

Assignment - 2

1) Apache live with Minimum 5 of live query language Command.

⇒ Apache live is a data Warehouse Infrastructure build on top of Hadoop for providing data summarization and hoc queries & Analysis of large data sets using SQL language Called HIVEQL. SQL Queries Over Petabytes of data using Hadoop & others The following features

- ① Tools to enable easy data Extraction, Transformation and loading.
- ② A Mechanism to Impose Structure On variety of data format
- ③ Query Execution via Map Reduce.

⇒ Hive Query Language Command:-

- ① To start Hive, Simply Enter the Hive Command

```
$ hive
```

(some Message may show up here)

```
hive>
```

- ② AS a simplest test, Create and drop table. Hive Must End with (;)

```
hive> CREATE TABLE Pokes (foo INT, bar STRING);
```

```
OK
```

Time taken : 1.075 Seconds.

- ③ hive> SHOW TABLES;

```
OK
```

```
Pokes
```

Time taken : 0.174 Seconds, fetched : 1 rows

```
hive> DROP TABLE Pokes;
```

```
OK
```

Time taken : 4.08 Seconds

④ A More detailed Example Can be developed using Web Server log file to summarize text Message types

```
hive > CREATE TABLE log (t1 string, t2 string, t3 string, t4 string, t5 string)
```

```
Row FORMAT DELIMITED FIELDS TERMINATED BY ',';
```

OK

Time Taken: 0.129 seconds

⑤ Load The data in this case from sample log file.

This file is available from Example code download.

```
hive > LOAD DATA LOCAL INPATH 'sample.log' OVERWRITE INTO TABLE LOG;
```

Loading data to table defaults. logs

Table Default. logs Stats: [numfiles=1, numRows=0, totalSize=99271, rawDataSize=0]

OK

Time taken: 0.953 seconds.

⑥ Explain Apache Sqoop Import and Export Methods With neat diagram

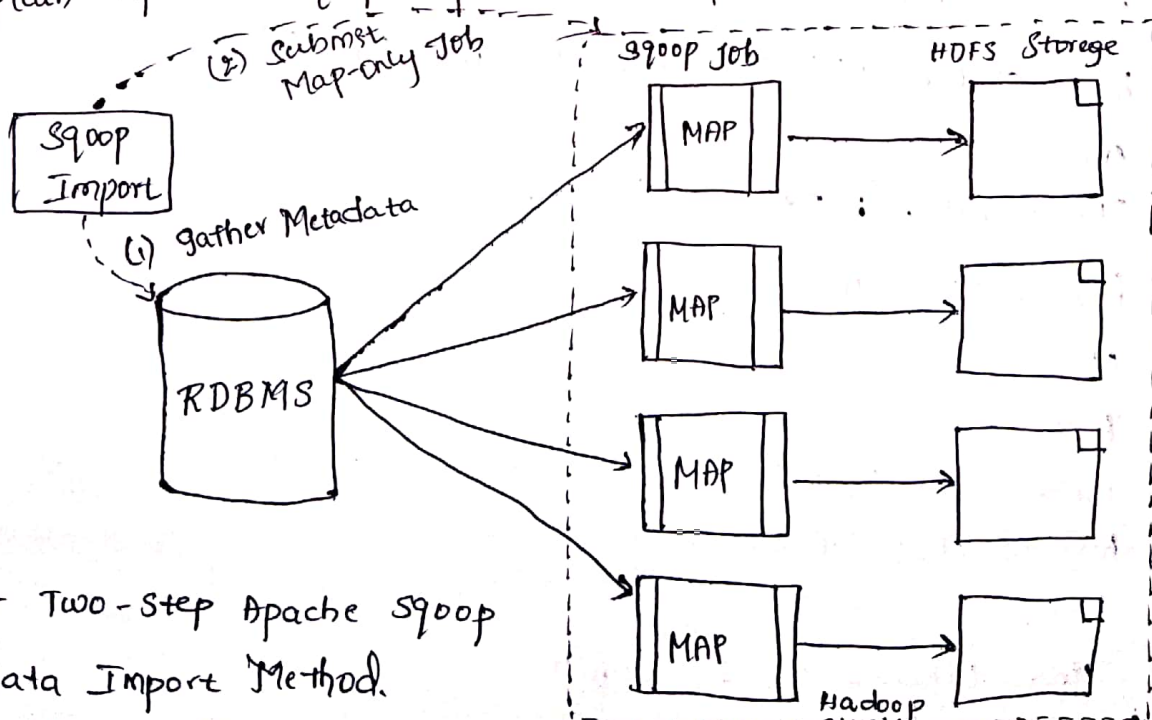


Fig 1:- Two-Step Apache Sqoop Data Import Method.

