

Data Warehouse

A data warehouse is a repository of information collected from multiple sources, stored under a unified schema, and that usually resides at a single site. Data warehouses are constructed via a process of data cleaning, data integration, data transformation, data loading, and periodic data refreshing. The data are stored to provide information from a historical perspective (such as from the past 5-10 years) and are typically summarized. For example: → Storing the details of each sales transaction, the data warehouse may store a summary of the transactions per item type for each store or summarized to a higher level, for each sales regions.

A data warehouse is usually modeled by a multidimensional database structure, where each dimension corresponds to an attribute or a set of attributes in the schema, and each cell stores the values of some aggregate measure, such as count or sales-amount.

* Usage

Health care, Banking, Retail, Data Mining.

* Need of Warehouse

- Large data store
- Smart store
- efficient retrieve
- Heterogeneous

* Million - H. Inmon

- Subject oriented
- Integrated
- Non-volatile
- time-variant

* Challenges face by design Warehouse

- Data quality
- Data Analytics
- User exception
- Cost
- Performance
- Quality Assurance.

* Framework of a data Warehouse

* Advantage of data Warehouse

- To clean data
- Query processing multiple options
- High query performance
- Local processing at source warehouse.
- Easy way of reporting access multiple system.

Framework

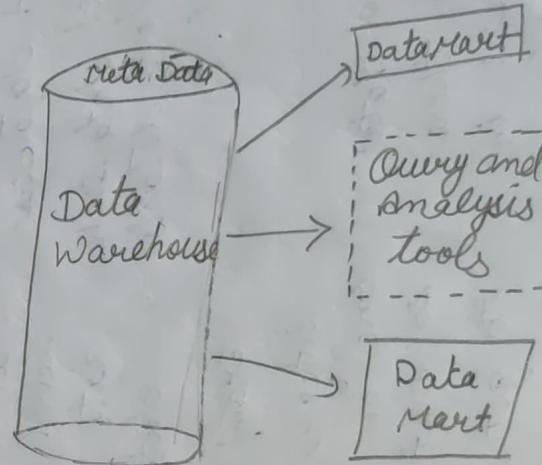
Data Source
in Chicago

Data source
New York.

Toronto

Vancouver

Clean
Integrate
Transform
Load
Refresh



DBMS V/S Data Warehouse

DBMS

- ① A common Database is based on operational or transactional processing.
- * A database stores currently and up-to-date data, which is used for daily operations.
- * A database is generally application specific.
- * Example - A database stores related data, such as the student details in a school.
- * Normalized data structure is there in a database in separate tables.
- * 100 MB to GB data store.
- * A Data warehouse is based on analytical processing.
- * A datawarehouse maintains historical data over time.
- * A datawarehouse is integrated generally at the organization level, by combining data from different database.
- Example :→ A data warehouse integrates the data from one or more databases, so that analysis be done to get result, such as the best performing school in a city.
- * Dynamic and quick analysis
- * 100 GB to KB of data is done

Introduction:-

a.i. Mining is defined as the procedure of extracting all or words we

Data Marts:^③ 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.

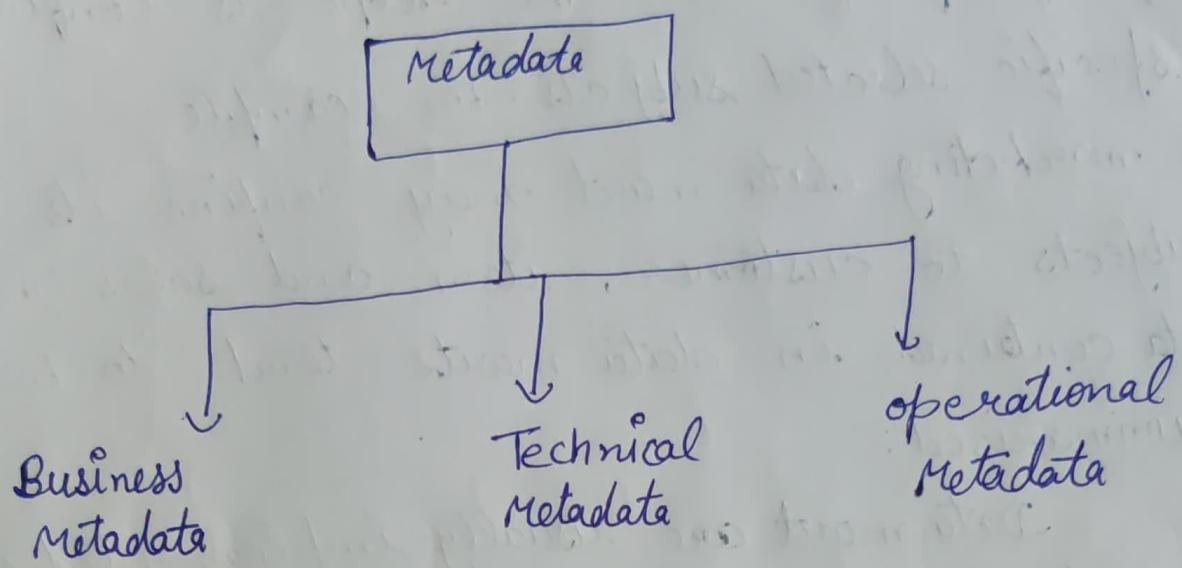
Data marts are usually implemented on low-cost departmental server that are Unix or Window based.

