

Unit - 1

Data management

- design data architecture and manage the data for analysis.
- understand various sources of data like sensor signals / GPS etc.
- Data management
- Data Quality (noise, outliers, missing values, duplicate data)
- data pre-processing & processing.

Data management : Design data architecture and manage the path for analysis

Data architecture design is like a detailed plan for how to handle data in a company, showing the steps for gathering, storing, accessing and using data. This plan helps keep data neat and well organized.

Data management adds to this by taking care of data from start to finish, including collecting, storing, arranging, maintaining and protecting it.

Together, these methods help a business by making sure data is correct, easy to get and dependable, which helps with making decisions and running the business smoothly.

Data Architecture Design

It is set of standards which are composed of certain policies, rules, models and standards which manages, what type of data is collected, from where it is collected, the arrangement of collected data, storing that data, utilizing and securing the data into the systems and data warehouses for

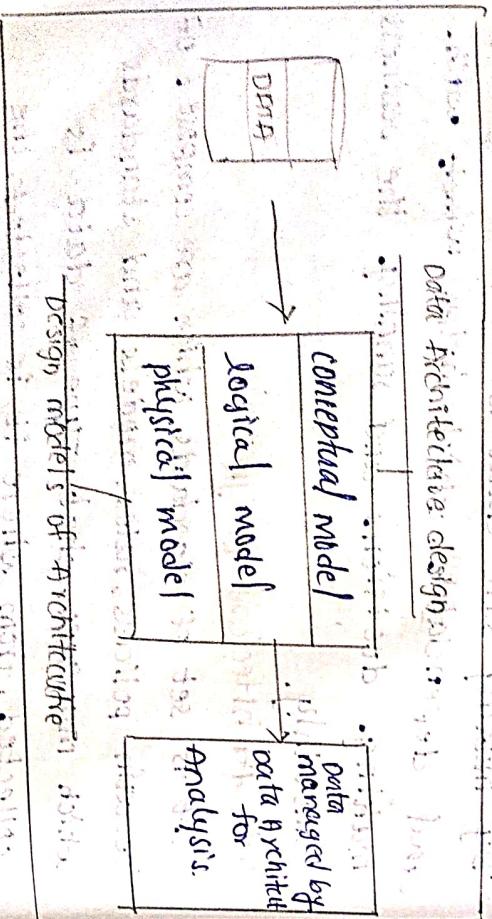
(or)

It is a strategic plan or blueprint that defines how data is collected, organized, stored and accessed, integrated, and secured within an organization. It serves as the foundation for all data-related processes.

Data architecture also describes the type of data structures applied to manage data and it provides an easy way for data preprocessing.

The data architecture is formed by dividing into 3 essential models and then are combined:

- Data architecture
- Database design
- Physical model



- Helps stakeholders (non-technical people) understand the system.
- Does not include technical details like data types or storage.

It serves as a blueprint that outlines major entities and their connections.

Ex:- In a hospital, the conceptual model shows relationships between patients, doctors and appointments.

2. Logical model

- Adds structure to the conceptual model

- Shows how data is organized logically (tables, XML, object classes)

- Developers use it to build databases.
- Contains tables, fields, primary keys, foreign keys, but still independent of actual software/hardware.

3. Physical model

The logical model provides a blueprint of how data is organized and how the databases should operate, translating complex systems designs into a readable form for technical developer.

e.g., patient table with columns like ID, Name, Age etc.

3. Physical model

- Actual implementation using specific technologies (SQL, Oracle etc.)

1. Conceptual model

- High level representation.

- Uses entity-relationship (ER) models to show how entities (like customer, orders) relate.

- Deals with indexes, partitions, storage paths and file formats.

- Used by DBA's (Database Administrators).

- Employee table stored in MySQL with indexing.

This model is crucial for database administrators who are responsible for the technical deployment and maintenance of the database systems.

A data architect is responsible for all design, creation, management, deployment of data architecture and defines how data is to be stored and retrieved. Other decisions are made by internal bodies.

Data Management

It is the end-to-end process of collecting, storing maintaining, securing and using data efficiently for analysis and business operations.

Key:-

- extracting data from sources
- storing data securely
- ensuring data accessibility, accuracy & consistency
- converting data into usable formats.
- maintaining security and reducing costs.

Responsibilities in Data management

- Accessibility & integrity :-
ensure authorized users can access accurate, consistent data.

- security
protect data using encryption, backups, access control.

- cost efficiency
use optimized tools to reduce storage / processing costs.

- support business decisions
provide reliable data for planning and strategy.

- handle Big Data :-
use tools like Hadoop, Scala, AWS, Tableau for massive datasets.

Relationship : Data Architecture vs Data management

Data Architecture

- framework for data flow execution of that framework.

- focuses on structure and policies

- focuses on operations and security.

- Built by data architects managed by IT / data team.

- Deals with how data deals with how data does should flow.

→ Both work together to ensure high-quality accessible, and secure data systems.

Benefits of both integrating

- Informed decision making

- Reduced Redundancy

- Automation

- Compliance

- Technology

Real world - Applications:-

- Banking sector

- Architecture: designs secure flow of transaction data management: ensures real time processing and backup of customer transactions.

- Health care

- Architecture: connects patient records across departments, management, ensures up-to-date records and compliance with laws like HIPAA.

Understand various sources of data like sensors/signals GPS etc.

To the digital era, data can come from various sources, and these can be:

- Structured (like databases)
- semi-structured (like XML, JSON)
- Unstructured (like audio, video, images)

Understanding the nature and behaviour of these sources is essential for building efficient systems for data collection, storage and analysis.

1. Sensor Data

Sensors are physical devices that detects and measure real-world conditions then convert them into digital data.

example of sensor data:
sensor type data captured applications

Temperature sensor Temperature ($^{\circ}\text{C}$ / $^{\circ}\text{F}$) Weather monitoring

Motion sensor Movement states smart homes, security systems.

Heart rate sensor Beats per minute (BPM) Health wearables

GPS etc. ECG monitoring.

- Real time: provides live updates

- Continuous: produces a stream of data.

- Time-series: each reading is linked with a timestamp.

use case:-

A smart home uses motion sensors to detect movement and turn on lights automatically.

2. Signal Data

Signal data refers to information derived from analog or digital signals, such as sound, electricity or electromagnetic waves.

examples:-

signal type

data captured

applications

audio signals

sound wave amplitude

voice assistants, music processing

ECG (electrocardiogram) hearts electrical activity heart disease detection.

EEG (electroencephalogram) brain wave patterns sleep research, neurology.

Radar / Reflected wave Navigation, ultrasound reflected pattern

Characteristics:-

- High Frequency: Captured many times per second.
- Time dependent: Values vary over time
- Requires processing: often transformed using techniques like:-
- FFT (Fast Fourier Transform)
- STFT (Short-time Fourier Transform)
- use case:- EEG signals help doctors analyze brain disorders like epilepsy.

