

⑤ Give a list of Built-in operators in Pig

Relational operators:

↳ allows you to transform data by sorting, grouping, joining, projecting and filtering

(i) LOAD

↳ load data from $\leftarrow \begin{matrix} \text{LFS} \\ \text{HDFS} \end{matrix} \rightarrow$ into pig relation

Ex: load1 = load '/' using PigStorage(',') as (——)

(ii) FOREACH

↳ generates data transformations based on columns of data
It is used to add or remove fields from a relation

Ex: for-each = foreach load1 generate url, id;

(iii) FILTER

↳ select tuples based on a condition

Ex: filter_group = filter load1 by id > 8;

(iv) JOIN

↳ perform a inner join of two or more relations based on common field values

Ex: join1 = join load1 by url, load2 by url;

(v) ORDERBY

order used to sort a relation based on one or more fields

Ex: load1 = ORDER load1 by url;

(vi) DISTINCT:

↳ removes duplicate tuples in a relation

Ex: load 3 = distinct load 1;

(vii) STORE:

↳ used to save results to the file system

Ex: Store load 3 into 'string';

(viii) GROUP:

↳ groups tuples together based on same key values

↳ Key field will be a tuple if the group key has more than one field

Ex: gc = group load 1 by url;

(ix) COGROUP:

↳ same as GROUP, used when multiple relations are involved

Ex: cgc = cgroup loading 1 by url, loading 2 by url;

(x) CROSS:

↳ used to compute cross product of two or more relations

Ex: ce = cross load 1, load 2

(xi) LIMIT:

↳ limits the number of output tuples

Ex: le = limit load 1 3;

(xii) SPLIT:

partition a relation into two or more based on some

expression

Ex: Split load 1 into x if ratio > 1, 1 if ratio <= 1;

Scanned by CamScanner

3) What are UDFs in Pig

Pig was designed to be extensible & customisable: all parts of Pig: loading, storing, filtering, grouping, joining can be altered by user defined functions. So If a function is not available you can write your own user defined function. UDFs can be written in Python, Java, JavaScript, Ruby etc

FilterUDFs are subclasses of EvalFunc which have abstract method `exec()` where we implement our code. To use the new function we first compile it and package it in a JAR file. Then we tell pig using REGISTER operator then we can invoke the function

An EvalUDF is a small step up from writing a filter function. To write a general eval udf you need to consider what the output's schema looks like. For scalars pig takes care but for complex types such as bags, tuples/maps. Pig needs more help you should implement `OutputSchema()`

— A LoadUDF

② Compare Apache Pig, MapReduce & MySQL

Pig

- ↳ Pig Latin is a procedural language, high level language
- ↳ platform for analyzing large data sets
- ↳ Easy and convenient to program uses less lines of code
- ↳ High level of Abstraction
- ↳ development effort is less, Code efficiency is relatively less

Hadoop MapReduce

- ↳ Compiled language
- ↳ Lower level of Abstraction
- ↳ More lines of code
- ↳ More development effort is needed
- ↳ Code efficiency is high when compared to Pig

SQL

- ↳ declarative language
- ↳ Used for transactional & analytical queries
- ↳ Schema is mandatory
- ↳ Relatively less lines of code

1) Diagnostic Operators

(i) DUMP:

↳ runs pig commands & displays the results on screen

Ex: dump load1

(ii) DESCRIBE:

↳ to review the schema of a particular relation

Ex: describe load1

(iii) ILLUSTRATE:

↳ to review how data is transformed through a sequence of Pig Latin statements

↳ but used to debug

Ex: illustrate load1;

(iv) EXPLAIN:

↳ prints the logical & physical plans.

Ex: explain load1;

② Describe CAP THEOREM

1. States that it is impossible for a distributed data store to simultaneously provide more than two out of the following

- 1) CONSISTENCY: Every read receives the most recent write or an error
- 2) AVAILABILITY: Every request a non-error response without a guarantee that it contains the most recent write
- 3) PARTITION TOLERANCE: The system continues to operate despite an arbitrary number of messages being dropped (delayed) by the n/w between nodes

ie) In the presence of a n/w partition, one has to choose b/w consistency & availability

④ What is the difference between NoSQL and SQL

SQL databases are primarily categorised as Relational databases. They have fixed (or) static or predefined schema. SQL databases display data in form of tables and are vertically Scalable. SQL databases are a powerful Structured Query language to define, manipulate data. Best suited for complex queries and not suitable for hierarchical data storage. Ex: MySQL, Oracle, MS-SQL etc

NoSQL

↳ categorised as Non-relational or distributed database

Systems.

↳ have dynamic schema

↳ Display data as a collection of key-pair values, documents, graph database

↳ horizontally scalable

↳ collection of documents are used to query the data. It is also called unstructured query language

↳ not so good for complex queries but best suited for hierarchical data storage

Ex: MongoDB, BigTable, Redis, Raven DB, Cassandra etc.