- fig 1 describes Sqoop data Import (HDFS) process. The data Import is done in 2 steps
- Step I :- Sqoop Examines the database together the Necessary Metadata for data to be Imported.
- Step II :- It is Map-only hadoop job that sqoop Submits to the cluster. This job does actual data transfer using Metadata captured in previous steps.
- Note Each Node doing Import must have access to Database.
- The Imported data are saved in HDFS directory
- Sqoop will use Database name for the directory where the files should be populated.

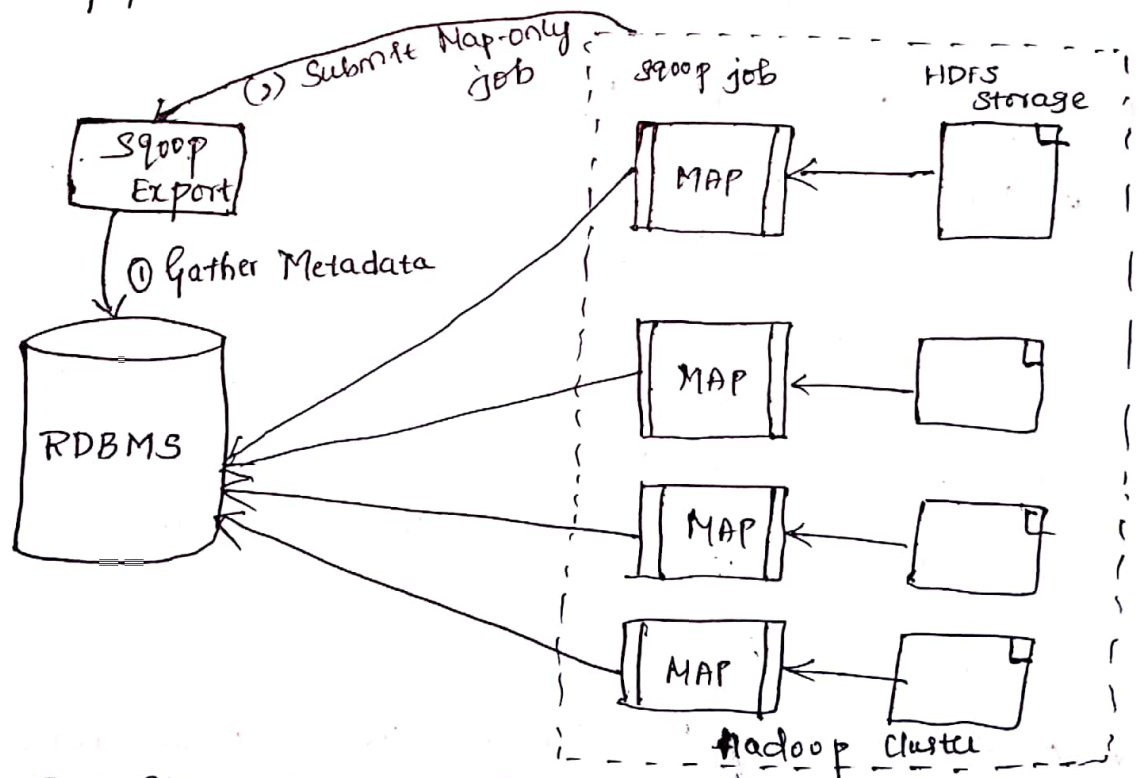


Fig. 2 :- Two Step sqoop data Export Method.

- ① The Export is done in 2 steps as shown in fig 2
- ② Step-1 is to examine database for Metadata
- ③ The Export step again use Map-only Hadoop job to write data to database
- ④ Sqoop divides its data set into splits, then each individual Map takes to push the splits to databases.
- ⑤ Again, This process assumes the Map tasks have access to database.

③ Explain YARN Application Framework.

⇒ The Yarn framework exists to Manage Application.

