

Unit III : NO SQL Databases.1.1 Introduction.1.1.1 History:-

- The term NOSQL was first used by Carlo Strozzi in the year 1998.
- He mentioned this name for his open source db sys. in which there was no provision of SQL query interface.
- In the early 2009, at conference held in USA, NOSQL was comes into picture and actually comes in practice.

1.1.2 Overview.

- NOSQL is not a RDBMS.
- NOSQL is specially designed for large amount of data stored in distributed environment.
- The important feature of NOSQL is, it is not bounded by the table schema restrictions like RDBMS.
- NOSQL generally avoids join operation.

1.1.3 Need:-

- In real time, data requirements are changed a lot. Data is easily available with Facebook, Google, Twitter and others.
- The data that includes user information, social graphs, geographic location data and other user-generated content.
- To make use of such abundant resource and data. It is necessary to work with a technology which can operate such data.

FJ - 208 Date \_\_\_\_\_  
Page \_\_\_\_\_

- SQL databases are not ideally designed to operate such data.

- NOSQL databases specially designed for operating huge amount of data.

1.2 Advantages..

1. Good resource scalability
2. Lower operational cost
3. Supports semi-structure data.
4. No static schema.
5. Supports distributed computing
6. Faster data processing.
7. No complicated relationships
8. Relatively simple data models.

1.3 disadvantages:-

1. Not a defined standard.
2. Limited query capabilities.

1.4 Companies working with NOSQL.

1. Google
2. Facebook
3. LinkedIn
4. McGraw-Hill

1.5 Four Types of NOSQL database.1.5.1 Key-value store databases:-

- This is very simple NOSQL db.
- It is specially designed for storing data in a schema-free data.
- Such data is stored in a form of data along with indexed key.

- this type is generally used when you need quick performance for basic create-read-update-delete operations and data is not connected.

Hamsterdb, riak, redis, oracle, memcached

### Example

- Storing and retrieving session information for a web page.
- Storing user profile and preferences
- Storing shopping cart data for ecommerce

### Limitations:-

- It may not work well for complex queries attempting to connect multiple relations of data.
- If data contains lot of many-to-many relationships, a key-value store is likely to show poor performance.

examples:- (can diff previous diagram)

- Cassandra, Azure Table Storage (ATS), Dynanodb.

ex of unstructured data for user records

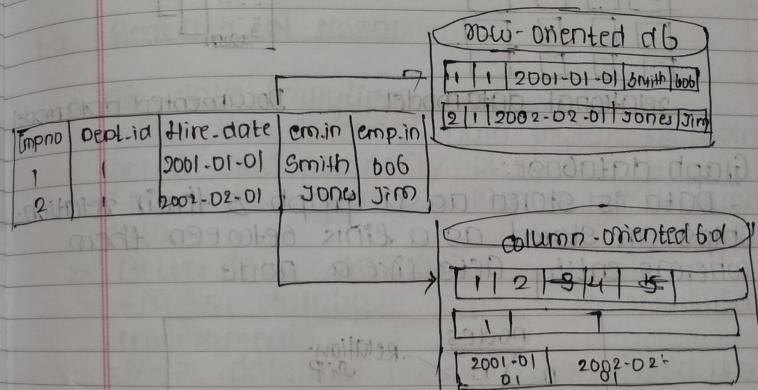
key 1 : ID : 123 [First name : Ganesh]

key 2 : email : abc@gmail.com [Age 37]

key 3 : facebookID : 12345 [Pass : xxx, fname : Max]

### 2. Column Store database:

- Instead of storing data in relational tuples (table rows), it is stored in cells grouped in columns.
- It offers very high performance and a highly scalable architecture.



### Example .

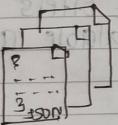
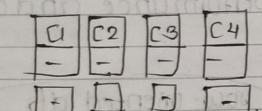
1. HBase
2. Big table
3. Hyper table.

### 3. Document database:-

- Document database works on concept of key-value stores where "documents" containing all of complex data.
- Every document contains a unique key, else to retrieve the documents.
- Key is used for storing, retrieving and managing document-oriented information.

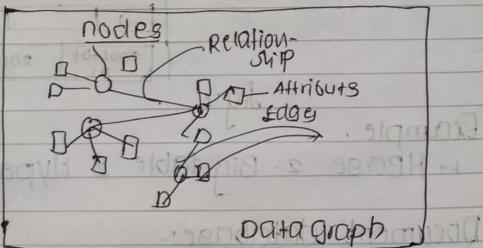
Known as semi-structured data.

- Examples:-
- Mongodb
  - CouchDB



#### 4. Graph databases:-

- Data is stored as a graph & their relation.
- Ships are stored as a link between them whereas entity acts like a node



Example:-

→ Neo4j , → Polyglot

#### 5. Comparison

Datamodel	Performance	Scalability	Flexibility
key-value store	High	High	High
column store	High	High	moderate
document store	High	variable(High)	high
graph database	Variable	Variable	High

#### 1.2 Benefits of NoSQL

##### 1. Big data Analytics:-

- Big data is one of main feature promotes growth and popularity of NoSQL
- NoSQL has good provision to handle such big data.

##### 2. Better data Availability:-

- NoSQL database works with distributed environments.
- NoSQL database environments should provide good availability across multiple data servers.
- NoSQL database supply high performance.

##### 3. Location Independence.

- NoSQL database can read and write database regardless of location of database operation.

#### 1.3 CAP theorem (Brewer's theorem)

CPA theorem states three basic requirement of NoSQL database to design a distributed architecture.

a) consistency:- database remain consistent state like before, even after execution of

b) Availability: - It indicates that NOSQL system is always available without any downtime  
c) partition tolerance :- this means that the system continues to function even the communication failure happens between servers i.e if one server fails, other server will take over.

there are many combinations of NOSQL rules.

#### 1. CA

- It is a single site cluster.
- All nodes are always in contact.
- Partitioning system can block the system.

