

* Data Warehouse :-

Data Warehousing provides Architecture & tools for business Executives to Systematically organize, Understand & Use their data to make Strategic decisions.

→ A data warehouse refers to a data repository that is maintained Separately from an organization's operational data-bases.

Data warehouse system allow for Integration of Variety of Application Systems.

→ William H. Inmon States "A data warehouse is a Subject-oriented, Integrated, Non-Volatile & High-Volume Collection of data to Support of Management's decision making process".

* Key features of Data Warehouse :- The key features of Data warehouse are as follows.

→ "Subject-oriented": Data warehouse is organized around Major Subject Such as Customers, Supplier, product & Sales. Rather than Concentrating on day-to-day operations & Transaction processing of Organization.

Data warehouse focuses on Modelling & Analysis of Data for Decision Making.

→ "Integrated": A data warehouse is usually Constructed by Integrating Multiple Heterogeneous Sources Such as Relational DB.

Flat files & online Transaction Records.

"Data Cleaning" & "Data Integration" techniques are applied to Ensure Consistency in Naming Conventions, Encoding Structures attributes, Measures & so on.

→ "Time - Variant" : Data are stored to provide information from an historic perspective. for ex 5-10 years.

Every key structure in data warehouse contains, either implicitly or explicitly a time element.

→ "Non-volatile" : A data warehouse is always a physically separate store of data transformed from the application data found in operational environment.

Due to this separation, A data warehouse does not need online transaction processing, Recovery & Concurrency Control Mechanisms.

→ The Data Warehouse allows "Knowledge Workers" such as Managers, Analysts & Executives to use the Warehouse to quickly & obtain an overview of the data, & make decisions based on info in Warehouse.

The construction of Data warehouse requires "Data Cleaning", "Data Integration" & "Data Consolidation".

→ Data Warehousing is very useful from the point of view of "Heterogeneous Database Integration".

In Traditional Database Approach to heterogeneous Database Integration is to build "Wrappers & Integrators" on top of Multiple Heterogeneous Databases.

→ When a query is posed to a client site, a Metadata dictionary is used to translate the query into queries appropriate for the individual heterogeneous sites involved.

* "Query-Driven Approach" Requires complex information filtering & Integration processes & Competes with local sites for processing resources.

"Traditional Database" Make use of Query-Driven Approach.

* "Update-Driven Approach" In which information from multiple heterogeneous sources is integrated in advance & stored in a warehouse for direct querying & analysis.

"Data Warehousing" makes use of Update-Driven Approach.

* Differences between Operational Database System & Data Warehouse

The Major task of online operational database system is to perform online transactions & query processing. These systems are called "Online Transaction processing (OLTP)".

For example, Most of day-to-day operations of an organization such as "purchasing", "Manufacturing", "Banking" etc.

-> "Online Analytical processing (OLAP)" serve users on knowledge workers in the role of data analysis & Decision Making.

Such systems can organize & present the data in various formats in order to accommodate the diverse needs of different users.

The Major distinguishing features of OLTP & OLAP are

* User & System Orientation :- An OLTP system is "Customer-oriented" and is used by clerks, clients, & IT professionals.

-> An OLAP system is "Market-oriented" and is used by knowledge workers such as Managers, Executives & Analysts.

* Data Contents :- OLTP system manages current data that, typically, are too detailed to be used.

→ OLAP System Manage large Amount historic data, provides facilities for Summarization & Aggregation & Store & Manage Information at different levels of granularity.

* Database Design :- OLTP usually adopts an Entity - Relationship (ER) data model & Appn oriented database design.

→ OLAP System adopts either "Star" or "Snowflake" data Models.

* View :- An OLTP System Focus Mainly on Current data within an Enterprise or department, without referring to historic data or data in different organization.

→ OLAP Systems Span Multiple Version of database Schema, due to Evolutionary process of an organization.

OLAP S/m deal with Info that originate from different organization, Integrating Info from different sources.

* Access :- OLTP System consist mainly of short atomic transa - ion so Concurrency Control & Recovery is required.

→ OLAP Systems are Mostly Read-only operations for any complex queries.

<u>Feature</u>	<u>OLTP</u>	<u>OLAP</u>
* <u>Characteristic</u>	* operational processing	* Informational processing
* <u>Orientation</u>	* Transaction	* Analytic
* <u>User</u>	* DBA, clerk, Database professional	* Knowledge Workers
* <u>Function</u>	* Day-to-day operation	* long-term for Decision Support
* <u>DB-Design</u>	* ER-Model, Appn oriented	* Star/Snowflake, Subject oriented
* <u>View</u>	* Detailed, Flat Relational	* Summarized & Multi-dimensional
* <u>Access</u>	* Read/Write	* Read
* <u>Focus</u>	* Data-in	* Information-out
* <u>Unit of work</u>	* Short, Simple transaction	* Complex query

* Operations	* Index / Hash	* Lots of Scan
* No of Records Accessed	* Tern	* Millions
* No of Users	* Thousands	* Hundreds
* DB Size	* GB to High-order GB	* \geq TB
* Priority	* High performance & Availability	* High flexibility, End-user Autonomy
* Metric	* Transaction Throughput	* Query throughput & Response time

* Data Warehousing : A Multitiered Architecture :-

Data warehouse often

adopts a 3-tier Architecture Such as

- * Bottom tier
- * Middle tier
- * Top tier

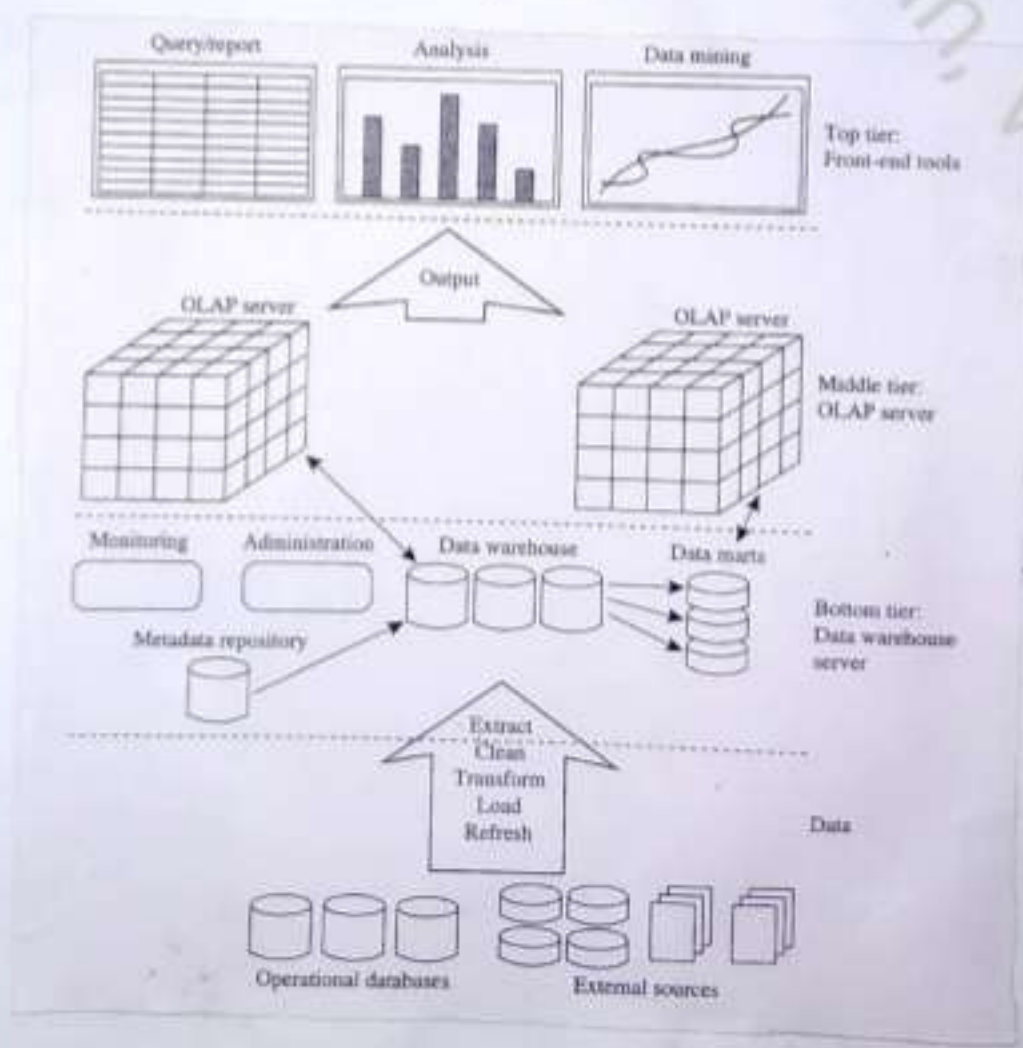


Fig:- Three - tier Datawarehouse Architecture

* Bottom tier :- Warehouse Database Server. It is almost a Relational Database & other External Sources.