- A yarn application Implements a specific function that runs on Hadoop.
- A yarn Application Involves 3 Components.
 - Client
 - Application Master (AM)
 - Container.

④ YARN client:-

- * Launching a new yarn application start with yarn client Communicated with Resource Manager to Create New yarn
- * Part of this process involves yarn client Informing the Resource Manager of Application Master's physical resource requirements.

⑤ YARN APPLICATION MASTER:-

- * Application Master is Master process of a YARN application
- * It doesn't perform any application specific work, as these functions are delegated to Container.
- * Once the application master is started it will periodically send heartbeats to the Resource Manager to affirm its health & to update the record of its resource demand.

⑥ YARN Container

- * A Container is an Application - Specific process that's Created by Node-Manager on behalf of an Application Master
- * At the fundamental level, a Container is Collection of Physical resource such as RAM, CPU wires & disks on the single Node.

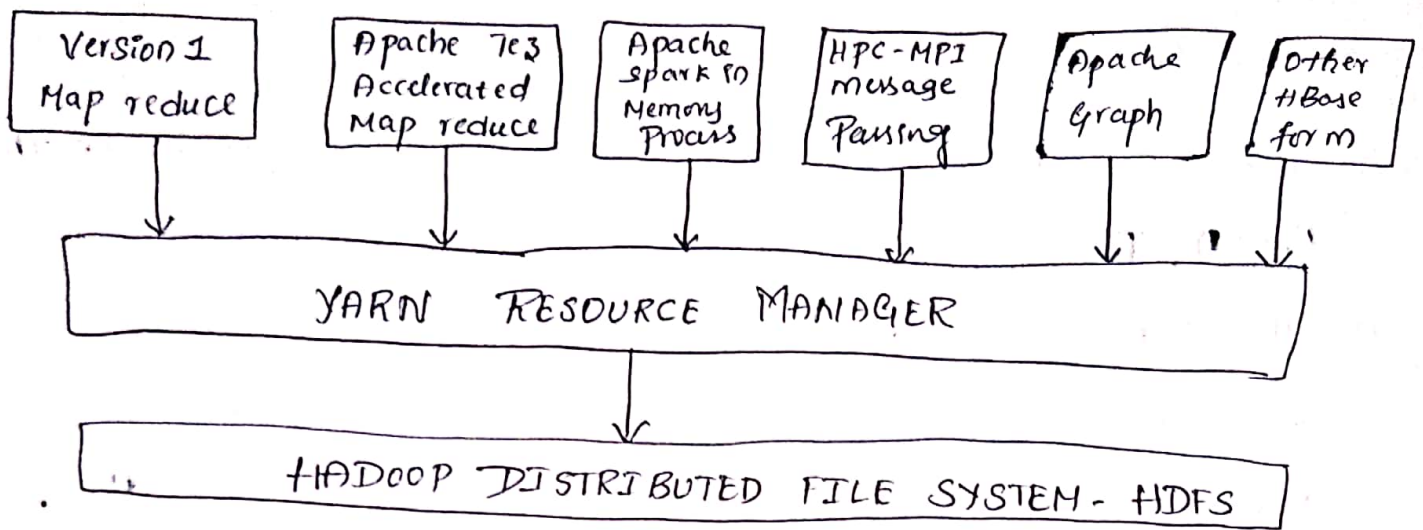


Fig : Example of Hadoop Version 2 Ecosystem.

④ Explain Apache Spark and Apache REEF

⊛ Apache Spark :-

- ① Spark was initially developed for an application in which keeping data in Memory Improve Performance Such as Algorithm, which are Common in Machine Learning.
- ② Spark differ from classic MapReduce in 3 Important ways:
 - ① Spark holds Intermediate results in Memory rather than Writing Them to disk.
 - ② Spark Supports more Than just MapReduce function.
 - ③ It greatly Expands the set of Possible Analyses That can be Executed Over HDFS data stores.
- ④ It also Provides API's in Scala, Java & Python.
- ⑤ The advantage of Porting & running Spark on Top of YARN is common resource Management and a single underlying file sy.

