

WEEK 04: APACHE HIVE BASICS-

- A. INTRODUCTION
- B. HIVE BASICS
- C. HIVE PRACTICAL
- D.

HIVE:

data warehouse tool to process structured, histor. data. (on top of HADOOP)

Tutorial
01
41 min.

A. INTRODUCTION

- Transactional & Analytical sys. ✓
- HIVE + HQL ✓
- Hive table, metadata ✓
- why hive over traditional database
- Transactional & Analytical Processing.
- Data warehousing
- Hive Architecture
- Hive on top of Hadoop
- working of Hive.

Transactional System

- day to day data / current data
- daily updates, inserts, delete
- ATM
- RDBMS (best suited)
- mysql, oracle
- monolytic system

Analytical sys.

- historical data / used for analysis
- no updates, only analysis / reads
-
- Dataware house (best)
- teradata
- distributed system

Hive is open source
data warehouse

(meant to solve analytical problems)

HIVE converts HQL to mapreduce job and will submit on the hadoop cluster.

utility : mapreduce works on java, insted of writing long java codes we can write hive queries in single lines and hive will do all further processing.

"Abstracting complexities from user"

* Hive table

2 parts of table.

1. Actual data (stored in HDFS)
(raw format).

2. Meta data (schema / format)
(stored in database).
(mysql).

* Why meta data is stored in database?

- 1). In HDFS updates & edits are difficult.
- 2). fast retrival of data in database (low latency)

* both are kept separately.

* metadata

emp-name (string)

city (string)

salary (int)

data

shiram banglore 100000

rahul delhi 100000

nitish delhi 20000

* hive metadata is stored in mysql inside database named "metastore".

Tutorial 02
Hive theory 1

- Transactional & Analytical processing
- data warehouse concepts.
- HQL, metastore
- Hive vs RDBMS
- HQL vs SQL.

Transactional Processing

- analyze individual entry
- recent and changing data
- updation of data
- real time access.
- single data source

example: delivery of any product by amazon.

(RDBMS)
mysql

Analytical Processing

- analyze large batch
- historical data.
- only reads data
- long jobs.
- multiple data source

revenue generated,
extraction of reports/
results.

(warehouse)
teradata

Data Warehouse

- Long running batch jobs
- optimized for read operations
- holds data from multiple sources
- holds data for long period.
- historical data.

Vertica, Teradata, Oracle, IBM.

Hive metastore

- stores metadata of all tables.
- maps files to table.
- holds schema for each table.

HIVE

v/s.

RDBMS

- large datasets
- parallel computation
- high latency
- read operation
- not ACID compliant
- HiveQL

- small datasets
- serial com.
- low
- write/update/read
- ACID compliant
- SQL

✱ cluster of machines

✱ low latency means, takes less time.

HQL

v/s

SQL

minimal index support
many builtin functions
only equi join
restricted subqueries

indexes allowed.
less b.i.f.
all joins. allowed.
many sub.

Tutorial 03

35 min
hive theory 2.

- ✱ data is by default stored in warehouse directory
/user/hive/warehouse.
- ✱ if creating a database
/user/hive/warehouse/trendytech.db
- ✱ files.
/user/hive/warehouse/trendytech.db/files.

Institute
Date
Page

* metadata is stored in metastore (rdbms (mysql))

Complex Data Types.

(ARRAY, MAP, STRUCT.)

- 1) ARRAY • homogenous collections.
- 2) MAP • data in form of key-values.
• unordered
- 3) STRUCT • class
• diff data types allowed.

Built In functions.

1. UDF user defined function
2. UDAF u.d. aggregate f.
3. UDTF u.d. table generating f.

UDF

- works on a single row
- outputs a single row
- trim(), concat(), length(), round(), floor()

UDAF

- works on multiple row
- outputs a single row
- count(*), sum(), avg()

UDTF

- single
- multiple
- explode(), posexplode()

→
* (flattens the data)
* (provides lateral view)

Set Operations on HIVE

- union (only operation of HIVE)
- minus (not supported in HIVE)
- Intersect (— P —)

union & union all

* How to use union instead of minus & intersect?

THE WHERE CLAUSE

- 1). IN / NOT IN
- 2). EXISTS / NOT EXISTS.

- (a) Select id from customers
where id in (1111, 3333, 5555);
- (b) Select id from customers where exists (
select cust_id from orders where
orders.cust_id == cust.id);

Tutorial 04
hive theory 3

VIEWS IN HIVE

view = virtual table.

" it is better to share few columns rather than all columns "

Benefits :- 1) security.

2). can change according to need.

3). create diff. logical path.

Normalization

* normalization is not preferred.

* denormalization is preferred

↓

Single table
