

5th March 2018

PIG (Cont...)

1. JSON Data Processing through PIG

Data.json <Key: Value>

{ "Website": "Amazon", "location": "California", "rating": 4,
"amount": 500000 }

-- --

JSONscript.pig

(same path pigscript)

initialData = LOAD 'home/gopal.krishna/Batch89/PIGPAC/
Data.json' USING

Jsonloader('website: chararray, location:
chararray, rating: int, amount: int');

filterData = FOREACH initialData GENERATE Website, location,
amount;

selectedData = FILTER filterData BY location;

sortedData = ORDER filterData INTO 'JSONFIL-DATA89';

STORE filterData INTO 'JSONSEL-DATA89';

STORE selectedData INTO 'JSONSORT-DATA89';

NOTE: - Not to refer any jar file support.
keys must be same as defined in JSON
keys, no change in case.

\$ ____/PIGPAC \$ pig -x local JSONscript.pig

2. SPECIAL JOINS IN PIG

inorder to improve the performance of JOIN
operations, below are the different joins available

- a) replicated join
- b) merge join
- c) skewed join

Note: - Merge and Replicated joins are included
in the Mapper phase.

skewed joins are included in Reducer phase.

a. Replicated join

Let's say:

example: part-m-00000
5 records

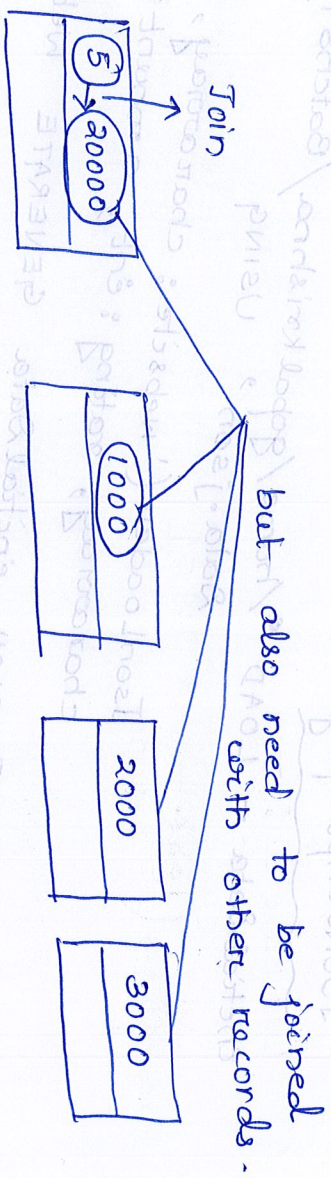
(Id: int)

85000 records

(Id: int)

We can perform the join, but the join operation suffers from performance problem.

example:-



For such case we use

cityAddress	Country
5 records	

cityAddress	Country	Capital	Website	Code
63456	records			

3. Embedded Mode of Pig

Steps for using Pig UDFs:-

a) Develop a class by extending from the base class of EvalFunc <string>

The package need to be imported as below:

org.apache.pig.EvalFunc

b) To write our own customizable logic we must override below methods:-

public String exec (Tuple input) throws IOException

c) We need to add required Pig jar file to the build path and create our own logic with business logic, which needs to be exported to the Hadoop environment where Pig is running.

4. Refer the same jar file in the main pig file using 'register' keyword.

→ register statement should be in first line of pig script.

** END **

6th March 2018 PIG UDFs (cont..)

A. Use case 1 :- Cisco Data log file ciscodata.log

Example :- (tab separated) 10101 | Network Related Issue | Sanjose |
2011-11-10 10:10:10
10201 | Router Related | Rose City | 2012-12-09 09:10:27

ciscoudfscript.pig

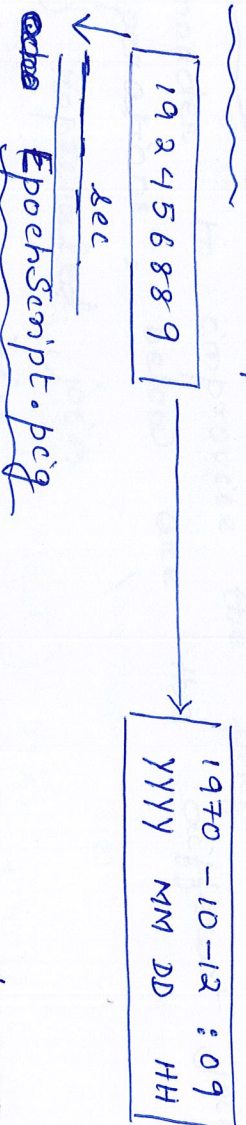
REGISTER /home/gopalkrishna/Batch88/PIGPRAC/PIG-CONCAT.jar;

A = LOAD 'ciscodata.log' USING PigStorage ('|') AS
(logid: int, logerron: chararray, logdate: chararray);
B = FOREACH A GENERATE logid, com.pig.concatenate.
CONCAT (logerron), logloc ;

STORE B INTO 'ciscoudfOUT89' ;

→ \$ pig -x local ciscoudfscript.pig <1

B. Use case 2 :- EpochTime



epoch EpochScript.pig
REGISTER /home/gopalkrishna/Batch88/PIGPRAC/PIG-EPOCH.jar ;
eneData = LOAD 'epochTimeData.log' AS
(timeData: chararray) ;

convertData = FOREACH eneData GENERATE timeData,
com.pig.epochusecase.EPOCH (timeData) ;

