

Experiment No. 07

Aim: Install and run Pig then write Pig Latin scripts to set group, join, project and filter your data.

7(a) : Install & run the pig on word count.

Phase - I Installation of Pig

- 1) Open Cloudera on virtual box
- 2) Open browser
- 3) Click on hive-web application developed by cloudera
- 4) Login with user id & password as cloudera
- 5) Click on Manage HDFS.
- 6) By default in user/Cloudera folder
- 7) Create folder as sec-B/Pig
- 8) Create file word count.txt save it
- 9) open terminal - write pig on command prompt to start pig.

Phase - II

- 1) working with Grunt shell.
- 2) Create word count application.
- 3) Execute word Count application.
- 4) Accessing HDFS from grunt shell.

Step 1: Start Grunt shell
 Open terminal & type pig

Create a file at /user/cloudera/Sec-B/Pig/
 wordcount.txt with following content:

I am learning Pig using Hadoop Exam
 I am learning spark using Hadoop Exam
 I am learning Java using Hadoop Exam.
 I am learning Hadoop using Hadoop Exam.

Step 2: Now load the file stored in HDFS (Space separated file)

input 1 = load / user / cloudera / Sec - B / Pig / wordcount.txt as (F1: chararray);

Dump input 1;
 I am learning pig using Hadoop Exam.
 I am spark using Hadoop Exam.
 I am learning Java using Hadoop Exam.

Step 3: Flatten the words in each line

words in Each line = For EACH input1 GENERATE
 flatten (TOKEN (F1)) as word;
 Dump words in Each line;

Step 4) Group the Name words.

Grouped words = group words /n Each line by word;

dump grouped words;

describe grouped words;

Step 5) Now do the word count

Counted words = For each grouped words generate group COUNT C words /n Each line);

dump counted words.

Pig Latin

Pig scripts can be linear workflow.

Pig scripts can have flattening like multiple data inputs arajoined (normalizing) and data splitting

Grunt: It is shell where we have been writing pig scripts. Generally production code will be written in a separate file. But while writing we want to test our scripts with test data, hence we will be using Grunt shell for prototyping our scripts.

Accessing HDFS You can use HDFS Commands inside grunt shells as > fs - ls

Output:

```
(21, { (4, Preethi Agarwal, 21, 9848022330, Pune), (1, Rajiv, Reddy,  
21, 9840223104, Hyderabad) } )  
(22, { (3, Rajesh Khanna, 22, 9848022339, Delhi), (3, Sidarth  
Bhattacharya, 22, 9848022339, Kolkata) } )  
(22, { 6, Aishwara . Mishra, 23, 984022335, Chennai )
```

Output:

```
(21, Pune) { (4, Preethi Agarwal, 21, 9848022330, Pune) } )  
(21, Hyderabad), { (1, Rajiv Reddy, 21, 984022337, Hyderabad) } )  
(22, Delhi), { (3, Rajesh Khanna, 22, 9840802231, Kolkata) } )
```

Output:

```
(1, Ramesh, 32, Ahmedabad, 2000, 1, Ramesh, 32, Ahmedabad, 2000)  
(2, Khilan, 25, Delhi, 1500, 2, Khilan, 25, Delhi, 1500)  
(3, Kaushik, 23, Kota, 2000, 3, Kaushik, 25, Kota, 2000 )  
(4, Chaitali, 25, Mumbai, 6500, 4, Chaitali, 25, Mumbai, 6500)
```

Experiment No 7 (b)

Objective: Write Pig Latin Scripts to sort, group, join, project and filter your data.

Apache Pig - Group operator

Student - details = LOAD '/user/cloudera/sec-B/Pig/'

Student - details.txt

using Pig Storage(',,') as (id: int, first name: chararray,
last name, chararray, age: int, phone: chararray,
city : chararray); group - data = GROUP Student -
details -> by age; dump group - data;

Grouping by multiple Columns

group - multiple = GROUP Student - details by (age, city),
dump group - multiple;

JOIN:

Customers = LOAD '/user/cloudera/sec-B/Pig/Customers.txt'
using Pig Storage (',,') as (id: int, name: chararray,
age: int, address: chararray, salary: int);

Output:

- (2, Khilan, 25, Delhi, 1500, 101, 2009-11-20, 00:00:00, 2, 1560)
- (3, Kaushik, 23, Kota, 2000, 100, 2009-10-08, 00:00:00, 3, 1500)
- (3, Kaushik, 23, Kota, 2000, 102, 2009-10-08, 00:00:00, 3, 3000)

Output:

- (1, Ram, 32, Ahmedabad, 2000, 1111)
- (2, Krishna, 31, Vundavan, 1500, 101, 2009-11-20, 00:00:00, 2, 1560)
- (3, Kaushik, 23, Kota, 2000, 102, 2009-10-08, 00:00:00, 3, 3000)
- (5, Hitesh, 42, Bhopal, 8500, 111),

Output:

- (2, Khilan, 25, Delhi, 1500, 101, 2009-11-20, 00:00:00, 2, 1560)
- (3, Kaushik, 23, Kota, 2000, 100, 2009-10-08, 00:00:00, 3, 1500)
- (4, Kautilya, 25, Mumbai, 6500, 103, 2008-05-22, 00:00:00, 4, 2060)

