

Dataware house

A dataware house is a powerful database model that significantly enhance the user's ability to quickly analyze large multidimensional data sets. It cleanse and organize data to allow users to make business decision based on facts. Hence, the data in data-warehouse must have strong analytic characteristics.

Dataware house is a field that has grown from the integration of a number of different technology and experience over the past decade. These experience have allowed the IT industries to identify the key problem that need to be solved.

Data can be classified into 3 categories.

- Reference and transaction data → Reference and transaction data originate from operational system and it normally kept in conventional database system. At a predetermind ~~fixed~~ periodicity. It will be loaded for refreshing the datawarehouse. Once it put in data warehouse then it can-not be modified.

- Derived data → Derived data on the other hand, is derived from reference and transaction data, based on certain rules and computations. It can always be

derived again on fresh or modified basis of derivation

- Denormalize data \Rightarrow Denormalized data, which is the basis for online analytical processing (OLAP) tools is prepared periodically but is directly based on detailed reference (transaction) data

* Data-warehouse characteristic

- It typically integrates \uparrow resources e.g. sales database from various regions/states/year
- It requires more historical data than generally maintained in operational database.
- It must be optimized for access to very large amount of data.
- It's mostly read-accessed and rarely write accessed.
- Data may be more coarse grained than in operational database.
- Data warehouse are maintained separately from operational database.

- It is based on client server architecture
- It provide multi-user support
- It is capable of handling dynamic sparse matrix
- It support unrestricted cross-dimensional operation
- It maintain transparency.

II Database vs Datawarehouse

Database	Data warehouse
It support operation process	It support analysis and performance reporting
Capture and maintain data	Explore the data
current state	multiple layers of history
Data is balanced within the scope of this one system	Data must be integrated and balanced from multiple systems
Data is updated when transaction occur	Data updated on scheduled process
Data verification is done when entry is done	Data verification occur after the fact
Size 100 mb to gb	100 mb to Tb

Application oriented

subject oriented

flat relational

multi-dimensional

• The compelling need of and purpose of Data Warehouse

The principle purpose of data warehouse is to provide information to business users for strategic decision-making. These users interact with data warehouse using front end tools, or by getting the required information through the information delivery system.

Need of data warehouse

- Rapid growth of data volume and variety in modern organization has created a need for a centralized repository for storing and managing data from disparate sources.
- Operational systems are often designed for specific tasks and are not well suited for complex analysis. A data warehouse provides a unified view of data from across the organization, making it easier to

perform queries and analysis.

- The need of timely and accurate information for decision making is more important than ever before. A dataware house can provide users with access to the data they need to make informed decisions.

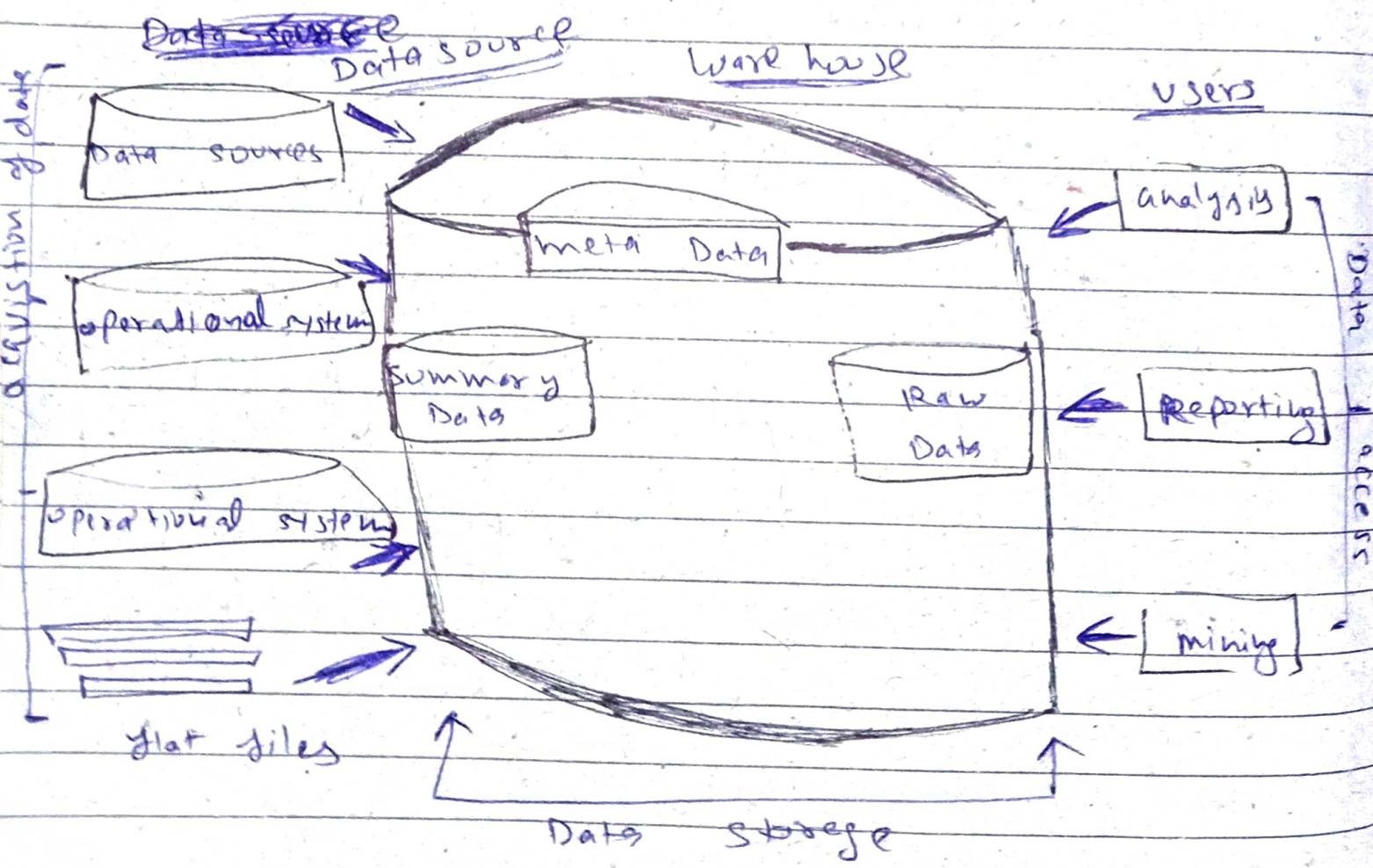
Purpose of data warehouse

- To provide a centralized repository for storing and managing the data from dispersed sources.
- To transform data into a consistent ~~informat~~ format that is easy to understand and analyze.
- To support the development of decision support system (DSS) and business intelligence.
- To enable data mining and machine learning initiatives.

Data warehouse architecture

The hardware and software and data resources

required to construct the data warehouse often depends upon the organization that want to construct it. The need and resource available, force the decision by organization regarding the architecture of particular data warehouse, there are many phases, that are common to all the data warehouse regardless of the organization or the design selected. The most common phases are acquisition of data, storage of data and data access.

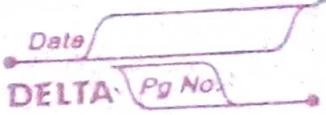


- acquisition of Data \Rightarrow All the data warehouse must have a source from where the data is acquired. Most of data in the warehouse is derived from the operational data of the organization. The required data is extracted, filtered & translated and integrated into data storage environment.
- storage of Data \Rightarrow the large amount of operational data that is historical in nature are defined, indexed and then partitioned allows for economic and efficient access.
- Data access \Rightarrow A number of data mining applications allow many users throughout the organization to retrieve, analyze, query and generate reports. The ability to access data is fundamental to the concept of data warehouse in the organization.

Data warehouse component

There are mainly six component

- (i) summarized data
- (ii) integration/transformation process
- (iii) Operational data store
- (iv) detail data
- (v) meta data
- (vi) archives



(i) Summarized data \Rightarrow The raw data generated by a transaction processing system may be too large to store online. However, many queries can be answered by just maintaining the ~~data base~~ summary data obtain by aggregation on a selection, rather than maintaining the entire relation.

Summary data is classified into two categories

(i) Light summarized data \Rightarrow obtained by aggregating detailed data. This represent data distilled from current detailed data. It is summarized according to some unit of time and always resides on disk.

(ii) Highly summarized data \Rightarrow This represent data distilled from slightly summarized data. It is more compact and easily accessible and resides on disk.

(iii) operational data-store \Rightarrow operational data-base are source data for the data-base. operational data-store is repository of

operational data

(iii) Integration / transformation system → The integration and transformation program convert the operational data that is application specific into enterprise data. The major function performed by these programs are follows

- Reformating, re-evaluation or changing key structure
- Adding time element
- Default value identification
- Providing logic to choose b/w multiple data store

(iv) Detailed data → Detailed is of two types older detail data or current detail data. The older detail data represent data whose not very recent may as old as ten year or longer.

The current detailed data represent data of a present nature and always has shorter time horizon than older detail data. It can be voluminous.

