

"Data Mining and warehousing"

Unit #1

Data Warehousing :-

Data warehousing is the process of constructing and using a data warehouse.

A Data warehouse is constructed by integrating data from multiple heterogeneous sources that support analytical reporting, structured and decision making.

Data warehousing involves data integration, data cleaning & data consolidations.

Need :-

• Consolidated Data :

Collects data from multiple sources into one place, making it easier for analysis.

• Faster Decision Making :

Provides timely and accurate insights for business decisions.

• Historical Analysis :

Allows businesses to analyze trends over time.

• Improved Data Quality :

Cleans and structures data, ensuring consistency.

Characteristics =>

- Subject - oriented
- Integrated
- Time - Variant
- Non - Volatile.

Functions =>

- Data Consolidation
- Data Analysis
- Data Mining
- OLAP

Types =>

There are 3 types of data warehouses :

① Enterprise Data Warehouse (EDW) =>

A Centralized warehouse Containing all the organizational data for reporting across departments.

② Data Mart :

A Smaller, focused subset of the data warehouse, catering to specific departments or business units.

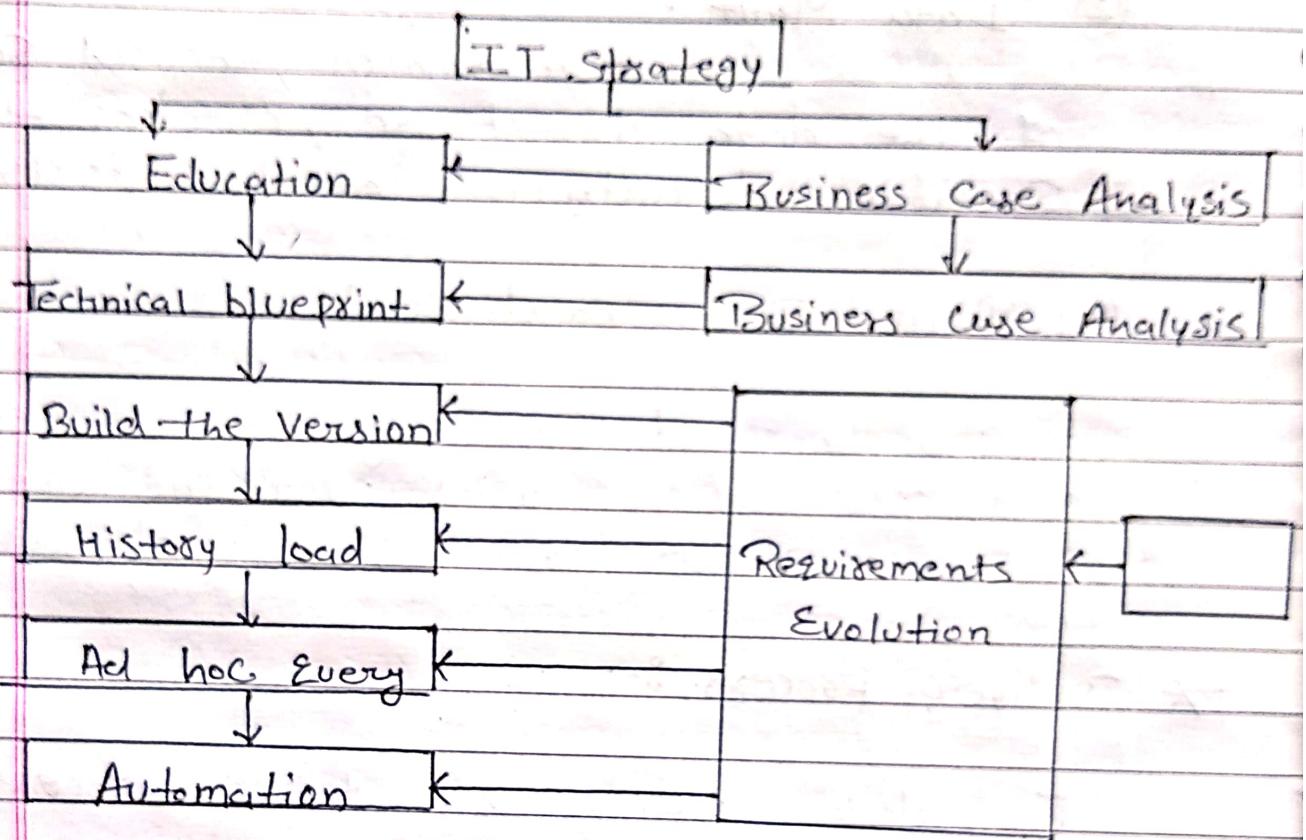
③ Virtual Data Warehouse :

A warehouse built using virtual views and real-time queries over the operational database rather than storing physical data.

* Delivery Process :-

The data warehousing delivery process involves the following steps :

- Extracting and loading data from different source systems.
- Cleaning and transforming the data into a form suitable for analysis.
- Generating aggregations from predefined definitions within the data warehouse.
- Determining aggregations to maintain Sys performance.
- Backing up, restoring & ~~archiving~~ archiving the data.



→ Data Warehouse Architecture :-

Data Warehouse Architecture uses a structured framework to manage and store data effectively.

There are two common approaches to constructing a data warehouse:

→ Top - Down Approach

→ Bottom - Up Approach

Components of Data Warehouse Architecture =>

A data warehouse architecture consists of several key components that work together to store, manage and analyze data.

* External Source :

External sources are where data originates. These sources provide a variety of data types, such as structured data, semi-structured data or unstructured data.

* Staging Area :

The staging area is a temporary space where raw data from external sources is validated and prepared before entering the data warehouse.

To handle this preparation effectively, ETL (Extract, Transform, Load) tools are used.

* Data Warehouse :

The data warehouse acts as the central repository for storing cleaned & organized data.

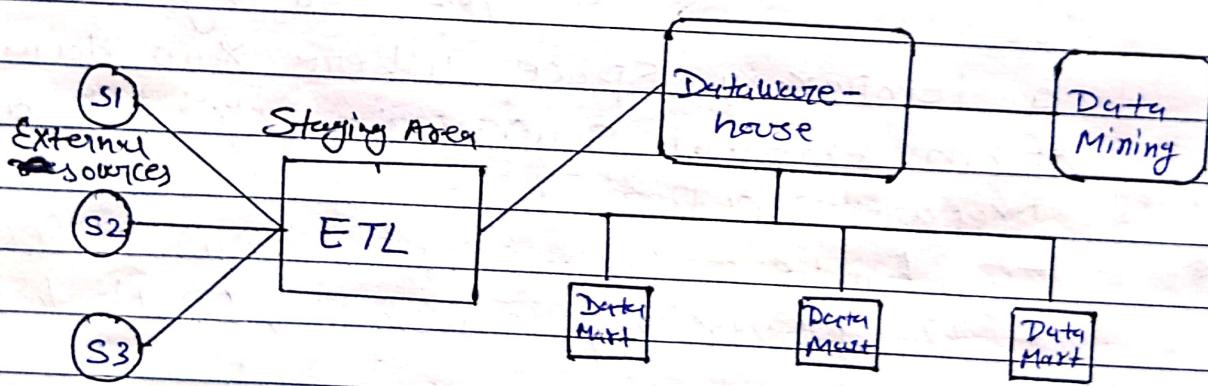
* Data Marts :

A data Mart is a subset of data warehouse that store data for a specific team or purpose, like sales or marketing.

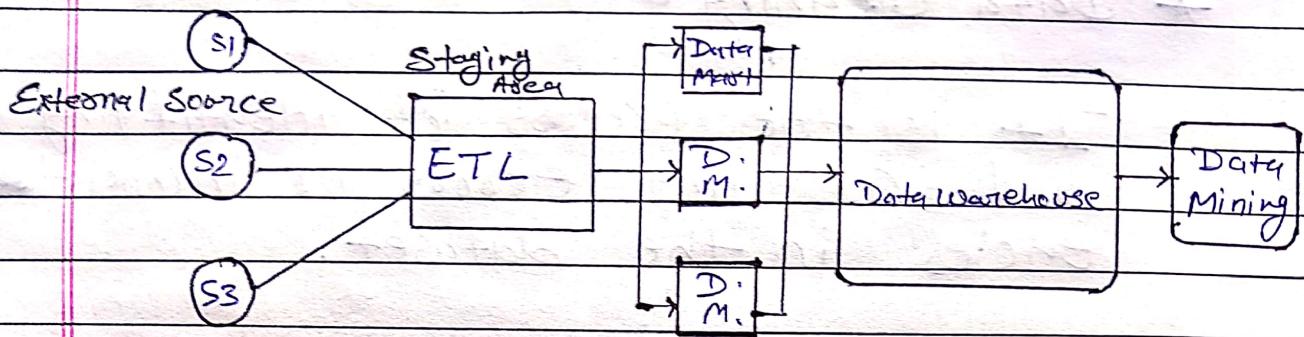
* Data Mining :

Data mining is the process of analyzing large datasets stored in the data warehouse to uncover meaningful patterns, trends, and insights.

