# Data Mining and Warehousing

## Assignment:1

(1) Explain snowflake schema for multidimensional database?

(2) Discuss the concept of OLAP? Why ROLAP is used in the rational database environment?

(3 ) Briefly describe 3-tier architecture?

(4) Discuss the facts of star schema ?

## UNIT I: Data Warehousing and Business Analysis

- Define **metadata** in the context of data warehousing.

- What is OLAP? List two types.

- Explain the role of **data transformation tools** in building a data warehouse.

- Describe **DBMS schemas** suitable for decision support systems.

- Illustrate how you would **design a star schema** for a retail company.

- Apply **data extraction** and transformation steps for integrating customer and sales data.

- Compare and contrast **ROLAP** and **MOLAP**.

- Analyze the relationship between **data cleaning** and **report accuracy**.

- Evaluate the effectiveness of a **centralized vs distributed warehouse** architecture.

- Justify the use of **dimensional modeling** in OLAP applications.

- Design a **multidimensional data model** for a healthcare organization.

- Develop a **reporting framework** using OLAP tools for educational analytics.

## UNIT I: Data Warehousing and Business Analysis

Data Warehousing and Business Analysis: - Data warehousing Components –Building a Data warehouse –Data Warehouse Architecture – DBMS Schemas for Decision Support – Data Extraction, Cleanup, and Transformation Tools –Metadata – reporting – Query tools and Applications – Online Analytical Processing (OLAP) – OLAP and Multidimensional Data Analysis.

## UNIT II: Data Mining

Data Mining: - Data Mining Functionalities – Data Preprocessing – Data Cleaning – Data Integration and Transformation – Data Reduction – Data Discretization and Concept Hierarchy Generation-

Architecture Of A Typical Data Mining Systems- Classification Of Data Mining Systems.

**UNIT III: Classification and Prediction**

Classification and Prediction: - Issues Regarding Classification and Prediction – Classification by Decision Tree Introduction – Bayesian Classification – Rule Based Classification – Classification by Back propagation Support Vector Machines – Associative Classification – Lazy Learners – Other Classification Methods – Prediction – Accuracy and Error Measures – Evaluating the Accuracy of a Classifier or Predictor – Ensemble Methods – Model Section.

**UNIT IV: Cluster Analysis**

Cluster Analysis: - Types of Data in Cluster Analysis – A Categorization of Major Clustering Methods – Partitioning Methods – Hierarchical methods – Density-Based Methods – Grid-Based Methods – Model-Based Clustering Methods – Clustering High-Dimensional Data – Constraint-Based Cluster Analysis – Outlier Analysis.

**UNIT V: Mining Object**

Mining Object, Spatial, Multimedia, Text and Web Data: Multidimensional Analysis and Descriptive Mining of Complex Data Objects – Spatial Data Mining – Multimedia Data Mining – Text Mining – Mining the World Wide Web.

# **UNIT-I**

## **Introduction to Data Warehouse:**

A data warehouse is a subject-oriented, integrated, time-variant and non-volatile collection of data in support of management's decision making process.

**Subject-Oriented**: A data warehouse can be used to analyze a particular subject area. For example, "sales" can be a particular subject.

**Integrated**: A data warehouse integrates data from multiple data sources. For example, source A and source B may have different ways of identifying a product, but in a data warehouse, there will be only a single way of identifying a product.

**Time-Variant**: Historical data is kept in a data warehouse. For example, one can retrieve data from 3 months, 6 months, 12 months, or even older data from a data warehouse. This contrasts with a transactions system, where often only the most recent data is kept. For

example, a transaction system may hold the most recent address of a customer, where a data warehouse can hold all addresses associated with a customer.