Back-End tools & utilities are used to perform Extraction, cleaning & Transformation (Ex:- to Merge Similar data from different Sources into a unified format).

→ The data are Extracted using Applu program Interface known as "gateways"

Example of gateways include ODBC (Open Database Connection) & OLEDB (Object Linking & Embedding Database) by Microsoft & JDBC.

→ This tier also contains a Metadata Repository, which stores Information about the Data Warehouse & Contents.

* Middle tier :- Consists of OLAP Server that is typically implemented using

→ "Relational OLAP (ROLAP)" :- Extended Relational DBMS that Maps Multidimensional data to Standard Relational-operation.

→ "Multidimensional OLAP (MOLAP)" :- Direct Implementation of Multidimensional data & operation.

* Top tier :- Consists of Front-End Client layer, which contains Query & reporting tools, Analysis tools & Data mining tools. Such as "Trend Analysis", "prediction" & so on.

* Data Warehouse Models :-

There are three Data Warehouse Models

* Enterprise Warehouse

* Data Mart

* Virtual Warehouse

* Enterprise Warehouse :- Collects all the Information about Subjects Spanning Entire organization.

It provides Corporate-wide data with Corporate-wide data - Integration.

→ Enterprise Warehouse Contains detailed data as well as Summarized data & can range in size from of few Gigabytes to hundreds of GB, Terabytes or beyond.

→ An Enterprise data Warehouse may be implemented on traditional Mainframes, Computer SuperServers or parallel Architectured -ware platforms.

It requires Extensive Business Modeling & may take years to Design.

* Data Mart :- Contains a subset of Corporate-wide data that is of value to Specific Group of users.
For Ex: A Marketing Data Mart may confine its Subject to Customers, Items & Sales.

→ Data Marts are usually implemented on "low-cost Department 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.

→ Data Marts can be Categorized as

* Independent → Sourced from operational Systems or External Info providers.

* Dependent → Sourced Directly from Enterprise Data Warehouse

* Virtual Warehouses :- Is a set of Views over operational databases for Efficient Query processing.

A Virtual Warehouse is Easy to build but requires Excess Capacity on operational database Servers.

→ For Ex: Air ticketing System, collects & displays business data relating to a Specific Moment in time, creating a Snapshot

of the condition of business at that Moment.

* Recommended Approach for Warehouse Development :-

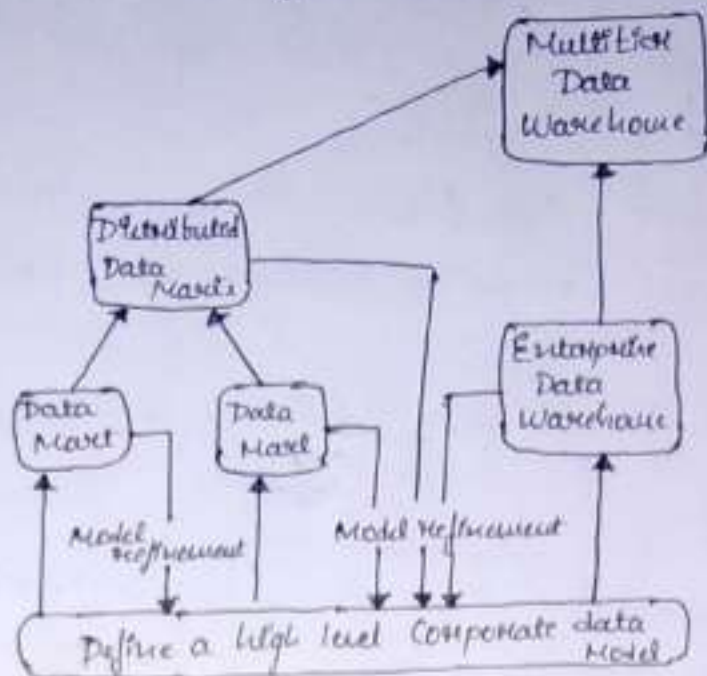


Fig:- A Recommended Approach for Warehouse Development

A Recommended Method for the development of data Warehouse Systems is to implement the Warehouse in an Incremental & Evolutionary Manner as shown in above figure.

- Firstly, a high-level Corporate data Model is defined within a short period of time, that provides a Corporate-wide view of data among different Subject & potential users.
- Second, Independent data Marts can be implemented in parallel with the Enterprise Warehouse based on the Corporate data Model.
- Third, Distributed data Marts can be constructed to integrate different data Marts via hub servers.

Finally a "MultiSite Data Warehouse" is constructed, where the Enterprise Warehouse is the Sole Custodian of all Warehouse data, which is then distributed to various dependent Data Marts.

* Data Cube : A Multidimensional Data Model

A Data Cube allows data to be Modelled & Viewed In Multiple Dimensions. It is defined by "Facts" & "Dimensions".

* Dimension → are the perspectives on Entities with respect to which an organization wants to keep records.

For Ex: "All Electronics" May Create a "Sales" data Warehouse in order to keep records of Stores Sales with respect to Dimensions "Time", "Item", "Branch" & "Location".

* Facts → are Numeric Measures, Quantities by which we want to analyze relationship between Dimensions.

For Ex: Facts for a Sales data Warehouse includes "Dollars-Sold", "Units-Sold", & "Amount-budgeted".

→ Consider a 2D Simple data Cube Example from "All Electronics" Sales data for Items Sold per Quarter in City Vancouver.

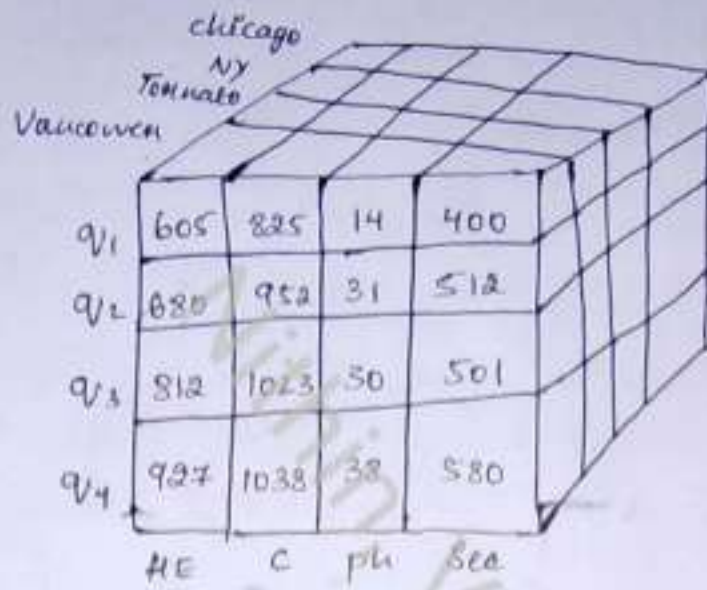
Location = "Vancouver"				
Time (Quarter)	Item (Type)			
	Home Ent	Computer	Phone	Security
Q ₁	605	825	14	400
Q ₂	680	952	31	512
Q ₃	812	1023	30	501
Q ₄	927	1038	38	580

In this 2-D Report, the Sales for Vancouver are shown w.r.t time dimension & Item.