① Top - Down Approach Archi. =>



② Bottom - UP =>



Data Pre-Processing :-

Data Pre-processing is the process of preparing raw data for analysis by cleaning and transforming it into a useable format.

- Goal is to improve the quality of data
- Helps in handling missing values, removing duplicates and normalizing data.
- Ensures the accuracy and consistency of the dataset.

Steps in Data Preprocessing :-

Some key steps in Data Preprocessing are Data Cleaning, Data Integration, Data Transformation and Data Reduction.

* Data Cleaning :-

It is the process of identifying and correcting errors or inconsistencies in the dataset.

It involves handling missing values, removing duplicates and correcting incorrect or outlier data to ensure the dataset is accurate and reliable.

Clean data is essential for effective analysis, as it improves the quality of results and enhances the performance of data models.

• Missing Values =>

This occurs when data is absent from a dataset. You can either ignore the row with missing data or fill the gaps manually, with the attribute mean or by using the most probable value.

• Noisy Data =>

It refers to irrelevant or incorrect data that is difficult for machines to interpret, often caused by errors in data collection or entry.

It can be handle in several ways:

- Binning Method
- Regression
- Clustering
- Removing Duplicates :

It involves identifying and eliminating repeated data entries to ensure accuracy and consistency in the dataset.

* Data Integration :-

It involves merging data from various sources into a single, unified dataset.

It can be challenging due to differences in data formats, structures and meanings. Techniques like record linkage and data fusion help in combining data efficiently, ensuring consistency & accuracy.

Record linkage :-

R.L. is the process of identifying & matching records from diff. datasets that refer to the same entity, even if they are represented differently.

It helps in Combining data from Various Sources by finding corresponding records based on Common identifiers.

Data Fusion :-

DF involves Combining data from multiple Sources to create a more Comprehensive and accurate dataset.

It integrates info. that may be inconsistent or incomplete from different sources, ensuring a unified & richer dataset for analysis.

* Data Transformation :-

It involves Converting data into a format suitable for analysis. Common techniques include normalization, which scales data to a common range.

Standardization, which adjusts data to have zero mean and unit variance.

and discretization, which converts continuous data into discrete categories.

These techniques help prepare the data for more accurate analysis.

- Data Normalization
- Discretization
- Data Aggregation
- Concept Hierarchy Generation

→ Data Reduction \Rightarrow

It reduces the dataset's size while maintaining key info. This can be done through feature selection, which chooses the most relevant features, and feature extraction, which transform the data into a lower dimensional space while preserving imp. details. It uses various reduction techniques such as.

→ Dimensionality Reduction (eg PCA)

→ Numerosity Reduction

→ Data Compression

→ Uses of Data Preprocessing :

- Data Warehousing
- Data Mining
- M.L.

- Data Science
- Web Mining
- Business Intelligence (BI)
- Deep Learning purpose.

Advantages :

- Improved data quality
- Better model performance
- Efficient data Analysis
- Enhanced Decision making

Disadvantages :

- Time - Consuming
- Resource - Intensive
- Potential Data Loss
- Complexity.

II. Data Warehousing Schemas :-

Schema :-

Schema is a logical description of the entire database. It includes the name and description of records of all record types including all associated data items and aggregates.

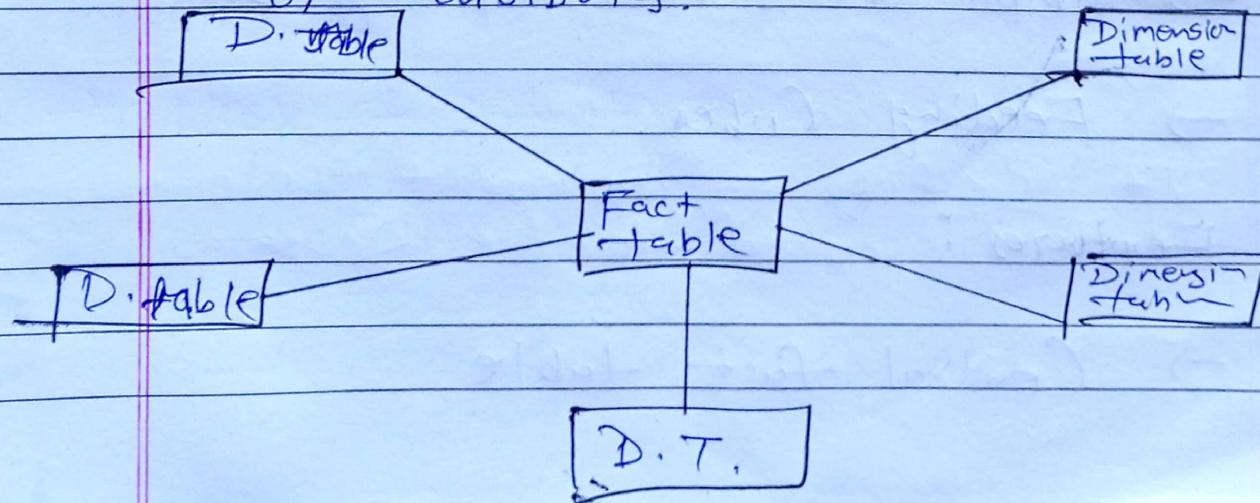
Much like a database, a data warehouse also requires to maintain a Schema.

A database uses relational model, while a data warehouse uses Star, Snowflake and fact Constellation Schema.

① Star Schema :-

→ Each dimension in a Star schema is represented with only one dimension table.

→ This dimension table contains the set of attributes.



- There is a fact-table at the center. It contains the keys to each of four dimensions.
- The fact-table also contains the attributes.
- This Schema is widely used to develop or build a data warehouse dimensional data marts.
- A Star Schema having multiple dimensions is termed as Centipede Schema.
- It is easy to handle a star Schema which have dimensions of few attributes.

Advantages :

- Simpler Queries
- Simplified Business Reporting Logic
- Feeding Cubes.

Features :

- Central fact-table

- Dimension tables
- Denormalized structure
- Simple queries
- Aggregated data
- Fast performance
- Easy to understand

ii) Snowflake Schema :-

A Snowflake Schema is a type of data modeling technique used in data warehousing to represent data in a structured way that is optimized for querying large amounts of data efficiently.