#### 2. CP

Some data may not be accessible always still it may be consistent or accurate.

#### 3. AP

- System is available under partitioning
- Some part of the data may be inconsistent.

#### 4. BASE Properties for NOSQL

BASE stands for  
B - Broadband  
A - Availability  
S - Soft state

Concurrent requests around the clock  
means field with entries might be  
inconsistent at points of read/write transaction  
between multiple sandwhiches simultaneously.  
In addition after a few write states

1.5 NOSQL Key value database: MongoDB.

#### NOSQL

- NOSQL means "NOT ONLY SQL", NOSQL database can use SQL like query concept.
- It includes all databases that is not a traditional relational RDBMS.
- NOSQL databases are more specialized for various types of data types.
- It is more efficient and better performing than RDBMS.

#### 2. MongoDB

- open source document dbms.
- one of the popular NOSQL db written in c++.
- High performance, high availability and automatic scaling are the important features.
- It has its own ad-hoc query language with rich features set.
- large number of edge cases can be easily handled.
- Mobile application, Content Management System, E-commerce, Gaming Application, Analytics, Advertising and Logging are the application areas of MongoDB.
- Multi-document transactions are not possible in MongoDB.
- It provides atomic operations on a single document.
- The biggest advantage of MongoDB is that it automatically uses a free memory available on the machine.

### 3. Mongo features:-

- High performance as compared to traditional SQL.
- Good support for embedded data models.
- Faster query processing using index support from embedded documents and arrays.
- Higher availability.
- Provides horizontal scalability.
- Provides automatic sharding techniques.

### 4. Collection & document:-

- Collection is equivalent to a table in RDBMS.
- A collection is a group of documents which exists within a single database.
- Collection is schemaless.
- Different fields can be created with different documents in a collection. But typically, all the documents in a single collection are of similar or related purpose.
- A document is set of key-value pairs. Documents have dynamic schema.

### RDBMS terminology with mongoDB terminology.

RDBMS	MongoDB
Database	Database
Table	Collection
Tuple / Row	document
Column	field
Table join	Embedded documents
primary key	Primary key (default key). id provided by mongoDB itself.

### 5. Document database:

- In MongoDB document is used to store data. Each document consists of field and values.
  - MongoDB documents and JSON objects are similar to each other.
  - Other documents arrays etc can be included in the field values.
- Documents have following advantages:-
- Due to embedded documents and arrays, joins are avoided. Hence less expensive.
  - Document uses dynamic schema.
  - In popular programming language documents are native data types.

### - A mongoDB document is as follows:-

{ movie : "Hanuman"

ticket : 100

Screen : 02

}

### 6. MongoDB Schema:-

- It uses dynamic schemas. Documents in a collection may have different set of fields. Due to this polymorphism can be easily used, without defining the structure.
- Adding new field and deleting existing field can be easily done.

- |  |                           |
|--|---------------------------|
| 1. Comparison MongoDB and RDBMS                            | RDBMS                     |
| 2. Schema less   | RDBMS                     |
| 3. Complex joins are avoided                               | • requires schema         |
| 4. document based db                                       | 2. complex joins are used |
| 5. it has its own ad-hoc query language with rich features | 3. table based db.        |
| 6. Easy to scale   | 4. RDBMS has SQL          |
|  | 5. NOT easy to scale.     |

1.6. Comparative study of RDBMS and NOSQL (SQL vs NOSQL)  
SQL dbases are RDBMS; whereas NOSQL database are non-relational database.

## 1. Data storage :-

- SQL databases stores data in a table which has relational database.

NoSQL databases are non-relational database based key-values.

graph databases or wide-column stores.

SQL data is stored in form of tables, with some rows.

- NoSQL data is stored as collection of key-value pair or documents or graph based data with no standard schema definitions.

## 2. Database Schema:-

- SQL databases have predefined schema which cannot be changed very frequently whereas NoSQL databases have dynamic schema which can be changed any time for unstructured data.

- SQL databases provides standard platform for running complex query.
  - NoSQL does not provide any standard environment for running complex query.
  - NoSQL are not as powerful as SQL query language.

Q1

- Structured Query Language
  - SQL is a declarative query language
  - SQL db works on ACID properties
    - Atomicity
    - Consistency
    - Isolation
    - Durability
  - Structured & organized data
  - Not Only SQL or Non-relational database.
  - This is Not a declarative query language.
  - NoSQL database follows the Brewers CAP theorem,
    - Consistency
    - Availability
    - Partition Tolerance
  - Unstructured and Unstable data.

# NOSQL

- Structured Query Language
  - SQL is a declarative query language.
  - SQL db works on ACID properties
    - Atomicity
    - Consistency
    - Isolation
    - Durability
  - Structured & organized data
  - Relational db is table based
  - Not Only SQL or Non-relational database.
  - This is Not a declarative query language.
  - NoSQL database follows the Brewers CAP theorem.
    - Consistency
    - Availability
    - Partition Tolerance
  - Unstructured and Unstable data.
  - Key-value pair storage column store, document store graph db.

5. Data and its relationships are stored in separate tables.

- light consistency
    - ex - MySQL, Oracle, MS SQL, SQLite, DB2.
  - eventual consistency
    - ex - MongoDB, big table, Neo4j, couchdb, cassandra, HBase.

## 17 Advantages of NoSQL

1. The growth of big data.
2. Continuous Availability of data.
3. Location Independence.
4. Modern transactional capabilities.
5. Flexible data models.
6. Better Architecture.
7. Analytics and Business Intelligence.

## 18 MapReduce

- Map reduce is a data processing paradigm for condensing large volume of data into useful aggregated results.
- MongoDB uses mapReduce command for map-reduce operations.
- generally used for processing large data sets.

### mapReduce command. (Syntax).

```
>db.collection.mapReduce(  
  function() {emit(key,value);},  
  function(key,values) {return reduceFunction();},  
  {  
    out: collection,  
    query: document,  
    sort: document,  
    limit: number  
  })
```

5.