8. GPS Data (Global positioning system)

GPS is a satellite based system that gives data about location, speed and movement.

e.g:-

Data type Description Application

latitude & longitude coordinates of a place Google maps, ride apps

speed & direction moment behaviour fleet vehicle tracking.

timestamps when data was recorded

Time-aware analysis (Route delivery).

Characteristics:-

- Spatio-temporal: Combines location + time
- Periodic or event driven: Data can be sent on schedule or during events.
- use case:- GPS in delivery drones tracks real-time routes and drop points.

4. Other sources of Data

source

Description

example

social media user-generated text/image / videos.

Tweets / Insta posts.

logs & clickstreams Tracks user activity on app/web

Login logs, click paths.

databases structured internal data

HR records, CRM systems.

DPS external sources providing

weatherinfo, stock prices.

drones multisensor data (videos, audio, GPS)

Surveillance delivery.

Data Quality (Noise, outliers, missing values, duplicate data)

It refers to the condition or state of data based on various attributes that determine its reliability,

accuracy, consistency, and suitability for a specific purpose.

High quality data is accurate, complete and timely ensuring it is free from errors, inconsistencies and bias.

Good data quality supports effective decision-making, analysis and modeling in various fields like business, research and technology.

key attributes of high-quality Data.

Attribute

- Data reflects the true value of what it represents.

Completeness

- No important information is missing.

Consistency

- No contradictions between datasets.

Timeliness

- Data is up to date when needed.

Reliability

- Can be trusted to make business or technical decisions.

Noise in Data

- Noise is irrelevant, distorted or corrupt data that masks the true pattern.

Noise refers to meaningless or corrupt data. It often includes typos, special characters, or unwanted symbols.

Sources:-

- Typing mistakes (e.g "R@3") - Data transmission errors
- Mal functioning sensors.

Ex:- A@3L, Anita, R@hul.

- Problem: "A@3L", "R@hul" contain special characters

→ Soln :- Remove noise using cleaning fn (e.g. regex, filtering)

Characteristics

- Reduce data readability

- Affects text analytics or machine learning models

- Often removed using data cleaning techniques

2. Missing values :-

A missing values is when no data is recorded for a particular field or variable.

Ex:-	Name	Age	Score
Anita	19	85	
Ravi	20	92	

Meena

Age missing for Ravi

Score missing for Meena.

Types of Missingness :-

1. MCAR (Missing completely at Random)

missing data is not related to any variable

e.g:- random nulls

Impact:- least bias.

2. MAR (Missing at random)

missing depends on other observed variables

e.g:- Income missing depending on age

Impact - moderate bias.

3. MNAR (Missing Not at Random)

missing depends on unseen factors

e.g:- Refusing to report income due to privacy

Impact:- Most bias.

Handling Techniques :-

- Fill with mean/median / mode

- Drop rows with too many nulls.

- Use predictive imputation (e.g KNN imputer)

3 Outlier

These are data points that are significantly different from others in the dataset. They can distort statistical analyses and machine learning models.

Outliers are extreme values that are far away from other data points.

e.g.: student

	score
Anita	85
Ravi	1000
Meena	78

→ Ravi's score is abnormally high—an outlier.

4. Duplicate Data.

Repeated rows or records in a dataset.

e.g.-

name	score
Ravi	92
Ravi	85

→ two rows for Ravi are identical → duplicate.

Why Remove Duplicates?

- saves storage space.

- improves model accuracy.

- ensures unbiased analysis.

DATA MANAGEMENT

Data management is the systematic process of collecting, storing, organizing, maintaining, securing, and preparing data for effective analysis and decision making.

key benefits:-

- ensures data quality (accurate, clean, up-to-date)

- enables quick and accurate analysis.

- helps meet data regulations like GDPR.

components of data management

1. Data collection

- source: sensors, APIs, forms, databases.

- goal: capture raw data in real-time or scheduled intervals.

e.g.- A health app collects heart rate data from smartwatches.

2. Data storage

- ~~structured~~ structured storage formats: SQL, NoSQL, CSV

- storage options :-

- cloud: AWS, Azure, Google Cloud

- local: Oracle, MySQL, Excel.

3. Data organization:-

- organizing data into tables, schemas, data warehouse, / lakes

- Helps improve retrieval speed & clarity.

e.g. separate tables for students, marks, attendance.

4. Data Cleaning

- removing noise, duplicates, fixing missing values.
- ensures data consistency and trustworthiness.

Ex:- Removing invalid timestamps in log data.

5. Data Integration:

- combining data from multiple sources
- resolves data soils and ensures completeness.

Ex:- Merging data from a CRM and an

Online shopping site.

6. Data Governance:

- defines policies and standards
- who can access data
- audit logs
- Data ownership & ethics.

Ex:- enforcing privacy rules like encryption or masked data views.

7. Data Preparation for Analysis

- Transforming data for easy visualization and

ML.

Ex:- calculating average sales per month.

DATA PREPROCESSING

It is the systematic process of collecting raw data and transforming it into usable meaningful information. This information can be then be analyzed, visualized, or used for decision-making.

- It's usually performed by data scientists or data teams using software tools and machine learning models.

Why is important :-

- converts unorganized raw data into structured insights
- make data understandable & actionable.
- helps in automated decision-making & real-time response.
- supports data storage, compliance and future use.

6 steps in Data Preprocessing

1. Data Collection

- gather raw data from data lakes, warehouses, sensors, apps, etc
- The quality of data depends on how reliable source is.

2. Data Preprocessing

- raw data is cleaned and organized
- remove noise - fix missing values - eliminate duplicates

3. Data Input

- clean data is entered into tools for

4. Data Preprocessing

- Algorithms or scripts process the data!

- Machine learning for pattern recognition

- Statistical operations.

5. Data Output & Interpretation

- The processed data is made readable & usable through:

Graphs, reports, Text summary

6. Data storage:-

Final results are stored in:

Databases, cloud storage, file

- must comply with data protection laws like

GDPR.

i. Batch Processing

Processes data in chunks at scheduled times

e.g.: Monthly salary processing.

2. Real-time processing

Instant data processing as it arrives.

e.g.: fraud detection in banks.

3. Multiprocessing

Use multiple CPUs to process data in parallel

e.g. Big data analytics.

4. Manual Processing

Done by humans;

slower used where precision is vital.

Challenges in Data Processing:

Organizations face several challenges when managing large volumes of data,

including:

- Quality issues

- Scalability constraints

- Integration complexity

- Regulatory compliance

3. Explain in detail about data preprocessing with an example.

Data preprocessing is the process of preparing raw data for analysis by cleaning and transforming it into a suitable format.