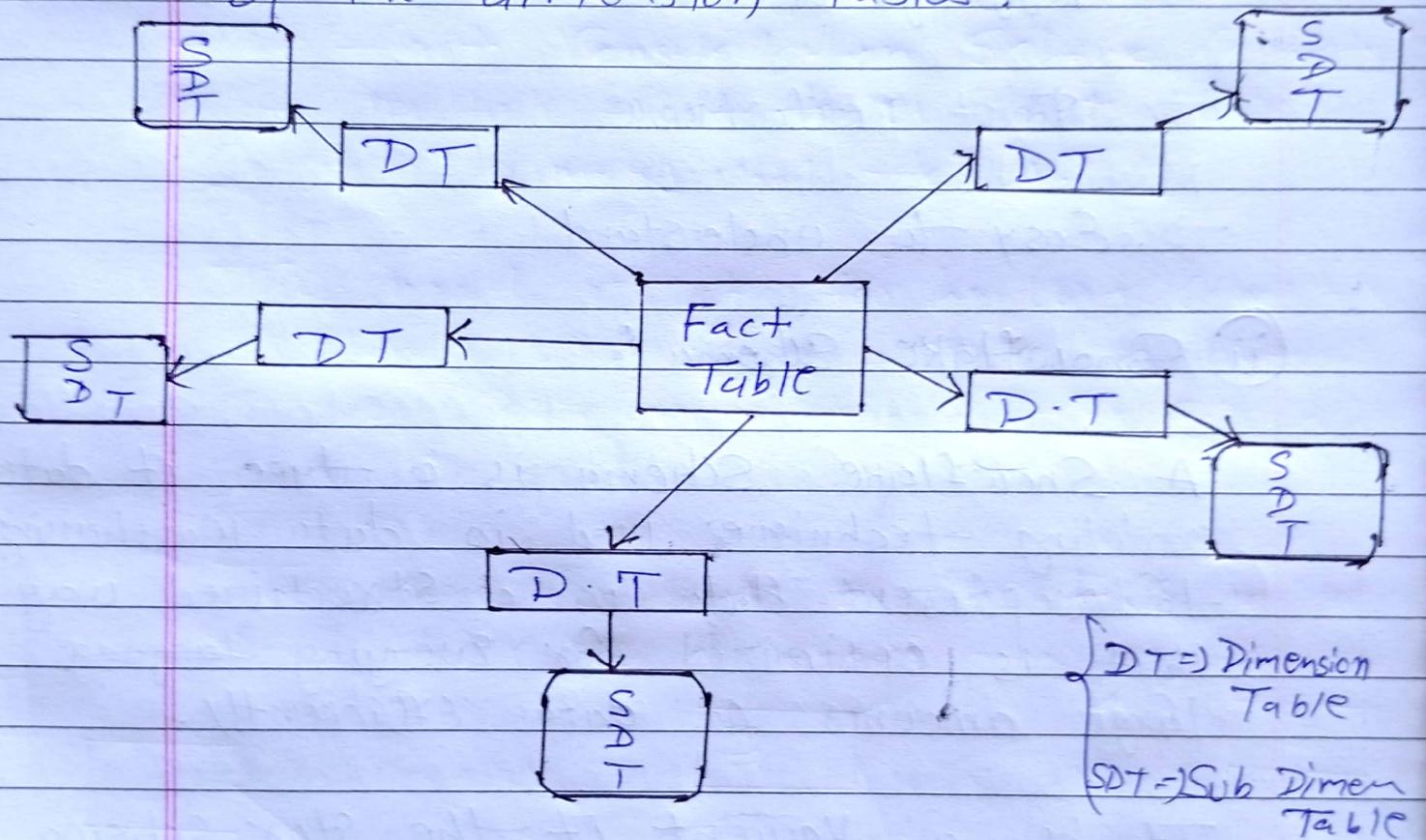
It is a variant of the Star Schema.

Here, the centralized fact table is connected to multiple related tables.

The Snowflake effect affects only the dimension tables & does not affect the fact tables.

The dimension tables are normalized into multiple related tables, creating a hierarchical or "Snowflake" structure.

The fact table is still located at the center of the schema, surrounded by the dimension tables.



Characteristics :-

- The Snowflake schema uses small disk space.
- It is easy to implement the dimension that is added to the schema.

- There are multiple tables, so performance is reduced.
- The dimension table consists of 100 or more sets of attributes.

Features :-

- Normalization
- Hierarchical structure
- Multiple levels
- Joins
- Scalability

Advantages :-

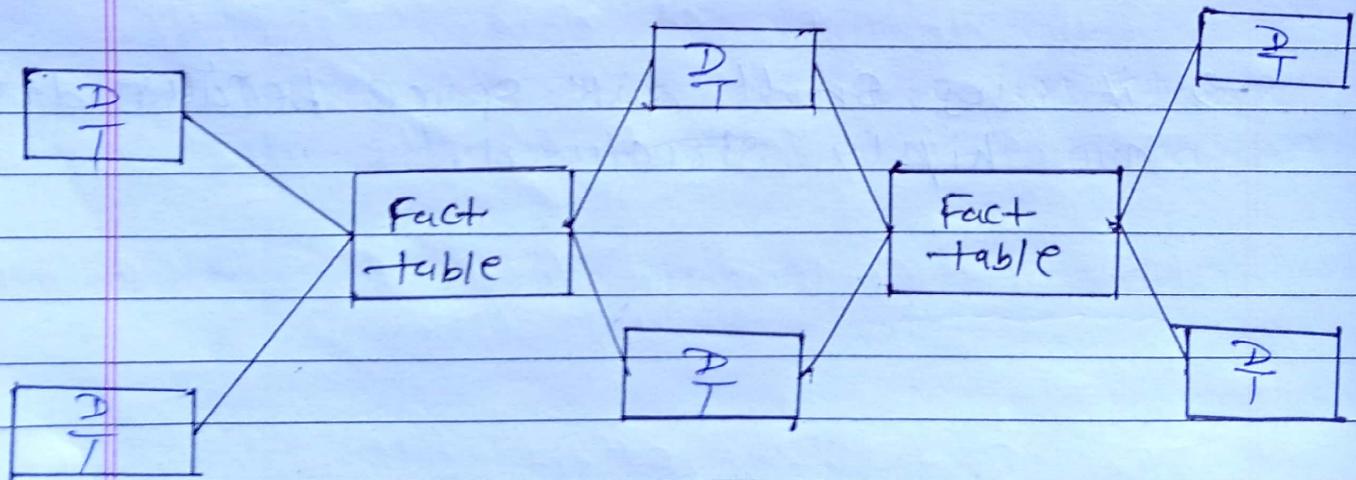
- It provides structured data which reduces the problem of data integrity
- it uses small disk space because data are highly structured.

iii) Fact Constellation :-

Fact Constellation Schema, also known as the Galaxy Schema, is an advanced data modeling technique used in designing data warehouse. Unlike simpler model like Star Schema and Snowflake Schema, the fact Constellation Schema consists of multiple fact tables that share common dimensional tables.

This model is ideal for handling complex systems and large-scale analytical queries, offering flexibility for BI & data mining.

The core components of the FCS include Fact table & dimension tables.



Benefits \Rightarrow

- \rightarrow Enhanced Query Performance
- \rightarrow Flexibility in Reporting & Analysis
- \rightarrow Improved Scalability
- \rightarrow Simplified Data Management
- \rightarrow Enhanced Data Consistency
- \rightarrow Support for Complex Analytical Query

Challenges \Rightarrow

- \rightarrow Increased Complexity in design
- \rightarrow Performance issues with Complex Queries
- \rightarrow Difficulty in maintaining Consistency
- \rightarrow Data Redundancy & Storage Overhead

Feature	Star Schema	Snowflake Schema
Structure	Central fact-table Connected to dimension tables	Fact-table connected to normalized dimension tables.
Data normalization	Denormalized dimension tables	Normalized dimension tables.
Performance	Faster query execution	Slower query performance due to fewer joins
Design Complexity	Simple & easy to understand	Complex design with multiple levels of normalization
Space Usage	Uses more storage due to denormalization	Use less storage due to normalization
Data Redundancy	Higher data redundancy	Lower data redundancy
Foreign Keys	fewer foreign keys	More foreign keys
Learning Curve	Easier to learn & implement	More complex to learn & implement

Pattern Warehousing :-

A pattern Warehouse is a specialized repo. that stores patterns identified through data mining processes in a persistent manner. Unlike traditional data warehouses, which focus on storing vast amounts of raw data for historical analysis, a pattern Warehouse consolidates refined and relevant patterns that carry significant semantic info, enabling quicker & more insightful decision making.

Key features \Rightarrow

- Stored Patterns for Later use
- Context - Based Org.
- Fast & Smart Searching
- Works with Data Mining

Compare \Rightarrow

\rightarrow Data Warehouse stores large amounts of raw data while pattern Wm stores summaries or patterns from that data.

→ DW good for reports & post data analysis
(while PW better for spotting trends
making decisions quickly.)

Application →

→ Business

→ Healthcare

→ Finance

→ Social Media

Conclusion ⇒

Pattern warehousing is like moving from just storing everything to storing only the useful patterns.

This helps business & org make faster, smarter decisions based on past behaviors & trends.

info. in more than two dimensions,
these multidimensional data by respective
Online analytical Processing (OLAP) systems

OLAP Analysis

Sales.

Popularity or how product placement impacts
QoS. Such as which color products are more
OLAP combines the attributes to answer
total sales value, in a different system.
The name of the items ordered a
collection of customer purchase data, such as
, size, cost, and location. The relevant QoS
The products it sells, such as color
for ex: a detailed sales data about all

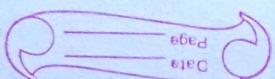
OLAP Combines & groups this data into
categories to provide actionable insights
for Strategic Planning.