The map-reduce function first queries the collection, then maps the result documents to emit key-value pairs, which is then reduced based on the keys that have multiple values.

In the above syntax

- map is a javascript function that reduces or groups all the documents having the same key.
- out specifies the location of the mapreduce query result.
- reduce is a javascript function that reduces or groups all the documents having the same key.
- query specifies the optional selection criteria for selecting documents.
- sort specifies the optional sort criteria.
- limit specifies the optional maximum number of documents to be returned.

example of mapReduce :-  
rdb.posts.mapReduce(  
 function() {emit(this.user\_id, 1);},  
 function(key,values) {return Array.sum(values)}

query: {status: "active"},  
out: "post\_total",  
sort: {}

output:-

```
{"result": "post_total",  
 "timeMillis": 9,  
 "counts": 8,  
 "input": 4,  
 "emit": 4,  
 "reduce": 2,  
 "output": 2}
```

1.9 HJVE

- HIVE is a datawarehouse infrastructure tool.
  - It processes structured data in HDFS. HIVE structures data into tables, rows, columns and partitions.
  - It resides on top of Hadoop.
  - It is used to summarize big data, analysis of big data.
  - It is suitable for online Analytical application processing
  - It supports ad hoc queries. It has its own SQL language called HiveQL or HQL.
  - Primitive datatypes like integers, floats, doubles and strings are supported by HIVE.
  - Associative arrays, lists, structs etc can be used.
  - Generic API and Deserialized API are used to store and retrieve data
  - HIVE is easy to scale and has faster processing

## 2.0 Architecture of HIVE

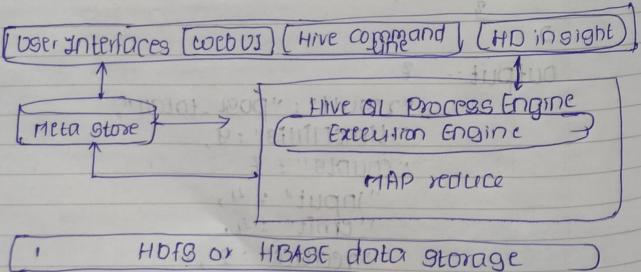
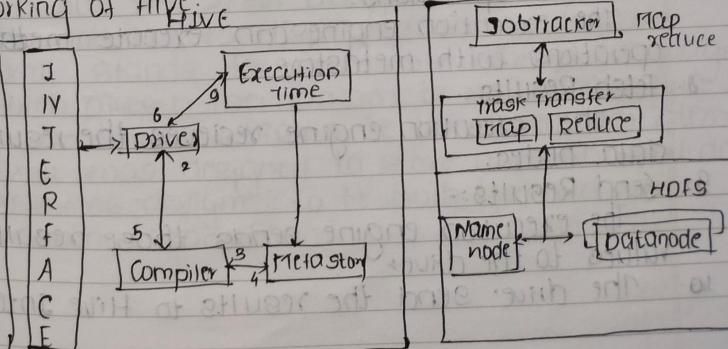


fig : Hive Architecture.

This information cannot be changed further as these are the standard components.

1. User Interface:- Hive supports Hive Web UI, Hive command line, and Hive TUI through which user can easily process queries.
  2. Meta Store - Hive stores meta data, schema etc. in respective database servers known as metastores.
  3. HiveQL Process Engine:- HiveQL is used as querying language to get information from metastore. It is an alternative to MapReduce Java program. HiveQL query can be written for MapReduce job.
  4. Execution Engine: Query processing and result generation is the job of execution engine. It is same as that of MapReduce results.
  5. HDFS or HBASE: Hadoop distributed file system or HBASE use the data storage techniques to store data into file system.

## 21. Working of Hive



## Sig Hive & Hadoop Communication

Command line or web UI sends query to JDBC or ODBC driver to execute.

1. Get plan:- With the help of query compiler driver checks the syntax and requirement of query.

#### 2. Get Metadata:-

The compiler sends metadata request to metastore for getting data.

#### 3. Send metadata:-

Metastore sends the required metadata as a response to the compiler.

#### 4. Send plan:-

The compiler checks the requirements and resends the plan to the driver. Thus the parsing and compiling of a query is complete.

5. Execute plan :- the driver sends the execute plan to the execution engine.

6. Execute Job :- The execution engine sends the job to JobTracker. JobTracker assigns it to TaskTracker.

#### 7. Metadata operations:-

The execution engine can execute metadata operations with metastore.

#### 8. Fetch Results:-

The execution engine receives the results from data nodes.

#### 9. Send Results:-

The execution engine sends those resultant values to the driver.

10. The driver sends the results to Hive interface.

## Hive Data models.

The Hive data models contain the following components:-

- Databases
- Tables
- Partitions
- Buckets or clusters

**Partitions**:- The table is divided into a smaller parts based on the value of a partition column. Then on these slices of data queries can be made for faster processing.

**Buckets**:- Buckets give extra structure to the data that may be used for efficient queries. Different data required for queries joined together. Thus queries can be evaluated quickly.

## 2.3

### Programming language - XML.

#### • Introduction:-

- XML is a markup language is very much like HTML but XML was designed to carry data and not to display data.

- XML stands for Extensible Markup Language which gives a mechanism to define structure a document which is to be transferred over internet.

- XML was designed to store and transport data.

- XML was designed to be both human and machine-readable.

- XML plays an important role in many different systems.
- XML is often used for distributing data over the Internet.
  - It is important (for all types of software developers) to have a good understanding of XML

2.5

XML does not do anything.

Notes

```
<notes>
  <to> Tove </to>
  <from> Jani </from>
  <heading> Reminder </heading>
  <body> don't forget me this weekend! </body>
</notes>
```

- The XML above is quite self-descriptive.
- It has sender and receiver information
  - It has a heading
  - It has a message body

But still, the XML above does not do anything.  
XML is just information wrapped in tags.