Apache REEF

- ① YARN's flexibility sometimes requires significant efforts on the part of application implementers.
 - ② The steps involved in writing a custom application on YARN include building your own application master, performing client and container management and handling aspects of faults, execution flow.
 - ③ It greatly expands the set of possible analyses that can be executed over HDFS data stores.
 - ④ It also provides APIs in Scala, Java & Python.
 - ⑤ REEF's design makes it suitable for both MapReduce & DAG-like execution as well as interactive and iterative computation.
- ⑤ Write a short note on Apache Ambari.
- ⇒ ① Managing a Hadoop installation by hand can be tedious and time consuming. In addition to keep configuration files synchronized across a cluster, starting, stopping and restarting Hadoop services & dependent services in right order is not a simple task.
- ② The Apache Ambari graphical management tool is designed to help you easily manage these and other Hadoop administration issues.
- ③ Along with being an installation tool Ambari can be used as a centralized point of administration for Hadoop cluster.

④ Using Ambari, The user can configure cluster services, monitor the status of cluster host or services, visualize hotspots by Service Metrics, start or stop service and add new hosts to the cluster.

⑤ All of these features infuse a high level of agility into the process of managing and monitoring a distributed computing environment.

⑥ Ambari also attempts to provide real-time reporting of important metrics.

⑥ Define Data Warehouse and write Design Consideration for DW.

⇒ A datawarehouse (DW) is an organized collection of integrated, subject-oriented databases designed to support decision support function.

- DW is organized at the right level of granularity to provide clean enterprise-wide data in standardized format for reports queries and analysis.

* Design Consideration for DW

The objective of DW is to provide business knowledge to support decision making. There are some requirements for good DW:-

① Subject Oriented :- To be effective, a DW should be designed around a subject domain i.e. to help solve a certain category of problem.

② Integrated :- The DW should include data from many functions that can shed light on particular subject area. Thus the organization can benefit from a comprehensive view of subject area.

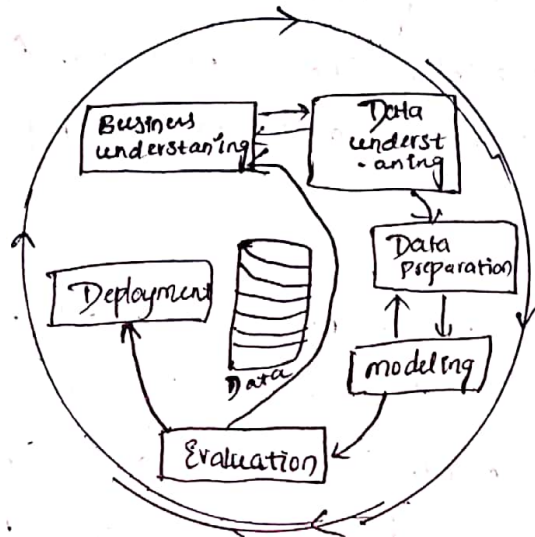
- ③ Time-Variant:- The data in DW should grow at daily or other choosen Interval. That allows latest Comparisons Over times.
- ④ Non-Volatile:- DW should be persistent, i.e it should not be Created on the fly from the Operational databases
- ⑤ Summarized:- DW contains rolled up data at right level of queries and analysis. The process of rolling up the data helps to Create Consistent for effective Comparisons.
- ⑥ Not Normalized:- DW often uses a star schema, which is rectangular central table, surrounded by some look-up tables. The single table view significantly enhances speed of queries.
- ⑦ Metadata:- Many of variables in databases are Computed from other Variables in Operational databases. The Method of Calculation for each variable should be effectively documented Every element in DW should be well-defined.
- ⑧ Near Real-time or Right-time:- DW should be updated in near real-time in many high transaction Volume industries. The cost of updating DW in real-time could be disclouring Through. Another Downside of real-time DW is Possibilities of InConsistencies in reports drawn just a few Minutes apart.

Q7 Explain CRISP-DM cycle with Neat Diagram.

⇒ Effective and Successful use of data Mining activity requires both business and technology skills.

- ① It also helps One imagine Possible relationship in data and Create hypotheses to test it.

- ① There are several best practices learned from the use of Data Mining techniques over a long period of time.
- ② The data Mining industry has proposed a cross-industry standard process for Data Mining.
- ③ It has Six Essential steps.