Meter, and internal system.
data sources, such as website, app., smart
obj's collect & store data from multiple
businesses, and internal system.

is soft. technology you can use to analyze
OLAP. Online analytical processing (OLAP)

OLAP :-

Unit :- 2



An OLAP Sys. works by Collecting, Organizing, Analyzing data using the following steps:

How does OLAP work?

- * OLAP Analytic tools
- * OLAP Cubes
- * OLAP database
- * OLAP Server
- * ETL tools
- * Data warehouse

This consists of the following steps:

← Time
← Location
← Product Type

dimensions:

Sales might consist of the following categories, multidimensional data for products measured how multiple characteristics involve columns and rows, but multidimensional data.

The MOLAP sys. stores pre-calculated data in the hypercube. Data engine needs use MOLAP because this type of OLAP technology provides fast analysis.

MOLAP involves creating a cube that represents multidimensional data from a data warehouse.

* Multidimensional OLAP (MOLAP) :

OLAP sys. operate in 3 main ways :

Types of OLAP :-

(3) Business analysts use OLAP tools to query & generate reports from the multidimensional data in the OLAP cube.

(2) Then, the ETL tools clean / aggregate, precalculate & store data in an OLAP cube according to the no of dimensions specified.

(1) The OLAP Server collects data from multiple data sources, including relational databases & data warehouses.



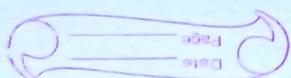
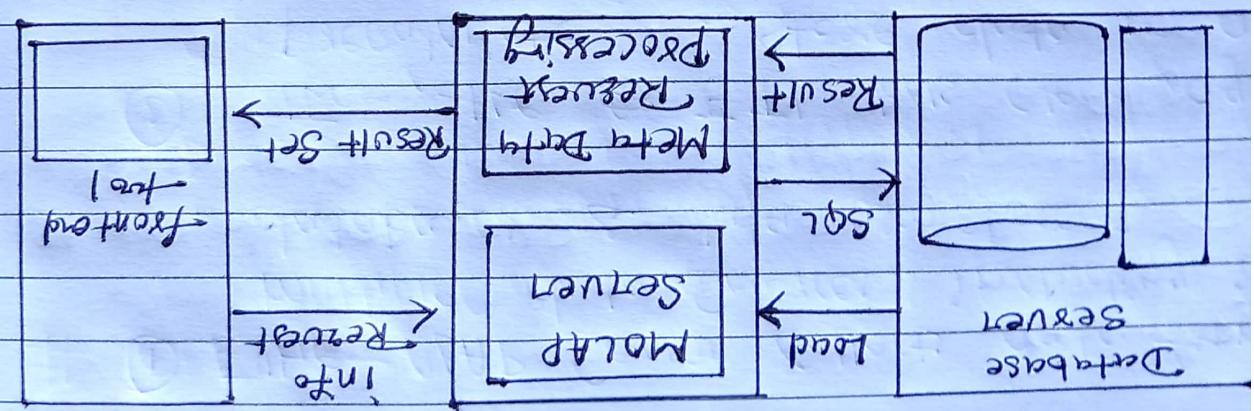
ROLAP is suitable for analyzing Extreme & detailed data. However, ROLAP has slow query performance compared to MOLAP.

In other words, data engineers use SQL queries to select for and retrieve specific information based on the required dimensions.

Instead of using a data cube, ROLAP allows data engineers to perform multi-dimensional analysis on a relational database.

* Relational OLAP (ROLAP) :-

MOLAP



if the roll up op.
Till down is the opposite
Till down =

summarizes the data for specific attribute.
In roll up, the OLAP sys.

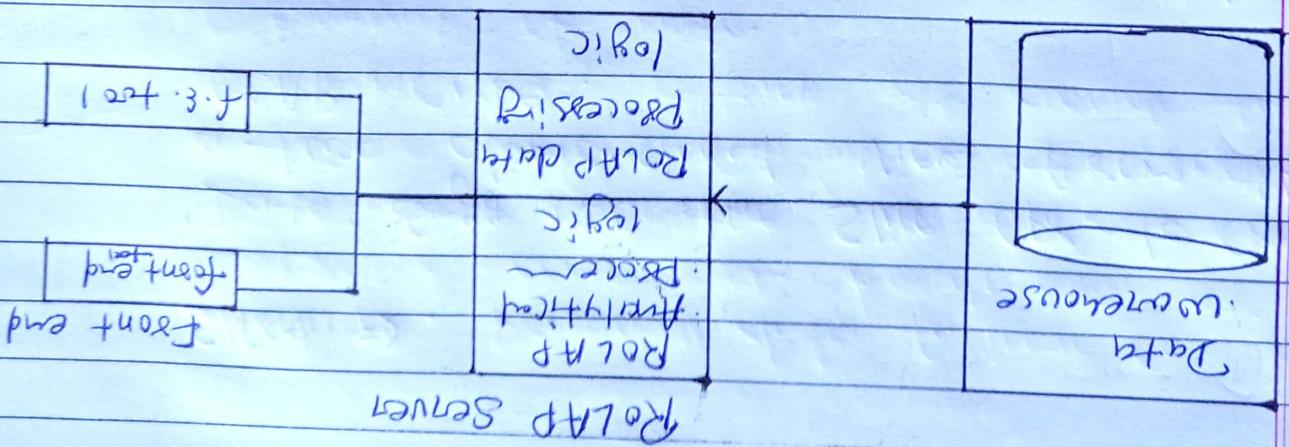
Roll up =

analytical op. with a multi OLAP cube
Business analysts perform several basic

OLAP Operations :-

info from relational database
from a data cube & extract detailed
to quickly retrieve analytical results
of both ends. HOLAP allows data entry.
MOLAP & ROLAP to provide the best
Hybrid online analytical processing combines

* Hybrid OLAP (HOLAP) :-



\Rightarrow Slice

By slicing the cube, different components
create a spread sheet - like table
consisting of products & chemicals for
a specific month.

for ex, a MOLAP cube consists of products / chemicals & month
according to products / chemicals & month

Data entry uses the slice app to create
OLAP cube.
a two dimensional view from the
cube dimensions

The goal of the backup is to make a copy of data that can be restored in the event of a primary data loss.

Backup & Recovery define the process of backing up records in the system that are needed for data loss. Tracing up data because it is applicable in case of data copying & archiving computer info., so that it is possible that data recovery because of a loss & setting up system that a user can recover easily.

Data warehouse Backup & Recovery:

- X-axis - location
- Y-axis - time
- Z-axis - product

Following definition:

Upon a plot, the DLP cube has the

- X-axis - product
- Y-axis - location
- Z-axis - time

Executive axes:

thus the following dimensions on the DLP cube for ex., a three-dimensional

The system has to be up
before the software engine is up

It is a backup that
is triggered by CPU.

The entire database
is packed up simultaneously. This
methodical order of data files (Central file),
is followed up sequentially. The entire database
is backed up sequentially.

This is the backup of the database —

File 1 is the header file, which contains the
information of index, as well as the address of
the memory address. File 2 is the first group of
records, which contains the first 10 records of the
group. File 3 is the second group of records, which
contains the next 10 records of the group. This
process continues until the last group of records
which contains the remaining records of the table.



↳ Simple Recovery Model

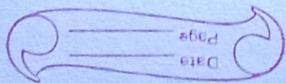
→ Bulk - Logged Recovery Model

→ Full Recovery Model

The model for a database can be changed after the database has been created.

The recovery model of a current database is inherited from the model database when the new database is generated.

Recovery is the phase of reconstituting a database after some element of a database has been hidden.



Used to Predict Continuously
Values, such as house prices or stock
prices. Common regression algos.
Include linear regression,
Polynomial regression,
SVM

① Regression:

SL is classified into two categories
of algo:

Type =

In SL the machine is trained on
a set of labeled data, which means
that the IP data is paired with
the detailed O/P. The machine then
~~predicts~~ learns to predict the O/P
for new I/P data.

Supervised learning is a Subcategory
of machine learning and AI that
involves training algos on labeled
datasets to classify data as predict-
outcomes accurately.