Someone must write a piece of software to send, receive, store, or display it.

Note

To: Tove  
From: Jani

Reminder

Don't forget me this weekend!

- 2.6. difference between XML and HTML
- XML - was designed to carry data - with focus on what data is.
  - HTML - was designed to display data - with focus on how data looks.
- XML tags are not predefined like HTML tags are

2.7

valid XML documents.

XML documents which satisfy given below conditions are called as XML documents.

a) conditions:-

- The doc. must be well formed
- Element used in the start & end tag pairs must follows the specified structure
- This structure is specified in a separate XML DTD (Document Type Definition) file or XML schema file

b) DTD Syntax :-

- Start with a name is given to the root tag of the document

Company

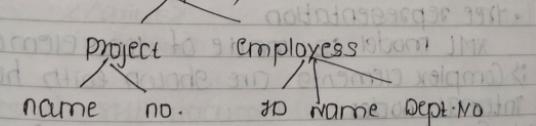


fig An XML DTD file called Company.

In the XML document the basic object is XML and it can be represented as hierarchical data model or tree data model.

## 2.9 Structuring element concepts used to construct an XML document

### a) Element:-

- An element is a group of tags data values that can contain character data, child elements or a mixture of both
- Elements can be two type simple and complex
- The elements are constructed from other elements by nesting them is called complex element.
- The elements contain data values are called as simple elements.

### b) Attributes:-

- Additional information that describes elements
- c) There are some additional concepts used in XML, such as entities, identifiers and references

### d) Tree representation

XML model is made of two elements

i) Complex elements are shown with help of internal nodes

ii) Simple elements are shown by leaf nodes.

Hence, XML model is called as a tree model or hierarchical model.

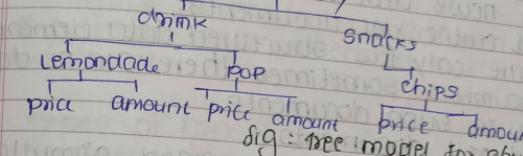
example:

as tree representation:-

Simple elements: <price>, <amount>

Complex elements: <drink>, <snacks> etc

3. Inventory



digraph tree model for above code.

### 6) Textual representation:-

- whenever value of the STANDALONE attribute in XML document is set to "YES" such XML documents is known as schemaless XML documents

ex <inventory>

<drinks>

<lemonade>

<prices> 2.50</prices>

</lemonade>

<drink>

</inventory>

### 3.1 Types of XML documents:-

#### a) Data - centric XML documents

- The document contains only small data items that follow a specific structure and hence it may be extracted from a structured database are called as data - centric XML documents.

- In order to exchange and display data over internet the document is formatted as XML documents.

b) document-centric XML documents.

The document contains large amounts of texts, such as news articles or books are referred as document centric XML documents.

- There are only few structured data elements in these documents. Sometimes there are no data elements in such documents.

c) Hybrid XML documents:-

- If both types of data are used simultaneously in document then it is referred as hybrid XML document.

- In such documents some part that contain structured data and other parts are predominantly unstructured.

### 3.2 XML- Structured data or Semi Structured data.

- Data centric XML documents sometimes considered as semi-structured data or sometimes considered as structured data.

Structured data if an XML document is written as per predefined XML schema or DTD.

- Semi structured if an XML follows documents that do not conform to any particular schema.

### 3.3 XML Document Type Definition (DTD).

- DTD describes the rules for analysing an XML documents.

- A DTD ensures the contents of an XML documents conform to expected rules the document should follow.

Types of DTD as per location of DTD.

a) Internal - DTD embedded in XML document.

b) External - DTD as a separate external file.

Types of DTD declarations.

a) Element type b) Attribute list or Entity d) Notation

a) Element type declarations:- defines one of kind of elements you can use that is, one of the tag types.

b) Attribute list declarations . identifies which elements have attributes they may have attribute values and optional attributes.

c) Entity declarations , = i>name ii> type iii> An optional delay va

### 3.4 Advantage of DTD:

- allow an XML parser to validate an XML document against its definition. This allows the online validation of an XML document before it is processed.

- It can extend capabilities of an XML document by allowing further processing file possibilities.

- DTD can serve as an automotive document for XML documents.

disadvantage of DTD.

- a> It is written in a different (non XML) syntax
- b> It has no supports for namespaces
- c> It only offers extremely limited data typing

Programming language - JSON.

~~Java~~

- JSON stands for JavaScript object Notation
- basically used for data transmission in the web applications
- is the extension of JS language
- can be used for storing and exchanging data
- can be used as an alternative to XML.
- is lightweight, text based data-interchange format. Text can be read and use as a data format by all programming languages.
- is self-describing and easy to understand
- has file name extension .json
- supports name-pair values, arrays, lists, vector, sequence etc.

JSON code

```
{ "name": "movies", "year": 2000, "list": [ { "name": "Tom and Jerry", "year": "2000" } ] }
```

JSON data types

i> Number:- defined as a double precision floating-point format.

- . Octal formats, hexadecimal formats, NaN, Infinity are not used.

ex : var jnum = (marks 69)

. Here jnum is the json object name given by user

ii> Boolean - used to indicate either true or false value

ex : var jnum = { name: 'Sachin', marks: 69, firstclass: true }

iii> null: null means empty, ex : var c = null;

iv> JSON value : may include string, boolean value integer, float null etc

ex : var z = 2;

var c = "Sachin";

var d = null;

v> string : sequence of characters &

ex : var jnum = { name: 'sachin' };

vi> Array : collection of similar data type.

vii> Object : uses name / value pairs.

All the keys should be different from each other

ex : {

"Screen": 0.5

"movie": "John Carter",

"ticket": 200

}

viii> whitespace : used between token pairs to make code more readable.

ex : var c = "Sachin"

for collecting & storing info. of a particular organization.