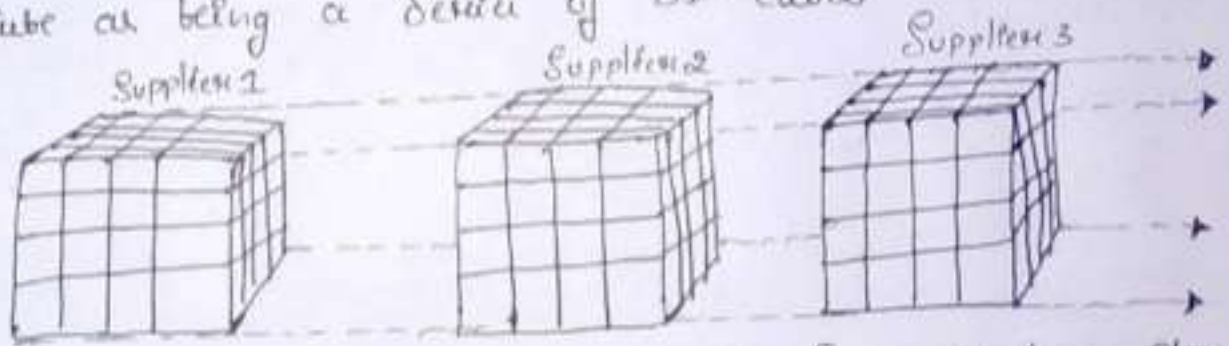
→ Suppose that we would like to view the Sales-data with 3-D Report the data according to Time, Item & Loc in as shown below

Time	Loc = "Chicago"				Loc = "New York"				Loc = "Toronto"				Loc = "Vancouver"			
	HE	C	Ph	Sec	HE	C	Ph	Sec	HE	C	Ph	Sec	HE	C	Ph	Sec
Q ₁	854	822	39	623	1087	720	64	620	1023	760	65	620	1000	610	65	60
Q ₂	943	890	64	698	1130	841	72	619	1014	820	66	619	1001	619	60	50
Q ₃	1052	924	59	789	1034	900	63	618	1025	823	67	720	1002	700	55	60
Q ₄	1129	992	63	876	1142	621	79	617	1026	824	68	723	1003	710	50	30

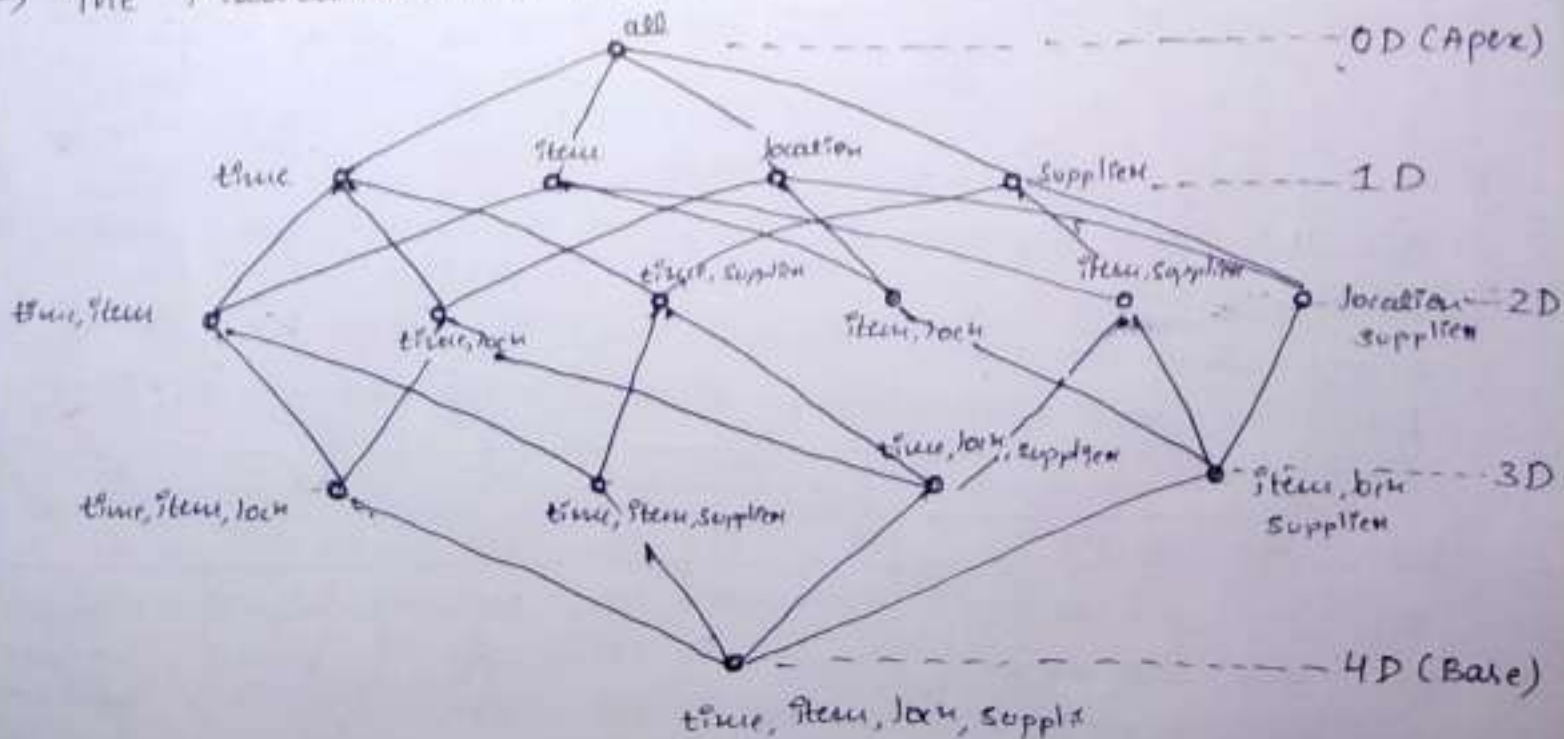
→ We may also represent the same data in the form of 3D data cube.



→ Suppose that we would like to view our sales data with addition of 4th Dimension such that Supplier. Viewing things in 4-D cube becomes complex. However, we can think of a 4-D cube as being a series of 3D-cubes.



→ The Multidimensional Data Model can be summarized as shown



The OD Cuboid, which holds the highest level of Summary
-ion is called "Apex Cuboid"

In our Example, the "total Sales" or "dollar Sold", Summary
-ed over all four dimensions.

* Star, Snowflake & Fact Constellation: Schema for Multidimensional Data Model

The Entity-Relationship data Model is commonly used to design the relational database, where database Schema consists of a set of Entities & the Relationships between them. Such data Model is appropriate for online Transaction processing.

→ The Most popular data Model for data Warehouse is a Multidimensional Model which can exist in the form of

- * Star Schema
- * Snowflake Schema
- * Fact Constellation Schema

* Star Schema :- is the most commonly used Multidimensional data Model, in which the data Warehouse contains

- * A large central table (Fact table) containing the bulk data with no redundancy.
- * A set of smaller attendant-tables (Dimension tables) one for each dimension.

The Schema Graph resembles a Star burst, with the Dimension tables displayed in radial pattern around the central Fact table.

→ A Star Schema for "all Products" is shown in the following figure.

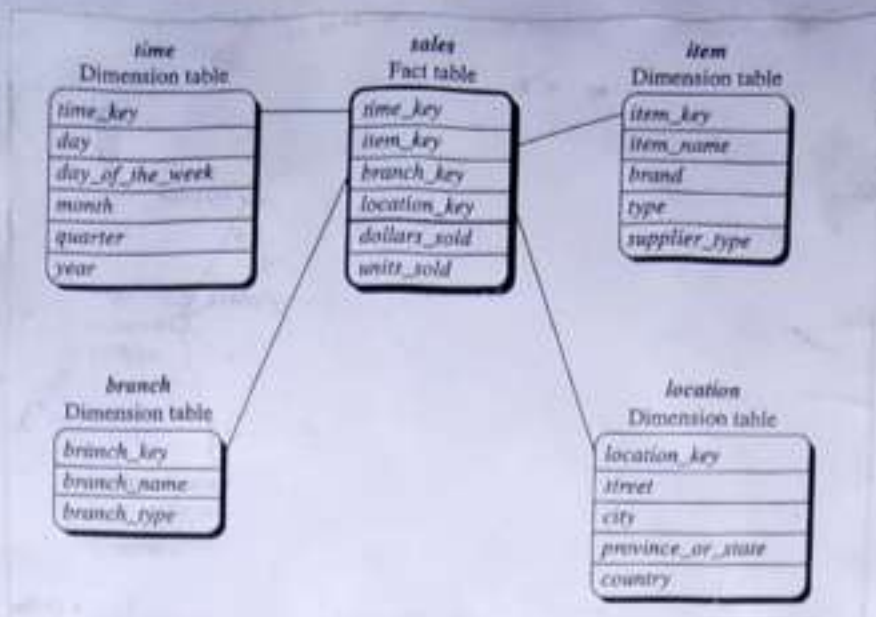


Fig:- Star Schema of Sales Data Warehouse

In the above Schema, Sales is considered along four Dimensions "Time, Item, branch & location". The Schema contains a Central Fact table for Sales that contains keys to each 4 Dimensions with two Measures "dollars sold" & "units sold".

To Minimize the Size of fact table, dimension-identifiers (time_key, item_key etc) are System generated identifiers.

→ In Star Schema, Each dimension is represented by only one table & Each table contains a Set of attributes.