(v) meta data ⇒ meta data is the data about data. meta data for data warehouse users

Date _____
DELTA P.O. No. _____

are part of the data warehouse itself and controls access and analysis of data warehouse contents. The meta data repository is the key data warehouse component. It contains both technical and business meta data. The technical data contains details about acquisition, processing, storage structure, data description, warehouse operation and maintenance functionality. The business meta data cover the relevant business rule and organizational details.

(vii) Archives → These contain old or historical data of significant interest and have value to enterprise. It is generally used for forecasting and trend analysis; thus these archives store old and meta data that describe the characteristics of old data.

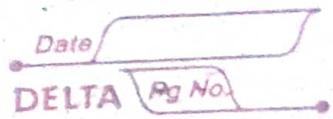
Data Mart

A data mart is a subset of a data warehouse that is designed to serve the specific data and analytic needs of a particular business unit, department, or group within an organization.

Data mart are smaller, more focused, and more specialized data storage or retrieval system compared to the larger and more comprehensive data warehouse.

Characteristic of Data mart

- focused data → Data mart contain data that is highly relevant to a particular business or user group. This data is typically a subset of data found in the organization's data warehouse.
- Department or business unit specific → Data marts are often created to cater needs of a specific department or business unit within an organization such as marketing, finance, sales or human resource.
- Data consolidation → In some case, data marts may store data from multiple source but only related to particular business area. Consolidation make it easier for user in that area to access data they need without the complexity of dealing with data from other part of organization.

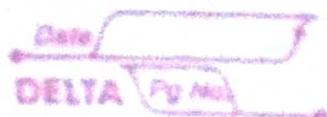


- Optimized query performance \Rightarrow Data marts are designed with query performance in mind. They are structured and indexed to facilitate quick retrieval of data, making them ideal for reporting and analysis.
- User friendly \Rightarrow Data marts are typically user friendly and tailored to the needs of specific user groups.
- Scalability \Rightarrow As data needs grow within a department or business unit, the data mart can be expanded and scaled accordingly.

A 3-tier Data warehouse architecture

There are mainly 3 types of Data warehouse architecture

- Single-tier architecture \Rightarrow The objective of a single tier architecture is to minimize the amount of data stored. This goal is to remove data redundancy. This architecture is not frequently used in practice.



• Two-tier architecture

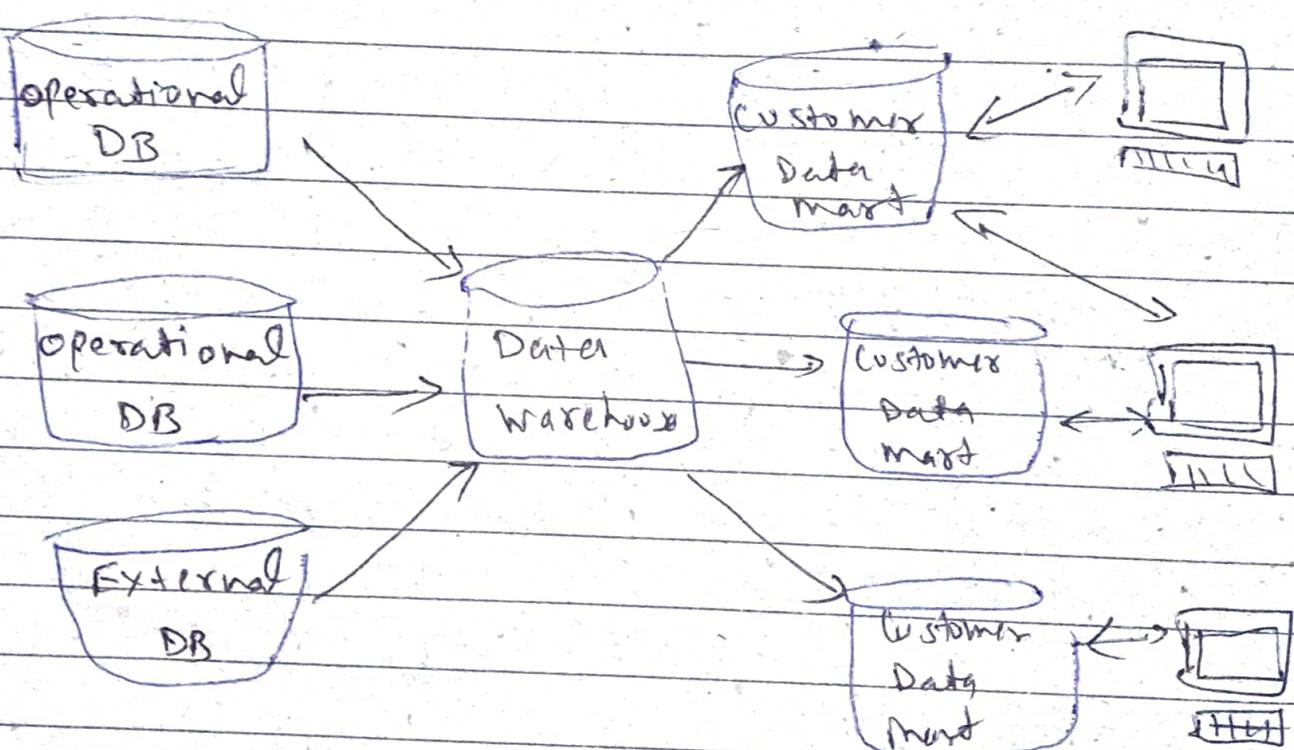
two-layer architecture separates physically available sources and data-warehouse. This architecture is not expandable and also not supporting a large number of end users. It also has connectivity problem because of network limitation.

• Three-tier architecture

This is most widely used architecture, it consists 3 layer top, middle, bottom

- ⇒ Bottom tier → The database of the Datawarehouse servers are the bottom tier.
- ⇒ is usually a relational database system. Data is cleaned, transformed, and loaded into this layer using back-end tools
- ⇒ Middle tier → The middle tier in Datawarehouse is an OLAP server which is implemented using either ROLAP or MOLAP model. For a user, this application tier present an abstracted view of database. This layer also acts as a mediator b/w the end-user and database.

③ Top tier → The Top tier is a front end client layer. Top tier is the tool and API that you connect and get data out from the data warehouses. It could be query tools, reporting tools, managed query tools, Analysis tools and data mining tools.



FNP
DDFP

Data Pre-Processing

Date _____
DELTA Pg No. _____

Data cleaning

Data cleaning is also known as data scrubbing. Data cleaning is a process which ensures the set of data is correct and accurate. Data accuracy and consistency, data integration is checked during data cleaning. Data cleaning can be applied for a set of records or multiple set of data which need to be merged.

Data cleaning is performed by reading all records in a set and verifying their accuracy. Typos and spelling errors are rectified. Mislabelled data if available is labelled and filed, missing or incomplete entry are completed. Unrecoverable records are purged, for not to take space and inefficient operations.

Data transformation

In data transformation process, data are transformed from one format to another format, that is more appropriate for data mining.

Some Data transformation strategies

- 1) Smoothing \Rightarrow It is the process of removing noise from Data
- 2) Aggregation \Rightarrow It is the process where summing or aggregation operations are applied to the Data
- 3) Generalization \Rightarrow In generalization low-lvl data are replaced with high-lvl data by using concept hierarchies climbing
- 4) Normalization \Rightarrow normalization scaled attribute data so as to fall within a small specified range such as 0.0 to 1.0
- 5) Attribute Construction \Rightarrow In attribute construction new attribute are constructed from given set of attribute

ETL (extract, Transform, Load)

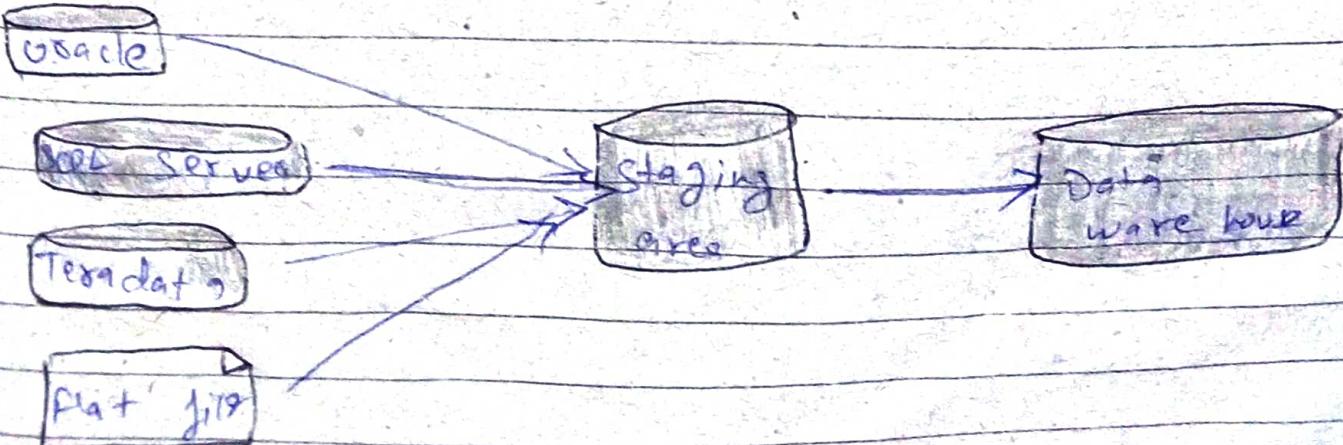
ETL is the process that defined as a process that extract the data from different RDBMS

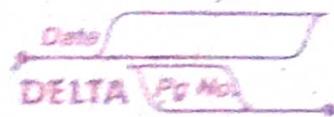
source system, then transform the data (like calculation, concatenation, etc.) and finally load the data into Data warehouse system.