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

Steps in Data Preprocessing:-

1. Data cleaning
2. Data integration
3. Data Transformation
4. Data Reduction.

1. 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 accurate and reliable.

a. Missing values:-

~ place of missing values, we can replace with "NA"
- sometimes replaced with most probable values
- missing values can be filled in 2 ways
manual/automatic

b. Noisy Data

Noisy data → inconsistent/error data
methods to handle

1. Binning method:-

The data is sorted into equal segments, and each segment is smoothed by replacing values with the mean or boundary values.

2. Regression

Numerical prediction of data

3. Clustering

Similar data items are grouped at one place

2. Data Integration

Merging or multiple heterogeneous source of data are combined into single dataset

- 2 types of data integration

a. Tight coupling

Data is combined together into a physical location.

b. Loose coupling

Only an interface is created and data is combined through the interface and also accessed through interface

- data remains in actual database only

3. Data Reduction

volume of data is reduced to make analysis easier

Methods:-

a. Dimensionality reduction:-

Reduces no. of s/p variables in the dataset
Bcz, large s/p variables → poor performance

"Data cube Aggregation
Data is combined to construct a datacube

c. Data compression:-

Reducing the size of data by encoding it in a more compact form, making it easier to store & process.

d. Numerosity Reduction:-

Reducing the no. data points by methods like Sampling.

4. Data Transformation

Data is transformed into appropriate form suitable for mining process.

a. Data Normalization:

The process of scaling data to a common range to ensure consistency across variables.

b. Discretization:

Converting continuous data into discrete categories for easier analysis.

c. Data Aggregation:-

Combining multiple data points into a summary form.

d. Concept Hierarchy Generation:-

Organizing data into a hierarchy of concepts to provide a higher-level view for better understanding & analysis.

Q. Explain about different tools used for Data Analysis

an example

Data analytics tools help in collecting, processing, analyzing, and visualizing data to gain insights.

Some most commonly used tools

1. Excel

- Best for small-scale data analysis, reporting and visualization.

Ex:- A sales manager uses excel to analyze monthly sales trends and create pivot tables.

2. Python

Used for data manipulation, statistical analysis and machine learning

Ex:- A data analyst uses Python Pandas library to clean a dataset of customer transactions and create a predictive model for customer churn.

3. R

Preferred for statistical computing & visualization

Ex:- A researcher uses R to analyze survey data and create regression models for predicting customer satisfaction.

4. SQL

Used for querying & managing large databases

Ex:- A bank uses SQL to extract customer transaction details & find outliers.

5. Tableau

A powerful tool for interactive data visualization and dashboards.

Ex:- A marketing team uses Tableau to create

a real-time dashboard of customer engagement metrics.

6. Power BI

Microsoft tool for business intelligence and data visualization.

Ex:- A company's HR department uses Power BI to track employee performance & attrition rates.

7. Apache Spark

Big data processing and real-time analytics.

Ex:- A streaming platform like Netflix uses Apache Spark to analyze user history

8. Hadoop

Storing & processing large-scale data across distributed systems.

Ex:- An e-commerce site like Amazon use Hadoop

9. Google Analytics

Website traffic analysis & digital marketing insights

10. KNIME

Data mining, machine learning & automation

Ex:- A pharmaceutical company uses KNIME

5. Explain Data modeling Techniques

Data modeling is the process of structuring and organizing data to make it useful for analysis.

1. Conceptual Data model

These are visual representation of an enterprise data elements and the connection b/w them.

Types of Data models

3 types of Data models

1. Conceptual Data model

- High level, abstract view of data
- Focuses on business concepts and rules
- used in the early project stage.

2. Logical Data model

- defines tables, columns & relationships.
- Independent of any database system
- Basis for database design

3. Physical Data model

- Actual database implementation
- Includes tables, constraints & keys
- Specific to database system (e.g MySQL Part)

Types of Data modeling

1. Hierarchical model:

- . Data is structured like a tree (One parent, multiple children)
- used less frequently today.

Ex:- A college (parent) has student & teacher (children)

Relational model :-

- Represents data in tables (row of col)
- Used in relational database management system.

Object - Oriented Data model.

- Data is stored as objects like in OOP
- Supports abstraction, inheritance & encapsulation

Network model

- Uses a graph structure with nodes (data) & edges (relationships)

Supports multiple parent-child relationships.

- ER Model (Entity - Relationship model)
- Use entities, attributes & relationships to model data.
- Visualized with ER diagrams.

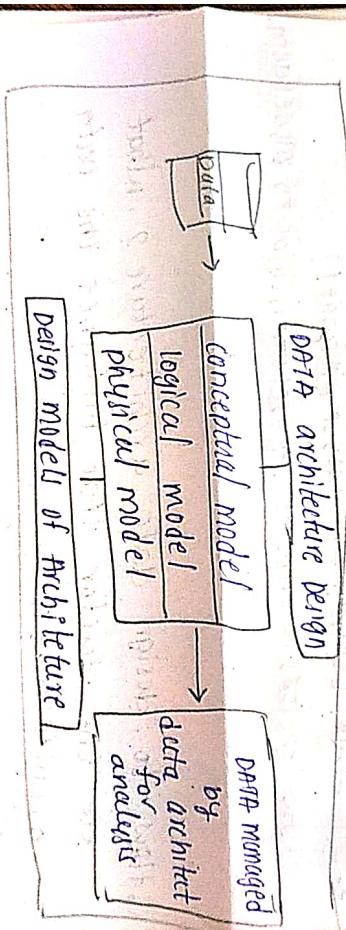
Q. How to design data architecture? What are the factors that influence the data architecture?

- Data architecture design is set of standards which are composed of certain policies, rules, models and standards.
- Which manages, what type of data is collected, from where it is collected, where storing that data for further analysis.
- Data architecture design is important for creating a vision of interactions occurring between data systems.

for example,

- If data architect wants to implement the architecture the visionary model of data interaction during the process can be achieved.

- Data architecture also describes the type of data structure applied to manage data and it provides an easy way for data preprocessing. The data architecture is formed by dividing into 3 essential models & then are combined.



i. conceptual model

It is a business model which uses ER [entity relationship] model for relation b/w entities and their attributes.

ii. logical model:

It is a model where problems are represented in the form of logic such as row, column, of data, class, XML tags & other DBMS techniques.

iii. physical model: It holds the data base design physical model holds the data base technology like which type of data base technology

will be suitable for suitable for architecture.

Factors that influence data architecture

1. Business Goals:-
Determines the data flow & structure based on
use cases
2. Data Volume & variety :-
Affect storage & processing decisions
3. Scalability Needs
Determines cloud vs on-premise architecture
4. Security & Compliance
Requires encryption, access control,
and legal compliance
5. Performance & Latency:-
Real-time processing needs vs batch
processing
6. Integration Requirements
7. Cost & Budget.