For ex- the "Location" Dimension table contains the attribute Set [location_key, Street, City, State, Country]. This constraint may introduce some Redundancy.

* Snowflake Schema :- Is a Variant of Star-Schema Model, where some dimension tables are Normalized, thereby splitting the data into additional tables.

→ The Major difference between the Star & Snowflake Schema Model is that the Dimension tables of the Snowflake Model may be kept in Normalized form to reduce redundancy. Such a table is Easy to Maintain & Save Storage Space.

→ The Snowflake Structure can reduce the Effectiveness of

Knowing, Since more joins will be needed to execute a query. The system performance may be adversely impacted. Hence, although the Snowflake Schema reduces redundancy, it is not as popular as Star Schema in Data Warehouse Design.

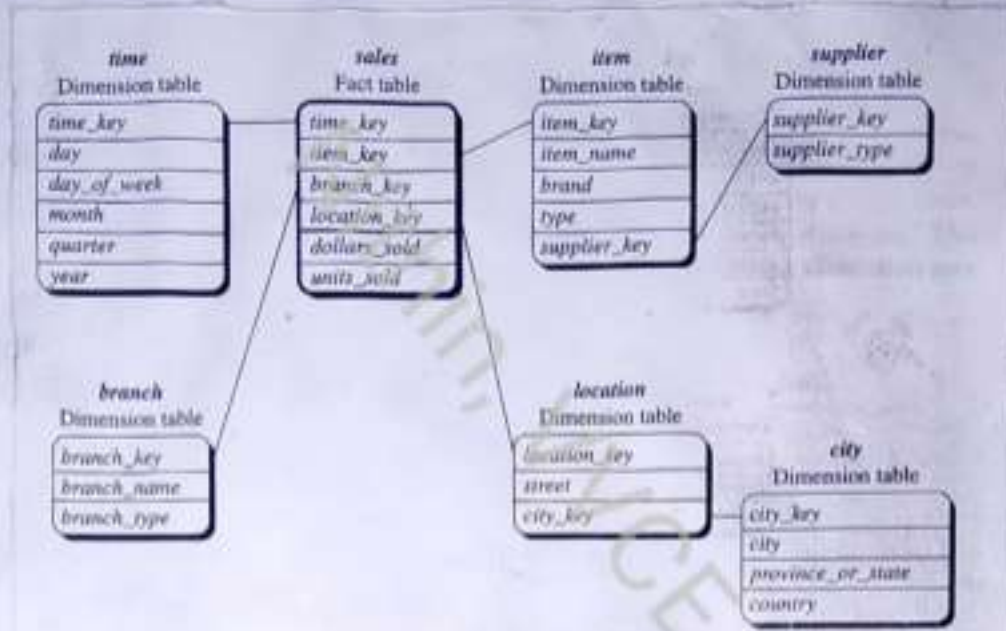


Fig:- Snowflake Schema of Sales

As shown in above Snowflake Schema for "All Electronics", the Sales fact table is identical to the Star Schema.

→ The main difference between Star & Snowflake Schema is Definition of Dimension table.

- * Dimension table of "Item" in Star Schema is Normalized in Snowflake Schema, resulting table form "Item" and "Supplier" Dimension tables.
- * Similarly Dimension table of "Location" in Star Schema is Normalized in Snowflake Schema, resulting table form "Location" & "City" Dimension tables.

* Fact Constellation :-

For each Star Schema or Snowflake Schema, it is possible to construct "Fact Constellation Schema".

→ The Fact Constellation Architecture contains Multiple Fact tables that share many Dimensions.

→ Consider the following Example of Fact Constellation

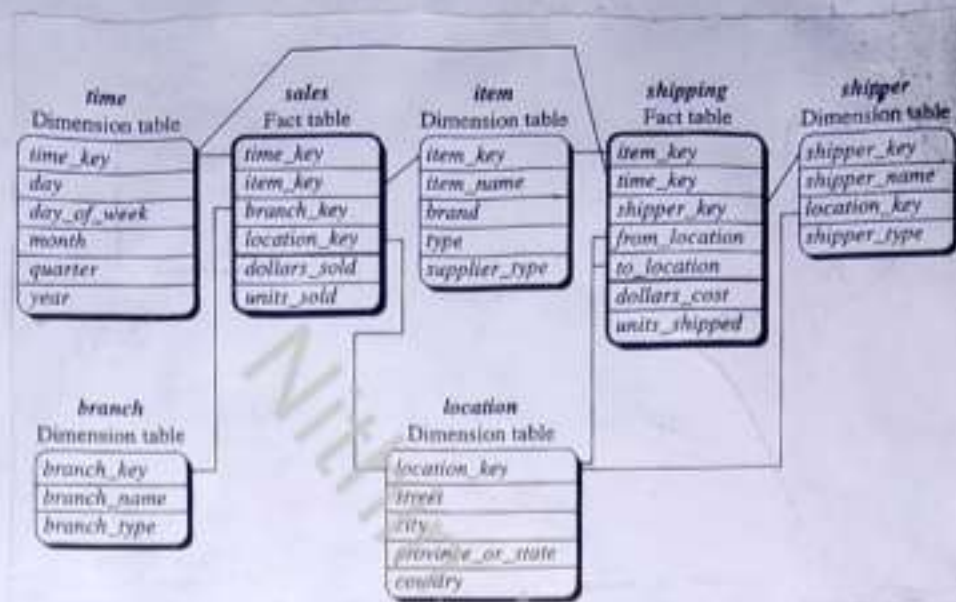


Fig. Fact Constellation Schema of Sales

The above Schema Specifies two Fact tables "Sales" & "Shipping"

→ The Shipping table consists of 5 Dimensions on keys & Two Measures

* Dimensions → Item_key, Time_key, Shipper_key, to_location
From_location

* Measures → Dollars_cost, Units_shipped.

A Fact Constellation Schema allows Dimension tables to be Shared between Fact tables.

→ A "Data Warehouse" Collects Information about Subjects that Span whole organization Such as Customers, Sales, Shipping Etc.

For Data Warehouse, the fact Constellation Schema is commonly used to Model Multiple Interrelated Subjects.

→ A "Data Mart" is a Subset of Data Warehouse that Focus on Selected Subjects.

For Data Mart, the Star or Snowflake Schema is commonly used since both are used to Model Single Subjects.

* Role of Concept Hierarchy

A "Concept Hierarchy" defines a Sequence

of Mapping from a Set of low-level concepts to High-level More general concepts.

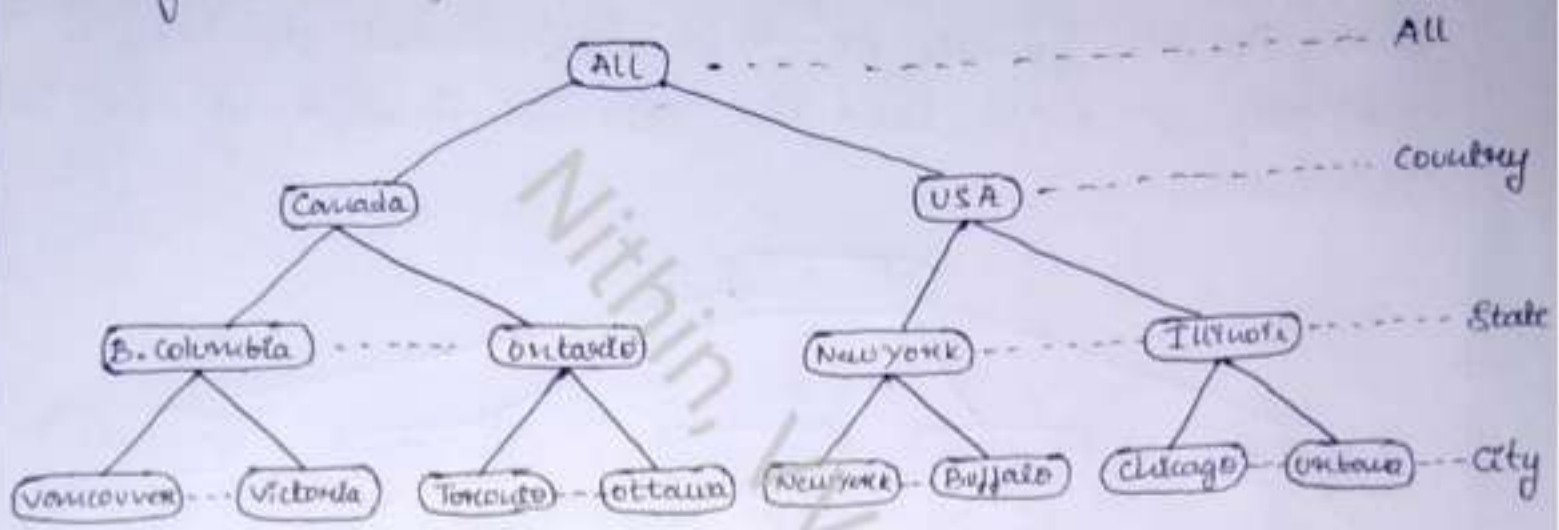


Fig: Concept Hierarchy of Location

Mapping a Set of low-level concepts (Cities) to High-level More general concepts (Countries).

