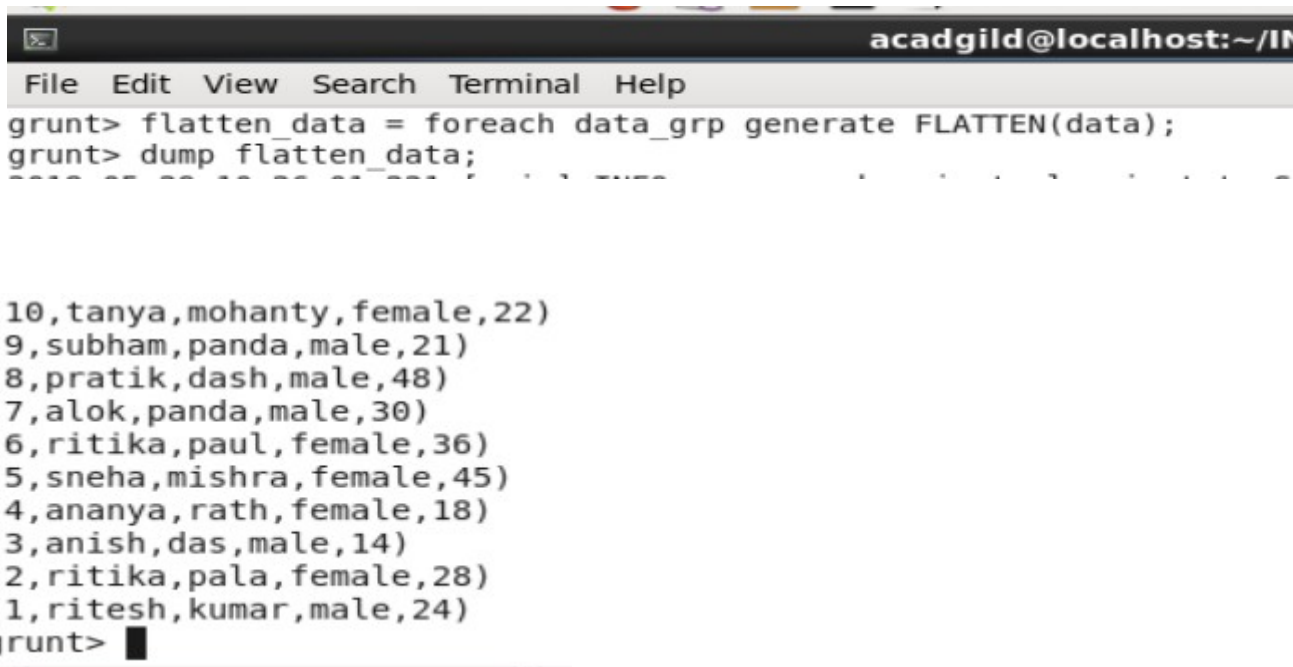O/P- DISTINCT is used to remove redundant tuples from a relation. Here , as there were no redundant tuples , the entire data was the output.

## 9. FLATTEN:-

**flatten_data =  foreach data_grp generate FLATTEN(data);**
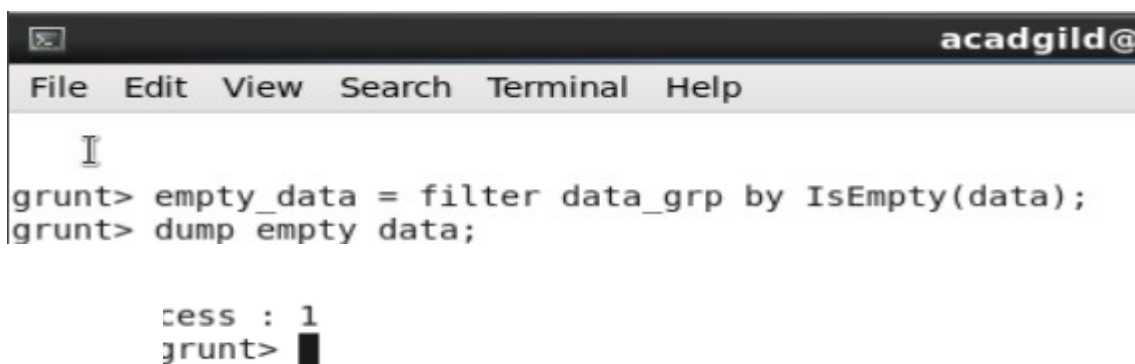**dump flatten-data;**



```
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;**



```
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.