Unit - II

Data Analytics

- Introduction to Analytics
- Introduction to Tools and environment
- Application of modeling in Business,
- Databases & Types of Data and variables
- Data modeling techniques.
- Missing Imputations etc
- Need for Business Modeling.

Data Analytics :-

- Data Analytics is a systematic process of converting raw data into valuable insights using tools, technologies and statistical techniques.
- These tools are used for collecting, cleaning, transformation and modeling.
- It supports decision-making and is crucial for business growth, risk management, and enhancing customer experience.
- Its applications span finance, agriculture, banking, retail, government and more.
- Process of Data Analytics:
1. Data collection:- Gathering raw data from various sources or subsets.
- 2. Data cleaning:- Correcting data quality issues like errors and duplicates, and white spaces which need to be corrected before moving to the next step.
- 3. Data analysis and interpretation:- Creating and testing analytical models using

tools like Python, excel, R, scala, and SQL

4. Data visualization:-

Using charts and graphs to identify patterns and insights.

Types of Data Analytics:-

1. Descriptive performance:- Summarizes past data to identify patterns and issues.
 - Used to compare past performance.
2. Real-time analytics:- Analyzes data as it is entered to track trends and competitors.
3. Diagnostic analytics:- Analyzes past data to find causes of anomalies using techniques like regression and correlation.
4. Predictive Analytics:- Uses current data and machine learning to predict future outcomes.
5. Prescriptive Analytics:- Suggest best solutions for decision making and automation.

- used in loan approvals, scheduling, pricing models.

Methods of Data Analytics

1. Qualitative Data Analytics:-

Analyzes non-numerical data like words, images and behaviours.

Some common qualitative methods are

- Narrative Analytics is used for working with data acquired from diaries, interviews and so on.

- Content analytics is used for analysis of verbal data and behaviour.
- Grounded theory is used to explain some given event by studying.

2. Quantitative Data Analytics

uses numerical data for statistical analysis.

- Hypothesis testing assesses the given hypothesis of the data set.

- Sample size determination is a method of taking a small sample from a large group of people and then analysing it.

- Average or mean of a subject is dividing the sum total nos in the list by the number of items present in that list.

Skills required to data analytics

These are multiple skills which are required to be a data analyst.

- Programming languages : R, python.

- Database language : SQL

- Knowledge of machine learning, probability statistics, data management, and data visualisation.

Data analytics jobs :-

- Entry-level roles

- Junior Data Analyst
- Associate Data Analyst.

- Experienced - level roles.

- Data Analyst
- Data Engineer
- Data Architect
- Data scientist
- Marketing Analyst
- Business Analyst.

Introduction to Analytics

Analytics is a systematic process of collecting, processing, and analyzing data to discover meaningful patterns, insights and trends. It involves using various tools, technologies and methodologies to transform raw data into actionable information that supports decision-making.

In essence, analytics empowers businesses and individuals to make informed decisions by leveraging data-driven insights.

And types are same as Data Analytics types.

Introduction to Tools and Environment

Analytics environment:-
Data Analytics environment refers to the complete ecosystem that enables data collection, processing, modeling, analysis and visualization.

It includes:

- Data storage systems
- Data processing platforms
- Programming environments
- Visualization tools
- Cloud and big data platforms.

This environment enables analysts to work with data at various stages - from acquisition and cleaning to modeling and presentation, and also helps data professionals carry out various analytical tasks efficiently and accurate.

Categories of Data Analytics Tools.

1. Programming languages
Used to write scripts, automate tasks, build model, and process data.
Ex:- Python, R, SQL, Java.

2. Data storage tools :
These tools are used to store, refine and manage large volumes of data in structured or unstructured formats.
e.g., MySQL, PostgreSQL for relational data and MongoDB, Hadoop, HDFS for non-relational & big data storage

8. Data processing tools

These tools help in cleaning, transforming & preparing data for analysis.
ex:- Python, R and pandas, spark, excel etc

4. Visualization tools

These tools converted processed data into visual forms like charts, graphs and dashboards.

Tools are :- Tableau, power BI, Matplotlib and seaborn.

5. Statistical and Machine learning tools

These are used to perform statistical tests, build predictive models, and discover patterns in data.

Ex:- scikit-learn, TensorFlow, SAS and IBM SPSS

6. Big data tools

These tools are designed to process, store and analyze very large data sets distributed across clusters.

Apache Hadoop handles batch processing, while Apache spark enables faster in-memory computation

7. Cloud platforms:

They provides scalable computing resources and services for storing and analyzing data.

In AWS, Google cloud & Microsoft Azure.

Application of modeling in Business

It refers to the use of mathematical, statistical and computational techniques to represent real-world business processes and scenarios.

These models enable businesses to simulate

- outcomes, predict trends, and optimize decision-making across various departments such as sales, finance, marketing, operations & HR.
- Modeling supports data-driven decisions,

helping organizations reduce risks, increase efficiency, and maintain a competitive edge in the market.

Major Applications :-

1. Sales forecasting

- predictive models are used to analyze historical sales' data.

- helps in forecasting demand, managing inventory and planning sales targets.

- enables organizations to avoid overproduction or stockout.

2. Customer segmentation:

- models like clustering and classification group customers by demographics, preferences or behaviour.

In DA, tools and environment form the foundation for effective data handling. A well-structured analytics environment ensures scalability, flexibility and accuracy in analytical tasks.

- Supports targeted marketing, personalization and improves customer retention.

3. Risk management

- Statistical and machine learning models assess financial credit, and operational risks.
- Used in credit scoring, fraud detection & insurance underwriting.

- minimizes losses by identifying high-risk activities.
- used in credit scoring, fraud detection & insurance underwriting.

4. Supply chain optimization

- Models like linear programming and simulation, optimize logistics, inventory and supplier decisions.
- Help reduce costs, shorten delivery times & improve efficiency.

5. Financial planning & Budgeting

- Forecasting and scenario analysis help in planning revenues, expenditures and investments.
- Models simulate best-case, worst case and average case financial scenarios.

Databases & Types of Data and Variables

A database is a structured collection of data stored electronically. It allows easy access, management, and updating of data using Database Management System (DBMS) such as MySQL, Oracle, MongoDB, etc.

Characteristics:

- Data organization - logical arrangement for quick access and modification.
- Efficiency - optimized performance for storage & retrieval.
- Scalability - can handle large, growing volumes of data.
- Security - ensures safe access & privacy.

7. Market Basket Analysis

- Association rule mining identifies product purchase patterns.
- Used in cross-selling, promotions, and retail shelf planning.