It's tempting to ~~think~~ think of creating a Data warehouse by simply extracting data from multiple source and loading into Database or Datawarehouse. This is ~~good~~ from the technical and require a complex ETL process. The ETL process require active input from various stakeholder including developer, analyst, tester etc..

ETL Process

ETL is a 3 step process





Step 1 Extraction

In this step data is extracted from source system into staging area. Transformation if ~~any~~ done in staging area so that performance of source system is not degraded. Also if corrupted data is copied directly from source into data-warehouse database, rollback will be challenge. Staging area give option to validate extracted data before it move to Data warehouse.

(3 Data extraction method)

- 1) Full extraction
- 2) Partial // without update notifyability
- 3) ————— with —————

Some validation are done during extraction

- Reconcile records with source data
- make sure that no spam/unwanted data loaded
- Data type check
- Remove all type of duplicate and fragmented data
- check weather all key are placed or not

Step 2 Transformation

Data extracted from source server is raw and not usable in its original form. Therefore it

needs to be cleansed, mapped and transformed. In fact this is the key step where ETL adds value and change data such as that insightful BI reports can be generated.

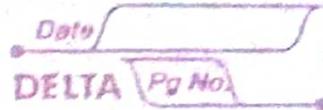
In this step, a set of function applied on extracted data. Data that does not require any transformation is called as direct move or pass through data.

Following are Data Integrity Problem

- 1) Different spelling of same person like John, Joun, Jon.
- 2) There are multiple way to denote company name like Cragle, Engle Inc.
- 3) In some Data required file was blank

Validation done during this stage,

- filtering - select only certain columns to load
- using rules and look-up table for Data standardization
- character set conversion
- Encoding handling
- Date threshold validation check



- Transposing Row and columns
- use look-up to merge data

Step 3] Loading

Loading data into target data warehouse database is the last step of ETL. In a typical Data warehouse, huge ~~increase~~ ~~in~~ volume of data need to be loaded in a ~~relatively~~ relatively short period. Hence load process should be optimized for performance.

In case of load failure, recovery mechanism should be configured to restart from the point of failure without data integrity loss. Data warehouse admin ~~need~~ ~~at~~ need to monitor, resume, cancel, load as per prevailing server performance.

TYPES of loading

- Initial load → Populating all the Data warehouse table
- Incremental load → applying ongoing change as we need Periodically
- full fresh → erasing the content of one more table and re-loading with fresh data

Load verification

- Ensure key field is never missing or null
- Test modeling view based on target table
- check combined value and calculated measure
- Data checks in dimension table as well as history table
- check BI report on the loaded fact and dimension table

ETL Tools

There are many ETL tools available in market
here are some most prominent one

marklogic

marklogic is data warehousing solution which make data integration is easier and faster using ~~any~~
of enterprise feature. It can query different type data like document, relationship and metadata

oracle

oracle is the industry leading database. It offer wide range of choice of Datawarehouse solution for both on-premises ~~in~~ and in cloud.

Amazon redshift

amazon redshift is Datawarehouse tool, it is

Date _____
DELTA Pg No. _____

Simple and cost effective to analyze all type
of data using Standard SQL and existing BI
tools. It allows running complex
queries against Petabytes of structured data

Defining the Business Requirements

* Dimensional Analysis

• capability to view their information over time

so the concept dimensional analysis became a method for defining data warehouse

Dimensional analysis is a technique used in data mining to identify patterns and relationship in data by analyzing the dimension of data. This is done by representing the data in a multidimensional space, where each dimension represent a different attribute of the data. The dimensions are then analyzed to identify patterns and relationships that may not be apparent from the data in its original form.

Dimensional data can be ~~explore~~ used to explore data from a variety of sources, including

- Transactional data → This type of data record events that occur over time, such as sale transaction, customer interaction and website visit
- Sensor data → This type of data is collected ~~by~~ from sensor that measure physical phenomena such as temperature, pressure, and humidity
- Text data → This type of data consist of written text, such as email, document and social media post

Dimensional analysis can be used to perform a variety of task including,

- clustering ⇒ This task involve grouping similar data point together
- classification ⇒ This task involve assigning data point to a predefined set of

Categories

- Anomaly detection \Rightarrow This task involve identifying ~~assigning~~ data points that are outliers or deviate from norm

* Multidimensional analysis

It is a data analysis ~~tech~~ technique that group data into multiple dimension, allowing for the exploration of complex relationships and pattern within data. It is often used in business intelligence (BI) and data mining application to gain insights into customer behaviors, sales, trending, market dynamics, and operational performance.

key characteristic of multidimensional analysis

- 1) Data organization. \Rightarrow MDA organize data into a multiple a multiple dimension array, where each dimension represent distinct attribute or category of the data. This structure allows for efficient analysis

- 2) Hierarchies → Dimensions can be organized into hierarchical relationship, reflecting Parent-child relationship b/w different lvl of granularity
 - This enable users to drill down or roll up through cluster

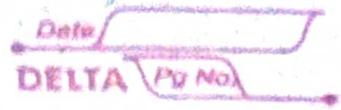
- 3) Measure → Numeric value, associated with data points are called measure. These measure represent the quantitative aspects of the data and are used to analyze trends, per羨 or more and relationship.

- 4) visualization techniques → multidimensional data is often visualized using techniques like Pivot table, charts and graph.

Application of MDA

used in:

- sales analysis
- customer segmentation
- Risk assessment
- operational efficiency



* Information Package

Information Packages are structured representation of specific business requirement or information needs in data mining. They serve as a bridge b/w abstracted business goals and concrete data mining tasks, ensuring that data mining efforts are aligned with overall business objectives.

Information Package are typically include a business question, business requirement, data source, data mining task and expected outcome. They offer several benefit, including enhanced clarity, improved communication, focused data mining efforts, and measurable outcomes.

* Requirement gathering

Good requirement start with good source. finding those quality source is an important task and fortunately done that takes few resources.

Example

- Customer
- Administrator
- Domain expert
- User
- Partner
- Industry analyst

Requirement gathering technique

After getting and identified these sources, there are number of technique that may be used to gather requirement. The following will describe various technique.

To go the requirement following technique are used

- Conduct a brainstorming session
- Interview user
- Send questionaries
- Work in target environment
- Examine suggestion and problem report
- Talk to support team.
- Conduct workshop
- Demonstrate prototype to stakeholders