- Many concept hierarchies are implicit within the database Schema
- Due to Space Constraint, not all of the hierarchy nodes are shown
- we can use Lattice Structures of attributes in warehouse dimension

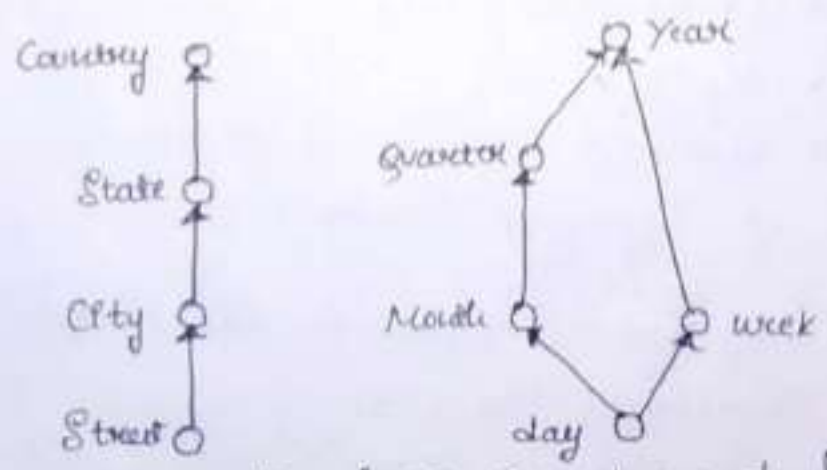


Fig: Lattice Structures of Location & Time

→ Concept Hierarchy may also be defined by discretizing or grouping values for a given dimension or attribute, resulting in a "Set-Grouping Hierarchy".

For the dimension "price" with interval (\$x --- \$y) denote the range from \$x (Exclusive) to \$y (Inclusive)

* Role of Concept Hierarchy

A "Concept Hierarchy" defines a Sequence

of Mapping from a Set of low-level concepts to High-level More general concepts.

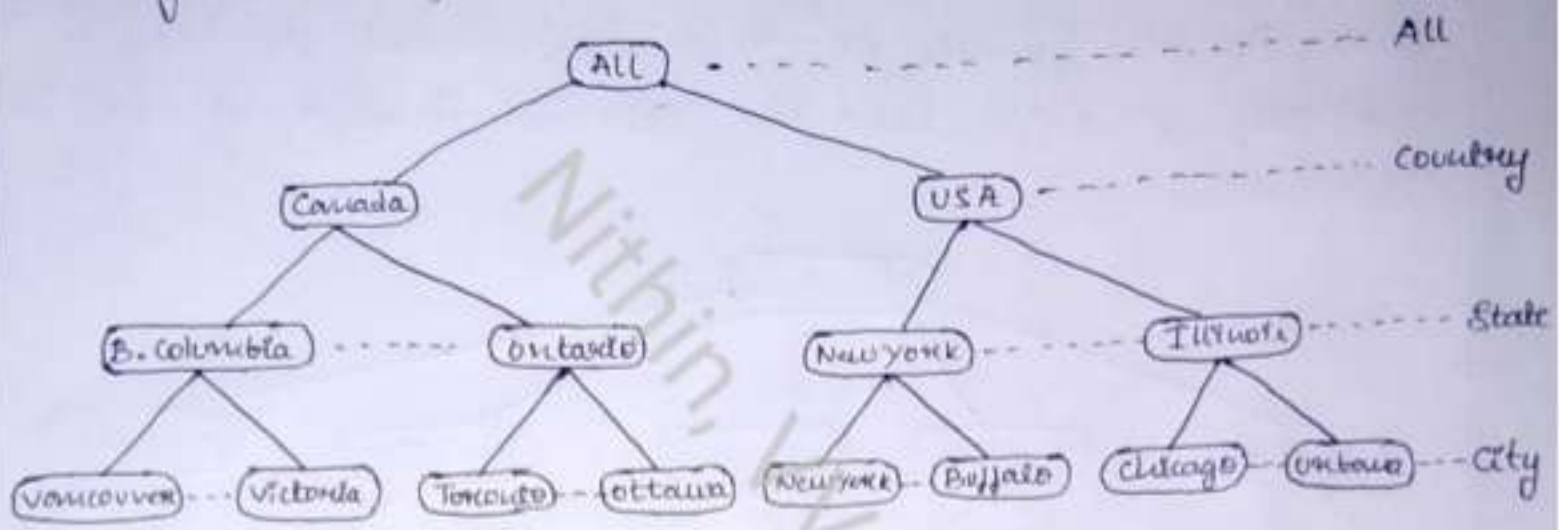


Fig: Concept Hierarchy of Location

Mapping a Set of low-level concepts (Cities) to High-level More general concepts (Countries).

- Many concept hierarchies are implicit within the database Schema
- Due to Space Constraint, not all of the hierarchy nodes are shown
- we can use Lattice Structures of attributes in warehouse dimension

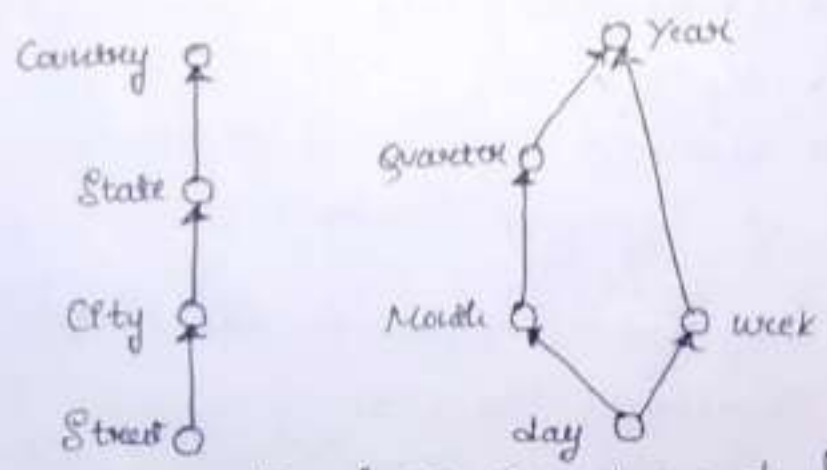


Fig:- Lattice Structures of Location & Time

→ Concept Hierarchy may also be defined by discretizing or grouping values for a given dimension or attribute, resulting in a "Set-Grouping Hierarchy".

For the dimension "price" with interval (\$x --- \$y) denote the range from \$x (Exclusive) to \$y (Inclusive)

→ There may be more than one concept hierarchy for a given attribute or dimension based on different user viewpoints.

Concept Hierarchies may be provided manually by System Users, Domain Experts or Knowledge Engineers or may be automatically generated based on Statistical Analysis of the Data Distribution.

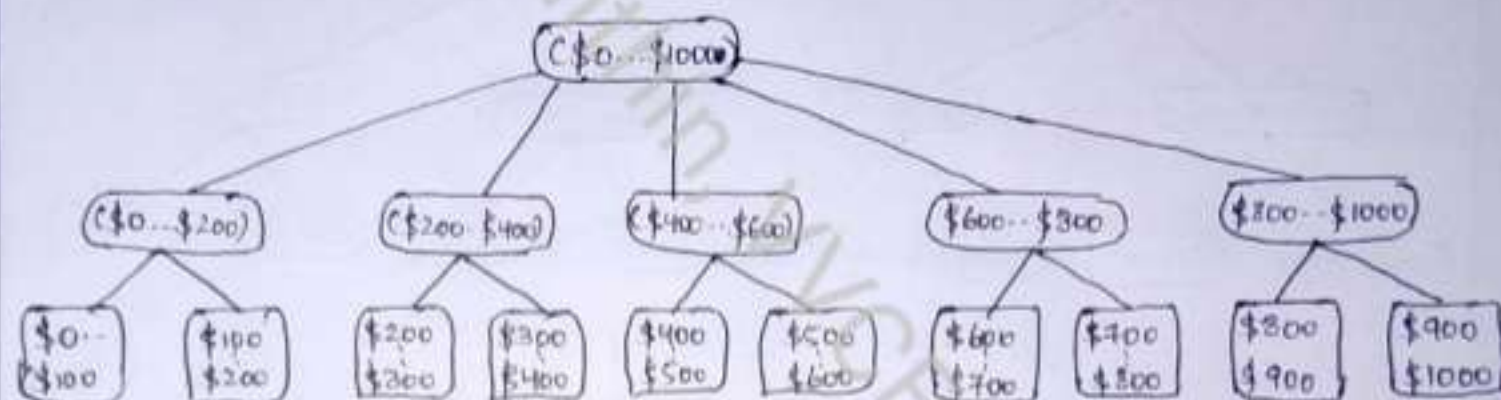


Fig- Set Grouping Hierarchy of price

* Measures & These Categorization & Computation :-

Measures can be

Organized into three categories such as

- * 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.