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

Metadata:- Metadata is simply defined as data about data. The data that is used to represent other data is known as metadata. For example:- The index of a book serves as

a metadata for the contents in the book.

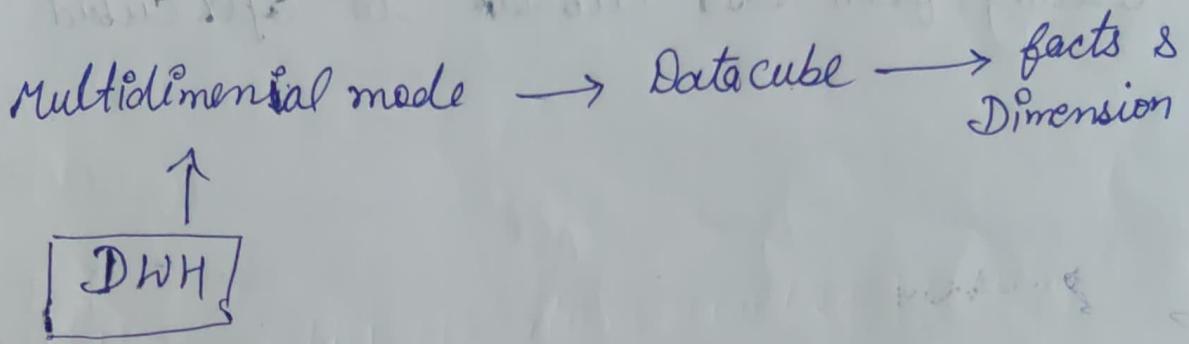
Categories of Metadata :-



- ① Business : - It has the data ownership information, business definition, and changing policies.
- ② Technical Metadata : - It includes database system names, table and column names and size, data types and allowed values. Technical Metadata also includes structural information such as primary and foreign key attributes and indices.
- ③ operational Metadata : - It includes currency of data and data lineage. Currency of data means whether the data is active, archived, or purged.

Introduction :-

Data warehouse and OLAP tool are based on multidimensional data model. In this model data is presented in the form of the data cube where as the data cube is define as represent of the data in multiple dimensions it is define in terms of fact & dimensions.



Data cube :-

When data is grouped or combined in multidimensional matrices called Data cubes.

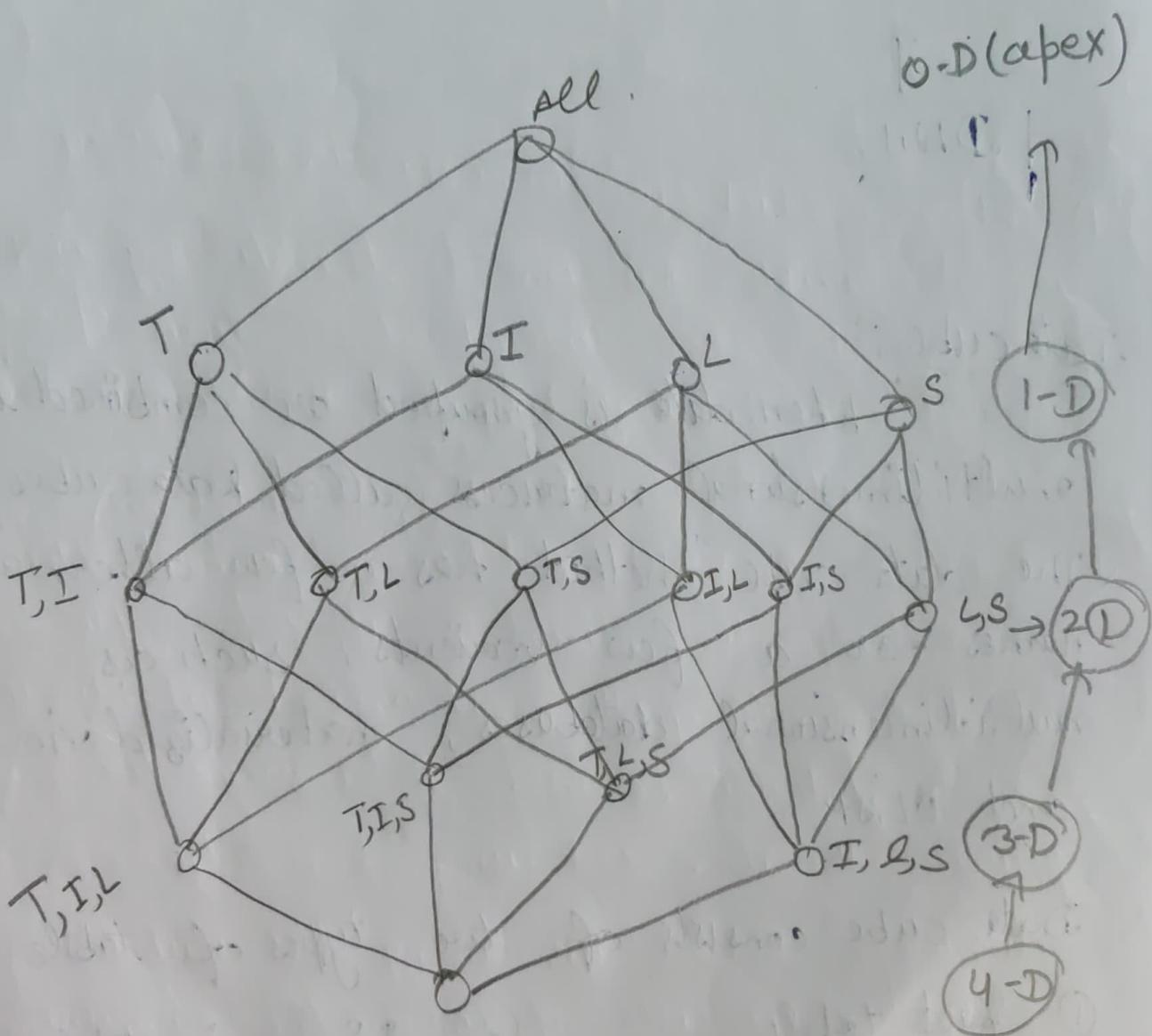
The data cube method has a few alternatives names or a few variants, such as multidimensional databases, materialized views, and OLAP