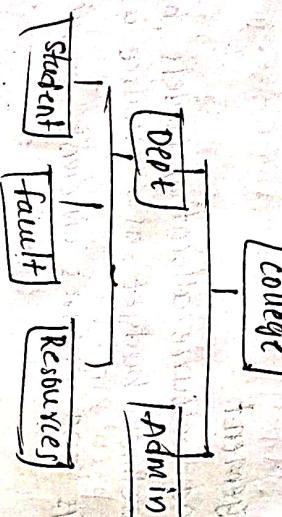
- Improves sales strategies by understanding customer buying behaviour.
- Improves sales strategies by understanding customer buying behaviour.

Types:-

Databases can be classified based on their structure, usage, storage methods, and intended applications.

1. Hierarchical databases

Tree-structured parent-child model



Ex:- DB2, IMS

Adv:- Simple and fast for structured data

disadv:- Rigid structure, limited flexibility.

2. Network database

→ Network databases builds on the hierarchical model but allows child records to link to multiple parent records, creating a web-like structure of interconnected data.

Ex:- IDBS - Integrated Data Store

Adv:- More flexible than hierarchical models

2. Handles complex relationships

Disadv:- Hard to design / manage

2. Object-oriented

Data stored as objects with attributes and methods.

Ex:- Berkeley DB

Adv:- Good for multimedia reusability
disadv:- Requires OOP knowledge.

4. Relational

Data in tables with rows (records) & columns (attributes)

Ex:- MySQL, PostgreSQL

Adv:- Simple, ACID (Atomicity, Consistency, Isolation, Durability)

disadv:- Difficult to scale, schema-dependent.

5. Cloud

A cloud database operates in a virtual environment hosted on cloud computing platforms. It

Ex:- AWS, Azure

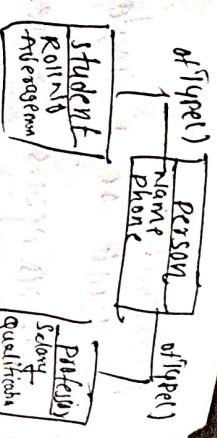
Adv:- Pay-as-you-go flexible
disadv:- Internet dependence, security risks.

6. Centralized

It is a database stored and managed at a single location, such as a central server or data center.

Ex:- Central servers

Adv:- Easy security control
disadv:- Limited performance at scale.



7. Personal :-

Meant for a single user lightweight

Ex:- MS Access, SQLite.

Adv:- Simple, low storage.

Disadv:- Not scalable, fewer features

8. NoSQL :-

Non-relational, supports flexible structures

(JSON, key value)

Ex:- Mongo DB

Adv:- Fast, scalable, unstructured

Disadv:- Limited GUI, weak backup tools.

Types of Data in Analytics

- Structured Data

Organized in a fixed format, usually in tables

Ex:- Excel files, SQL databases.

- Unstructured Data

No predefined format. Requires more processing

Ex:- emails, PDFs, videos.

- Semi-structured Data

Partially organized data with tags or markers

Ex:- JSON, XML

- Time Series Data

Data points indexed in time order

Ex:- stock prices, tem. pressure logs.

- Spatial Data :-

Data that represents locations & shapes

Types of Variables:-

Variables represent measurable characteristics or properties of data. They are key to statistical and machine learning analysis.

a) Nominal categories without any order
Ex:- Gender (Male, Female), colors.

b) Ordinal categories with a meaningful order
Ex:- education level (High school, BE)

c) Interval:- Numerical data without a true zero.
Ex:- Temperature in Celsius.

d) Ratio:- Numerical data with true zero.
Ex:- Age, Height, weight, sales.

Classification of Variables

a. Quantitative (Numerical) Variables

These represent measurable quantities

(Countable & Discrete Variable): can take only specific values

Ex:- No. of students, cars.

e) Continuous Variables:-

can take any value within a range(measurable)

Ex:- height, weight, temperature.

b. Qualitative (Categorical) Variables:-

These represent categories or labels

Nominal:- No order

Ordinal: Ordered categories.

- databases form the backbone of modern data storage and management

- Understanding the types of data (structured, unstructured) and types of variables (categorical, numerical) is crucial for effective data analysis.

- Proper classification helps in choosing the right tools, techniques and visualization methods in data analytics.

Data modeling techniques

Data modeling:-

It is the process of structuring, organizing and representing data to define how it is stored, connected and accessed in databases or analytics systems.

It helps analysts and businesses understand:

- what data is collected.
- how data elements relate
- how to efficiently store and retrieve data.

Modeling Techniques

1. Entity - Relationship (ER) modeling

It is one of the most common techniques used to represent data. It's concerned with defining 3 key elements:

- entities (objects or things within the system)

- relationships

(how these entities interact with each other)

- attributes (properties of the entities)

The ER model provides a clear, visual representation of how data is structured to help map the connections b/w different data points.

Ex:-

E-commerce store

customers	
int	CustomerID
str	Name
float	Email

orders	
int	Order-ID
int	CustomerID
float	Total

products	
int	Product-ID
float	Price

order_products	
int	Order-ID
int	Product-ID

↳ contains

↳ included in

2. Hierarchical modeling

It is like a tree structure

The data is stored in parent and child records, each of which may have a set of attributes



1

[student]

{college}

{faculty}

2

3

4

3. Relational Modeling

- stores data in tables (relations) with rows and columns
- based on keys (primary, foreign) and constraints
- supports SQL-based systems like MySQL.

PostgreSQL

Each table contains information relevant to a single logical entity, like a store for example and the link b/w these tables is represented by a relationship.

Product		Sales		Store	
ID	Name	ID	Product ID	ID	Name
1	Washing machine	1	1	1	Electronics Mart
2	UK Flatscreen TV	2	2	2	Furniture Mart
3	Apple Mac	3	3	3	Quantum Mart

4. Dimensional modeling

The dimensional data model is also a variation of the relational model.

- designed for data warehousing & OLAP systems
- focuses on facts (quantitative data) and dimensions (contextual data)
- used to build star schema & snowflake scheme

5. Object-oriented modeling

- represents data using objects
- supports inheritance, encapsulation & polymorphism
- used in object-oriented databases and modern programming.

6. Conceptual, logical & physical models
- **Conceptual** :- High focuses level, focuses on what data exists and its relationships (ER diagrams)
 - **Logical** :- Defines how data is organized logically (tables, fields)
 - **Physical** - specifies how data is stored physically on system (indexes, partitions)

Using appropriate data modeling techniques ensures efficiency, clarity and scalability in both operational and analytical systems. It acts as a bridge b/w real-world data needs and technical implementation.

In Missing Imputations

What is missing data?

Missing data refers to the absence of a data value in a data set where it should exist.

- Impacts:-
- leads to inaccurate analysis & biased results
 - affects model performance in ML
 - can cause misleading interpretation.

