

## Ex No: 5.a

## Pig Latin scripts to sort, group

### AIM:

To write a script for sorting and grouping of data.

Student data:

Assume we have a file **student\_data.txt** in HDFS with the following content.

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

Step 1:

Load and store the student data in HDFS .

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

The **ORDER BY** operator is used to display the contents of a relation in a sorted order based on one or more fields.

```
grunt> Relation_name2 = ORDER Relatin_name1 BY (ASC|DESC);
```

Verify the relation **order\_by\_data** using the **DUMP** operator as shown below.

```
grunt> Dump order_by_data;
```

## Output

It will produce the following output, displaying the contents of the relation **order\_by\_data**.

```
(8,Bharathi,Nambiayar,24,9848022333,Chennai)  
(7,Komal,Nayak,24,9848022334,trivendram)  
(6,Archana,Mishra,23,9848022335,Chennai)  
(5,Trupthi,Mohanthy,23,9848022336,Bhuwaneshwar)
```

```
(3,Rajesh,Khanna,22,9848022339,Delhi)
(2,siddarth,Battacharya,22,9848022338,Kolkata)
(4,Preethi,Agarwal,21,9848022330,Pune)
(1,Rajiv,Reddy,21,9848022337,Hyderabad)
```

The **GROUP** operator is used to group the data in one or more relations. It collects the data having the same key.

Given below is the syntax of the **group** operator.

Now, let us group the records/tuples in the relation by age as shown below.

```
grunt> group_data = GROUP student_details by age;
```

Verify the relation **group\_data** using the **DUMP** operator as shown below.

```
grunt> Dump group_data;
```

**Output:**

```
(21,{(4,Preethi,Agarwal,21,9848022330,Pune),(1,Rajiv,Reddy,21,9848022337,Hyderabad)})
(22,{(3,Rajesh,Khanna,22,9848022339,Delhi),(2,siddarth,Battacharya,22,9848022338,Kolkata)})
(23,{(6,Archana,Mishra,23,9848022335,Chennai),(5,Trupathi,Mohanthy,23,9848022336,Bhuwaneshwar)})
(24,{(8,Bharathi,Nambiayar,24,9848022333,Chennai),(7,Komal,Nayak,24,9848022334,trivendram)})
```

**Ex No: 5.b**

Pig Latin scripts to project, and filter your data.

**AIM:**

To write a script to performing project and filtering.

The **FILTER** operator is used to select the required tuples from a relation based on a condition.

Given below is the syntax of the **FILTER** operator.

```
grunt> Relation2_name = FILTER Relation1_name BY (condition);
```

**student\_details.txt**

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

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

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

Let us now use the Filter operator to get the details of the students who belong to the city Chennai.

```
filter_data = FILTER student_details BY city == 'Chennai';
```

**Verification**

Verify the relation **filter\_data** using the **DUMP** operator as shown below.

```
grunt> Dump filter_data;
```

## Output

It will produce the following output, displaying the contents of the relation **filter\_data** as follows.

```
(6,Archana,Mishra,23,9848022335,Chennai)
(8,Bharathi,Nambiayar,24,9848022333,Chennai)
```