Suppose the data are partitioned into 'n' sets, we apply the Function which results 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 Entire Data Set in Distributed Manner.

→ Some of good examples of Distributive Aggregate Functions are

"Count()", "Min()", & "Max()".

* Algebraic :- An Aggregate Function is Algebraic, if it can be computed by an Algebraic function with M Arguments, each of which is obtained by applying a distributive Aggregate Function.

-tion. → A Measure is "Algebraic" if it is obtained by Applying an Algebraic Aggregate Function.

For Ex:- "Avg()" can be computed by Sum()/count(), where both Sum() & count() are distributive Aggregate Functions.

* Holistic :- An Aggregate Function is Holistic, if there is no constant bound.

There does not exist any Algebraic function that characterizes the computation. Such Measures are called "Holistic".

→ Some of good Examples of Holistic Functions include "Median", "Mode()" & "Rank()".

More large data cube Applications require Efficient computation of distributive & Algebraic Measures. Many Efficient techniques for this exist.

In contrast, it is difficult to compute Holistic Measures Efficiently.

* Typical OLAP operations :-

Number of operations may be

applied to Data-Cubes are

- * Roll-up
- * Drill-Down
- * Slice & Dice
- * Pivot or Rotate

Consider the following Example of "All Electronics" Sales Data Cube which

Contains the Dimensions [Location, Time & Item].

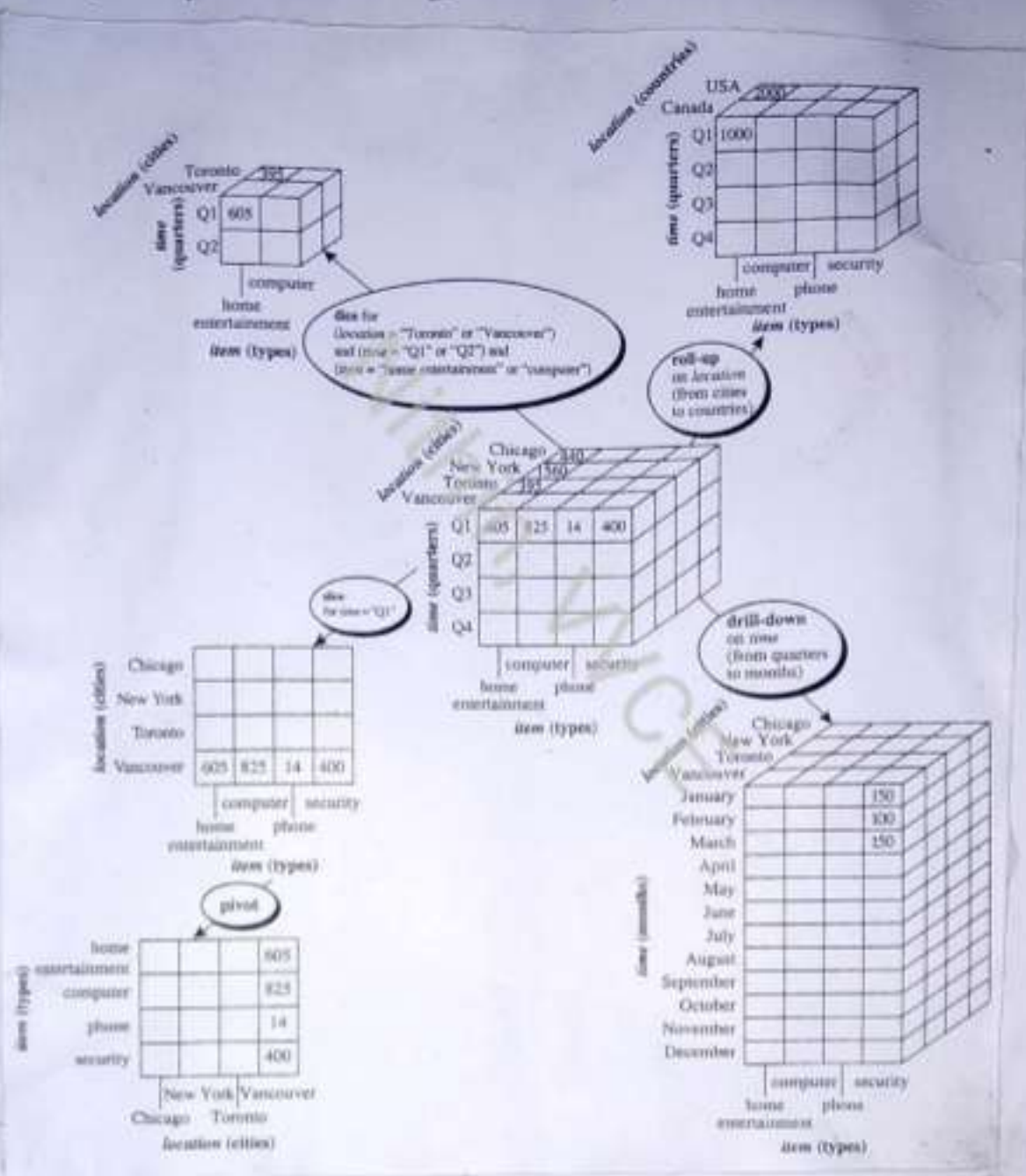


Fig: Typical OLAP operations

* Roll-up :- Is like "Zooming-out" on the Data Cube, It is required when the user needs further Abstraction or Less Detail.

→ This operation performs further Aggregation on the Data. For ex. Roll-up on location from "cities" to "countries"

* Drill-Down :- Is like "Zooming-In" on the Data Cube, i.e. Reverse of roll-up.

→ It is an appropriate operation when the user needs further detail or when user wants to partition more finely on

Wants to focus on some particular values of certain dimension.

→ For ex, if the time in data cube is represented in quarters then it can be drilled down or zoomed-in the time to month wise.

* Slice & Dice :- Slice & Dice are operations for browsing the data in the cube. The terms refer to the ability to look at information from different viewpoints.

* A "Slice" is a subset of cube corresponding to a single value for one or more members of dimension.

* A "Dice" operation is similar to slice but dicing does not involve reducing the no of dimension. A dice is obtained by performing a selection on two or more dimension.

* Pivot or Rotate :- Pivot operation is used when the user wishes to re-orient the view of data cube. It may involve swapping the row & column, or moving one of row dimension into column dimension.

* Data Mining :-

Data Mining is a technology that blends traditional Data analysis Method with Sophisticated Algorithms for processing large volume of data.

→ It has also opened up exciting opportunities for Exploring & Analyzing new types of data & for analyzing diff types of data in new ways.

"Data Mining is a process of automatically discovering useful Information in Large Data Repositories".

* Data Mining & Knowledge Discovery :- Data Mining is an integral part of Knowledge Discovery in Database (KDD), which is a overall process of converting raw data into useful info.

→ This process consists of a series of Transition Steps from "Data preprocessing" to "post processing of Data Mining Results".