Types of missing data.

1. MCAR (missing completely at random)

missing values occur independently of any data value

Ex:- A survey respondent accidentally skips a question

2. MAR (missing at random)

Missingness is related to observed data but not the missing value itself.

Ex:- People with higher income level tend to skip the salary question.

3. MNAR (missing not at random)

missing values depend on the variable itself.

Ex:- Individuals with very low income purposely leave the income field blank.

Missing Imputation:-

It is the process of filling in or replacing missing values in a dataset to ensure that the data is complete and usable for analysis or modeling.

- Imputation techniques

1. Deletion methods:

- pairwise deletion:-

remove the entire row if any value is missing

- pairwise deletion: removes all available data without removing the row entirely useful in correction analysis.

2. mean / median / mode Imputation

- replace missing values with the mean, median or mode

- easy but may reduce variance and introduce bias
e.g:- student age

A	20	mean = $(20+22)/2$
B	22	= 21
C		$B = 21$

3. forward/backward fill:-

→ forward fill (ffill): use the last known value

→ backward fill (bfill): use the next known value.

4. K-nearest Neighbors (kNN) imputation

Missing values are filled using the average value of the k most similar rows.

e.g:- If a row is missing age, kNN finds k rows with similar income and education and averages their ages.

5. Regression Imputation

predicts the missing value using regression models based on other variables.

Ex:- predict missing income using age and education.

6. Multiple Imputation:-

- generates multiple datasets with estimated values and averages the result.

e.g:- Used when the dataset has 20-30% missing values, uncertainty must be preserved.

Missing imputations are essential for preparing clean and usable data. Depending on the nature of the data & missingness, the appropriate technique must be selected to preserve data integrity, reduce bias, and improve model accuracy.

Need for business modeling

what is business modeling

Business modeling refers to the process of creating abstract representation of business processes, operations and data flows using mathematical, statistical, or computational tools.

These models help stakeholders to agree, understand, analyze and improve business processes, and act as the foundation for effective decision-making, strategic planning, and IT alignment.

why is needed

1. Clarity in business processes
- Provides a visual representation of workflows of activities
- Helps stakeholders understand process interconnections.
- Aids in identifying areas for process improvement.

2. Improved Decision-making

- Offers data-driven insights through simulations and "what-if" analysis.

- Helps managers take informed decisions backed by analytical models.

3. Alignment b/w Business & IT

- Acts as a bridge b/w business teams & technical developers.

- Ensures that IT aligns with organizational goals & work flows.

4. Process Optimization and Efficiency

- Identifies redundancies, bottlenecks, & inefficiencies in work flows.

- Facilitates workflow redesign to improve productivity & reduce costs.

5. Strategic Planning & Forecasting

- Models break down complex problems into manageable components.

- Used in "what-if" analyses to test customer for different strategies.

6. Communication & Documentation

- A business model acts as a common language across departments.

- Helps in onboarding, training & maintaining consistency in operations.

7. Improved customer understanding

- models analyze customer behaviors, preferences, and feedback.
- enables targeted marketing, product recommendation and better user experience.
- e.g. e-commerce platforms model customer buying patterns.

8. Competitive Advantage:-

- modeling helps businesses adapt faster, make smarter moves, to stay ahead of competitors
- transforms raw data into strategic insights.

Real time example:-

Amazon - e-commerce & Marketplace Model

Business modeling is essential for organizations to analyze operations, optimize workflows and innovate with confidence. It improves communication, enables strategic planning and drives data-backed decision-making.

In a competitive environment, businesses that leverage models gain a clear advantage in terms of efficiency, adaptability of customer satisfaction.

Business modeling concepts used:-

- process optimization:- Amazon uses logistics & inventory models to streamline its global supply chain
- customer segmentation
- Revenue simulation:-
models forecast earnings from Subscriptions vs 3rd party sales.

Why it fits? :-

Amazon's model is a perfect case of multi channel business modeling - blending retail, platform, logistics & subscriptions, all aligned via structured data models for decision making & forecasting.

Unit-III

Regression

- Concepts, Blue property assumptions,
- Least square estimation
- Variable Rationalization
- Model building etc.

Logistic Regression:-

- Model Theory
- Model fit statistics
- Model construction
- Analytics applications to various Business Domains etc.

Regression analysis

- It is a statistical technique used to examine the relationship between dependent and independent variables.
- It determines how changes in the independent variable(s) influence the dependent variable, helping to predict outcomes, identify trends and evaluate causal relationships.
- Widely used in fields like business, economics, healthcare and social sciences, regression analysis provides a robust framework for data-driven decision-making.
- It helps in:
 - predicting future values
 - understanding variable relationships.
 - quantifying the impact of predictors.

Methods of Regression Analysis

1. Ordinary least squares (OLS)

It minimizes the sum of squared errors between

actual and predicted values

$$\text{minimize } \sum (y_i - \hat{y}_i)^2$$

use case:

- most commonly used method in simple and multiple linear regression.
- Suitable when assumptions like linearity & homoscedasticity are met.

Q. Maximum likelihood estimation(MLE)

- It choose parameter values that maximize the likelihood that the observed data occurred

$$\hat{\theta} = \arg\max_{\theta} L(\theta | x)$$

Use case:-

commonly in logistic regression and generalized linear models.

3. Ridge and Lasso Regression

- Both are regularization techniques
- Ridge adds L_2 penalty \rightarrow shrinks all coefficients
- Lasso adds L_1 penalty \rightarrow can set some coefficient to zero (feature selection)

$$\text{Ridge: } \sum (y_i - \hat{y}_i)^2 + \lambda \sum \beta_j^2$$

$$\text{Lasso: } \sum (y_i - \hat{y}_i)^2 + \lambda \sum |\beta_j|$$

use case:-

- Useful when dealing with multicollinearity or high-dimensional data
- Prevents overfitting.

4. Stepwise Regression:-

- Adds or removes predictor systematically to improve the model.

Use Case:-
Used in exploratory data analysis to identify significant predictors.

5. Bayesian Regression:-

Incorporates prior beliefs and data evidence into regression modeling

use case:- forecasting

useful in exploratory data analysis to

identify significant predictors.

Types of Regression Analysis:

1. Simple Linear Regression

- Assume that there is only one independent variable x . If the relationship b/w x (independent variable) and y (dependent or output variable) is modeled by the relation,

$$y = \alpha + bx$$

then the regression model is called a linear regression model.

use case:- predicting housing prices based on square footage

2. Multiple regression

- Assume that there are multiple independent variables say x_1, x_2, \dots, x_n . If the relationship one dependent and 2 or more independent variables

$$y = \alpha_0 + \alpha_1 * x_1 + \alpha_2 * x_2 + \dots + \alpha_n * x_n$$

use :- analyzing how advertising budget, product price and seasonality affect sales revenue.