Data cube consists of two types of table

① fact table

② Dimension table

Based on different number of Dimensional or different subsets of data can be created from data cube. Such subset are known as data cuboids. The arrangement of all possible cuboids in the manner is known as data lattice. Lattice consists of different level of summarization starting from base cuboid to apex cuboid.



Time, item, location, supplier

✓ Introduction :-

Data Mining is defined as the procedure of extracting information from huge sets of data. In other words we can say that data mining is mining knowledge from data.

- There is a huge amount of data available in the information industry. This data is of no use until it is converted into useful information. It is necessary to analyze this huge amount of data and extract useful information from it.
- Extraction of information is not only the process we need to perform; Data mining also involves other processes such as data cleaning, Data integration, Data transformation, Pattern Evaluation and Data presentation. Once all these processes are over, we would be able to use this information in many applications such as
Market Analysis (M)
Fraud Detection (F)
Customer Retention (C)
Production Control (P)
Science Exploration (S)

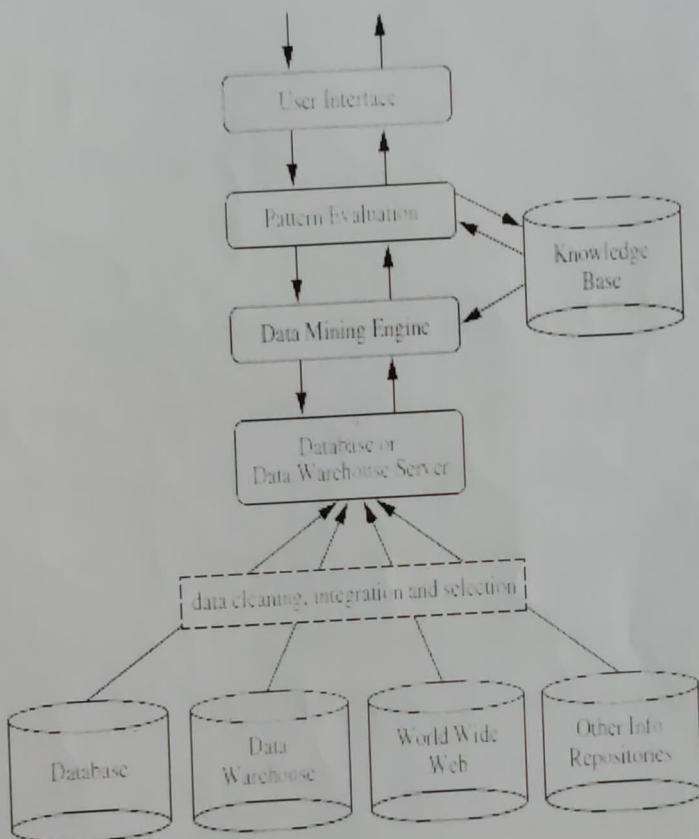
* WHAT IS DATA MINING:-

Data Mining is defined as extracting information from huge sets of data. The information or knowledge extracted can be used for many applications.

- **Summarization** – providing a more compact representation of the data set, including visualization and report generation.

1.4 Architecture of Data Mining

A typical data mining system may have the following major components.



1. Knowledge Base:

This is the domain knowledge that is used to guide the search or evaluate the interestingness of resulting patterns. Such knowledge can include concept hierarchies,

used to organize attributes or attribute values into different levels of abstraction. Knowledge such as user beliefs, which can be used to assess a pattern's interestingness based on its unexpectedness, may also be included. Other examples of domain knowledge are additional interestingness constraints or thresholds, and metadata (e.g., describing data from multiple heterogeneous sources).

2. Data Mining Engine:

This is essential to the data mining system and ideally consists of a set of functional modules for tasks such as characterization, association and correlation analysis, classification, prediction, cluster analysis, outlier analysis, and evolution analysis.

3. Pattern Evaluation Module:

This component typically employs interestingness measures interacts with the data mining modules so as to focus the search toward interesting patterns. It may use interestingness thresholds to filter out discovered patterns. Alternatively, the pattern evaluation module may be integrated with the mining module, depending on the implementation of the datamining method used. For efficient data mining, it is highly recommended to push the evaluation of pattern interestingness as deep as possible into the mining process so as to confine the search to only the interesting patterns.