"Supervised Learning"

Unit = 4

There are 3 common types of data marts:

Types =

Data marts are small in scale & scope, typically holding relevant data for a specific group of users, such as sales, marketing or finance.

It provides a simplified & targeted view of data, address specific reporting & analytical needs.

A data mart is a specialized subset of data warehouse, designed to focus on a specific functional area of department within an organization.

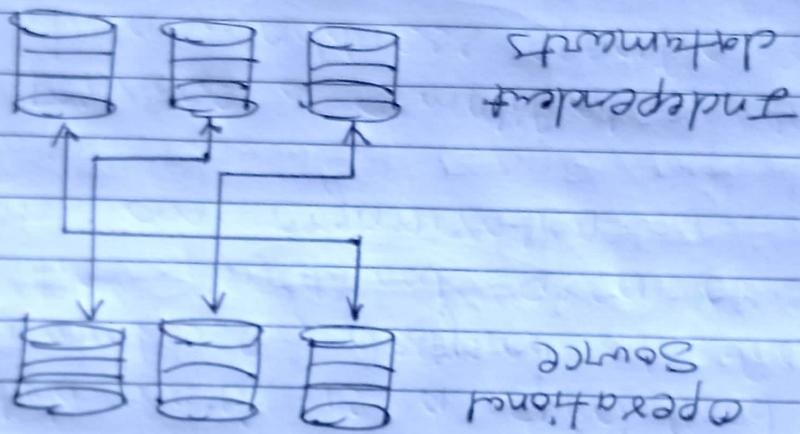
Q: What is the imp of Data Marts in DW?

⑥ Classification is used to predict categorical values, such as whether an email is spam or not. Common classification algos include Logistic regression, SVM, Decision Trees, Random Forests & Naive Bayes.

Dependent data units benefit from the data integration, quality and consistency provided by the data warehouse.

Generates directly from a data warehouse, it classifies data to meet the needs of a specific industry.

② Dependent Data

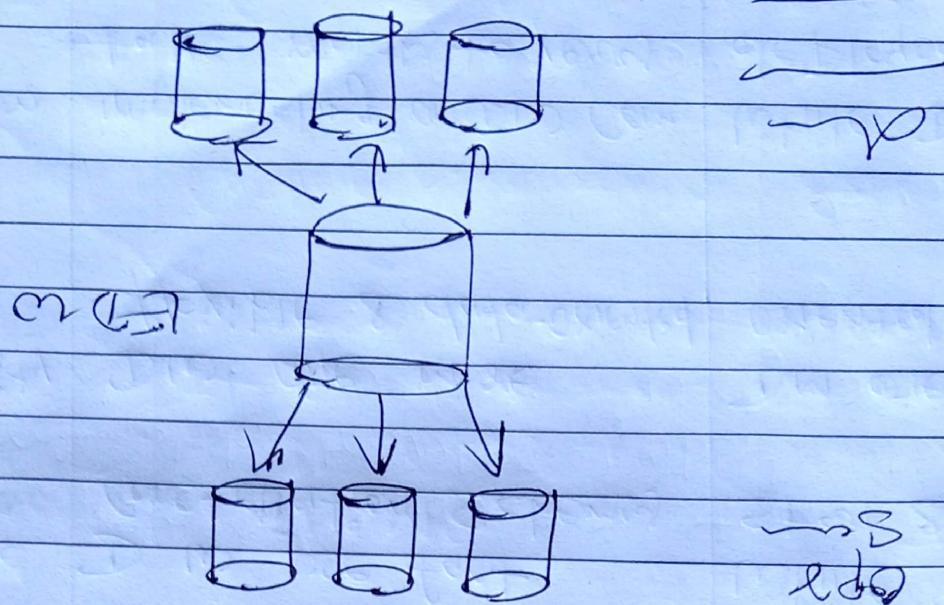


But many result in data redundancy and after flexibility & agility.

Creates and maintains separate specific business units at departmental level. From the data warehouse, it satisfies the particular needs of a specific business unit at department level.

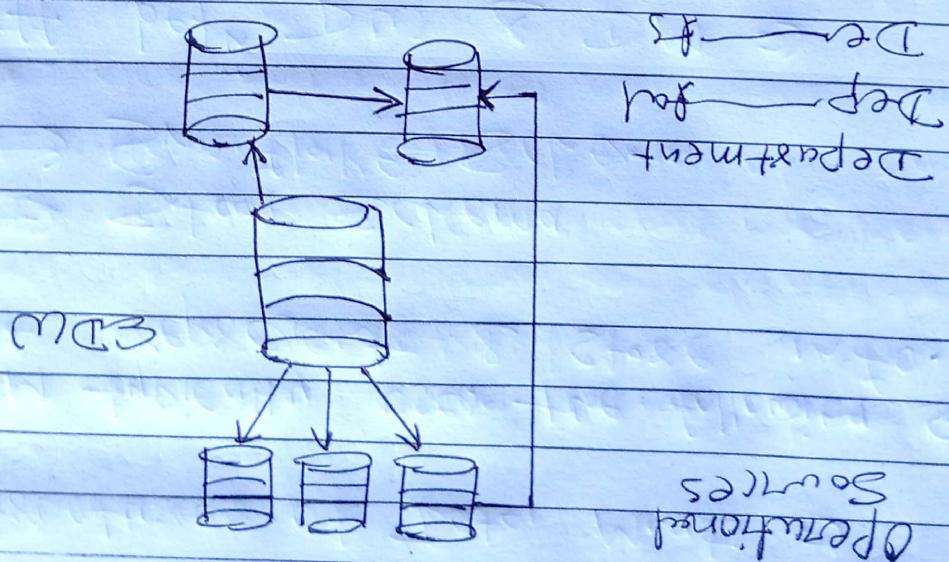
① Independent Data Mart





Combines Compartments of both Index & data
data musts If uses the centralized
& Consistent while incorporation
addition data sources specific to a
business unit or department

③ Hybrid T-M :-



Structures of Data Marts =>

DM typically uses the following structure
to represent & store info. :

→ Star Schema
→ Snowflake Schema

Diff b/w DW & DM =>

Feature	D.W.	D.M.
Scope	DW are large, centralized sys that integrate data from various sources within an org.	While DM are smaller, decentralized sys focused on specific business areas.
Data Integ. & Normalization	D W use fact constellation Schema	While DM use star & snowflake
Flexibility	DW are more flexible & data-oriented, oriented to few fields	DM are project-oriented & less flexible
Implementation Time	implementing a DW can take months to years.	While DM can be deployed in few minutes.

Advantages ⇒

→ Optimized Performance

→ Foster Query Response

→ Empower Business users

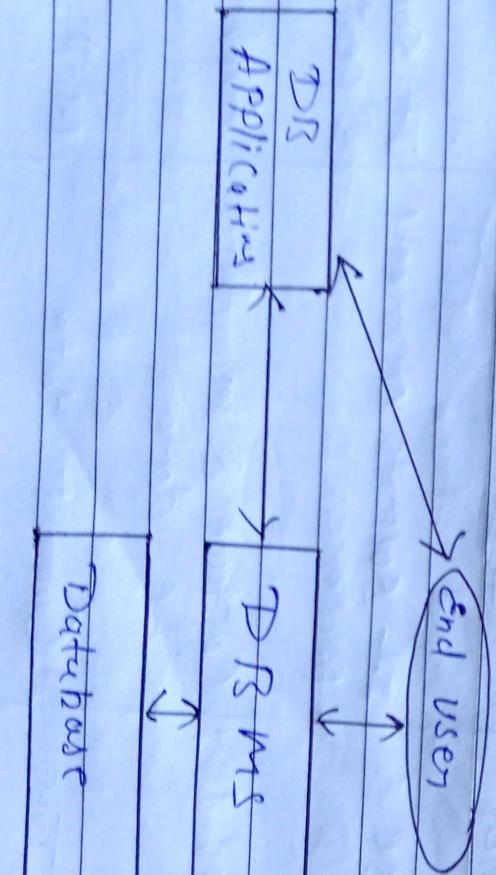
DBMS vs DW ⇒

DBMS :

DBMS is used in the traditional way of storing & retrieving data.

The major task of a database system is to perform query processing.

These sys are generally referred to as Online Transaction Processing systems. These sys are used in the day-to-day operations of any organization.