① Business Understanding:-

The Most important step in data mining is asking right business question. A question is good one if answering it would lead to large Payoffs for Organization, financially and otherwise. There should be strong support for DM Project, which means that Project aligns well with business strategy. Thinking Outside the box is important, both in terms of a Proposed model as well as data sets available are required.

② Data Understanding:-

A related important step is to Understand the data available for mining. One need to be Imaginative in looking for many elements of data through many sources in helping address of hypotheses to solve a Problem. Without relevant data, the hypothesis cannot be tested.

③ Data Preparation :-

The data should be clean, relevant & high quality. It's important to assemble a team that has mixture of technical and business skills. Data cleaning can take 60-70% of time in Data Mining Project.

④ Modeling :- This is actual task of running many Algorithms using the available data to discover if hypotheses are supported. A hosts of Modeling tools and algorithm should be used. A tool could be tried with different Options such as running different decision tree Algorithm.

⑤ Model Evaluation :- One should not accept what the data says at first. It is better to triangulate the analysis by applying multiple data mining Techniques and Conducting many what if scenarios, to build confidence in the Solⁿ. One should Evaluate and Improve the Model's predictive accuracy with More test-data.

⑥ Dissemination and rollout :- It is Important that Data Mining solution is presented to key stakeholders, and is deployed in the Organization. Otherwise the project will be wasted of time and will be setback for establishing and supporting a data-based decision process culture in the Organization. The model should be eventually embedded in the Organization's business processes.

⑦ Explain 5 Important data Mining Techniques.

(a) Decision Trees :- This are most important DM techniques for many reasons :

- (a) Decision tree are easy to understand and use by analyst, as well as Executives, They also Show high predictive accuracy.
- (b) Decision tree Select most relevant Variable automatically Out of all The available Variables for decision making.
- (c) Decision tree are tolerant of data availability issues & do not require much data preparation from users.
- (d) Even Non-linear relationship can be handled well by decision tree.

⑧ Regression :-

- (a) Most popular Statistical Data Mining Technique.
- (b) The goal is to derive smooth well-defined Curve to best the data.
- (c) It can be used to Model & predict The Energy Consumption as function of daily temperature.
- (d) Applying Non-linear regression Equation will fit data very well with high Accuracy.
- (e) The Accuracy of regression Model depends entirely upon the data-sets used and Not at all on the Algorithm or tools used.

⑨ Artificial Neural Network:-

- (a) ANN is Sophisticated DM technique from Artificial Intelligence in Computer Science.
- (b) A decision task may be processed by just one neuron & result may be Communicated Spd.
- (c) The Neural N/w can be trained by making a decision over again with many data points.
- (d) The Neural N/w can be learned Enough & begin to match the predictive Accuracy of Human Expert or Alternative Classification Tech.
- (e) ANNs require a lot of data to train it to develop from Predictive Ability.

⑨ Define Data Visualization and Explain types of charts.

⇒ Data Visualization is the art and science of Making data easy to understand and Consume for The end User.

* Types of Charts

(i) line graph:-

It Shows data as a series of Point Connected by Straight line segment. If Mining with time-series data, time is usually shown on x-axis. Multiple Variable Shows on y-axis to Compare The graph lines of all variables.

(ii) Scatter plot:-

It Helps to reveal the relationship b/w 2 Variables. Unlike in a line graph, There are no line segments Connecting the Points.

(iii) Bar graph:-

A Bar graph Shows Thin colorful rectangular bars with their length being Proportional to values represented the bars can be Plotted vertically or Horizontally.

(iv) Stacked Bar graph:-

Values of Multiple Variable are Stacked up one on top of other to tell an Interesting story. Bar can be Classified such as total height of every bar is equal can be relative composition of each bar.

⑤ Histograms :

These are like bar graphs, except that they are useful in showing data frequencies or data values on classes Numerical values.

⑥ Pie chart :

Pie chart shows distribution of variable such as sales by region.

⑦ Box chart :-

Shows distribution of variables. The box shows middle half of values while whiskers on both sides extend to extreme value.

⑧ Bubble graph :-

This is interesting way of displaying multiple dimension in one chart. It is variant of scatter plot with many data points marked on 2 dimensions.

⑨ Dials :

Charts are like speed dial in car, that shows whether the variable value is in low range, medium or high range.

⑩ geographical Data Map :-

Useful map to denote statistics

⑪ Pictographs :-

One can use picture represent data.