3. Logistic Regression:-

- used when the dependent variable is binary (e.g. yes/no, success/failure). $P(y=1) = \frac{e^{\alpha_0 + \alpha_1 x}}{1 + e^{\alpha_0 + \alpha_1 x}}$
- Predicting whether a customer will buy a product

based on demographic features.

4. Polynomial regression.

Assumes that there is only one dependent variable x . If the relationship b/w independent variables x and dependent or dep variable y is modeled by the relationship

$$y = a_0 + a_1 * x + a_2 * x^2 + \dots + a_n * x^n$$

for some positive integer $n \geq 1$, then we have

a polynomial regression.

BLUE property of estimators (Gauss-Markov Theorem):

The term BLUE stands for

Best Linear Unbiased Estimator.

Variance tells us how much the data is spread out.

Variance = How far values are from the average (mean)

What is an estimator

In statistics an estimator is a rule for

calculating an estimate of given quantity based on observed data.

Ex: x follows a normal distribution, but we don't know the parameters of our distribution, namely mean (μ) and variance (σ^2)

To estimate the unknowns the usual procedure is to draw a random sample of size ' n ' and use the sample data to estimate parameters

Ols is a method used in linear regression to draw the best possible straight line through a set of data points.

- It helps us predict something based on something else.

Ols chooses the line that minimizes the error or it calculates the square of each error, adds them all, and finds the line where this total squared error is smallest.

Ordinary least squared = least total of squared errors.

What is BLUE

In linear regression, the term BLUE refers to the properties of an estimator that is

BEST LINEAR UNBIASED ESTIMATOR

These properties are guaranteed for ordinary least squares (OLS) estimators if certain assumptions are met. This is formalized in Gauss-Markov theorem.

meaning of each term:-

Best - Has the minimum variance among all linear biased estimators.

Linear - Is a linear function of the observed data.

Unbiased - On average, the estimator gives the true value of the parameter.

Estimator - A rule or formula used to calculate an estimate from sample data.

Gauss-Markov Theorem

These states that:-

Under certain assumptions, the OLS estimator is Best Linear Unbiased Estimator (BLUE) of the regression coefficients.

Assumptions required for BLUE

1. Linearity :- The model is linear in parameters

$$Y = \beta_0 + \beta_1 X + \epsilon$$

2. Random sampling :- Data is collected via random sampling.

3. No perfect multicollinearity :-

Independent variables are not perfectly correlated.

4. Zero mean of errors :- $E[\epsilon] = 0$, meaning errors average to zero.

5. Homoscedasticity :- Constant variance of error terms across observations.

6. No autocorrelation :- Errors are not correlated with each other.

Why BLUE is important?

- ensures that OLS estimates are efficient and accurate.

- forms the foundation of linear regression analysis

- If assumptions are violated, estimates may still be unbiased but not best.

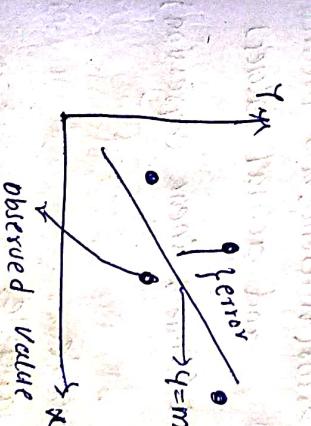
The BLUE property ensures that OLS regression estimators are statistically optimal when the Gauss-Markov assumptions hold true. It gives confidence that predictions are reliable, unbiased and have minimum variance - a key strength in data modeling and analytics.

Least Square Estimation

It is a method used in linear regression to find the best fitting straight line through a set of data points.

It works by :-

minimizing the sum of squares of the errors (the differences b/w actual values and predicted values).



least squares method is used to derive a general linear equation between two variables. When the value of the dependent and independent variables they are represented as x and y , coordinates in 2D cartesian coordinate system.

The least squares method is a popular statistical technique used in

- regression analysis
- predictive modeling.

Its main purpose is to find the line of

best fit that minimizes the sum of squared differences between the actual data points and the values predicted by the line.

Why is it used:-

- To predict a dependent variable (y) from an independent variable (x)
- To draw a straight line that fits the data as closely as possible.
- To reduce the total error (difference b/w predicted and actual values)

Basic equation of the line (Simple linear regression)

$$y = mx + c$$

where y :- Dependent variable

x :- Independent variable.

m :- Slope of the line

c :- Intercept.

Least squares formulas

- slope (m)

$$m = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$

$$\text{Intercept } (c): \\ c = \frac{\Sigma y - m(\Sigma x)}{n}$$

These formulas help us calculate the best fitting line that minimizes the sum of squared errors.

Steps in the least square method

1. Identify x and y (independent and dependent variables)
2. Create a scatter plot of the data points.
3. Calculate the mean of x and y
4. Use formulas to calculate the slope m and intercept c .

$$m = \frac{\Sigma (x - \bar{x})(y - \bar{y})}{\Sigma (x - \bar{x})^2}$$

5. write the equation of the best fit line as

$$y = mx + c$$

Example :-

$$(1,3) (2,4) (4,8) (6,10) (8, 15)$$

Step 1 :- Find the mean \bar{x} and \bar{y}

$$\bar{x} = \frac{1+2+4+6+8}{5} = 4.2$$

$$\bar{y} = \frac{3+4+8+10+15}{5} = 8$$

Step 2 :- Use the formula

x_i	y_i	$x - \bar{x}^i$	$y - \bar{y}^i$	$(x - \bar{x}^i)(y - \bar{y}^i)$	$(x - \bar{x}^i)^2$
1	3	3.02	5	16.0	10.04
2	4	2.02	4	8.8	4.04
4	8	0.02	6	0.0	0.04
6	10	-1.02	-2	-2.04	3.6
8	15	-3.02	-7	-21.04	9.24
				<u>55.0</u>	<u>32.08</u>

Step 3 :- Calculate slope m

$$m = \frac{\sum (x - \bar{x}^i)(y - \bar{y}^i)}{\sum (x - \bar{x}^i)^2} = \frac{55}{32.08} = 1.68$$

Step 4. Calculate Intercept c

$$y_{\text{line}} = y - mx = 8 - (1.68 \times 4.02) = 0.94$$

Final Line :-

$$y = 1.68x + 0.94$$

Graph Interpretation :-

- The red dots represent the actual data points
- The straight red line shows the best fit line
- This line can now be used to predict y values for any x.

Limitations :-

- Not suitable for non-linear data
- Gets affected by outliers
- Assumes equal variance and normal distribution of errors.

The least squares method helps you find the best line through data points by minimizing total squared error between actual and predicted values. It is essential for linear regression and accurate prediction in business, science and engineering.