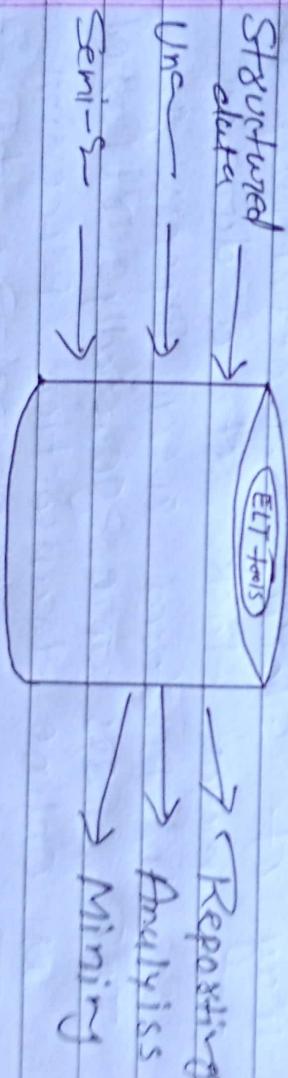


Data Warehouse :-

Dw is the place where huge amount of data is stored.

It is meant for users or knowledge workers in the role of data analysis & decision making.

These systems are referred as Online analytical processing



Features	DBMS	DW
Purpose	It supports operational processes	It supports analysis & performance reporting
Data Handling	Capture & maintain the data	Explore the data
Data Type	Current data	Multiple years & historical

Data Verification	D.V. occurs when entry is done	D.V. occurs after the fact.
Data Size	100 MB to GB	100 GB to TB
Data Model	ER based	Snowflake/Snowflake
Orientation	App. oriented	Subject oriented
Data Specificity	Primitive & highly detailed	Summarized & Consolidated
Storage Structure	Flat relational	Multidimensional
# Knowledge discovery	o	
Knowledge discovery refers to the complete process of uncovering valuable knowledge from large datasets. It starts with selection of relevant data, followed by preprocessing to clean & organize it, transformation to prepare it for analysis, data mining to uncover patterns & relationships, & concludes with the evolution & interpretation of results,		

Ultimately producing Valuable knowledge or insights.

KD is widely utilized in fields like
ML, Pattern Recognition, Statistics, AI
2. Data Visualization.

Here is the list of Steps involved
in the KD process :-

① Data Cleaning =>

In this step, the noise
& inconsistent data is removed.

② Data Integration =>

In this step,
multiple data sources are combined.

③ Data Selection =>

In this step, data
relevant to the analysis task are
selected from the database.

④ Data Transformation =>

In this step,
data is transformed or consolidated into
forms appropriate for mining by
performing Summary or aggregation
operations.

⑤ Data Mining =>

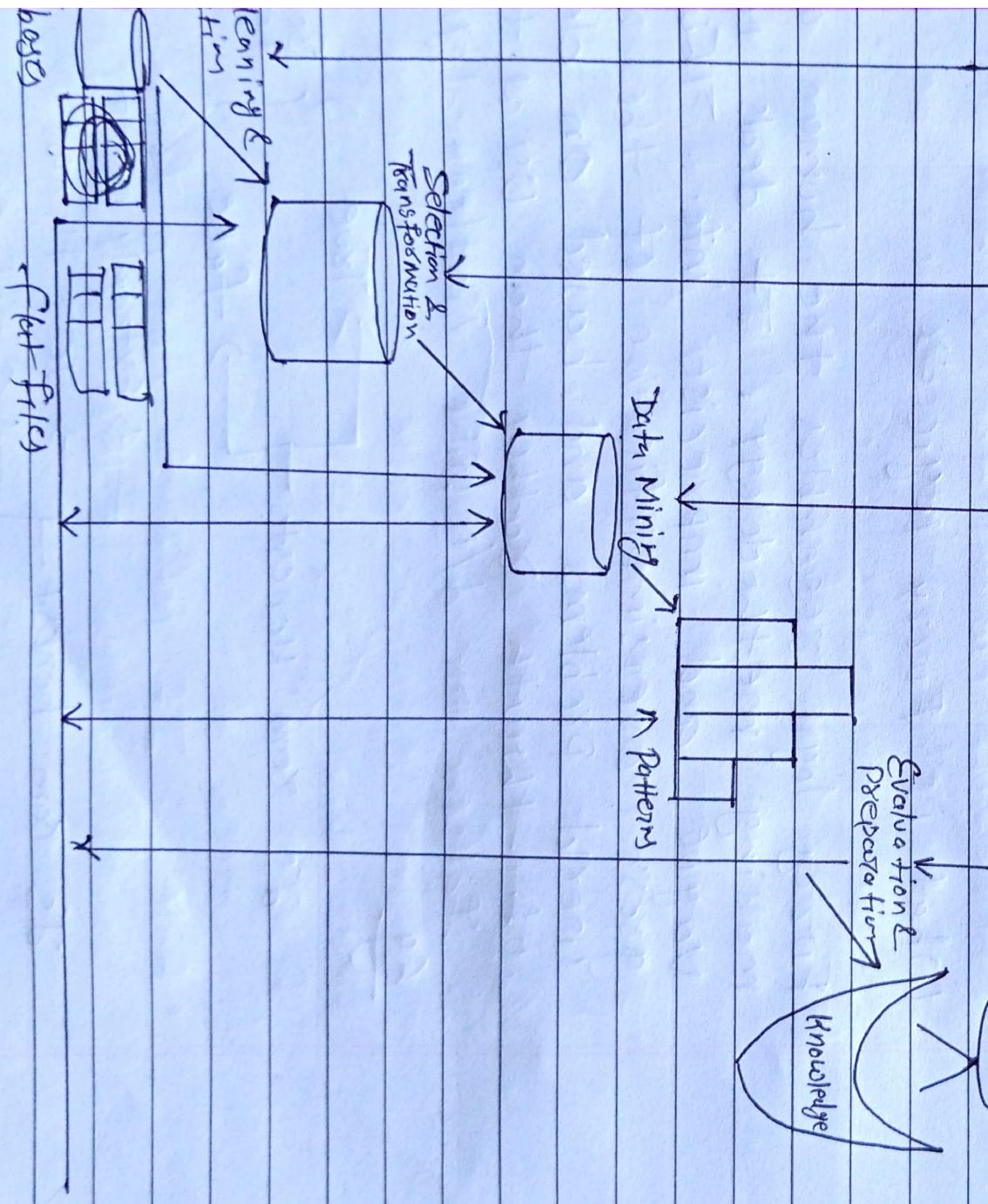
In this step, intelligent methods are applied in order to extract data patterns.

⑥ Pattern Evaluation =>

In this step, data patterns are evaluated.

⑦ Knowledge presentation =>

In this step, knowledge is represented.