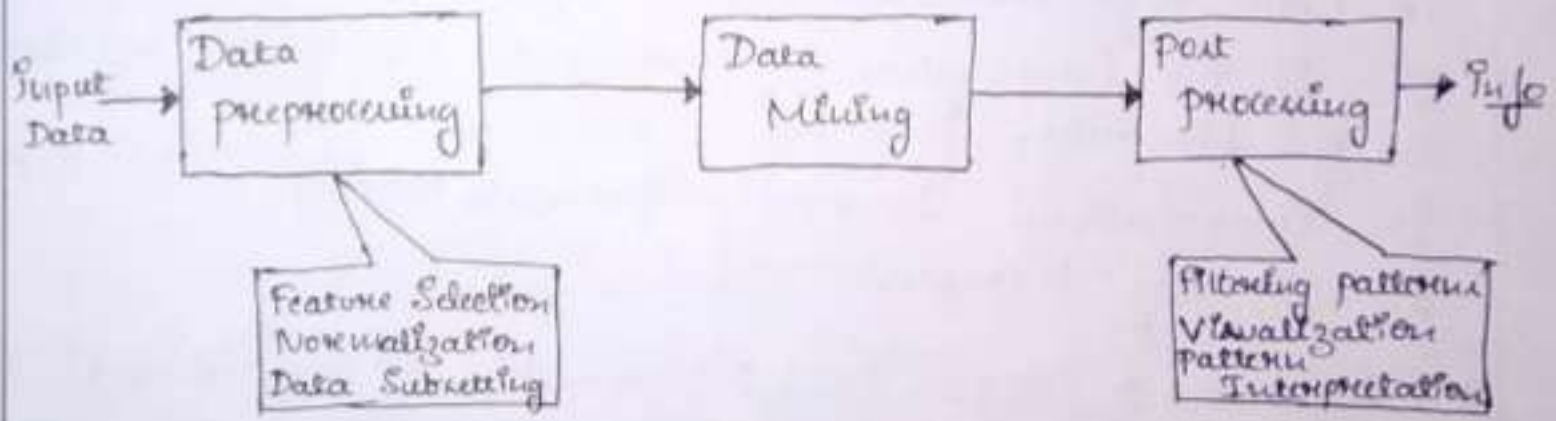


Fig 1

The Input data can be Stored in a Variety of file formats such as C flat files, Spread Sheets or Relational tables.

→ The purpose of preprocessing is to Transform the raw Input data into an appropriate format for subsequent Analysis.

The Steps Involved in Data preprocessing are

- * Finding Data from Multiple Sources.

* Cleaning Data to Remove noise.

* Selecting Records & features to data Mining, this is the most time consuming step in knowledge discovery process.

→ "post processing" steps that ensure that only valid & useful results are incorporated into the Decision Support System.

* Challenges :-

The motivating challenges of Data Mining are.

* Scalability :- Advances in Data generation & collection of data sets with sizes of Giga bytes, Terabytes or even Petabytes are becoming common.

→ If the data Mining Algorithms are to handle these massive data sets, then they must be Scalable.

* High Dimensionality :- It is now common to encounter data sets with hundreds or thousands of attributes.

→ Data sets with temporal or spatial components also tend to have high dimensionality.

Ex:- If the temporal measurements are taken repeatedly no of dimensions increase.

→ The Computational Complexity increases rapidly with the dimensionality increase.

* Heterogeneous & Complex Data :- Traditional Data Analysis Methods often deal with data sets containing attributes of the same type, either continuous or categorical.

→ The role of Data Mining in Business, Science, Medicine & other fields have grown, so has the need for techniques that handle heterogeneous attributes.

* Data Ownership & Distribution :- Sometimes the Data Model for Analysis is not stored in one location or owned by one org. Instead data is geographically distributed.

→ The key challenges faced by distributed data Mining Algorithms include

- * How to reduce the Amount of communication required to perform Distributed computation.
- * How to efficiently Consolidate the Data Mining Results obtained from Multiple Sources.
- * How to Address Data Security Issues.

* Non-Traditional Analysis :- Traditional Statistical Approach is based on a "hypothesize-and-Test" paradigm.

→ Current Data Analysis tasks often requires the Renovation & Evolution of thousands of hypothesis & consequently the development of some data Mining techniques has been motivated. It is a non-traditional Approach.

* Data Mining Tasks :- The Data Mining tasks can be classified generally into two types based on what a specific task tries to achieve.

- * Predictive Tasks
- * Descriptive Tasks.

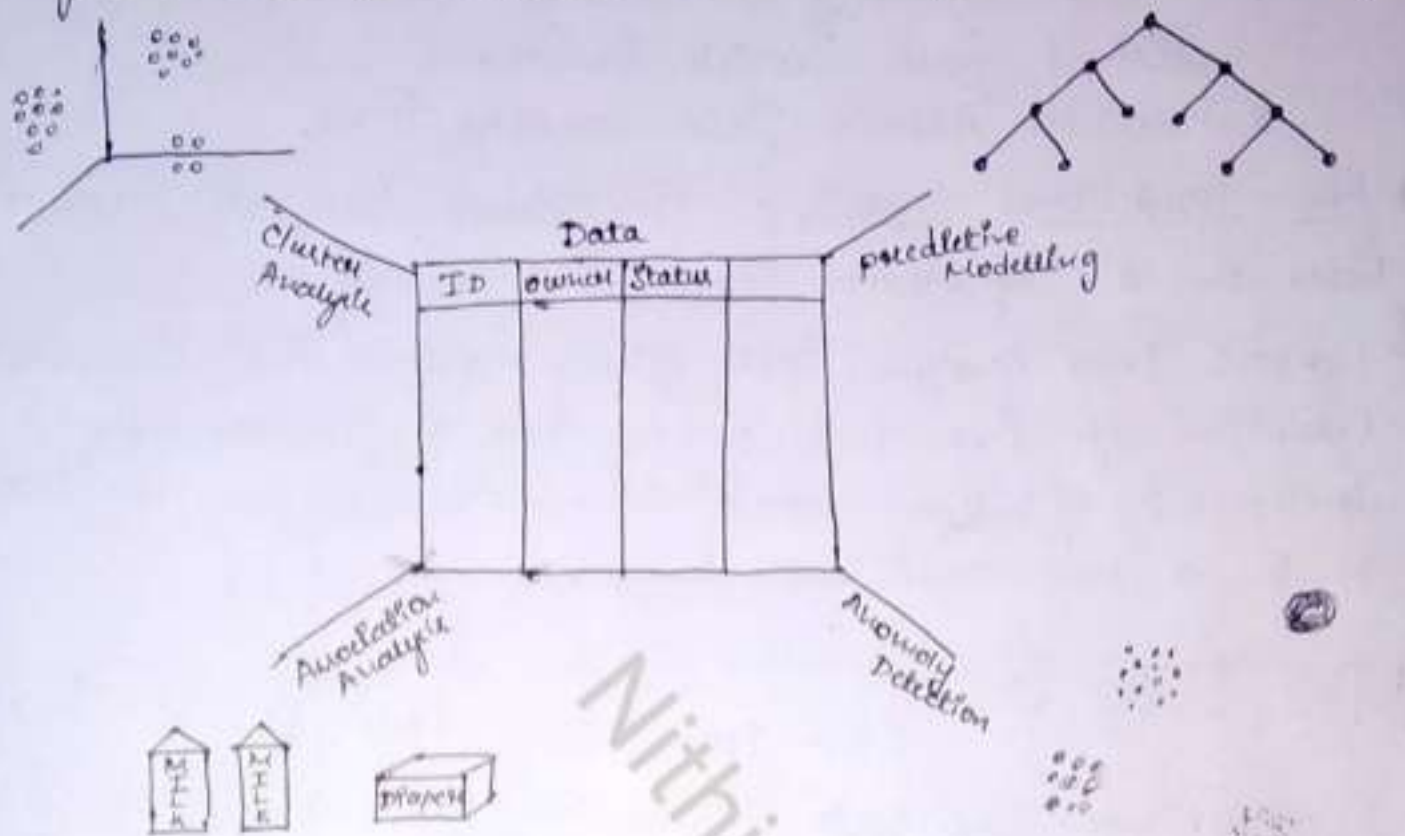
There are a number of Data Mining tasks such as Classification, Prediction, Time-Series Analysis, Association, Clustering, Summarization etc. All these tasks are either predictive tasks or Descriptive tasks.

* Predictive Tasks → Usually from the available data set that is helpful in predicting unknown or future values of another data set of interest.

A Medical practitioner trying to diagnose a disease based on the Medical test results of a patient can be considered as a predictive Data Mining task.

* Descriptive Tasks → Usually find data describing patterns & come up with new, significant information from the available Data Set.

A retailer trying to identify products that are purchased together can be considered as a descriptive data Mining Task.



* predictive Modelling :- Is a process that uses Data Mining & probability to forecast outcomes.

→ There are two types of predictive Modelling

- * Classification → used for Discrete Variables
- * Regression → used for Continuous Variables.

The goal of both tasks is to learn a Model, in which Error between the predicted & True values is low.