**Non-volatile**: Once data is in the data warehouse, it will not change. So, historical data in a data warehouse should never be altered.
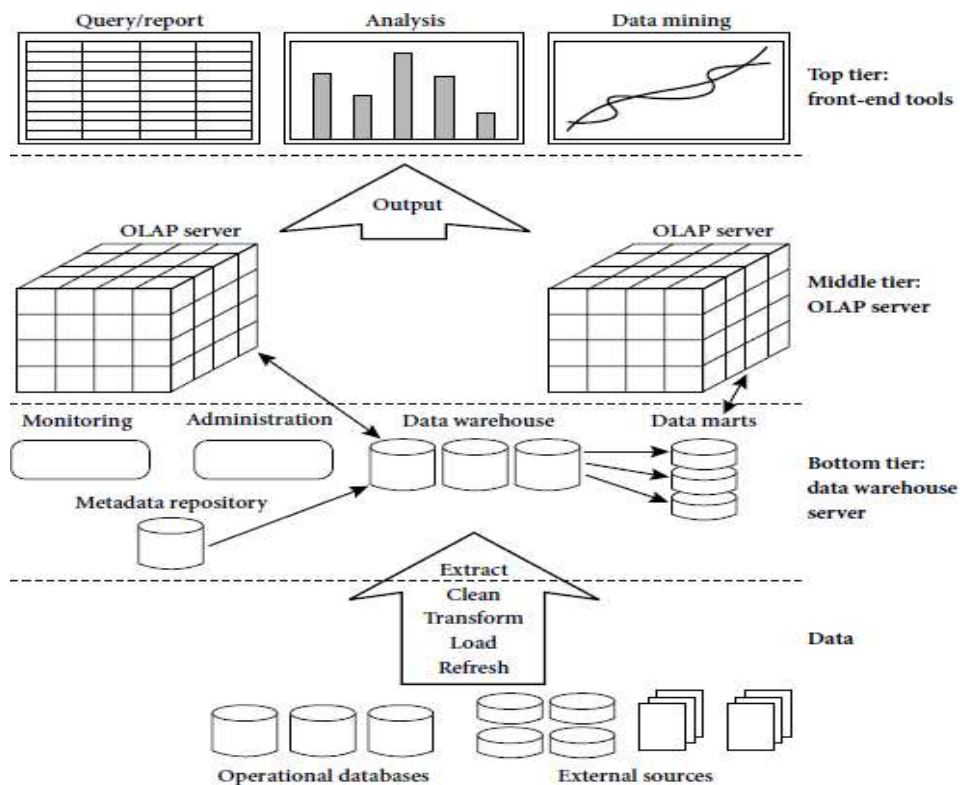
**Data Warehouse Design Process:**

A data warehouse can be built using a top-down approach, a bottom-up approach, or a combination of both.

- The top-down approach starts with the overall design and planning. It is useful in cases where the technology is mature and well known, and where the business problems that must be solved are clear and well understood.

- The bottom-up approach starts with experiments and prototypes. This is useful in the early stage of business modeling and technology development. It allows an organization to move forward at considerably less expense and to evaluate the benefits of the technology before making significant commitments.

- In the combined approach, an organization can exploit the planned and strategic nature of the top-down approach while retaining the rapid implementation and opportunistic application of the bottom-up approach.

The warehouse design process consists of the following steps:

- Choose a business process to model, for example, orders, invoices, shipments, inventory, account administration, sales, or the general ledger. If the business process is organizational and involves multiple complex object collections, a data warehouse model should be followed. However, if the process is departmental and focuses on the analysis of one kind of business process, a data mart model should be chosen.

- Choose the grain of the business process. The grain is the fundamental, atomic level of data to be represented in the fact table for this process, for example, individual transactions, individual daily snapshots, and so on.

- Choose the dimensions that will apply to each fact table record. Typical dimensions are time, item, customer, supplier, warehouse, transaction type, and status.

- Choose the measures that will populate each fact table record. Typical measures are numeric additive quantities like dollars sold and units sold.

# A Three Tier Data Warehouse Architecture:

## Tier-1:

The bottom tier is a warehouse database server that is almost always a relational database system. Back-end tools and utilities are used to feed data into the bottom tier from operational databases or other external sources (such as customer profile information provided by external consultants). These tools and utilities perform data extraction, cleaning, and transformation (e.g., to merge similar data from different sources into a unified format), as well as load and refresh functions to update the data warehouse. The data are extracted using application program interfaces known as gateways. A gateway is supported by the underlying DBMS and allows client programs to generate SQL code to be executed at a server. Examples of gateways include ODBC (Open Database Connection) and OLEDB (Open Linking and Embedding for Databases) by Microsoft and JDBC (Java Database Connection). This tier also contains a metadata repository, which stores information about the data warehouse and its contents.