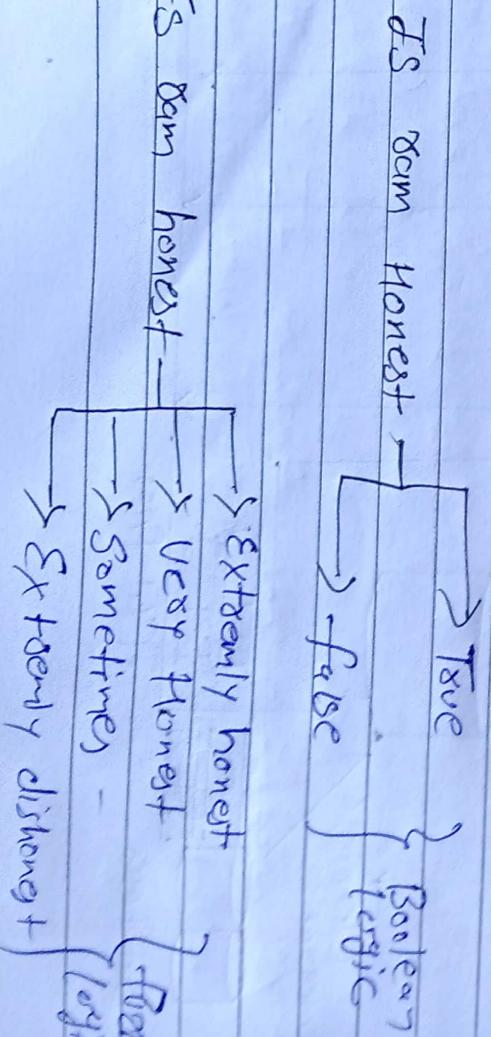
Fuzzy logic & Sets \Rightarrow

Fuzzy logic :-

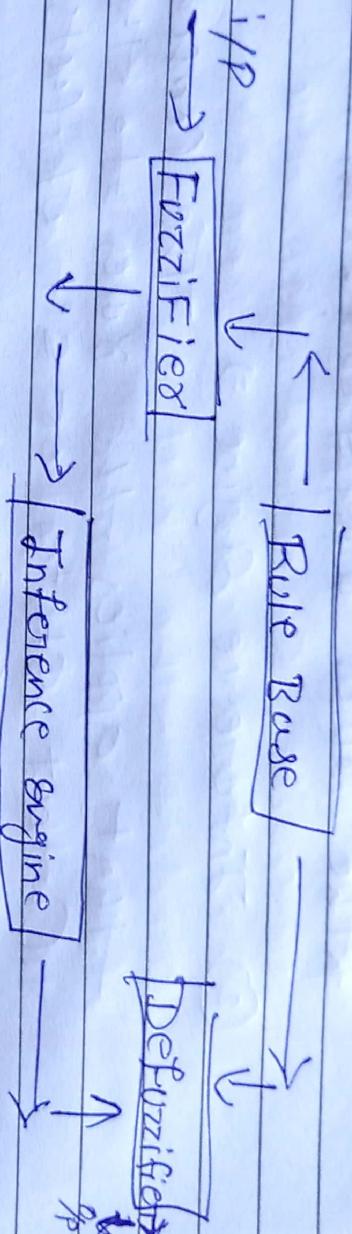
The word Fuzzy refers to things which are not clear or are vague. Any event / process / or function that is changing continuously can not always be defined as either true or false, which means that we need to define such activities in a fuzzy manner.

What is Fuzzy logic :-

Fuzzy logic resembles the human decision making methodology. It deals with vague and imprecise info. This is gross oversimplification of the real world problems and based on degrees of truth rather than usual true/false or 1/0 like Boolean logic.
Eg



Architecture :



Rule based Algo. =>

A rule based algo is a sys that uses a set of specific rules to reach decisions or conclusions.

These rules typically follow an "IF Condition THEN action" structure.

for eg, a rule might state, "if the temp is above 30°C THEN issue a heat warning."

This transparency makes such algo easy to understand & implement, especially in field like AI & data mining.

How rule based algo work ?

① Knowledge Base \Rightarrow

This is where all the rules are stored.

② Inference Engine \Rightarrow

The Component that applies the rules to the current State or set of facts to derive conclusions.

③ Working memory \Rightarrow

Holds the current facts & data being processed by the algo.

④ User interface \Rightarrow

Allows users to interact with the sys, inputting data & receiving outputs.

Eg :-

Consider a simple medical diagnosis system using a rule based algo :-

Rules :-

- Rule 1 : IF patient has a fever AND Cough THEN consider flu.

- Rule 2 : IF Patient has a sore throat AND fever THEN Consider strep throat.

Applications =)

- Expert Sys
- Decision support Sys
- Control Sys

Advantages =)

- Transparency
- Consistency

Disadvantages =)

- Rigid
- Maintenance

KNN :-

K - Nearest Neighbors (KNN) is a simple, versatile, Supervised ML algo used for classification & regression that predict the label or value of a data point based on the 'K' NN in the dataset.

Overview of KNN =>

- > S.L. algo applicable in both classification & regression
- > KNN classifies a data point based on the majority class among its KNN or predicts a value using the avg of the K neighbors values.
- > KNN does not build a model upfront; it stores the training data & performs calculation at the time of prediction, making it a non-parametric method.

How KNN works =>

① Select K :

Choose the no of neighbors (K) to consider.

② Calculate distance :

Compute the distance b/w if p data point & all points in the training dataset using distance metrics such as :

- Euclidean Distance
- Manhattan
- Minkowski

③ Find nearest neighbors :

Sort the distance from the query point to every other point & select the top K closest points.

④ Predict C/P :

- For classification, assign the class label that is most common among the K neighbors.
- For regression, return the avg value of the K neighbors.

How to choose right K Value :

Choosing the optimal K is critical ; techniques include :

- Cross Validation
- Elbow method

Advantages \Rightarrow

Simplicity
flexibility
Adaptability

Disadvantages \Rightarrow

Computationally intensive
Sensitive to outliers
Curse of dimensionality

Applications \Rightarrow

Recommendation Sys
Image Classification
Medical Diagnosis.

#

Data Types used in clustering?

Main Data Types in clustering \Rightarrow

1. Numerical Data \Rightarrow

Numerical Quantitative data consists of
Continuous values that can be measured
& ordered.

Clustering algo's like K-mean &
Gaussian mixture models (GMM) typically
work well with num data.

② Categorical data:

Categorical or qualitative data refers to variables that represent distinct categories or groups.

Algo's like k-nodes can be used for Clustering Categorical data.

- * Types of Clustering algo's & data user =>
- * Centroid - Based Clustering primarily for numerical data
- * Density - — Suitable for both
- * Distribution — Primarily segmented data
- * Hierarchical — works with both numerical depending on linkage method
- * Fuzzy — Mostly numerical data, but can be used with categorical data where membership probabilities apply.

Apriori algo :-

The Apriori algo, proposed by Rakesh Agrawal and Ramakrishnan Srikant in 1994, is widely employed in data mining to uncover patterns of items that frequently occur together in transactional databases.

This particularly effective for opp. like market basket analysis, where businesses aim to ~~to~~ understand consumer buying behaviour & optimize product placement.

Key Concepts :-

① Itemsets \Rightarrow

An itemset is a collection of one or more items

② Frequent itemsets \Rightarrow

These are itemsets that appear in the dataset with a frequency above a set threshold known as support.

(3) Support \Rightarrow

This measures how often an itemset occurs in the dataset

$$\text{Support}(X) = \frac{\text{No of transactions containing } X}{\text{Total no of transactions}}$$

(4) Confidence \Rightarrow

Confidence measures the likelihood that an item y is purchased when item x is purchased.

$$\text{Confidence}(X \rightarrow Y) = \frac{\text{Support}(X \cup Y)}{\text{Support}(X)}$$

(5) Lift \Rightarrow

$$\text{Lift}(X \rightarrow Y) = \frac{\text{Confidence}(X \rightarrow Y)}{\text{Support}(Y)}$$

How its work :

- * Set Parameters
- * Find Frequent 1 - itemsets
- * Generate Candidate itemsets
- * Prune Candidates
- * Repeat for Higher - Order itemsets
- * Generate association rule

Applications :-

Markets ~~Basket~~ Basket Analysis
Recommend sys
Inventory man
Fraud detection

Advantage :-

The algo is easy to implement & understand, providing interpretable results in the form of association rules.

Limitations :-

The Apriori algo can be computationally expensive especially with large datasets, due to the multiple database scans required.

BIRCH :-

BIRCH is a Scalable and efficient clustering technique that enables effective data grouping in applications where datasets are too large to fit into memory.

It is designed to summarize extensive datasets into smaller clusters, which can subsequently be clustered further.



by other algo's, enhancing performance & reducing processing time.

Key Components \Rightarrow

① Clustering Feature (CF) :

The CF is a Summary of a cluster that Comprises three Statistics:

- $N \Rightarrow$ The no of Points in the cluster.
- $LS \Rightarrow$ The linear sum of the points,
- $SS \Rightarrow$ The squared sum of the points.

② CF Tree :

The CF-Tree is hierat Chical Tree Structure where Each node Represents a cluster.

It enables efficient storage of clustering features, allowing the algo to summarize large datasets without needing to store all individual data points

Advantages :

- Scalability
- Memory Efficiency
- Incremental Clustering

Limitations of BIRCH

Sensitivity to parameters

Assumption of spherical clusters

Distance metric limitations.

Data Mining Tasks :-

Data Mining Task Can be Classify into two Categories :

1. Descriptive Data mining :-

This Category focuses on finding human-interpretable patterns describing the data.

It aims to Summarize, Visualize and Understand the inherent characteristics of dataset.

Common task within this category include:

- Classification \Rightarrow

- Predicting if an Email is spam or not based on specific features.