4. User interface:

This module communicates between users and the data mining system, allowing the user to interact with the system by specifying a data mining query or task, providing information to help focus the search, and performing exploratory datamining based on the intermediate data mining results. In addition, this component allows the user to browse database and data warehouse schemas or data structures, evaluate mined patterns, and visualize the patterns in different forms.

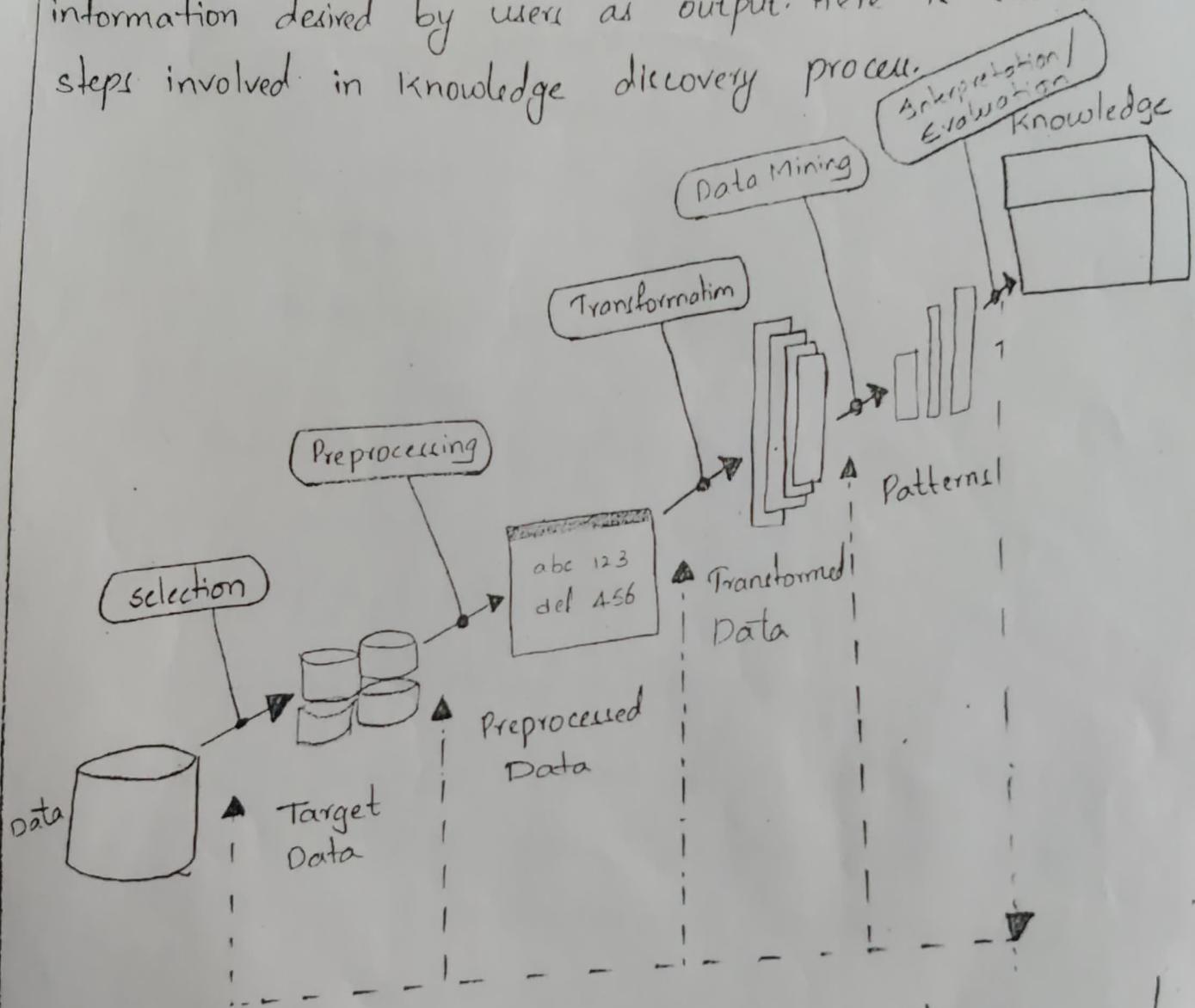
1.3 Tasks of Data Mining

Data mining involves six common classes of tasks:

- **Anomaly detection (Outlier/change/deviation detection)** – The identification of unusual data records, that might be interesting or data errors that require further investigation.
- **Association rule learning (Dependency modelling)** – Searches for relationships between variables. For example a supermarket might gather data on customer purchasing habits. Using association rule learning, the supermarket can determine which products are frequently bought together and use this information for marketing purposes. This is sometimes referred to as market basket analysis.
- **Clustering** – is the task of discovering groups and structures in the data that are in some way or another "similar", without using known structures in the data.
- **Classification** – is the task of generalizing known structure to apply to new data. For example, an e-mail program might attempt to classify an e-mail as "legitimate" or as "spam".
- **Regression** – attempts to find a function which models the data with the least error.