**Tier-2:**

The middle tier is an OLAP server that is typically implemented using either a relational

OLAP (ROLAP) model or a multidimensional OLAP.

● OLAP model is an extended relational DBMS thatmaps operations on
multidimensional data to standard relational operations.

● A multidimensional OLAP (MOLAP) model, that is, a special-purpose server that
directly implements multidimensional data and operations.

**Tier-3:**

The top tier is a front-end client layer, which contains query and reporting tools, analysis

tools, and/or data mining tools (e.g., trend analysis, prediction, and so on).

# Data Warehouse Models:

There are three data warehouse models.

● **Enterprise warehouse:**

● An enterprise warehouse collects all of the information about subjects
spanning the entire organization.

● It provides corporate-wide data integration, usually from one or more
operational systems or external information providers, and is cross-functional in scope.

● It typically contains detailed data as well as summarized data, and can range in size
from

a few gigabytes to hundreds of gigabytes, terabytes, or beyond.

● An enterprise data warehouse may be implemented on traditional
mainframes, computer super servers, or parallel architecture platforms. It requires
extensive business modeling and may take years to design and build.

● **Data mart:**

● A data mart contains a subset of corporate-wide data that is of value to a
specific group of users. The scope is confined to specific selected subjects. For example, a
marketing data mart may confine its subjects to customer, item, and sales. The data
contained in data marts tend to be summarized.

● Data marts are usually implemented on low-cost departmental servers that are
UNIX/LINUX- or Windows-based. The implementation cycle of a data mart is more likely to
be measured in weeks rather than months or years. However, it may involve complex

integration in the long run if its design and planning were not enterprise-wide.

● Depending on the source of data, data marts can be categorized as independent more dependent. Independent data marts are sourced from data captured from one or more operational systems or external information providers, or from data generated locally within a particular department or geographic area. Dependent data marts are source directly from enterprise data warehouses.

● **Virtual warehouse:**

● A virtual warehouse is a set of views over operational databases. For efficient query processing, only some of the possible summary views may be materialized.

● A virtual warehouse is easy to build but requires excess capacity on operational database servers.

## Meta Data Repository:

Metadata are data about data. When used in a data warehouse, metadata are the data that define warehouse objects. Metadata are created for the data names and definitions of the given warehouse. Additional metadata are created and captured for time stamping any extracted data, the source of the extracted data, and missing fields that have been added by data cleaning or integration processes.

A metadata repository should contain the following:
● A description of the structure of the data warehouse, which includes the warehouse schema, view, dimensions, hierarchies, and derived data definitions, as well as data mart locations and contents.
●Operational metadata, which include data lineage (history of migrated data and the sequence of transformations applied to it), currency of data (active, archived, or purged), and monitoring information (warehouse usage statistics, error reports, and audit trails).

● The algorithms used for summarization, which include measure and dimension definition algorithms, data on granularity, partitions, subject areas, aggregation, summarization, and predefined queries and reports.
● The mapping from the operational environment to the data warehouse, which includes source databases and their contents, gateway descriptions, data partitions, data extraction, cleaning, transformation rules and defaults, data refresh and purging rules, and security (user authorization and access control).
● Data related to system performance, which include indices and profiles that

improve data access and retrieval performance, in addition to rules for the timing and scheduling of refresh, update, and replication cycles.