UNIT 2

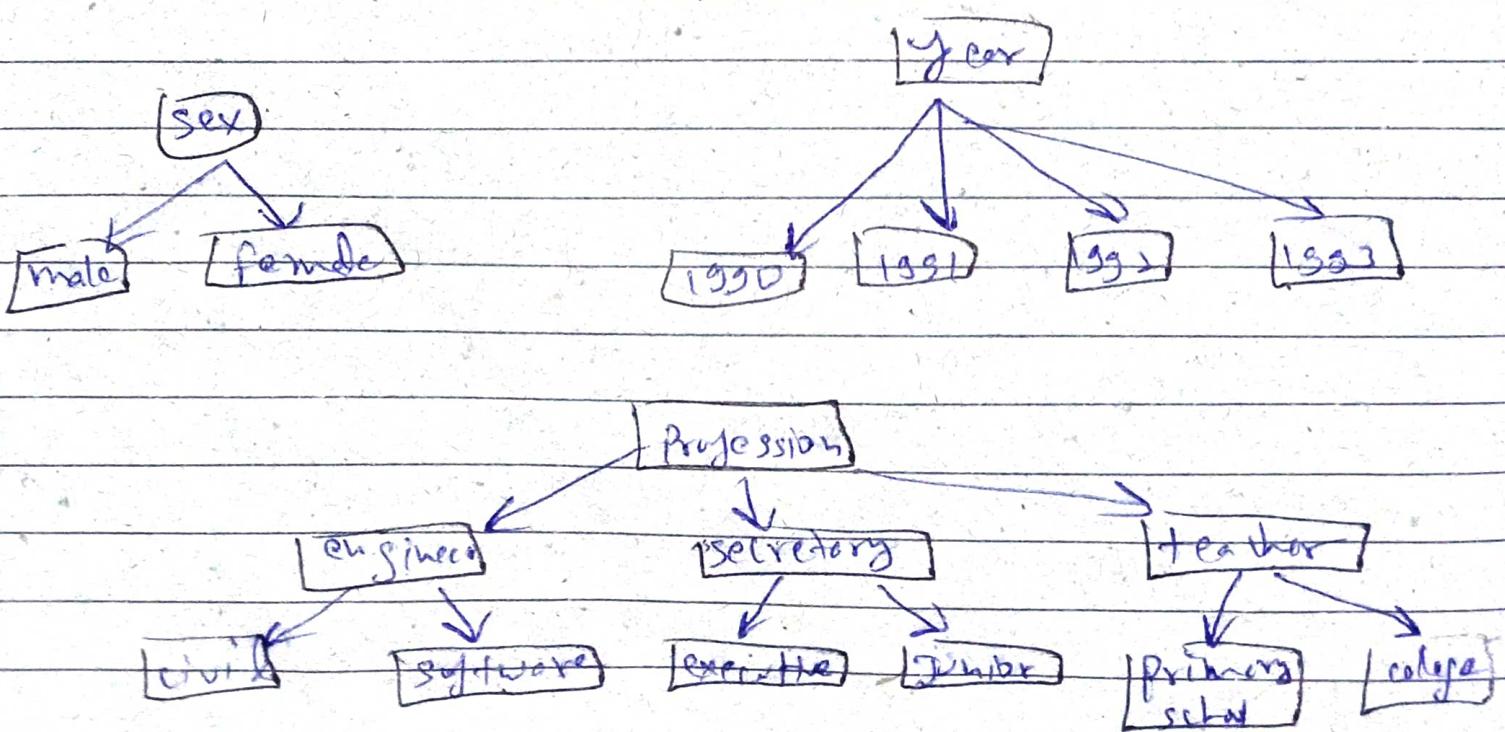
Date _____
DELTA Pg No. _____

Dimension modeling

The notion of a dimension provide a lot of semantic information, especially about the hierarchical relationship b/w its elements. It is important to note that dimension modeling is a special technique for structuring data around business concepts. Unlike FR modeling, which describe entities and relationships, dimension modeling structures the numeric measure and the dimension. The dimension schema can represent the details eg the dimensional modeling.

The following figure show the dimension modeling.

Diagram



The Dimension hierarchy helps us view the multi-dimensional data in several different data cube representation.

Steps to create dimension modeling

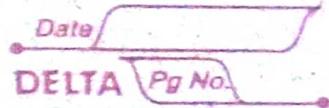
- Step 1 \Rightarrow Identifying the business objective

The first step is to identifying the business objective. Sales, HR, Marketing, etc are some examples as per need of the organization. Since it is the most important step of Data modeling, the selection of business objective also depends on quality of data available for that process.

- Step 2 \Rightarrow Identifying granularity.

Granularity is the lowest level of information stored in the table. The label of detail for business problem and its solution is described by grain.

- Step 3 \Rightarrow Identifying Dimension and its attribute. Dimension are object of things. Dimension



categorize and describe data warehouse facts and measure in a way that support meaningful answer to business question. A data-warehouse organizes descriptive attribute as column in dimensional table for example the Date dimension may contain date like a year, a month, a date

Step 4

- Identifying fact building Schema

The measurable data is held by fact table. Most of the fact table rows are numerical value like price or cost etc.

- Step 5 Building of Schema

We implement the Dimension model in this step. A schema is a database structure. There are two popular schemas

→ Star Schema

→ Snow-flake schema

Identifying Business objective

↓
Identify granularity

↓
Identify Dimension and attribute

↓
Identify fact

↓
Build Schema

from Requirement to Data design effects

Good "Business Intelligence (BI)", allows your organization to query data obtained from trusted sources and use it gain a competitive edge in your industry. The first step to achieving effective BI is a well designed warehouse.

- Requirement Gathering → Identifying long-term business strategy alongside current and technical requirement

- Analyze long-term user needs, reporting to management, and plan for hardware development, testing, implementation, and user training
- Develop a disaster recovery plan for quick response to threats
- Physical Environment Setup → Establish separate physical environment for development, testing, and production to ensure smooth transitions and data integrity
- Data modeling → Define how data structure will be accessed, connected, processed and stored
- Identify data sources, ensuring knowledge of data locations and availability for full project execution
- ETL (extract, transform, load) → identify data source during data modeling to reduce ETL development time.

- prioritize optimized load speeds without sacrificing quality to prevent poor performance
- OLAP cube design → obtain specification ~~specification~~ and measure during the design process
- optimize the OLAP cube generation partly to prevent performance issues after the data warehouse goes live
- From End development → Ensure secure access to data from any device
- choose a tool that allows backend modification and provides a user-friendly GUI for report customization

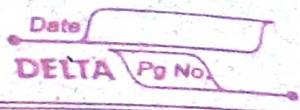
multidimensional Date model

multi-dimensional date model stores data in the date cube. mostly, data warehousing supports two or three-dimensional cubes.

A date cube allows data to be viewed in multiple dimensions. Dimension ~~are~~ entities with respect to which an organization wants to keep records. for example. In store sales record, dimensions allows the store to keep track things like monthly sales of items and the branches and location.

A multidimensional database helps to provide data-related answer to complex business queries quickly and accurately.

Data warehouse and online analytic processing (OLAP) tools are based on a multidimensional data model. OLAP in data warehousing enable users to view data from different angles and dimension



What is multidimensional schemes?

Multidimensional scheme is especially designed to model data warehouse systems. The schemes are designed to address the unique needs of very large databases designed for the analytical purpose (OLAP).

Types of Data warehouse Scheme

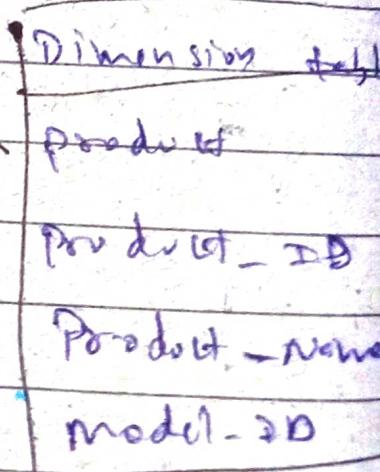
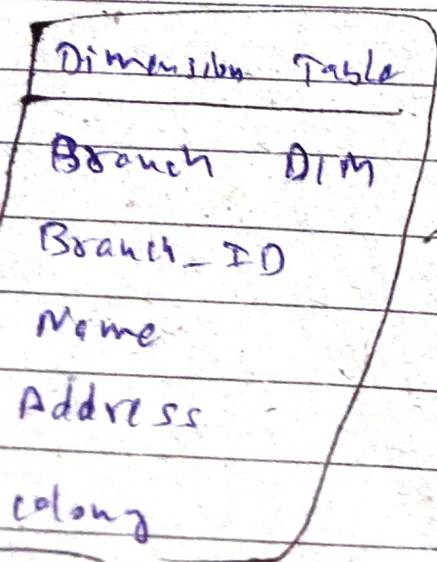
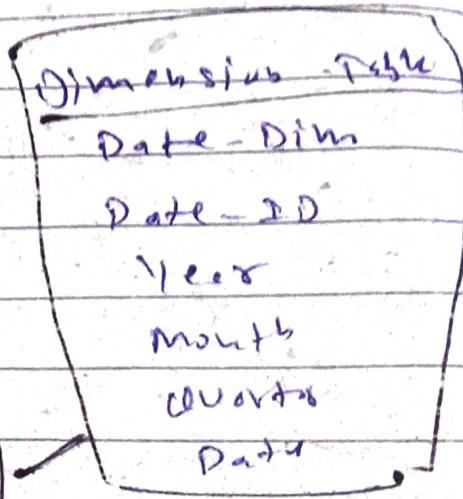
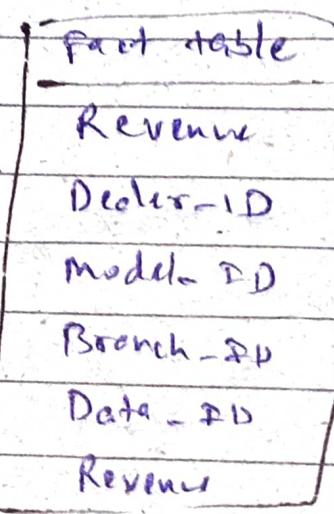
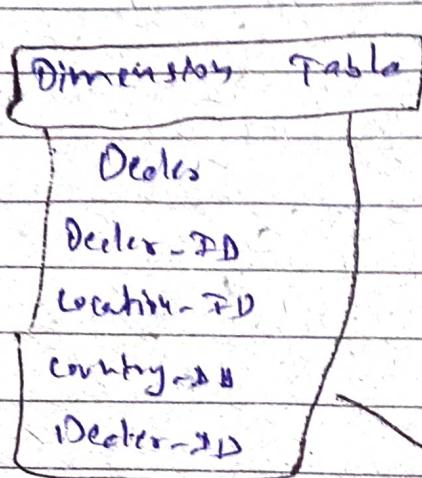
- Star Schema
- Snow-flake Schema
- Galaxy Schema

Star Schema

The Star Schema is the simplest type of Data warehouse scheme. It is known as Star Schema as its structure resembles a star. In the Star schema, the center of the

DATA
DELTA Po No.