* KDD:-

KDD stands for Knowledge Discovery in Database. Data mining is an essential step in the process of knowledge discovery. There are seven different stages in KDD process. This process takes raw data as input and provides useful information desired by user as output. Here is the list of steps involved in knowledge discovery process.



- Preprocessing of database consists of Data cleaning and Data integration.
- KDD is an iterative process.

Data Warehouse Implementation

It is defined by number of cube or cuboid occupy the memory space as the warehouse. content huge volumes of data is become necessary to implement data-cuboid for fast data analysis and efficient query processing.

1. Efficient Computation of data cubes
 2. The compute cube operation
 3. Indexing. OLAP Data (Modeling. OLAP data)
 - 4) Efficient processing. OLAP queries.
 - 5) OLAP Server architecture.
- 1.) Efficient computation of data cubes.
- i) Data cube can be viewed as a lattice of cuboids.
The bottom-most cuboids is the base cuboids.
The top-most cuboids (apex) contains only one cell.
 - ii) Materialization of data cube:-
 - Full materialization refers to the computation of all cuboids in a data cube lattice.
 - Partial materialization refers to the selective computation of a subset of the cuboid cells in the lattice. Iceberg cubes and shell fragments are examples of partial materialization.

② Partition Rule Mining →

- i) Iceberg cube is a data cube that stores only those cube cells that have an aggregate value some minimum support threshold
- ii) Shell fragments of a data cube, only some cuboids involving a small number of dimensions are computed, and queries on additional combinations of the dimensions can be computed on-the-fly.

All

A

B

C

AB

AC

BC

ABC

iceberg cube.

→ No materialization :- Do not precompute any of the nonbase cuboids.

③ Indexing OLAP Data →

Indexing is another approach which can help efficient data searching.

Two types of indexing :-

i) Bitmap ii) Join

i) Bitmap :-

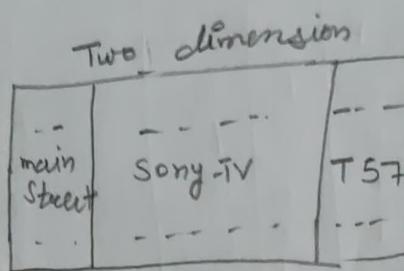
Bitmap indexing allows you to reduce the response time of the bigger SQL queries.

Join Index table

Location	Sales-key
mainstreet	T 57
"	T 238
"	T 884
---	---

Join index table

Item	Sales-key
Sony-TV	T 57
Sony-TV	T 459
---	---



④ OLAP queries:-

- Determine which operations should be perform on the available cuboids:- This involves transforming any selection, projection, roll-up, drill-down operations specified in the query into correspond SQL and OLAP operations.
- Determine to which materialized cuboid the relevant operations should be applied:- This involved identifying all of the materialized cuboids that may potentially be used to answer the query. pruning the above set using knowledge of dominance relationships among the subcubes, estimating the costs of using the remaining materialized subcubes, and selecting the subcube with the least cost.

The bitmap indexing method is popular in OLAP products because it allows quick searching in data cubes. The bitmap index is an alternative representation of the RID (Record-ID).

RID	item	city
R ₁	H	D
R ₂	C	D
R ₃	P	D
R ₄	S	D
R ₅	H	M
R ₆	C	M
R ₇	P	M
R ₈	S	M

Item bitmap				
RID	H	C	P	S.
R ₁	1	0	0	0
R ₂	0	1	0	0
R ₃	0	0	1	0
R ₄	0	0	0	1
R ₅	1	0	0	0
R ₆	0	1	0	0
R ₇	0	0	1	0
R ₈	0	0	0	1

City		
RID	D	M
R ₁	1	0
R ₂	1	0
R ₃	1	0
R ₄	1	0
R ₅	0	1
R ₆	0	1
R ₇	0	1
R ₈	0	1

H → Home entertainment

D → Delhi

C → computer

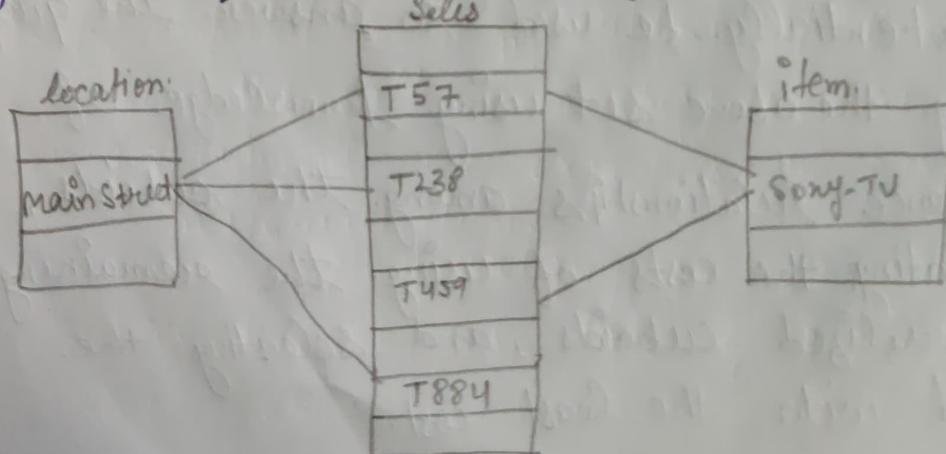
M → Mumbai

P → Phone

S → Security

Join Indexing:-

This indexing is popular used in relational query process. It registers the join row to different table for faster data access. Join indexing registers the joinable rows of two relations from a relational database.