- goal of using a dw is to have an efficient way of managing information and analyzing data.

Data warehouse Components:-

1. Source Data Components

- In DW Source data is coming from various heterogeneous data sources grouped into four broad categories.

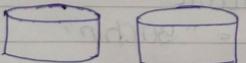
a) Production data

- this data is coming from the various heterogeneous (different) operational system of organization

- Based on data requirements in the data warehouse,

we may choose segments of data from the different operational systems.

ex:- in case of ABC choppers stop all sales data can be considered as production data

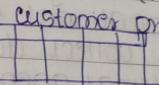


Sig: representation of production data

b) Internal data:-

- In any organization will have various documents like customer profiles, registers, employee data and sometimes even financial databases.

- data is generated by internal organizational operations.  
ex - In big organizations employee's personal data may be treated as internal data, or in sales database customer profile is an internal data



Sig: representation of internal data

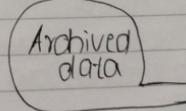
c) Archived data

- the main task of operational systems is to run the current business operations.

- In every operational system, we will periodically fetch the old data and store it as archived files.

- As per corporate requirements, we will decide how often and which portions of the operational databases are archived for storage.

ex: Daily transaction will be archived during night time in case of daily sales database



Sig: Representation of archive data

## Datawarehouse Architecture

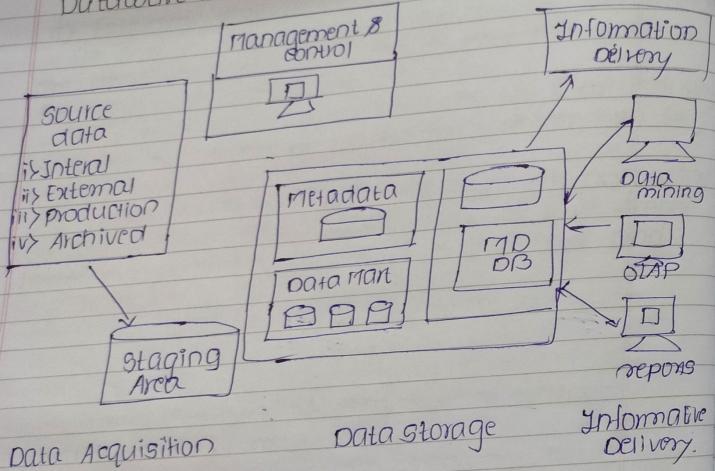


fig: Architectural components in the 5 major areas

- It is primarily based on the business processes of a business enterprise taking into consideration the data consolidation across the business enterprise with adequate security, data modelling and organization, extent of query requirements.

- Meta data management and application, warehouse staging area planning for optimum bandwidth utilization and full technology implementation.

- 2 main areas of data warehouse:-
- a) Data acquisition
- b) Data storage
- c) Information delivery.

- 3 Building blocks of the data warehouse.
- a) Source data
- b) Data staging
- c) Data storage
- d) Information delivery
- e) Metadata
- f) Management & control.

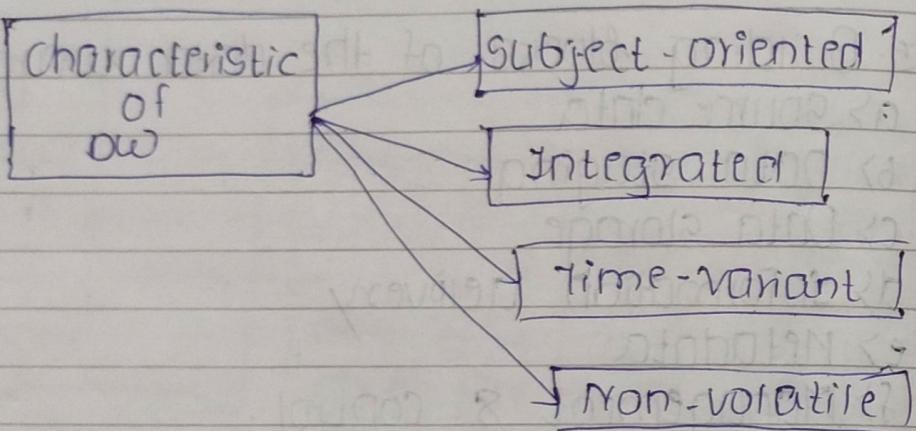
### 4 DW Architecture:-

- In order to set up this information delivery system, we need different building blocks.
- These building blocks are amalgamated together in the most optimal way to serve intended target.
- Architecture, in the context of an organization's dw efforts, is a conceptualization of how the dw is built.
- DW relates all components (which has definite functions and provides specific services together) so as to make fully functional data warehouse.
- Architecture is the proper arrangement of the components.
- We can build at a data warehouse with software and hardware components.
- To suit the organizational requirements, we need to arrange these building blocks in a manner to maximum benefit.