Star can have one fact tables and
and number of associated dimension
tables. It is also known as Star
Join schema and is optimized for
querying large data sets.



for example, as you can see in the above - given image the fact table is at the center which contains key to every dimension table like Deal-ID, Model-ID, Product-ID, Branch-ID and other attributes like units sold and revenue.

Characteristics of star schema

- Every dimension in a star schema is represented with the only one-dimensional table.
- The dimension table should contain the set of attribute.
- The dimension table is joined to the fact table using a foreign key.
- The dimension table are not joined to the fact table using a foreign key.
- The dimension table are not ~~either~~

Date/
DELTA
PO NO.

Joined to each other,

- Fact table would contain key and measure

- The schema is widely supported by BI tools

• Snowflake Schema

We notice that star schema consist of a single fact table and a single denormalized dimension table for each dimension of the multi-

dimensional data model, to support attribute hierarchy, the dimension table can be normalized to create snowflake schema. As shown

fact schema consist of single fact table and multiple dimension table. Like star schema, each tuple of fact table consist

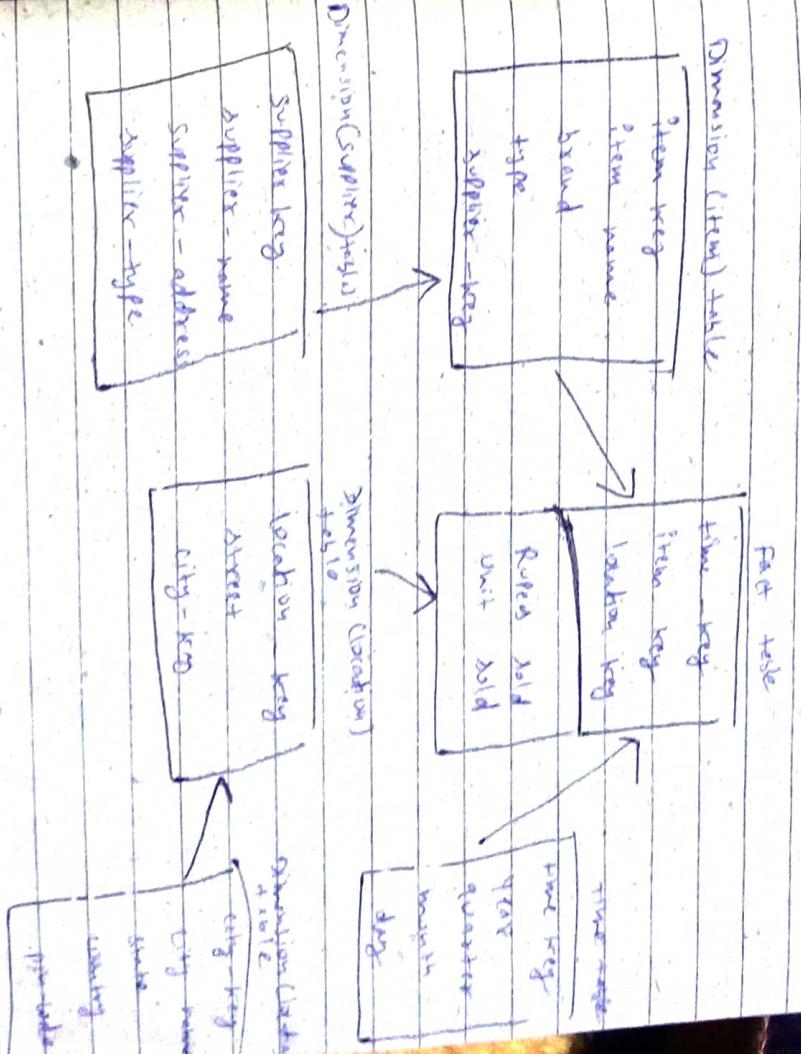
of a (foreign key) pointing to each of

the dimension table that provide it's

multidimensional co-ordinates. It also stores numerical values (non-dimensional attribute) for those co-ordinate. Dimension table

- fact constellation

most often, there may be a need to have more than one fact table and these are called fact constellation.

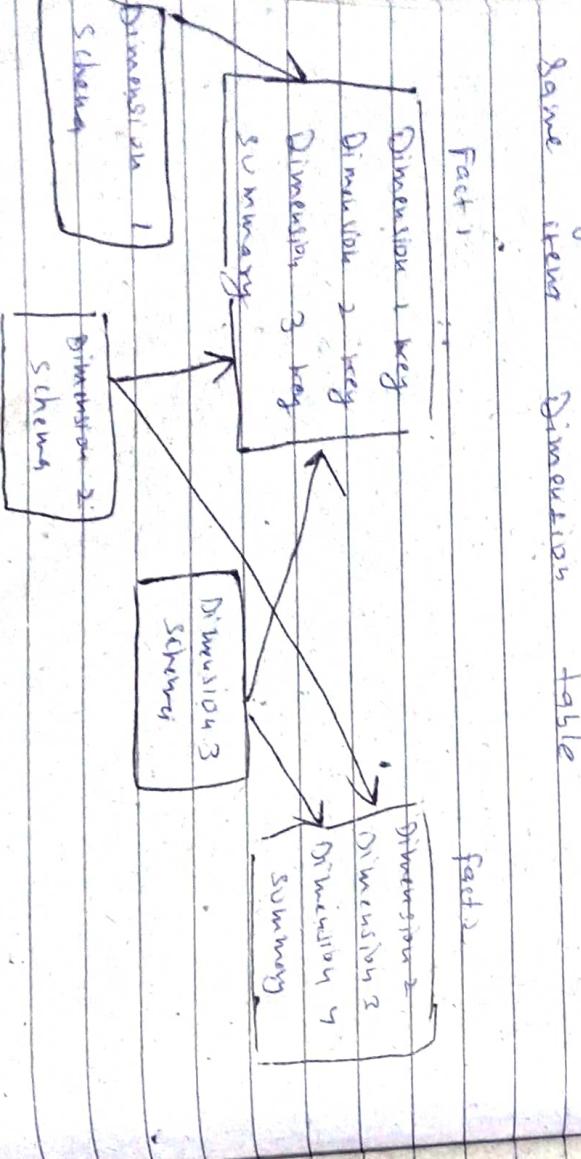


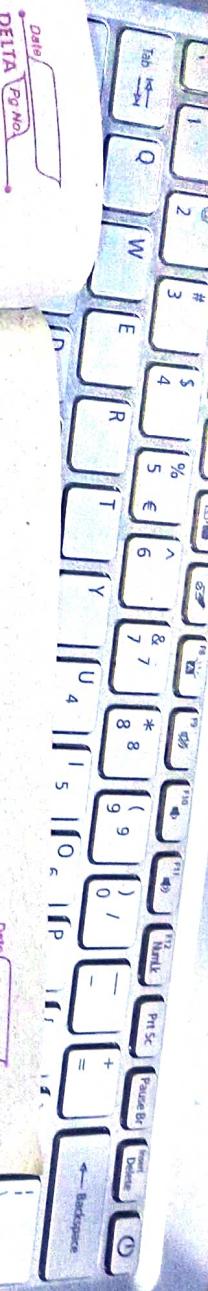


- A fact constillation is a kind of dimension where we have more than one fact table sharing among them. Some dimension tables are also called dimension tables. It is also called fact delivery system. For example let's assume what DCCW Electronic would like to have another fact table for supply and delivery. It may contain five dimensions or key, time, item, delivery-agent, origin, destination along with the numeric measure as the number of unit supplied and those lost by delivery. It can be seen that both fact table can share the same items Dimension table.

Key characteristic and operation

- a) Multidimensional Data model \Rightarrow One system organizes data into multidimensional cubes. Each dimension represent a specific aspects of data, and the intersection of dimensions contain aggregated data point.
- b) Cube and measure \Rightarrow The core composed of OLAP cube. Data cube. Measure represent the numeric data being aggregated and dimension categorized and organized.





These measure

- 3) drilling, drill and pivoting \Rightarrow OLAP allows user to "slice" and "dice" through the data cubes to view subset of information. Allowing "pivot" selecting a single dimension to view a single layer of cube, while drilling involve selecting specific elements from multiple dimensions, pivoting involves rotating the cube to view it from a different perspective.

- 2) ROLAP (Relational OLAP)
 - ROLAP system stores data in a multidimensional cube format. This means that data is pre-aggregated and organized into a structure where each cell in the cube represents the unique intersection of various dimensions. And data stored in this representation is very efficient in data retrieval and its fast and response to user queries very well and it will handle large data.

- 2) Drill down and Roll up \Rightarrow User can drill down into more detailed data or roll up to view higher level summary. This flexibility allows for detailed analysis at a granular level or broader overview.