OLAP Server Architectures.

- * **Relational OLAP (ROLAP)** :→ These are the intermediate servers that stand in between a relational back-end server and front-end tools. They use a relational or extended-relational DBMS to store and manage warehouse data and OLAP middleware to support missing pieces. ROLAP servers include optimization for each DBMS back-end, implementation of aggregation navigation logic, and additional tools and services. ROLAP technology tends to have greater scalability than MOLAP technology.
- * **Multidimensional OLAP** :→ These servers support multidimensional views of data through array-based multidimensional storage engines. They map multidimensional views directly to data cube array structure. Many OLAP servers adopt a two-level storage representation to handle sparse and dense data sets. The dense subcubes are identified and stored as array structure, while the ~~sparse~~ sparse subcubes employ compression technology for efficient storage utilization.
- * **Hybrid OLAP** :- The hybrid OLAP approach combines ROLAP and MOLAP technology, benefitting from the greater scalability of ROLAP and the faster computation of MOLAP. For example, a MOLAP server may allow large volumes of detail data to be stored in a relational database, while aggregations are kept in a separate MOLAP store.
- * **Specialized SQL Servers** :- To meet the growing demand of OLAP processing in relational databases, some relational and data warehousing firms implement specialized SQL servers which provide advanced query language and query processing support for SQL queries over star and snowflake schemas in a read-only environment.

Data-Warehouse Back-End Tools :-

- Data extraction:- which typically gathers data from multiple, heterogeneous, and external sources.
- Data cleaning:- which detects errors in the data and rectifies them when possible.
- Data transformation:- which converts data from legacy or host format to warehouse format.
- Load:- which sorts, summarizes, consolidates, computes views, checks integrity, and builds indices and partitions.
- Refresh:- which propagates the updates from the data sources to the warehouse.

Association Rule Mining :-

Association rule learning is a type of unsupervised learning technique that checks for the dependency of one data item on another data item and maps accordingly so that it can be more profitable. It tries to find some interesting relations or associations among the variables of datasets. The association rule learning is one of the very important concept of machine learning, and it is employed in market based analysis, Web usage mining, continuous production etc.

For example, if a customer buys bread he most likely can also buy butter, eggs, or milk. So these products are stored within a shelf or mostly nearby.

Association rule learning can be divided into three type of algorithms:-

- 1) Apriori
2. Eclat
- 3 F-P Growth Algorithm

Spiriori Algorithm It uses of association rules & it is designed to work on datasets or databases that contain transactions.

This algo. uses a BPS & Hash tree to calculate the itemset associations efficiently.

Frequent itemset — FI are those items whose support is greater than threshold value.

Eg:- $A = \{1, 2, 3, 4, 5\}$ $B = \{2, 3, 7\}$

2,3 are FI values.

Step-I Determine support of itemsets & Select min support & confidence.

Step-II Take all supports in transaction with higher support value than min or selected support value.

Step-III Find all rules of these sets that have higher confidence value than threshold or min confidence

Step-IV Sort the rules as decreasing order of lift.

Apriori Algorithm:-

It uses frequent itemsets to generate association rules & it is designed to work on datasets or databases that contain transactions.

This algo. uses a BFS & Hash tree to calculate the itemset associations efficiently.

Frequent itemset :- FI are those items whose support is greater than threshold value.

Eg:- $A = \{1, 2, 3, 4, 5\}$ $B = \{2, 3, 7\}$

2,3 are FI values.

Step-I Determine support of itemsets & Select min support & confidence.

Step-II Take all supports in transaction with higher support value than min or selected support value.

Step-III Find all rules of these sets that have higher confidence value than threshold or min confidence

Step-IV Sort the rules as decreasing order of lift.

Eg:-

Apriori Algo.

②

Tid.	Itemsets
T ₁	A, B
T ₂	B, D
T ₃	B, C
T ₄	A, B, D
T ₅	A, C
T ₆	B, C
T ₇	A, C
T ₈	A, B, C, E
T ₉	A, B, C

Given :-

min support = 2.

min confidence = 50%.

C → customer

FI → Frequent itemset

Step ① Calculating C₁ and F₁

Itemset	Support count
A	6
B	7
C	6
D	2
E	1

→ F₁

Itemset	SC
A	6
B	7
C	6
D	2

Step ② Calculating C₂ and F₂.

Itemset	SC
(A, B)	4
(A, C)	4
(A, D)	1
(B, C)	4
(B, D)	2
(C, D)	0

→ F₂

Itemset	SC
(A, B)	4
(B, C)	4
(A, C)	4
(B, D)	2

Step ③ Candidate Generation $C_3 \& F_3$

	SC
(A, B, C)	2
(B, C, D)	0
(A, C, D)	0
(A, B, D)	1

F_3

Itemset	SC
(A, B, C)	2

Step 4) Finding Association rules for subsets:-

→ We will calculate the confidence using

$$\frac{\text{Sup}(A \cap B)}{A} \rightarrow \frac{\text{Sup}((A \cap B) \cap C)}{\text{Sup}(A \cap B)} \cdot \text{Sup} = \frac{A \rightarrow B}{A}$$

→ After calculating we will exclude rules that have less ~~than~~ confidence than min threshold (50%).