# Datawarehouse - characteristic, Advantages & Disadvantages.

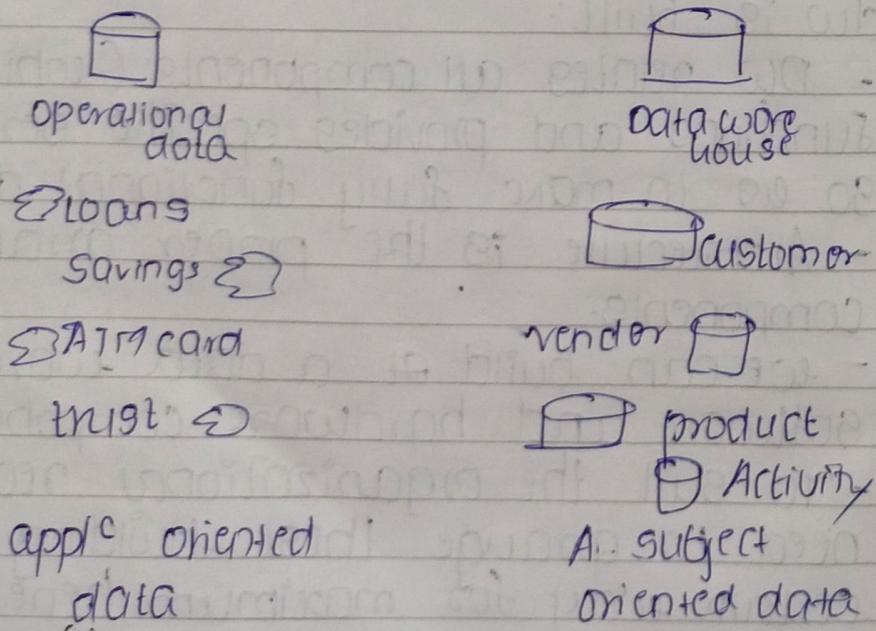
## i. characteristic

### a. Subject - Oriented data



### a. Subject - Oriented

- DW is organized around subject such as sales, product, customer etc.
- it focuses on modeling & analysis of data for decision makers like managers, CEO higher hierarchy in organizations



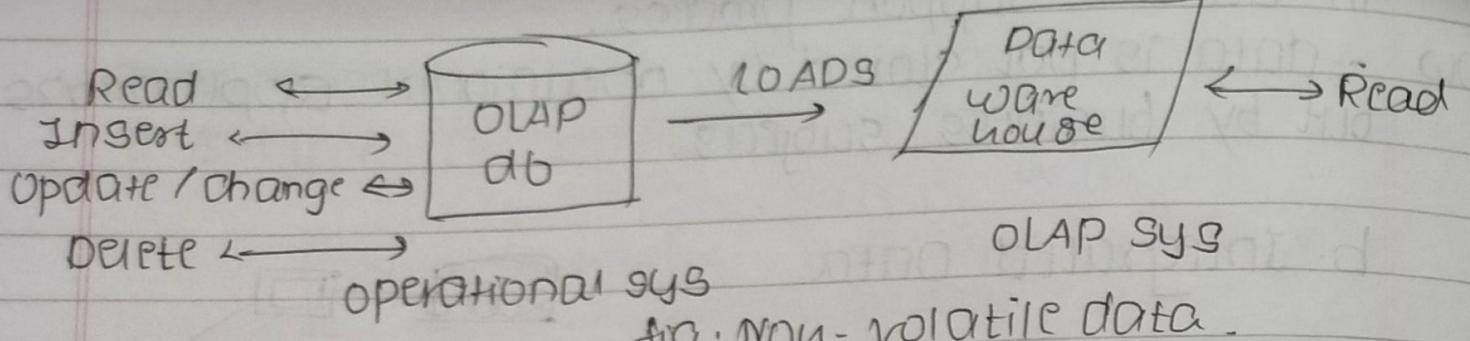


fig: Non-volatile data.

The dw pt is a physically generated data storage, which is transformed from the source Operational RDBMS.

The operational updates of data do not occur in dw. ie update, insert & delete operations are not performed.

- It usually requires two procedures in data accessing : initial loading of data and access to data.

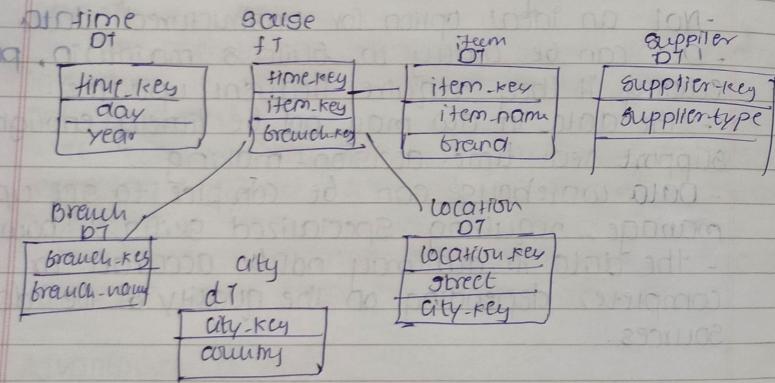
- Non volatile defines that once entered into the warehouse, and data should not change.

### Advantages of DW.

- DW allows business user to quickly access critical data from some sources all in one place.
- DW provides consistent information on various cross functional activities.
- It is also supporting ad-hoc reporting & query.
- DW helps to integrate many sources of data to reduce stress on the production system.
- DW helps to integrate many source of help to

- dimension table containing the set of attributes
- There is a fact table in center (fig\*)
- It contains the key to each of four dimensions
- The fact table also contains the attributes, namely dollar sold & units sold

## 2) Snowflake schema



- Some dimension tables in the snowflake schema are normalized
- normalization splits up the data into additional tables
- unlike star schema, the dimensions tables in snowflake schema are normalized.
- for ex, the item dimension table in star schema is normalized. → split into two dimension tables, namely item & supplier

Supplier key is linked to the supplier D1, the supplier D1 contains the attributes supplier key and supplier-type.

OLAP and types of OLAP Architecture  
a) ROLAP b) MOLAP c) HOLAP.

online Analytical processing (OLAP) Server is based on the multidimensional data model.

It allows manager, and analysis to get an insight of the information through fast, consistent, and interactive access to information.

### Types of OLAP servers.

1. Relational OLAP (ROLAP)
2. Multidimensional OLAP (MOLAP)
3. Hybrid OLAP (HOLAP)

### 1. ROLAP.

- ROLAP servers are placed between relational back-end server and client front-end tools.
- To store and manage warehouse data, ROLAP uses relational or extended-relational DBMS.

ROLAP includes:-

- Implementation of aggregation navigation logic.
- Optimization for each DBMS back end.
- Additional tools and services.

## 2. MOLAP:-

MOLAP uses array-based multidimensional storage engines for multidimensional view of data.

- With multidimensional data stores, the storage utilization may be low if the data set is sparse.
- Many MOLAP servers use two levels of data storage representation to handle dense and sparse data sets.