Output:

- (1, Ram, 32, Ahmedabad, 2000, 111)
- (2, Khilan, 25, Delhi, 1500, 101, 2009-11-20, 00:00:00, 2, 1560)
- (3, Kaushik, 23, Kota, 2000, 102, 2009-10-08, 00:00:00, 3, 3000)
- (5, Krishna, 24, Vundavan, 5000, 105, 2009-10-07, 00:00:00, 4, 4600)

POORNIMA

Orders = LOAD '/user/cloudera/sec-B/Pig/orders.out' using Pig Storage (',') as (id:int, data:chararray);
Customer - id:int amount:int);

def Train:

Customer1 = LOAD '/user/cloudera/sec-B/Pig/customer.out' using Pig Storage (',') as (id:int, name:chararray, age:int, address:chararray, subarray,salary:int);

Customer2 = LOAD '/user/cloudera/sec-B/Pig/customer.out' using Pig Storage (',') as (id:int, name:chararray, age:int, address:chararray, salary:int);

Customer3 = JOIN Customer1 By id, Customer2 By id;

dump Customer3;

Inner JOIN :

Customer_orders = JOIN Customer By id, Orders By customer-id;

dump Customer_orders;

Output

(8, Bharatti, Namidiyan, 24, 98480, 22333, Chennai)

Outer JOIN :

left - Outer JOIN

Outer - left = JOIN Customer By id LEFT Outer , orders
By customer - id ;
Dump outer - left ;

Right Outer JOIN :

Outer - right = JOIN Customers By id RIGHT , orders
By customer - id ;
Dump outer - right .

Full Outer JOIN

Outer - Full = JOIN Customer By id Full Outer ,
orders By customer - id ;
Dump outer - full ;

Apache - Pig - filter operator

Student - details = LOAD '/user/cloudera/sec-B/pig /
Student - details.txt ' using Pig Storage ('.') as
id : int , first name : chararray , last name : chararray , age :
int , phone : chararray , city : (chararray) ; filter . data
= FILTER Student - details By city = 'chennai' ;
Dump filter - data ;

Experiment No. 08

Objective: Install and run Hive then use hive to create, alter & drop database, tasks, views, function and indexes.

Hive is pre-installed in Cloudera framework

- 1) open terminal - run command `hive --enter`
- 2) In hive mode new
`hive >`
- 3) To create database
`hive > CREATE DATABASE user;`
 OR

Time Taken : 4.478 seconds

- 4) Shows database in hive

`hive > show database;`

OR

default
user

Time Taken = 0.77 seconds , Fetched : 2 rows(s).

- 5) Drop database :-

`hive > drop database if exists user;`

OR

Time Taken : 0.751 seconds.

6) Create Table:-

hive > CREATE TABLE IF NOT EXISTS employee
 (id: int, name: string, salary: string, department: string).

> COMMENT 'EMPLOYEE - details'

> ROW FORMAT DELIMITED

> FIELDS TERMINATED BY '\t'

> LINES TERMINATED BY '\n'

> STORED AS TEXTFILE;

OK

Time Taken : 2.316 seconds

7) Show Tables:-

hive > Show Table;

OK

emp

Time Taken : 0.77 seconds, Fetched : 1 row(s)

8) Create txt file in /home/cloudera/hive / sample.txt
 - write data in file.9) Load data local in PATH '/Home/cloudera/hive /
 sample / test /' OVER WRITE INTO TABLE emp;

Writing data to table default.emp

Table default.emp Stats: [num files = 1, num rows = 0, Total size = 20. Max datysize = 0]

OK

Time Taken = 1.409 seconds

13) Hive - Alter Table

The statement takes any of the following syntax based to what attributes we wish to modify in a txt file:

ALTER TABLE name RENAME To new-name

ALTER TABLE name ADD COLUMNS (col-space [col-space])

ALTER TABLE name DROP [COLUMN] Column_name

ALTER TABLE name CHANGE (Column name new-name New type)

ALTER TABLE name REPLACE COLUMNS (col-space [col-space]);

14) The following query rename the table from employee to emp

Hive > ALTER TABLE employee RENAME to emp;

OR

Time Taken : 0.23 seconds

Hive > show tables;

OR

OK

P

Patients

Prescription

Time Taken : 0.05 seconds , fetched : 4 rows (s)

(12) Change Statement

Change column name :-

Nine > ALTER TABLE emp CHANGE name pname
: string ;
OR

(13) Change type of column :-

Nine > ALTER TABLE emp CHANGE salary salary
Double ;
OR

Time Taken : 0.269 seconds

(14) To add new column in table :

This query add a column named dept to the
emp table

Nine > ALTER TABLE emp ADD COLUMN dept
STRING COMMENT 'Department name' ;
OR

Time Taken : 0.281 seconds

POORNIMA

(15.)

Replace Statement:

This query deletes all the columns from the emp table & replaces it with emp and name columns:

hive > ALTER TABLE emp RERLACE COLUMNS (eid :int , empid : int , ename ; STRING , name: STRING) ;

(16)

Drop Table Statement :-

This query drops a table named "employee".

hive > DROP TABLE IF EXISTS emp ;