- Clustering ⇒
Clustering is the task of grouping similar data objects together based on their attributes.
 - Association Rule Learning ⇒
 - This technique identify interesting relationships b/w variables in large datasets, often used in market basket analysis.
 - Summarization ⇒
 - It Can Summarize key patterns, trends or statistics in a more compact form for decision making.
 - Outlier detection ⇒
 - This task identifies data points that deviate significantly from the majority of the data.
- (2) Predictive Data Mining :
- predictive tasks aim to develop a model that predicts future outcome based on existing data.
 - This tasks typically use historical data

To drive predictions. Key types :-

- Regression Analysis =>

This technique is used to predict a quantitative response based on the relationships b/w independent variables and the dependent variable.

A common app is predicting Sale based on advertising Spend.

- Time Series Analysis =>

Technique analyze trends over time for forecasting future values like Stock market predictions.

- Forecasting =>

Similar to segm but often used for broader prediction contexts forecasting helps predict unknown values within a dataset & it is commonly used in various industries for demand or supply forecasting.

DBSCAN :-

DBSCAN is an unsupervised ML algo used for Clustering tasks. It groups together points that are closely packed in the feature space while marking points in low-density regions as noise.

DBSCAN is particularly useful in Scenarios where the no of clusters is not pre-defined & is capable of discovering Clusters of arbitrary shapes.

Key Concepts

- Core Points \Rightarrow Points that have at least a minimum no of other points within a specified distance (ϵ).
- Border Points \Rightarrow Points that are within ϵ distance of a core point but do not have enough neighbours to be considered core points themselves.
- Noise Points \Rightarrow Points that are neither core nor border points; they don't belong to any cluster.

How its work \Rightarrow

② Parameters Selection :

Set two parameters, epsilon (ϵ) for the neighborhood radius & minpts for the mini. no of points required to form a dense region.

③ Cluster formation :

\rightarrow for each point in the dataset, DBSCAN retrieves the neighborhood within ϵ distance.

\rightarrow If enough neighbors (minpts) are found, a new cluster is initiated; otherwise, the point is labeled as noise

\rightarrow The algo recursively expands the cluster by adding all density-connected points to the cluster until no new points can be added.

Advantages :

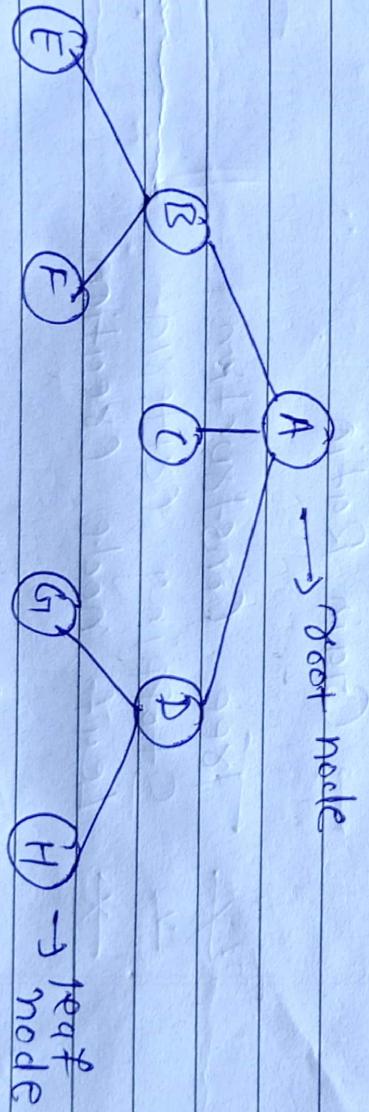
- * No need for pre-defined no of clusters
- * Flexibility in cluster shape
- * Noise Handling

Decision Tree

- * Sensitive to Parameters
- * Difficulty with Varying Densities
- * Performance on High Dimensions

Decision Tree induction algo :-

A decision tree is a structure that includes a root node, branches, and leaf nodes. Each internal node denotes a test on an attribute, each branch denotes the outcome of a test, and each leaf node holds a class label. The topmost node in the tree is the root node.



The benefits of having a decision tree are as follows -

- * It does not require any domain knowledge
- * It is easy to comprehend
- * The learning & classification steps of a decision tree are simple & fast

A machine researcher named J. Ross Quinlan in 1980 developed a decision tree algo known as ID3 (Iterative Dichotomiser).

In this algo, there is no backtracking, the trees are constructed in a top down recursive divide-and-conquer manner.

How its work :-

- * Data Preparation
- * Attribute Selection
 - Info Gain
 - Gini Index
 - Gain Ratio
- * Tree Construction
- * Stopping Criteria
- * Leaf Node Creation

Pruning :-

Pruning is important to prevent overfitting, where the model becomes too complex & performs poorly on unseen data.

There are two main types of pruning:

Pre-pruning \Rightarrow Stopped the tree growth early based on certain criteria

post \Rightarrow first creating the full tree & then removing branches that have little significance \rightarrow improving accuracy on validation data.

~~Numpy & Pandas~~ Advantages \Rightarrow

- Easy to Understand
- Handles Various data
- Flexible Application

Disadvantages \Rightarrow

- Prone to overfitting
- Sensitive to Data Changes
- Tries Toward Certain Attributes.

#

Neural Network based algo ?

Neural Network Classification involves using various neural network archi. To automatically categorize ips into Specific Classes.

This technique has proven successful in domains such as image recognition, NLP & medical diagnosis, where complex data needs sophisticated handling.

Common neural network archi. :-

- * CNN
- * RNN
- * ENN
- * GAN etc

#

FP - Growth :-

The Frequent Pattern Growth (FP-Growth) algo is widely used in data mining to discover frequent patterns or itemsets in transactional dataset. Such as market basket analysis. It improves upon the Apriori algo by minimizing the no of database scans and eliminating the need for candidate generation, which can be computationally expensive.

Key Concepts :

- Frequent Itemsets \Rightarrow

Groups of items that appear together frequently in transactions.

- Support \Rightarrow

The frequency of an itemset in the dataset. It is defined as the proportion of transactions that contain the itemset.

- FP - Tree \Rightarrow

A compressed tree structure that retains the ~~essential~~ essential info of the dataset, allowing for efficient frequent pattern mining.

How it works :

Step 1 : Build the FP - Tree

- Scan the data
- Sort items
- Construct the Tree

Step 2 : Mine the FP - Tree for frequent patterns

- Start from the bottom of the tree
- & identify frequent itemsets up to the root

→ for each item in the tree, Create a Conditional Pattern base, which consists of the paths leading to that item.

→ Construct a new Conditional FP-tree from the Conditional pattern base to find associated frequent pattern

Advantages :

- * Efficiency
- * No Candidate Generation
- * Scalability

Limitation :

- * Complex Implementation
- * Memory usage

Naïve Bayes :-

* Naïve Bayes is a family of probabilistic classifiers that apply Bayes' Theorem, assuming Strong independence between the features.

This means that the presence of one feature does not affect the presence of another, simplifying the computation of probabilities.

The algo computes the posterior prob. of a class given a set of features & selects the class with the highest probability. The formula is :

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Where,

$P(Y|X)$ is the posterior prob,
 $P(X|Y)$ is the likelihood,
 $P(Y)$ is the prior prob.,
 $P(X)$ is the evidence

Types :-

→ Gaussian Naïve Bayes

→ Multinomial Naïve Bayes

- Bernoulli
- Complement

Advantages :

- Simplicity
- Speed
- Performance

Disadvantages :

- Independence Assumption
- Zero Probability Problem

Application :

- Text Classification
- Medical Diagnosis
- Recommendation Sys
- Weather Prediction

Category	OLAP	OLTP(online Transaction Processing)
Definition	It is well known as online database query management Sys.	It is well known as an online database modifying Sys.
Data Source	Consists of historical data from various DR	Consists of only operational current data.
Method Used	It makes use of a data warehouse.	It makes use of a Standard DRMS.
Application	If Master is Subject oriented. Used for data mining, Analytics, Decision making etc	It is app. Oriented. Used for business tasks.
Normalized	In an OLAP database, tables are not normalized.	In an OLTP database, tables are normalized.
Nature	The process is focused on Audience on the Customer	— — — the market
Operations	Only read & query write ops.	Both read & write ops.
DB design	Design with a focus on the subject	Design that is focused on the app.