Rules	Support	Confidence.
$A \cap B \rightarrow C$	2	$\text{Sup}((A \cap B) \cap C) / \text{Sup}(A \cap B) = 2/4 = 50\%$
$B \cap C \rightarrow A$	2	$\text{Sup}((B \cap C) \cap A) / \text{Sup}(B \cap C) = 2/4 = 50\%$
$A \cap C \rightarrow B$	2	$\text{Sup}((A \cap C) \cap B) / \text{Sup}(A \cap C) = 2/4 = 50\%$
$C \rightarrow A \cap B$	2	$\text{Sup}((C \cap (A \cap B)) / \text{Sup}(C) = 2/5 = 40\%$
$A \rightarrow B \cap C$	2	$\text{Sup}((A \cap (B \cap C)) / \text{Sup}(A) = 2/6 = 33.33\%$
$B \rightarrow A \cap C$	2	$\text{Sup}((B \cap (A \cap C)) / \text{Sup}(B) = 2/7 = 28.57\%$

Frequent Pattern Growth Algo! → It is improved version of Apriori Algo. It represents the database in the form of tree structure that is known as frequent pattern or tree.

Tid	Item
T ₁	E, K, M, N, O, Y
T ₂	D, E, K, N, O, Y
T ₃	A, E, K, M
T ₄	C, K, M, U, Y
T ₅	C, E, I, K, O, O

min Support = 3.

②

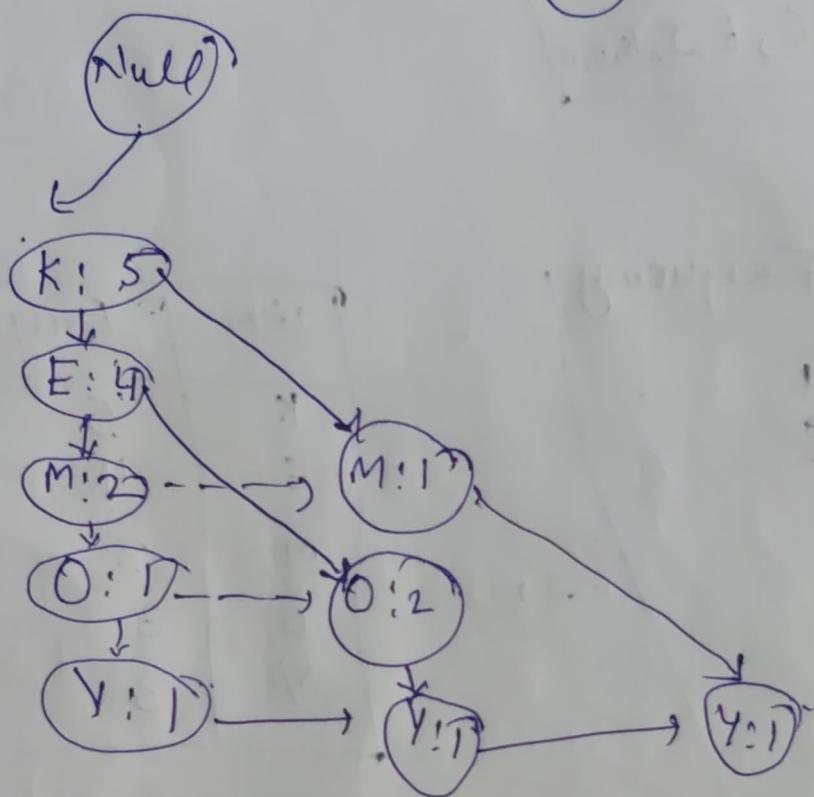
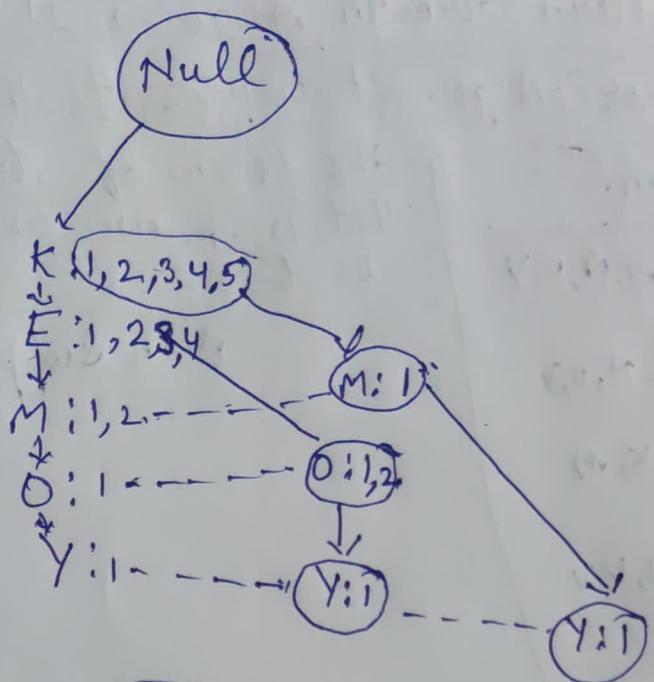
Item	Frequency
A	1
C	2
D	1
E	4
I	1
K	5
M	3
N	2
O	1
U	3
Y	3