## 3. HOLAP:-

Hybrid MOLAP is a combination of both ROLAP and MOLAP.

It offers higher scalability of ROLAP and faster communication of MOLAP.

- HOLAP servers allow to store the large data volumes of detailed information.
- The aggregations are stored separately in MOLAP store.

## ⑤

### OLAP operations:-

- a) roll up
- b) Drill down or slice/dice.
- c) pivot (rotate).

#### rollup

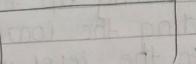
rollup performs aggregation on a data cube in any of the following ways:-

- By climbing up a concept hierarchy for a dimension.
- By dimension reduction.

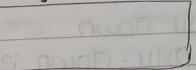
### 3. Slice:-

The slice operation select one particular dimension from a given cube and provides a new sub-cube.

Here slice is performed for "dimension 'time' using the criterion time = "Q1"



Slice for time  
= "Q1"



- It will form a new sub-cube by selecting one or more dimensions

### 4. Dice:-

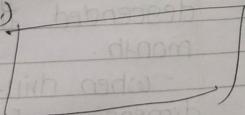
Dice selects two or more dimensions from a give cube and provides a new sub-cube

The dice operation on the cube based following selection

Criteria involves 3 dimensions

one for location

- [location = "Toronto" or "Vancouver"]
- [Time = "Q1" or "Q2"]
- [Item = "Mobile" or "Modern"]

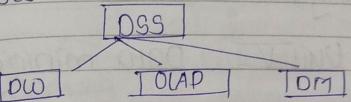


4. Pivot:- it is known as rotation. It rotates the data axes in view in order to provide an alternative presentation of data.

### 6. Decision Support System:-

DSS is a kind of computerized information systems that support decision making activities.

- DSS provides each kind of decision information as well as solution of commercial question for the enterprises.



• There are two kinds of DSS system

1. thinks so long as the system has certain supports to the decisions is DSS.

2. thinks DSS is interactive computer-based system which can help the decision-maker using data and model to solve non-structure questions.

• The majority is not DSS. bcoz most of them do not help to solve non-structure problems

• Some are not interactive:- some databases become the model stores house and they are not entire

DW:- DW does not like traditional database which faces the service level. It mainly faces the high level application and carries on the decision support. Thus, in fact, it is expansion of DSS database.

interactive way and can promptly purpose information from mutative data and not too complete data the movement which relates with enterprise management.

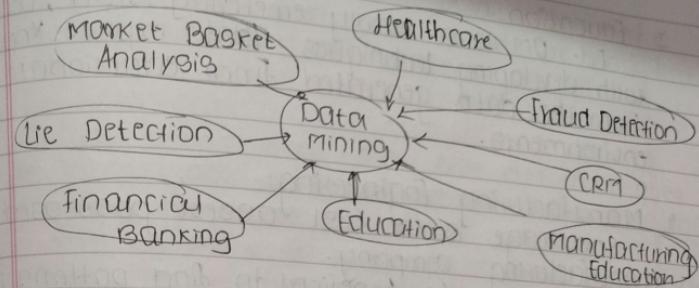
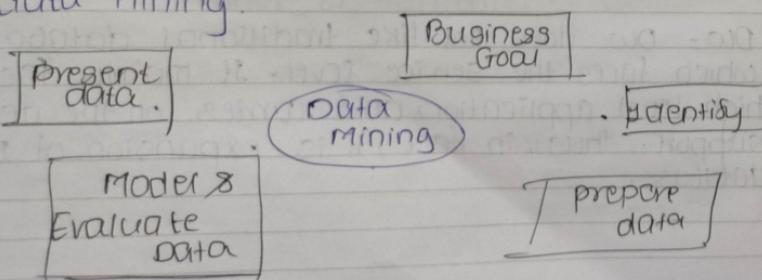
## Unit IV Data mining

### 1. Data mining & it's application.

- DM is widely used in diverse areas.
- There are a number of commercial data mining system available today and yet there are many challenges in the field.

DM helps banks work with credit rating and anti-fraud systems, analyzing customer financial data, purchasing transactions, and card transactions.

[The process of extracting information to identify patterns, trends, and useful data that would allow the business to take the data driven decision from huge sets of data is called data mining.



### 1. Healthcare:-

DM in healthcare has excellent potential to improve the health system.

- It uses data and analytics for better insights and to identify best practices that will enhance health care services and reduce costs.

- Analysts use DM approaches such as ML, multi-dimentional db, Data visualization, soft computing & statistics.

- DM also enables healthcare insures to recognise fraud & abuse.

### 2. Market Basket Analysis.

- MBA is a modeling method based on hypothesis.

- If we buy a specific group of products, then are more likely to buy another group products.

- Using a different analytical comparison of results between various stores, between custom in different demographic groups can be done.

### 3. Education:-

- Ed. DM is a newly emerging field, concerned with developing techniques that explore knowledge from the data generated from educational environments.

### 4. Manufacturing Engineering:-

Knowledge is the best asset possessed by a manufacturing company.

- DM tools can be beneficial to find patterns in a complex manufacturing process.
- It can also be used to forecast the product development period, cost, and expectations among the other tasks.

### 5. Customer Relationship management:-

CRM is all about obtaining & holding customers also enhancing customer loyalty & implementing customer-oriented strategies.

To get a decent relationship with the customer, a business organization needs to collect data and analyze the data.

With DM technologies, the collected data can be used for analytics.

### 6. Fraud detection:-

- Billions of dollars are lost to the action of frauds.

- An ideal fraud detection system should protect the data of all the users.