- Business metadata, which include business terms and definitions, data ownership information, and charging policies.
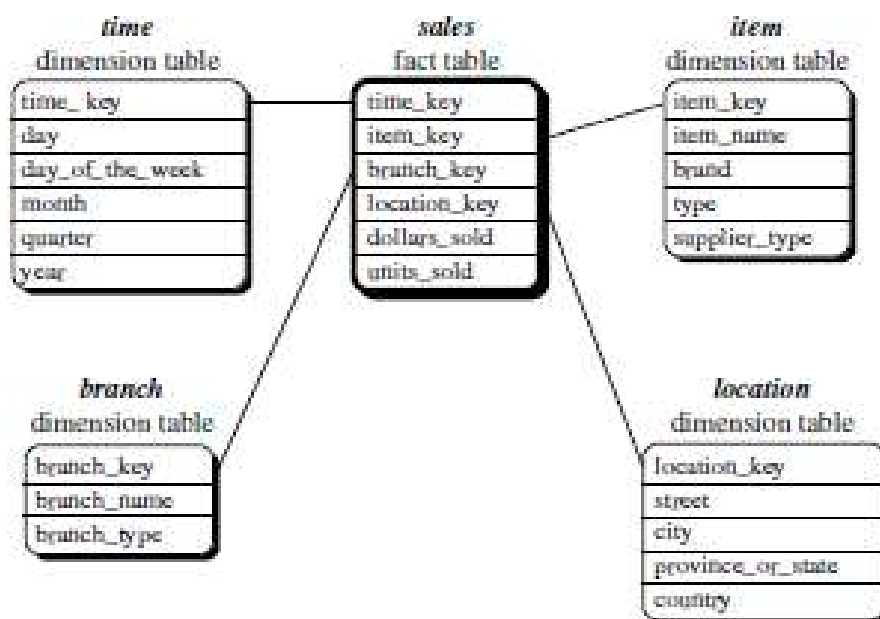
# Schema Design:

Stars, Snowflakes, and Fact Constellations: Schemas for Multidimensional Databases The entity- relationship data model is commonly used in the design of relational databases, where a database schema consists of a set of entities and the relationships between them. Such a data model is appropriate for on- line transaction processing. A data warehouse, however, requires a concise, subject-oriented schema that facilitates on-line data analysis. The most popular data model for a data warehouse is a multidimensional model. Such a model can exist in the form of a star schema, a snowflake schema, or a fact constellation schema. Let's look at each of these schema types. Star schema: The most common modeling paradigm is the star schema, in which the data warehouse contains (1) a large central table (fact table) containing the bulk of the data, with no redundancy, and (2) a set of smaller attendant tables (dimension tables), one for each dimension. The schema graph resembles a starburst, with the dimension tables displayed in a radial pattern around the central fact table.

**Star schema:**

A star schema for AllElectronics sales is shown in Figure. Sales are considered along four dimensions, namely,time, item, branch, and location. The schema contains a central fact table for sales that contains keys to each of the four dimensions, along with two measures: dollars sold and units sold. To minimize the size of the fact table, dimension identifiers (such as time key and item key) are system-generated identifiers. Notice that in the star schema, each dimension is represented by only one table, and each table contains a set of attributes. For example, the location dimension table contains the attribute set

{location key, street, city, province or state, country}. This constraint may introduce some redundancy.

For example, "Vancouver" and "Victoria" are both cities in the Canadian province of British Columbia. Entries for such cities in the location dimension table will create redundancy among the attributes province or state and country, that is, (..., Vancouver, British Columbia, Canada) and (..., Victoria, British Columbia, Canada). Moreover, the attributes within a dimension table may form either a hierarchy (total order) or a lattice (partial order).
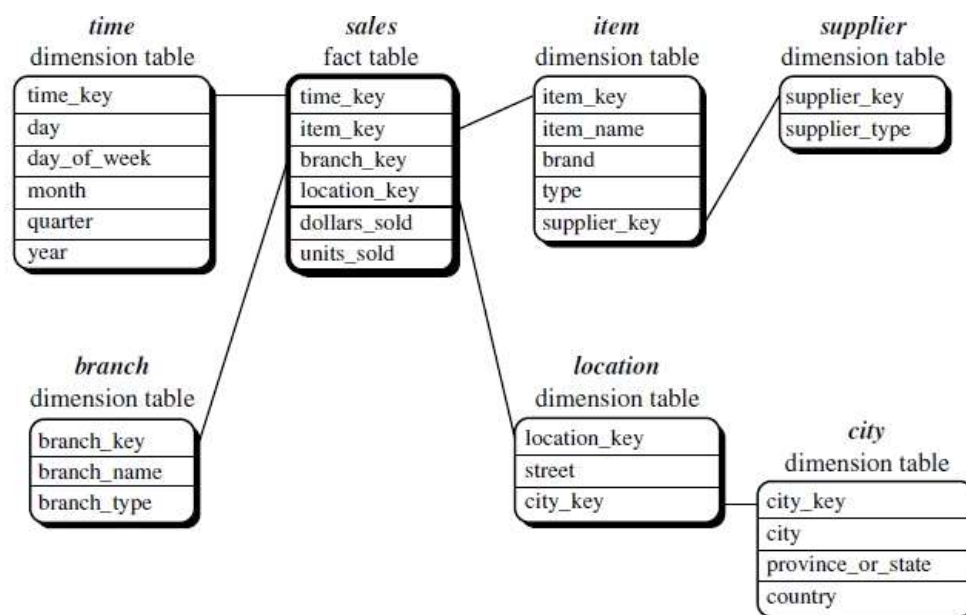