- 2) Drill down and Roll up \Rightarrow User can drill down into more detailed data or roll up to view higher level summary. This flexibility allows for detailed analysis at a granular level or broader overview.

Solving modern business problems such as market analysis and financial forecasting requires query centric database schema

Date _____
DELTA Pg No. _____

Date _____
DELTA Pg No. _____

1) Complex business reporting

What are array - oriented and multidimensional
in nature. These business problem are characterized by the need to retrieve large
number of records from very large data sets (hundred of gigabytes) and even
multi-dimensional nature of the problem. It is designed to address it
live key driver of OLAP

2) Fast query response

OLAP system in particular store pre-aggregated data in a tree format leading to faster query response times. This responsiveness is crucial for interactive and ad-hoc querying.

3) User friendly interface

OLAP systems provide a user friendly interface with feature like drag and drop, point and click, and visualisation tool working it easier for non technical user.

4) Flexibility

OLAP facilitates complex reporting by allowing user to slice, dice, drill down and pivot through data cubes enabling the creation of detailed and comprehensive reports.

5) Strategic planning

Enter

it

Pg No. _____

Delta

Date _____
DELTA PO NO. _____

Date _____
DELTA PO NO. _____

OLAP features

Alerts allow for even more complex
and detailed analysis

• Multidimensional conceptual view

• Transparency

• Reusability

• Consistent reporting performance

• Great server architecture

• Granule dimension

• Sparse multidim. handling

• Multi-user support

• Flexible reporting interface

• Scalable

key parts

• Dimension → Each axis in a hypercube
represent a different dimension

• Cells → The intersection of dimensions

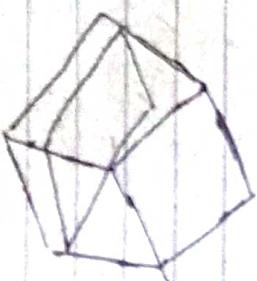
• Hypercube contains data points or cell
each cells represent a unique combination
of values from dimension

• Analysis → Hypercube enable users to ana-

lyze data across multiple dimensions simultaneously

Hypercube

It is refer to multidimensional cube in
with more than 3 dimension. In
traditional three - dimensional OLAP cubes
data is organized along three axes
representing different dimensions. A hyper-
cube extends this concept to include
additional dimension, creating a structure



OLAP server / engine

An OLAP server or OLAP engine is a soft component responsible for managing and facilitating online analytical processing operation. It serves as the backbone system that handles the storage, retrieval, and processing of multidimensional data to support analytic queries. The OLAP server plays crucial role in providing a fast and efficient environment for user to interact with and analyze large dataset.

rolling up

② Integration with BI tools → OLAP Server one

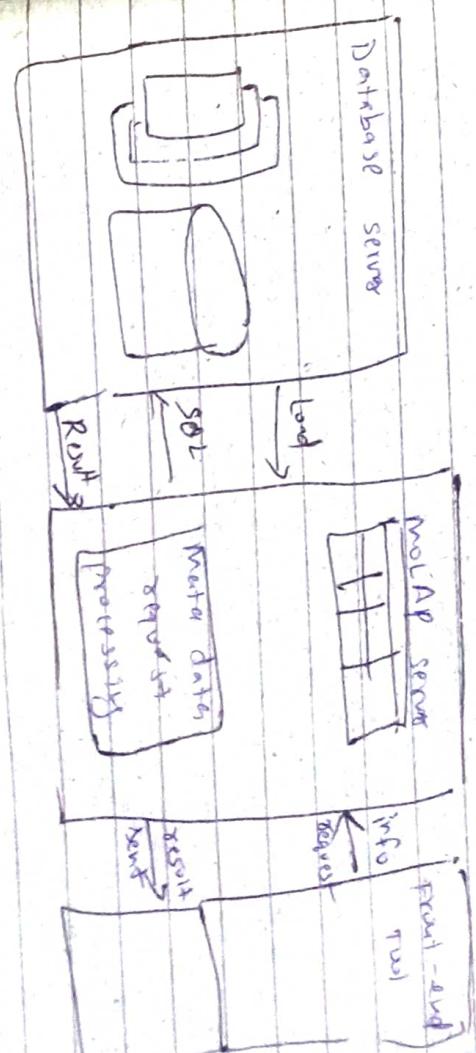
Commonly integrate with Business intelligence tool like BI, provide user interface to interact with OLAP server.

③ Scalability → OLAP servers need to scale efficiently to handle increasing volume of data and user queries

MOLAP architecture

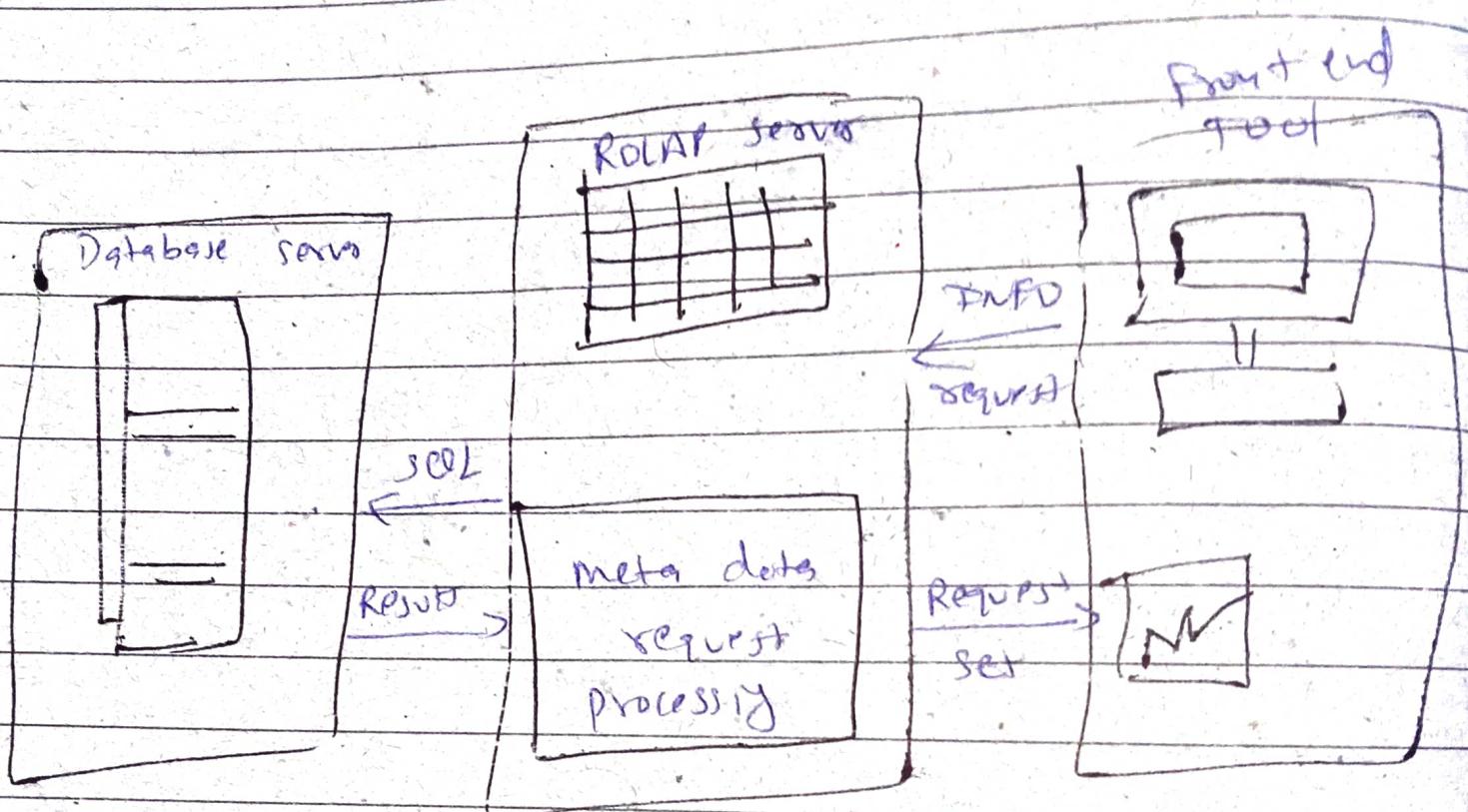
④ Multidimensional storage → The OLAP server store data in a multi-dimensional format, typically organized in cube.