→ For Example, a Model can predict the Income of an Employee based on Education, Experience & other demographic factors like place of Stay, gender etc.

* Anomaly Detection :- Is a task of identifying observations whose characteristics are significantly different from the rest of data.

It is also called as "outlier Detection".

→ A good Anomaly detector must have a high detection rate & a low false alarm rate.

The good anomaly detection algo is to discover real anomalies & avoid false labelling.

→ Typically the anomalous items will translate to some kind of problem such as bank fraud, structural defect, medical problem etc.

* Association Analysis :- Association discovers the association or connection among a set of items. Association identifies the relationships among objects.

→ Association Analysis is used for commodity management, advertising, direct marketing etc.

→ For example, a retailer can identify the products that normally customers purchase together or even find the customers who respond to the promotion of same kind of products.

* Clustering Analysis :- Is used to identify data objects that are similar to one another.

→ The similarity can be decided based on a no of factors like purchase behavior, responsiveness, geographic location & so on.

→ For example, An insurance company can cluster its customers based on age, residence, income etc.

* Types of Data :-

Data set can be often be viewed as a collection of data objects. Other names for a data object are "record", "point", "vector", "pattern", "Event", "case", "sample" or "observation".

Data can be classified into two broad types such as

* Qualitative

* Quantitative (numeric)

* Qualitative Data :- arise when the observations fall into separate distinct categories.

→ For Ex:- colour of eyes - blue, green, brown
Exam result - Fail, pass

Such data are inherently "discrete", in that there are finite number of possible categories into which each observation may fall.

* Quantitative Data :- or numerical data arise when the observations are "counts" or "measurements".

→ The data are said to be "discrete" if the measurements are integers (Ex No of people in a house) & "continuous" if the measurements can take on any value, usually within some range (Ex weight).

* Attribute of Measurement of Data :- An "Attribute" is a property or characteristic of an object that may vary, either from one object to another or from one time to another.

→ For Ex:- Eye colour varies from person to person, while the temperature of an object varies over time.

* Measurement Scale is a rule (function) that associates a numerical or symbolic value with an attribute of an object.

For Ex:- we count the no of chairs in a room to see if there will be enough to seat all the people coming to a meeting.

* Different Types of Attributes :- A useful way to Specify the type of an attribute is to Identify the properties of numbers that correspond to underlying properties of the attribute.

→ Distinctness = \neq (Nominal) :- The values of a nominal attribute are just different names

Nominal values provide only Enough info to distinguish one object from another (=, \neq)

For Ex: Zip code, Employee ID

→ Order $< \leq > \& \geq$ (Ordinal) :- The values of an ordinal attribute provide Enough info to order objects

For Ex: Grades, Street Numbers

→ Interval + $\& -$ (Addition) :- For Interval attributes, the differences between values are meaningful, i.e. a unit of measurement exists (+, -)

For Ex: Calendar dates.

→ Multiplication * $\& /$ (Ratio) :- For ratio variables both differences & ratios are meaningful (*, /)

For Ex: Counts, Age, Mass Etc.

* General Characteristics of Data Sets :- There are 3 characteristics that apply to many data sets & have a significant impact on the data Mining techniques.

* Dimensionality

* Sparsity

* Resolution

→ Dimensionality :- refers to how many attributes a dataset has.

For Ex, Healthcare data is notorious for having vast amounts of variables (Blood pressure, weight, cholesterol level)

This data could be represented in spreadsheet, with one column representing each dimension.

→ Sparsity : Is one in which a relatively high percentage of the variable cells do not contain actual data.

For ex: a new variable dimensioned by MONTH for which you do not have data for part months.

→ Resolution : It is frequently possible to obtain data at different levels of resolution, & often the properties of the data are different at different resolutions.

For ex: the surface of Earth seems very uneven at a resolution of few meters, but is relatively smooth at a resolution of tens of kilometers.

* Types of Data Sets :- There are many types of data sets & as the field of Data Mining develops & matures, a greater variety of data sets become available for analysis.

→ Record Data : Is a collection of records (data objects), each of which consists of a fixed set of data fields (attributes).

There is no explicit relationship among records or data fields, & every record has the same set of attributes.

→ Record data is usually stored either in flat files or in relational databases.

→ Different types of record data are

* Transaction Data → where each record (Transaction) involves a set of items.

* Data Matrix → collection of data with same fixed set of numeric attributes.

* Document-Term Matrix → where each term is a component (attribute) of the vector.

→ Graph-Based Data : A graph can sometimes be a convenient & powerful representation of data.

→ There are two diff types of Graph-Based Data

* Data with relationship among objects → Relationship among objects frequently convey important information.

Ex: Web pages in World Wide Web

* Data with objects that are Graphs → If objects contain subobjects that have relationships, then such objects frequently represented as Graphs.

→ Ordered Data: In one in which the attributes have relationship that involve order in time or space

→ Different types of Ordered Data are

* Sequential Data → an Extension of Record data, where each record has a time associated with it.

Ex: Retail transaction data with time of transaction.

* Time Series Data → is a special type of Sequential Data in which each record is a time series.

a series of measurements taken over time.

* Spatial Data → Some objects have spatial attributes, such as position or areas, as well as other types of attributes.

Ex: Weather Data.

* Data Quality :-

Today's real-world databases are highly susceptible to noisy, missing & inconsistent data due to their typically huge size & their likely origin from multiple heterogeneous sources

-1-

→ The Data Mining Applications focuses on

* the detection & correction of Data quality problems

* Use of Algorithms that can tolerate poor Data quality

"Detection" & "Correction" is often called "Data cleaning".

* Measurement & Data Collection Errors - refers to any problem resulting from the Measurement process.

A common problem is that the value recorded differs from the true value to some extent.

→ The term "data collection error" refers to errors such as omitting data objects or attribute values or inappropriately including data object.

* Noise and Artifacts - Noise is the random component of a Measurement Error. It may involve the distortion of a value or addition of spurious objects.



Fig: True Series



Fig: True Series with Noise

→ Data Errors may be the result of a more deterministic phenomenon. Such deterministic distortion of the data are often referred to as "Artifacts".

Ex: A Streak in same place on a set of photographs.

* Precision, Bias & Accuracy - "Precision" is a closeness of repeated measurements to one another.

"Bias" is a systematic variation of measurements from the quantity being measured.

→ Precision is often measured by the standard deviation of a set of values, while bias is measured by taking the difference between the mean of the set of values & quantity being measured.

→ It is common to use the general term, "Accuracy" to refer to the degree of measurement error in data.
Accuracy depends on precision & bias.

* Outliers :- An outlier is an observation point that is distant from other observations (Data).

An outlier may be due to variability in the measurement or it may indicate experimental error.

→ The latter are sometimes, outliers will be excluded from the data set.

An outlier can cause serious problems in analysis.

* Missing Values :- Imagine you need to analyze an Enterprise data, you note that many tuples have no recorded value for several attributes. In such case, fill the missing values by following methods

1. Ignore the tuple → Usually done when class label is missing, this method is very ineffective.
2. Fill the missing value manually → It's very time consuming & not feasible.
3. Use a global constant to fill missing value → Replace all missing values by some constant such as "unknown" or "co". This method is simple, it's not fool proof.
4. Use a measure of central tendency for the attribute (mean or median) to fill missing values.
5. Use the most probable values to fill the missing values → This may be determined with regression, interface based tools using Bayesian or Decision tree.

* Inconsistent Values :- Data can contain inconsistent values.

Some types of inconsistent values are easy to detect.

For ex :- A person's height should not be negative

Once an inconsistency has been detected, it is sometimes possible

to correct the data. The correction of an inconsistency requires additional or redundant information.

* Duplicate Data :- A data set may include data objects that are duplicates, or almost duplicates of one another.

→ To detect & eliminate such duplicates, two main issues must be addressed

* If there are two objects that actually represent a single object with different values, & these inconsistent values must be resolved

* Care to be taken to avoid accidentally combining data-objects that are similar, but not duplicates

The "Deduplication" is often used to deal with duplicate data.

* Data preprocessing :-

is a data mining technique that involves transforming raw data into an understandable format.

Real-world data is often incomplete, inconsistent, and/or lacking in certain behaviors or trends, & is likely to contain many errors.

Data preprocessing is a process method of resolving such issues.

→ Some of the most important approaches for data preprocessing are

- * Aggregation
- * Sampling
- * Dimensionality reduction
- * Feature Subset Selection
- * Feature Creation
- * Discretization & Binarization
- * Variable Transformation

These approaches fall into two categories, selecting data objects & attributes for analysis or creating/changing the attributes.