Star schema of a data warehouse for sales.

**Snowflake schema.:**

A snowflake schema for AllElectronics sales is given in Figure Here, the sales fact table is identical to that of the star schema in Figure . The main difference between the two schemas is in the definition of dimension tables.

The single dimension table for item in the star schema is normalized in the snowflake schema, resulting in new item and supplier tables. For example, the item dimension table now contains the attributes item key, item name, brand, type, and supplier key, where supplier key is linked to the supplier dimension table, containing supplier key and supplier type information. Similarly, the single dimension table for location in the star schema can be normalized into two new tables: location and city. The city key in the new location table links to the city dimension. Notice that further normalization can be performed on province or state and country in the snowflake schema



Snowflake schema of a data warehouse for sales.
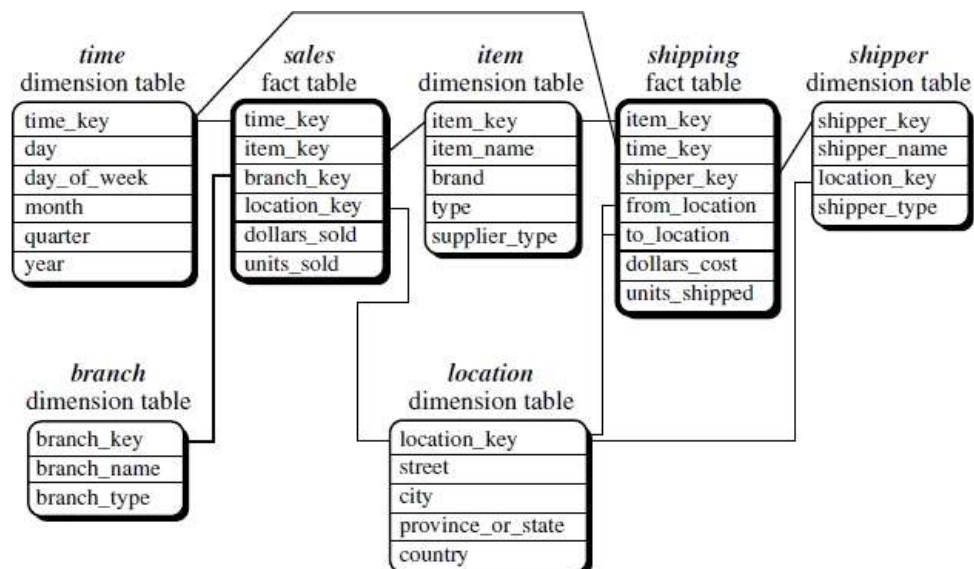
**Fact constellation.**

A fact constellation schema is shown in Figure. This schema specifies two fact tables, sales

and shipping. The sales table definition is identical to that of the star schema . The shipping table has five dimensions, or keys: item key, time key, shipper key, from location, and to location, and two measures: dollars cost and units shipped.

A fact constellation schema allows dimension tables to be shared between fact tables. For example, the dimensions tables for time, item, and location are shared between both the sales and shipping fact tables.

In data warehousing, there is a distinction between a data warehouse and a data mart.

A data warehouse collects information about subjects that span the entire organization, such as customers, items, sales, assets, and personnel, and thus its scope is enterprise-wide. For data warehouses, the fact constellation schema is commonly used, since it can model multiple, interrelated subjects. A data mart, on the other hand, is a department subset of the data warehouse that focuses on selected subjects, and thus its scope is department wide. For data marts, the star or snowflake schema are commonly used, since both are geared toward modeling single subjects, although the star schema is more popular and efficient.



