| | |
|---|---|
| clear | exec [-param Parameter-name = perm-value] [-param |
| help | or run |
| history | Parameter-name = param_value] script. |
| quit | run or exec / ~ .pig / diff exec & run |
| Execution in batch mode | run-commands in history. |

$ pig -x local simple-script.pig
$ pig -x mapReduce simple-script.pig

A.

**operator** | Case when Then else end
--- | ---
+ |
− |
* |
/ division | Case $2%2$
% modulus | when 0 Then
 | "even"
**Comparsion** | when 1 then
== | "odd"
!= matches: | 
>= Pattern mach. | end.
<= |
> |
< |

**loading (or) reading data.**

relationname → name to store data.

load "filepath" using function as schema
hdfs/local                    self load
based on                      function
pig execution
mode

BinStorage
JsonLoader
schema = (col1, col2)       PigStorage
           (or)             TextLoader
(col1: datatype, col2: dtype)

without schema, colours      $01, $02, $03, $04, ---

Dump A→statement (relation)
└→ execute (run) pig latin statement &
display output

<u>Storing data:</u>                        -- single line comment
                                /* */ → multiline comment
Stor Relation-name into (pathtostore) [using
                                            function];
&
Stor Relatio A into "└─── (⌣) " using
                                Pig Storage.


Dignostic Operator:
└→ load → Just load data to specified relation
 └→ to very load we us Dignostic Operaton.


4 different types of dignostic operators:-
   * Dump Operator
   * Describe operator.
   * Explanation Operator
   * Illustration Operator.
Dump Operator → run labin statement & give.
                    or display output on screen.
                    * used for debugging purpose.
Dump relation-name;


Describe operator → view schema of relation.
Describe relation name:
Explain operator → show(or) display logical,
                    physical, Mapreduce execution
                    plans of a relation.

illustrate operation → give step by step execution
of sequence of statements

## Group Operator

Group: group the data in one (or) more state-
ments.

Grouped-data = GROUP relation-name BY age;
Group multi-col data,
                                                            feild.
                    ↳ = GROUP relation-name BY
                                        (age, city);

to group all columns
group by all columns

group-all = GROUP relation name All;

---

## Cogroup

⌐ Same like group operator
└→ can Group ___work on→ 1 relation
         cogroup ————→ 2 (or) more relations

CoGroup-data = Cogroup A1 by c1, A2 by c2;

( ~1  { — }, { (?)(?) } )

# Join

Self-Join
Inner-Join
Outer-Join — (left Join, right Join, full Join)

2 relations · · · customers = (cid:int, name:chararray,
$\qquad\qquad\qquad$ age:int, address:chararray,
$\qquad\qquad\qquad$ salary:int);

$\qquad$ orders = (oid:int, date:charray,
$\qquad\qquad\qquad$ customer_id:int, amount:int)

used when we want to compare data
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ with itself.

<u>Self Join</u> : Join table itself.

$\quad$ * create 2 relations with same data, schema.

$\qquad\qquad\qquad$ Syntax

$ relation_3 $\quad$ Join Relation-1 by key, Relation2 by key
$\qquad\qquad$ Join customer-1 by id, customer-2 by id.

$\qquad\qquad\qquad\qquad\qquad$ both same bot just renamed

$\qquad\qquad$ equiJoin

Inner join :-

Customer_order = Join Customer by Id, orders by
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ customer_id.

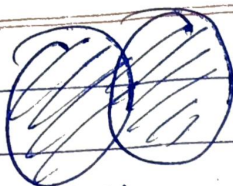# Outer Join:
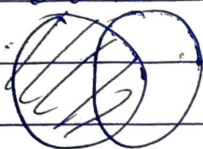


left       Right       Full
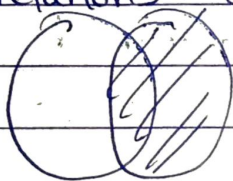
return all rows atleast one relation.

## leftouter

Relation3 = Join customers by ID left outer,
  orders by customer-id;



## Right outer

Relation3 = Join customers by ID Right outer,
  orders by customer FD



Full outer Join → return rows that match atleast one relation.

Relation = Join customers by ID Full outer,
  orders by columnID

## Join operation based on multiple keys:

syntax
       employee        Key1, Key2

Relation3 = Join customer by (ID, JobID),
  employee-contact by (ID, JobID);

Join
or
cross operator :- ___ cross product of two or more relation

↳ combine each row of one table &with another row of another table.

ex

Car model
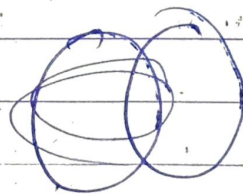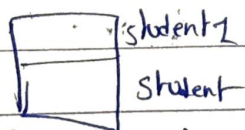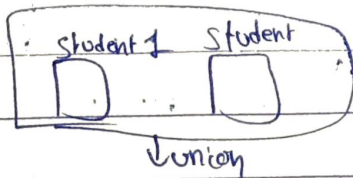
1    camry
2    Corolla
3    Prius

|  | Cormodel | Colorname |
|---|---|---|
| 1 | Camry | Black |
| 2 | corolla | Red |
| 3 | prius | sliver |
| 1 |  |  |
| 1 |  |  |
|  |  |  |

crossJoin →

Color_name

1    Black
2    Red
3    Sliver

c

Relation3 = Cross    Relation1 , Relation2

_combining

union operator
merge two or more relations



Student 1    Student

∪union

student1
Student

Student 1 = (id, firstname, lastname, phone, city) - 50 reco
Student 2 = (id, firstname, lastname, phone, city) - 50
s = union Student 1, Student 2; → 100 records

$\square \rightarrow \overset{D}{\underset{D}{}}$

split → divide the relation into two or max.

splitting into 1.

split renation into new relation-name if $age71$

var value

split inpterelation into output-relation1 if condition1,
output-relation2 if condition2;

filter → to select required tuples ~~based on~~ from a
relation based on condition.

Relation2 = Filter Relation_1-name by condition;

Distinct → remove the duplicate tuples from ~~tuples~~ relation

relation2 = Distinct Relation1;

for each → generate specific data transformations
based on column data,

Relation-2 = Foreach Relation1 Generate columnname;

to select id, age, city - foreachstudent
Relation3 = Foreach student details Generate id, age, city;

order by (→ sort by ~~columns~~ based on. one or more fields)

order relation-name by column (ASC|DESC);

age DESC;

limit → select Top n num of Tuples

limit_data = LIMIT student_details 4;
└ Top 4 tuples

Converting datatype
date2 - int
(char Array) date1
✓

EVAL (evaluate) functions:-

Avg
ForEach relation Generate
AvgC ) j

AVG ( ~~~ ) → preciding use group all or group by

CONCAT ( expres firstnam, lastname).
↳ 2 or more expressions

Count (Age) — not include nullvalues.

Count_STAR ( ) - Include nullvalug

student_details
· Student_group_all = Group student_details All;

studentcount = FOREACH student_group_all Generate
COUNT_STAR(student_details.
gpa}

DIFF ( expl, exp2 ) - compare bags
expl tuple1 = exp2 tuple2 means {} - empty bag
≠ return both

AVG
CONCAT.
COUNT
IN
MAX
MIN
SIZE
SUM
TOKENIZE