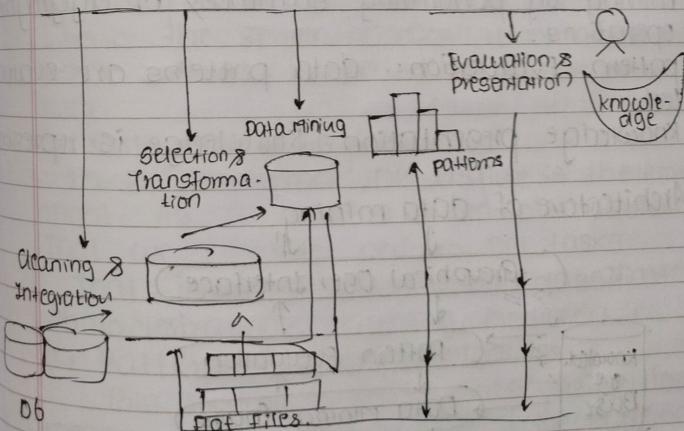
Supervised methods consist of a collection of sample records, and these records are classified as fraudulent or non-fraudulent.

### 7. Lie Detection:-

Law enforcement may use dm techniques to investigate Offenses, monitor suspected long communication etc,

This technique include text mining also, and it seeks meaningful patterns in data, which is usually meaningful pattern in data, which is usually unstructured text.

### ⑥ KDD process with diagram.



Date \_\_\_\_\_  
Page \_\_\_\_\_

DM Architecture system contains too many components. That is a data source, data warehouse, server, data mining engine, and knowledge base.

#### a. DATA SOURCE.

There are so many documents present. That is db, data warehouse, www.

- That are the actual source of data.
- Sometimes, data may reside even in plain text files or spreadsheet.
- www or the Internet is another big source of data.

#### b. DB or DW <sup>Warehouse</sup> Server.

- The db server contains the actual data that is ready to be processed.
- Hence the server handles retrieving the relevant data. That is based on DM request of the user.

#### c. Data mining Engine.

In DM system DM engine is the core component. As it consists a no. of modules that we used to perform DM tasks.

- That includes association, classification, characterization, clustering, prediction, etc.

#### d. Pattern Evaluation Modules.

This module is mainly responsible for the measure of interestingness of the patterns. For this, we use a threshold value. Also, it interacts with the DM engine.

That's main focus is to search towards interesting patterns.

## e. Graphical user interface.

- we use this interface to communicate between the user & the DM system.
- Also, this module helps the user use the system easily & efficiently.
- they don't know the real complexity of the process.

When the user specifies a query, this module interacts with the DM system. thus, displays the result in an easily understandable manner

## f. Knowledge Base:-

- In whole DM process, the knowledge base is beneficial.
- we use it to guiding the search for the result patterns.
- The KB might even contains user beliefs and data from user experience.
- That can be useful in the process of DM.
- The pattern evaluation module interacts with the KB.

## DM in business intelligence

Business intelligence & DM differ in a few core ways. Namely, in purpose, volumes & results.

The purpose of BI is to convert data into useful information for executive.

BI tracks key performance indicators & presents data in a way that encourages

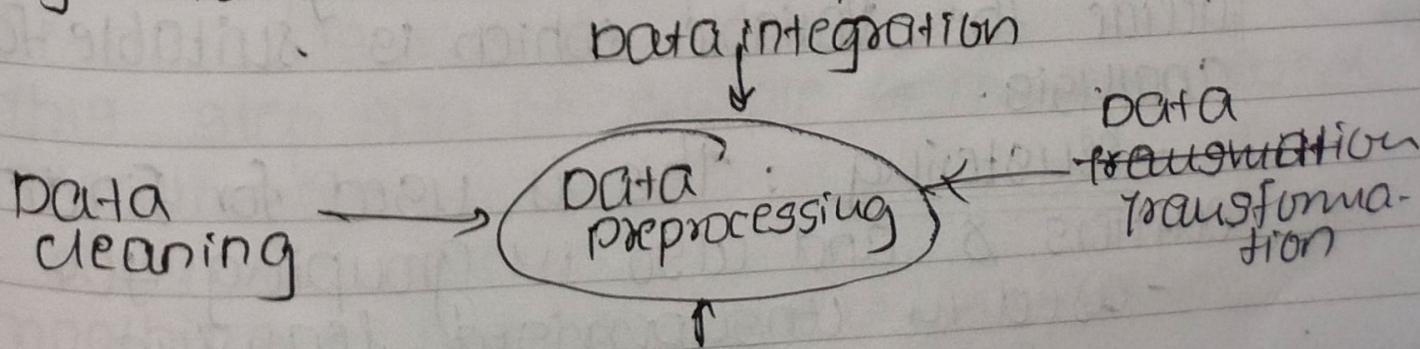
- It is also an important step in DM as we cannot work with raw data.

The quality of the data should be checked before applying ML or DM algor.

Preprocessing of data is mainly to check the data quality.

- The quality can be check by following
- Accuracy :- To check whether the data entered is correct or not.
- Completeness :- To check whether the data is available or not recorded.
- Consistency :- To check whether the same data is kept in all the places that do or do not match.
- Timeliness :- The data should be updated correctly.
- Believability :- The data should be trustable.
- Interpretability :- The understandability of the data.

Major tasks in Data Preprocessing



Data cleaning :- is the process to remove incorrect data, incomplete data & inaccurate data from the datasets, and it also replaces the missing values.

\* noisy :- generally means random error or containing unnecessary data pts.  
Some methods to handle noisy data.

\* Binning - This method is to smooth or handle noisy data.

- There are 3 methods for smoothing by bin mean method -

- Smoothing by bin mean method.
- the values in the bin are replaced by the mean value of the bin.

- Smoothing by bin boundary :- the using minimum and maximum values of the bin values are taken and the values are replaced by the closest boundary value  
~~sometimes~~ sometimes by the bin median - the values in the bin are replaced by the median.

regression :

This is used to smooth the data & will help to handle data when unnecessary data is present.

- for analysis purpose regression helps to decide the variable which is suitable for our analysis.

- clustering : This is used for finding the clusters & and also in grouping the data used in unsupervised learning.