Item	Frequency
K	5
E	4
M	3
O	3
Y	3

③

Tid	Items	Needed - Itemset
T ₁	E, K, M, N, O, Y	K, E, M, O, Y
T ₂	D, E, K, N, O, Y	K, E, O, Y
T ₃	A, E, K, M	K, E, M
T ₄	C, K, M, U, Y	K, M, Y
T ₅	C, E, I, K, O, O	K, E, O



Item.	Conditional Pattern Base	Conditioned Freq.
Y.	{K, E, M, O: 1}, {K, E, O: 1}, {K, M, Y: 1}	3K: 33
O	{K, E, M: 1}, {K, E: 2}	3K, E: 33
M.	{K, E: 2}, {K: 1}	3K: 33
E	{K: 4}	3K: 43

Items.

Frequent Pattern Generated

(4)

Y

$\{<K, Y: 3>\}$

O

$\{<K, O: 3>\}, \{<E, O: 3>\}, \{<K, E, O: 3>\}$

M

$\{<K, M: 3>\}$

E

$\{<K, E: 3>\}$

K

—

$K \rightarrow Y, Y \rightarrow K$

$K \rightarrow O, E \rightarrow O, O \rightarrow K, O \rightarrow E, E \rightarrow K, K \rightarrow E$

$K \rightarrow M, M \not\rightarrow K$.

$K \rightarrow E, E \rightarrow K$.

* Eclat Algorithm

Eclat Algorithm stands for Equivalence class clustering and bottom-up lattice traversal.

It is one of popular methods of Association rule mining algorithm.

Transaction id	Bread	Butter	Milk	Coke	Jam
T ₁	1	1	0	0	1
T ₂	0	1	0	1	0
T ₃	0	1	1	0	0
T ₄	1	1	0	1	0
T ₅	1	0	1	0	0
T ₆	0	1	1	0	0
T ₇	1	0	1	0	0
T ₈	1	1	1	0	1

$K=1$, Minimum Support = 2.

$K=2$, MS = 2.

Item	Tidset
Bread	$T_1, T_4, T_5, T_7, T_8, T_9$
Butter	$T_1, T_2, T_3, T_4, T_6, T_8, T_9$
Milk	$T_3, T_5, T_6, T_7, T_8, T_9$
Cookie	T_2, T_4
Jam	T_1, T_8

Item	Tidset
(B, B)	T_1, T_4, T_8, T_9
(B, M)	T_5, T_7, T_8, T_9
(B, C) \times	T_4
(B, J)	T_1, T_8
(B, U, M)	T_3, T_6, T_8, T_9
(B, U, C)	T_2, T_4
(B, U, J)	T_1, T_8
(M, C) \times	0
(M, J) \times	T_8
(C, J) \times	0

$K=3$, MS = 2.

Item	Tidset
B, BU, M	T_8, T_9
B, BU, J	T_8, T_1
B, BU, C \times	T_4
B, M, J \times	T_8

$K=4$.

Item	Tidset
B, BU, M, J	T_8

Advantages:

- **Memory Requirements:** → ECLAT Algo. uses a Depth-first Search Approach, it uses less memory than Apriori.
- **Speed:** → The ECLAT Algo. is typically faster than Apriori Algo.
- **Number of Computations:** → The ECLAT Algo. does not involve the repeated scanning of the data to compute the individual support values.

Decision Tree :-

Decision tree algorithm falls under the category of supervised learning. It can be used to solve both regression and classification problems. It is a tree that helps us in decision making purpose. The decision tree creates classification or regression models as a tree structure. Decision tree uses the tree representation to solve the problem. Decision tree contains ~~tree~~ 3 types of nodes.

① Root node:- It is the top most node in the tree. Data which is inside the node is known as attribute.

② Internal node:- Each internal node denotes a test on attribute. Nodes which are in between root node and leaf nodes are called as internal nodes.

③ Leaf node:- We call last nodes as leaf nodes. They represents output is class label.

- Advantages :-
- ① It does not require any domain knowledge.
 - ② Classification steps of decision tree are simple and fast.
 - ③ Missing values in data does not effect output.
 - ④ A decision tree model is automatic and does not require a standardization of data.

key factors:- Building a decision tree is all about discovery. attributes that return the highest data gain

Entropy:-

Entropy refers to a common way to measure impurity. In the decision tree, it measures the impurity in dataset.

Information Gain:-

Information Gain refers to the ~~define~~ decline in entropy after the dataset is split. It is also called as entropy reduction.

Example:-

Day.	Weather.	Temperature	Humidity	Wind.	Play.
1	Sunny.	Hot	High	Weak	No
2	Cloudy.	Hot	High.	Weak	Yes.
3	Cloudy	mild.	High	Strong.	Yes
4	Rainy.	mild	High.	Strong	No
5	Sunny.	mild.	Normal	Strong	Yes.
6	Rainy.	Cool	Normal	Strong	No
7	Rainy.	mild.	High	Strong	Yes
8	Sunny.	Hot	High	Weak.	No
9	Cloudy	Hot	Normal	Strong	Yes
10	Rainy.	mild.	High	Strong	No