5 Fact constellation schema of a data warehouse for sales and shipping.

**Measures: Their Categorization and Computation:**

"How are measures computed?" To answer this question, we first study how measures can be categorized.1 Note that a multidimensional point in the data cube space can be defined by a set of dimension-value pairs, for example, htime = "Q1", location = "Vancouver",item = "computer"i. A data cube measure is a numerical function that can be evaluated at each point in the data cube space. A measure value is computed for a given point by aggregating the data corresponding to the respective dimension-value pairs defining the

given point.

Measures can be organized into three categories (i.e., distributive, algebraic, holistic), based on the kind of aggregate functions used.

Distributive: An aggregate function is distributive if it can be computed in a distributed manner as follows. Suppose the data are partitioned into n sets.We apply the function to each partition, resulting in n aggregate values. If the result derived by applying the function to the n aggregate values is the same as that derived by applying the function to the entire data set (without partitioning), the function can be computed in

a distributed manner. For example, count() can be computed for a data cube by first partitioning the cube into a set of subcubes, computing count() for each subcube, and then summing up the counts obtained for each

subcube. Hence, count() is a distributive aggregate function. For the same reason, sum(), min(), and max() are distributive aggregate functions. A measure is distributive if it is obtained by applying a distributive aggregate function. Distributive measures can be computed efficiently because they can be computed in a distributive manner.

# OLAP(Online analytical Processing):

- OLAP is an approach to answering multi-dimensional analytical (MDA) queries swiftly.

- OLAP is part of the broader category of business intelligence, which also encompasses relational database, report writing and data mining.

- OLAP tools enable users to analyze multidimensional data interactively from multiple perspectives.

OLAP consists of three basic analytical operations:

- Consolidation (Roll-Up)

- Drill-Down

- Slicing And Dicing

- Consolidation involves the aggregation of data that can be accumulated and computed in one or more dimensions. For example, all sales offices are rolled up to the sales

department or sales division to anticipate sales trends.

- The drill-down is a technique that allows users to navigate through the details. For instance, users can view the sales by individual products that make up a region's sales.

- Slicing and dicing is a feature whereby users can take out (slicing) a specific set of data of the OLAP cube and view (dicing) the slices from different viewpoints.

## Types of OLAP:

- **Relational OLAP (ROLAP):**

- ROLAP works directly with relational databases. The base data and the dimension tables are stored as relational tables and new tables are created to hold the aggregated information. It depends on a specialized schema design.

- This methodology relies on manipulating the data stored in the relational database to

give the appearance of traditional OLAP's slicing and dicing functionality. In essence, each

action of slicing and dicing is equivalent to adding a "WHERE" clause in the SQL statement.

ROLAP tools do not use pre-calculated data cubes but instead pose the query to the

standard relational database and its tables in order to bring back the data required to

answer

the question.

- ROLAP tools feature the ability to ask any question because the methodology does not limit to the contents of a cube. ROLAP also has the ability to drill down to the lowest level of detail in the database.

- **Multidimensional OLAP (MOLAP):**

- MOLAP is the 'classic' form of OLAP and is sometimes referred to as just OLAP.

- MOLAP stores this data in an optimized multi-dimensional array storage, rather than in a relational database. Therefore it requires the pre-computation and storage of information in the cube - the operation known as processing.
- MOLAP tools generally utilize a pre-calculated data set referred to as a data cube. The data cube contains all the possible answers to a given range of questions.
- MOLAP tools have a very fast response time and the ability to quickly write back data into the data set.

- **Hybrid OLAP (HOLAP):**

- There is no clear agreement across the industry as to what constitutes Hybrid OLAP, except that a database will divide data between relational and specialized storage.

- For example, for some vendors, a HOLAP database will use relational tables to hold the larger quantities of detailed data, and use specialized storage for at least some aspects of the smaller quantities of more-aggregate or less-detailed data.

- HOLAP addresses the shortcomings of MOLAP and ROLAP by combining the capabilities of both approaches.

- HOLAP tools can utilize both pre-calculated cubes and relational data sources.