⑤ Data aggregation → OLAP server pre-aggregate data at different level of granularity to optimize query performance



⑥ Query Processing → OLAP servers process complex query such as slicing, dicing, drill down and

ROLAP architecture

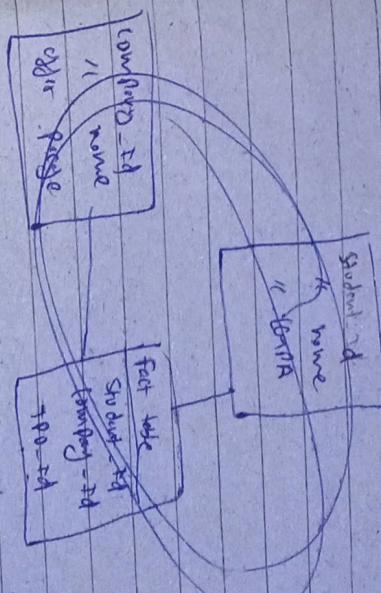


Star Schema vs Snowflake

Example of Inflex or fact constellation

Star

- Fact table and dimension fact table and dimension tables are contained within fact table and dimension are contained
- Star schema is simpler - it is more difficult to understand.
- Star schema uses less space
- Space
- It takes less time for execution of query
- Its design is very simple but design is complex.
- It is simple to understand and stand alone
- Does number of foreign key number of keys with data redundancy



Dimension table

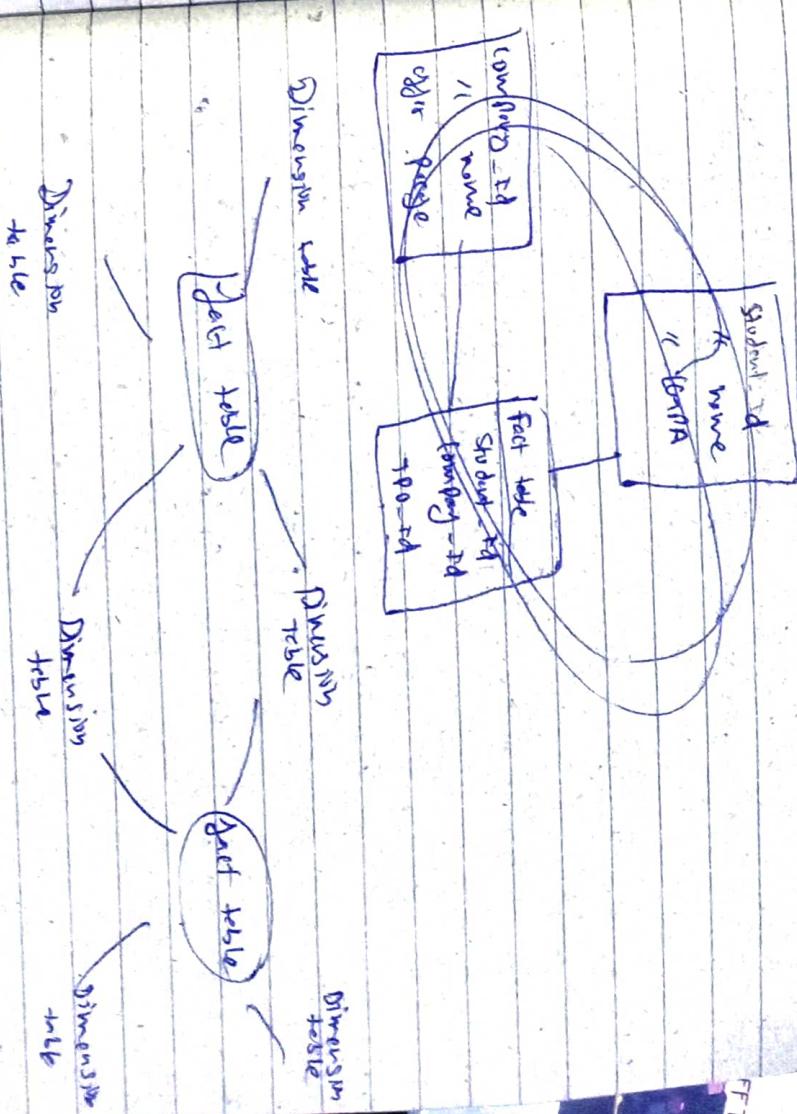
Date _____
DELTA Pg No. _____

Date _____
DELTA Pg No. _____

Star Schema vs Snowflake

Example of Olap or fact constellation

- fact table and dimension fact table and dimension tables are contained within dimension are contained
- star schema is top-down model
- use less space
- star schema use more space
- less time for execution of query
- design is very simple
- it is difficult to understand
- hand
- less number of foreign key
- high data redundancy



MOLAP vs ROLAP

MOLAP

- ROLAP
- full form
- used for limited volume
- access is slow
- success of ROLAP slow

- Data store in relational database
- Data store in multidimensional array
- Table

- Data is fetched from MOLAP database
- sparse matrix is used
- complicated SQL queries are used

Delta /
DELTA Pg No.

Date /
DELTA Pg No.

ff EIS

An executive Information System (EIS) is a decision support system (DSS) it used to assist senior executive in the decision making process. It does this by providing easy access to important data needed to achieve strategic goal in an organization. An EIS normally features graphically display on an easy to use interface.

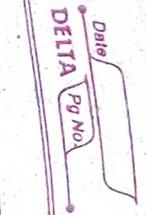
EIS can be used many different types of organization to monitor enterprise performance as well as to identify opportunities and problems.

~~#~~ Codds rule for OLAP

In 1993, Dr. E.F. Codd's originate 12 rule as the basis for selecting OLAP tool.

These 12 rules are :

- 1) Multidimensional conceptual view of database
- 2) Concept of transparency
- 3) Concept of accessibility
- 4) Consistent reporting performance
- 5) Client server architecture
- 6) Generic dimensionality
- 7) Multi-user support
- 8) Dynamic matrix sparse matrix handling
- 9) Unrestricted cross dimensional operation
- 10) Intuitive data manipulation
- 11) Flexible reporting
- 12) unlimited dimension and aggregation



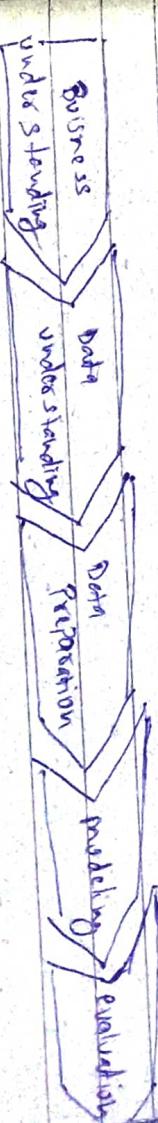
Types of DATA

Data mining can be performed on following types of data.

- Relational database
- Data warehouse
- Advance DB and information repositories
- Object - oriented and object - relational databases
- Transactional and spatial database
- Heterogeneous and legacy database
- Multimedia and streaming database
- Text database
- Text mining and web mining

The insight derived via Data-mining can be used for marketing, fraud detection and scientific discovery, etc.

Data mining Implementation process



Data mining is also called as knowledge discovery, knowledge extraction, data pattern analysis, information harvesting, etc.

Deployment

D) Business understanding and data-mining plan

In this phase, business and data-mining plan are established.

- First need to understand business and client objective. You need to define what your client want.

• Take stock of the current data-mining scenario, resource, assumption, constraints, factors and their significant factor into your analysis and assessment.

- Using business objective and current scenario, define your data mining goals

3) Data Preparation

- A good data mining plan is very detailed and should be developed to accomplish both business and data mining goals

⇒ Data understanding

- In this phase, sanity check on data is performed to check whether it's appropriate for data mining goals

- first data is collected from multiple data source available in the organization.
- These data sources may include multiple databases, flat either or data cubes.
- next step is to search for properties of acquired data. A good way to explain the data is to answer data mining question using query, reporting and visualization tools.
- Based on result of query, the data quality should be ascertained.
- The data preparation consume about 30% of the time of the project.
- The data from different source should be selected, cleaned, transformed, formatted, anonymize and constructed.
- Data cleaning is process to "clean" the

- Data by smoothing noisy data and filling data by missing values
- Data transformation operation would contribute to data transformation process.
- The success of mining process.
- Smoothing \Rightarrow Help to remove noise from data.
- Aggregation \Rightarrow Aggregation operation is applied to the data example \Rightarrow The weekly sales data is aggregated to calculate the monthly and yearly total.
- Normalization \Rightarrow In this step, low-level data is reported by high-level concepts with help of concept hierarchy.
- Standardization \Rightarrow Standardization is performed when the attribute scales are scaled up or scaled down.

④ Modeling

- In this phase, mathematical models are used to determine pattern.
- Based on business objective, suitable modeling techniques should be selected for prepared data set.
- Create a scenario to test check the quality and validity of model.

- Run the model on prepared dataset

⑤ Evaluation

- This phase, patterns are evaluated against the business objective.
- Result generated by the data mining model should be evaluated against business objective.
- Training business understanding is an iterative process.
- In fact, while understanding how business requirement may be raised because of data mining.

⑥ Deployment

\Rightarrow Plan your data-mining discoveries to everyday business operations.

- The knowledge or information discovered during data mining process should be made easy to understand for non-technical methods.
- A detailed deployment plan for skipping maintenance and monitoring of data mining discoveries is created.

⑦ Modeling

- In this phase, mathematical models are used to determine pattern.

- Based on business objective, suitable modeling techniques should be selected for prepared data set.

- Create a scenario to test check the quality and validity of model.

