

Exercise-3: Perform Sort, Group, Join, Project and Filter Operation in PIG

Note: You need to submit the complete program as part of your assignment.

Report Template

Name: SWETHA K
Roll Number : 235229143
Class : II M.Sc Data Science
Subject : Big Data Analytics Lab
Date : 17-08-2024

Step-1: Download the data files "employee data.txt" and "dept data"

Step-2: Share the data files to the Cloudera platform

hdfs dfs -mkdir /user/cloudera/employee_data

hdfs dfs -mkdir /user/cloudera/dept_data

Step-3: Move the datafiles into theHDFS

hdfs dfs -put employee_data.txt /user/cloudera/employee_data/

hdfs dfs -put dept_data.txt /user/cloudera/dept_data/

Step-4: Enter into the PIGcommand prompt Open Terminal in Cloudera Type "pig" and press Enter

Step- 5: Load Employee Data
employee_data = LOAD
'employee_data.txt' USING
PigStorage(',') AS (emp_id:int,
name:chararray, dept_id:int);

Step-6: Display the Employee Data
DUMP employee_data;

```

(1,Human Resources,)
(2,Finance,)
(3,Marketing,)
(4,IT,)
(5,Sales,)
(6,Customer Service,)
(7,Research & Dev,)
(8,Operations,)
(9,Legal,)
(10,Administration,)
(1,akkim,)
(2,anu,)
(3,ashraf,)

```

Step-7: Load Department Data

**dept data = LOAD 'dept_data.txt' USING PigStorage(',') AS
(dept_id:int, dept_name:chararray);**

step-8: Display the Department data

DUMP dept_data;

```

(1,Human Resources)
(2,Finance)
(3,Marketing)
(4,IT)
(5,Sales)
(6,Customer Service)
(7,Research & Dev)
(8,Operations)
(9,Legal)
(10,Administration)

```

Step-9: **Join employee** data with dept data on dept id

joined_data = JOIN employee_data BY dept_id, dept_data BY

dept_id;

Step-10: Display the **Joined** data

DUMP joined_data;

```
(,Kavitha,4,4,IT)
(,Divya,4,4,IT)
(,Ramesh,5,5,Sales)
(,Mohan,5,5,Sales)
(,Geetha,5,5,Sales)
(,Balaji,5,5,Sales)
(,Karthik,5,5,Sales)
(,Vidhya,6,6,Customer Service)
(,Shanthi,6,6,Customer Service)
(,Santhosh,6,6,Customer Service)
(,Rajini,6,6,Customer Service)
(,Selvi,6,6,Customer Service)
(,Mala,7,7,Research & Dev)
(,Vikram,7,7,Research & Dev)
(,Saravanan,7,7,Research & Dev)
(,Suresh,7,7,Research & Dev)
(,Aravind,7,7,Research & Dev)
(,Ranjani,8,8,Operations)
(,Lakshmi,8,8,Operations)
(,Anitha,8,8,Operations)
(,Ram,8,8,Operations)
(,Priya,8,8,Operations)
(,Krishna,9,9,Legal)
(,Gopal,9,9,Legal)
(,Karthik,9,9,Legal)
(,Bala,9,9,Legal)
(,Nithya,9,9,Legal)
(,Padma,10,10,Administration)
(,Latha,10,10,Administration)
(,Vasanth,10,10,Administration)
(,Revathi,10,10,Administration)
(,Priya,10,10,Administration)
```

Step-11:Project fields: **employee** name and department **name** **projected_data =**
FOREACH **joined_data** **GENERATE** **employee_data::name** **AS** **emp_name**,
dept_data::dept_name **AS** **dept_name**;

Step-12:Display the **projected** data

DUMP **projected_data**;

```
(Kavitha,IT)
(Divya,IT)
(Ramesh,Sales)
(Mohan,Sales)
(Geetha,Sales)
(Balaji,Sales)
(Karthik,Sales)
(Vidhya,Customer Service)
(Shanthi,Customer Service)
(Santhosh,Customer Service)
(Rajini,Customer Service)
(Selvi,Customer Service)
(Mala,Research & Dev)
(Vikram,Research & Dev)
(Saravanan,Research & Dev)
(Suresh,Research & Dev)
(Aravind,Research & Dev)
(Ranjani,Operations)
(Lakshmi,Operations)
(Anitha,Operations)
(Ram,Operations)
(Priya,Operations)
(Krishna,Legal)
(Gopal,Legal)
(Karthik,Legal)
(Bala,Legal)
(Nithya,Legal)
(Padma,Administration)
(Latha,Administration)
(Vasanth,Administration)
(Revathi,Administration)
(Priya,Administration)
```

Step-13:Group by department
grouped_by_dept = GROUP projected_data BY

dept_name;

Step-14: Display the grouped by dept

DUMP grouped_by_dept;

```
(IT, {(Divya, IT), (Kavitha, IT), (Aishwarya, IT), (Sindhu, IT), (Murali, IT)})
(Legal, {(Nithya, Legal), (Bala, Legal), (Karthik, Legal), (Gopal, Legal), (Krishna, Legal)})
(Sales, {(Ramesh, Sales), (Karthik, Sales), (Balaji, Sales), (Geetha, Sales), (Mohan, Sales)})
(Finance, {(Saranya, Finance), (Meena, Finance), (Radha, Finance), (Mani, Finance), (Suganya, Finance)})
(Marketing, {(Kumar, Marketing), (Vinoth, Marketing), (Prakash, Marketing), (Vijay, Marketing), (Sudha, Marketing)})
(Operations, {(Priya, Operations), (Ram, Operations), (Anitha, Operations), (Lakshmi, Operations), (Ranjani, Operations)})
(Administration, {(Priya, Administration), (Revathi, Administration), (Vasanth, Administration), (Latha, Administration), (Padma, Administration)})
(Research & Dev, {(Mala, Research & Dev), (Aravind, Research & Dev), (Suresh, Research & Dev), (Saravanan, Research & Dev), (Vikram, Research & Dev)})
(Human Resources, {(Kannan, Human Resources), (Rajesh, Human Resources), (Srinivasan, Human Resources), (Ganesh, Human Resources), (Arun, Human Resources)})
(Customer Service, {(Vidhya, Customer Service), (Selvi, Customer Service), (Rajini, Customer Service), (Santhosh, Customer Service), (Shanthi, Customer Service)})
```

Step-15: Count the **number of employees** per department

**count_per_dept = FOREACH grouped_by_dept GENERATE group AS dept_name,
COUNT(projected_data) AS employee_count;**

Step-16: Order the results by department name (ascending)

ordered_by_dept_name = ORDER count_per_dept BY dept_name;

Step-17: Display the ordered by dept name

DUMP ordered_by_dept_name;

```
1 - Total input paths to process : 1
(Administration,5)
(Customer Service,5)
(Finance,5)
(Human Resources,5)
(IT,5)
(Legal,5)
(Marketing,5)
(Operations,5)
(Research & Dev,5)
(Sales,5)
```

Step-18: Filter data to get employees in the Sales department

sales_employees = FILTER projected_data BY dept_name == 'Sales';

Step-19: Order Sales employees by name (ascending)

ordered_sales_employees = ORDER sales_employees BY emp_name;

Step-20: Display the Ordered sales employees

DUMP ordered_sales_employees;

```
(Balaji, Sales)
(Geetha, Sales)
(Karthik, Sales)
(Mohan, Sales)
(Ramesh, Sales)
```