Data mining techniques

- ④ Classification \Rightarrow This analysis is used to retrieve important and relevant information about data, and measure. This data mining method help to classify data in different classes.
- ⑤ Clustering \Rightarrow Data analysis is a data mining technique to identify data that are like each other. This process help to understand the difference and similarities between the data.
- ⑥ Regression \Rightarrow Regression method is the data mining method of identifying and analyzing the relationship between two variable. It is used to identify the likelihood of a specific variable, given the presence of other variable.
- ⑦ Association rule \Rightarrow This data mining technique help to find the association b/w two or more items. It discover a hidden pattern in the data set.
- ⑧ Outlier detection technique \Rightarrow This type of technique refer to observation of data item in the data set which do not match an expected pattern or

Expected behavior, this technique can be used in various domain such as intrusion, detection, fraud or fault detection etc.

- ⑨ Sequential patterns \Rightarrow This data mining technique help to discover or identify similar patterns or trend in transaction data for certain period.

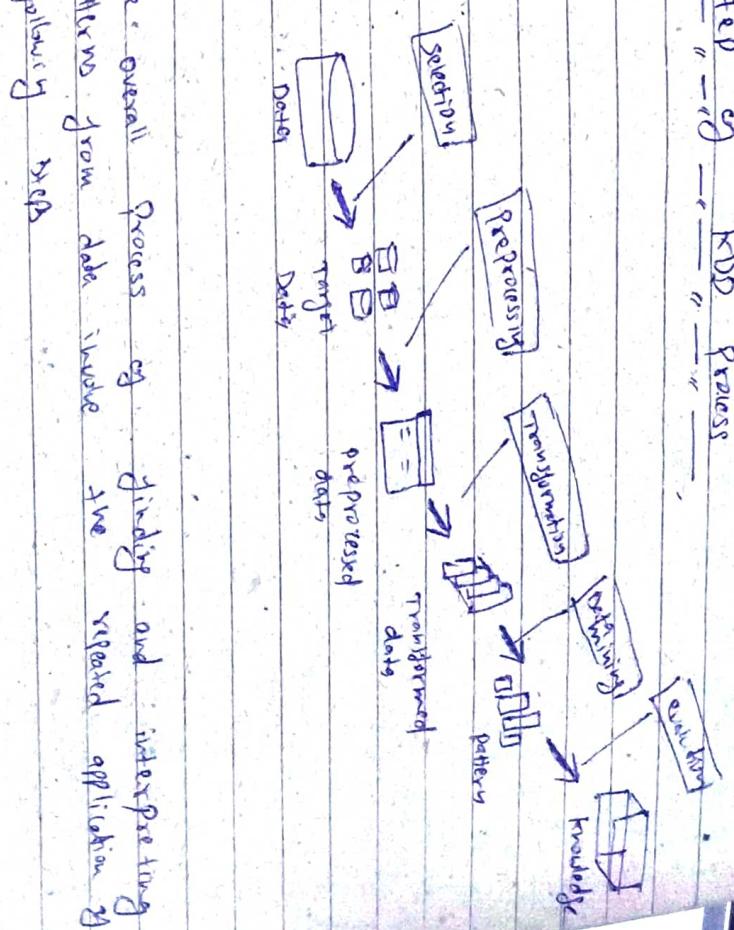
- ⑩ Prediction \Rightarrow It has used in combination of the other data mining technique like trend, sequential pattern, clustering, classification etc. It analyze past event or instance in a right sequence of predicting future event.

Challenges for implementation of Data mining

- Skilled experts are needed to formulate the data mining queries.
- Overfitting, due to small size of database, ~~and sometimes are different to fitting model may not fit in future~~ fitting model may not fit in future
- Data-mining need large database which sometime are difficult to manage.
- Business practice may need to be modified to determine to use the information uncovered.

- If the data set is not diverse, data mining res. will may not be accurate.

Step "n" - KDD Process



The term "knowledge discovery in database" (KDD) refer to the broad process of finding knowledge in data, and emphasizes the "high-level" application of particular data mining methods. It is of interest to researcher in machine learning, pattern recognition, database statistic, artificial intelligence, knowledge acquisition, expert system and data visualization.

The overall process of finding and interpreting patterns from data involve the repeated application of following step:

- 1) Identifying goal of KDD process is to extract knowledge from data in context of large database.
- 2) Doing this by using determining method to extract what is needed knowledge, all according to the specification of measure and threshold, using a database along with any required preprocessing, subsampling and fragmentation of database.

- 3) Data cleaning and Preprocessing.



- Removing of noise and outliers
- collecting necessary information to model or account for noise
- strategies for handling missing data field
- Accounting for time sequence information and known usage
- Applying useful feature to represent the data depending on goal of task
- using dimensionality reduction or data selection method to reduce the effective number of variables under consideration or to find invariant representation for data
- choosing data mining task
- Decide whether the goal of KDD process is classification, regression, clustering etc.
- choosing data mining algorithm
- selecting method to be used for searching patterns in data.
- Deciding which model and parameters may be appropriate.



- matching of particular data mining method with overall criteria of KDD process

7) Data mining

- Searching for pattern of interest in a particular representation of data in a set of such representation as classification rules or trees, regression, clustering and so on

- a) Just extracting mined patterns
- b) Consolidating discovered knowledge

The Business context of data mining

Why does an organization have to practise data mining when it does not bring impact to their business? To predict marketing, the marketing manager should identify the segment of the population who is most likely to respond to your product. Identifying these segments of population involve understanding the overall population and deploying the right techniques to classify the population. Likewise in predictive modeling, there



are several way to interact with customer using customer using different channel. These include direct marketing, print advertisement, telemarketing, radio, and so on. It is only through data mining, that an analyst would conclude which is optimal channel for sending communication for customer.

Process improvement through Data mining

Process improvement through data mining involve using data mining technique to discover pattern, trends, and insights within business processes. By identifying areas for improvement, organization can make better decision to enhance efficiency, reduce cost, and optimize performance.

How data mining can be applied to process improvement:

- 1) Hypothesis generation as researcher can use data mining to explore large dataset and generate hypothesis or research questions base on observed pattern and trend.
- 2) Data collection and integration as gather relevant data from various source to create a comprehensive dataset for analysis.

- ② Data integration \Rightarrow Integrate data from various resources and create a comprehensive data set.
- ③ Data mining \Rightarrow Identify and handle missing or erroneous data and manage data to ensure consistency.
- ④ Data mining as research tool
- ⑤ Data mining serve as powerful research tool across various domain enabling researcher to extract valuable patterns, trends and insights from large dataset.
- Here: How data mining can be applied as research tool
- ⑥ Hypothesis generation as researcher can use data mining to explore large dataset and generate hypothesis or research questions base on observed pattern and trend.
- ⑦ Data collection and integration as gather relevant data from various source to create a comprehensive dataset for analysis.



- Aggregate and preprocess the data to ensure high quality and consistency
- Regression analysis \Rightarrow employ regression analysis to model relationship b/w variable, allowing researchers to predict and understand the impact of one variable on another

- Data mining information can help to
 - \rightarrow increase return on investment
 - \rightarrow improve CRM and market analysis
 - \rightarrow reduce marketing costs
- There is multiple benefit by using data mining in marketing.

- Clustering \Rightarrow use clustering algorithm to group similar data points together, providing insights into the inherent structure of data.

- \rightarrow predict future trend
- \rightarrow understand customer purchase behavior
- \rightarrow help with decision making
- \rightarrow improve company revenue and lower cost
- \rightarrow market basket analysis

Data mining in marketing

Data mining play crucial role in marketing by helping businesses analyze large volume of data to extract valuable insights, patterns and trends.

- Data mining technique for marketing
- Knowledge base marketing

Data mining technology allows to better profile about whom customers and make smart marketing decisions.

- There are 3 major area of application

Date _____
DELTA Pg No. _____

Date _____
DELTA Pg No. _____

- The customer profiling system can analyse the customer purchase behaviour by data mining for knowledge based marketing.
- The frequency of purchase time the customer can know how many times he has visited the store by this product or visit the store.
- The deviation analysis give a marketer a good capability to query changes that occurred as a result of recent price changes or promotion.
- The trend analysis can determine the trend in sales, cost and profits by products or market in order to achieve the highest amount of sale.

- 1) **Market Basket Marketing**
 - The customer profiling companies can know how many times the customer buy this product or visit the store.
 - The deviation analysis give a marketer a good capability to query changes that occurred as a result of recent price changes or promotion.
 - The trend analysis can determine the trend in sales, cost and profits by products or market in order to achieve the highest amount of sale.
- 2) **Market Basket Marketing**
 - The customer profiling system can analyse the frequency of purchase time the customer can know how many times he has visited the store by this product or visit the store.
 - The deviation analysis give a marketer a good capability to query changes that occurred as a result of recent price changes or promotion.
 - The trend analysis can determine the trend in sales, cost and profits by products or market in order to achieve the highest amount of sale.

- Determine what customer purchase together and sales tactic using customers data already available to company.
- Improve the effectiveness of marketing and sales tactic using customers data.