Variable Rationalization

It is the process of identifying, selecting and transforming variables in a dataset to ensure that only the most relevant and meaningful variables are used in the analytical model.

- Improve the accuracy of a model
- Reduce noise and redundancy
- make the model easier to understand and faster to compute.

Why it is Needed :-

- When building a regression model or any analytical model:
- You may have too many variables
 - Some variables may be irrelevant, correlated or repetitive
 - keeping only the meaningful variables leads to better performance.

Steps in Variable Rationalization:-

Example :-

1. Understanding Data Context:-

- Review variables based on domain knowledge.
- Identify which variables are logically important

2. Feature selection

choose only the most important variables for the model.

Methods:-

• Correlation analysis (remove highly correlated variables)

• P value (in regression, if $P > 0.05$, variable may be dropped)

• Forward / Backward Stepwise Selection

3. Feature transformation

Transform variables to make them more useful.

e.g:- Taking logarithms of large nos

• Creating polynomial features for non-linear patterns.

• Normalizing / scaling variable to same range.

4. Remove redundant or noisy variables.

- Remove variable that:

• Repeat the same information as other

• Contain too many missing values.

• Add randomness, but not meaningful

5. Domain-Based selection

Sometimes, your knowledge of the business or domain helps decide which variables are meaningful and which are not.

After rationalization, you keep

• Hours studied

• Remaining are not required.

• Attendance (Attendance required)

Tools:-

• PCA (Principle Component Analysis)

• Variance Inflation Factor (VIF)

Advantages:-

- Reduce the model complexity

- Improves prediction accuracy

- Faster computation

- reduce overfitting

Disadvantages:-

- can be time consuming

- risk of removing useful variables

- may introduce bias

- requires expertise

Summary:-

Variable Rationalization is the process of selecting the most relevant variables and removing or transforming those that are unnecessary or harmful to model accuracy.

It improves performance, simplifies models and helps make better predictions.

Model Building

It is the process of creating a predictive or analytical model using data. The goal is to understand relationships between variables or predict outcomes accurately.

Model building is the process of

- preparing data.
- choosing a model
- training it to learn from the data.
- evaluating its performance.
- and using it for predictions or analysis.

The main goal is to create a model that can understand patterns in training data and predict accurately on new/unseen data.

Key Objectives

- Achieve high accuracy on training & testing dataset
- Ensure the model generalizes well, not just memorizes data.
- make the model suitable for real-world use.

Steps in Model Building:-

1. Split the Dataset

To avoid overfitting, we split the data into:

- Training dataset - used to train the model
- Testing dataset - used to evaluate the model's performance.

Usually done in a 75% vs 25% or 80:20 ratio using python's train-test-split.

Code:-

```
x = data['x', 'y']  
y = np.random.rand(1000)  
train_x, test_x, train_y, test_y =  
train_test_split(x, y, test_size=0.2)
```

2. Scale the Data

Why?

- prevent features with big values from dominating small-scale feature.
- make models like KNN and Gradient descent accurate and faster.

→ use min max scale x to scale data to range [0, 1]

from sklearn.preprocessing import MinMaxScaler

```
scaler = MinMaxScaler()  
train_x_scaled = scaler.fit_transform(train_x)  
train_x_scaled = scaler.transform(test_x)
```

3. choose the model

depending on the problem

Regression - linear regression, Decision tree
Classification - logistic regression, KNN, SVM

here we are using regression problem so we use decision tree.

from sklearn.tree import DecisionTreeRegressor

reg = DecisionTreeRegressor(min_samples_split=4, max_leaf_nodes=10)

reg.fit(X_train_scaled, y_train)

4. make predictions

5. evaluate the model

- (how we measure how well the model predicated using Mean squared error (MSE))
—A lower MSE means better performance.

6. visualize the model

- Advantages :-
—easy to interpret
—can handle both numerical & categorical data.
—visualize structure make decisions transparent.
Limitations:- —can overfit if not pruned or limited
—sensitive to small changes in data.

Summary :-

Model building is the process of training a machine to learn patterns from data and make predictions. It includes data splitting, scaling, training, testing, evaluation & visualization.

Logistic Regression

Logistic regression is a statistical technique

used to predict the probability of a binary outcome used to:

e.g Yes/No, 0/1, true/false).

Unlike linear regression which predicts a continuous value, logistic regression is used when the target variable is categorical.

Real-life example

- will a customer buy a product? (Yes or No)
—Is the email spam? (spam or Not spam)
—will a student pass the exam? (pass or fail)

Model theory of logistic regression

It is a classification technique used to predict the probability of a binary outcome e.g (Yes/No)
—Why not use linear regression?
It predicts any value, even below 0 or above 1 which is not valid for probability.
—logistic regression fits.

It uses the sigmoid function to keep the output between 0 and 1.

Sigmoid function (S-shaped curve)

$$\sigma(z) = \frac{1}{1+e^{-z}}$$
 where $z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$

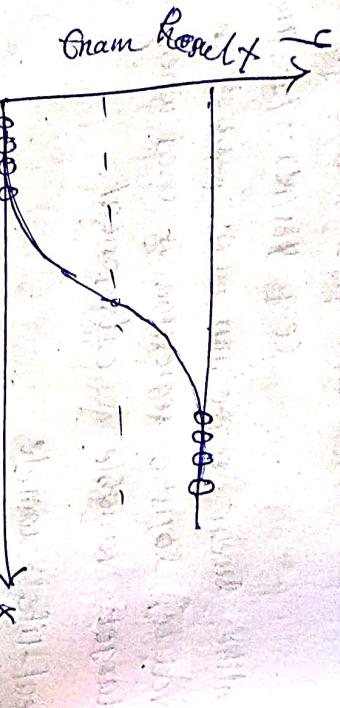
Output: A probability between 0 and 1

- If $\sigma(z) > 0.5$, predict class 1

- If $\sigma(z) < 0.5$, predict class 0

$$\alpha = -4 + 1 \times 3 = -1$$

$$\sigma(-1) = \frac{1}{1 + e^{-(z)}} = 0.27 \Rightarrow \text{predict class 0}$$



for 5 hours:-

$$z = -4 + 1 \times 5 = 1$$

$$\sigma(z) = \frac{1}{1 + e^{-1}} \approx 0.73$$

73.01%. → predict pass(1) probability of passing.

Here β_0 & β_1 are not calculated manually like linear regression.

And in β_0 (Intercept) and β_1 (slope).

here they use logit function.

Advantages :-

- Simple and easy to implement
- Predict probability (0 to 1)
- Efficient for small datasets.
- less prone to overfitting

Ex:- predict if a student passes based on hours studied.

Hours (x) Passed (y)

1	0
2	0
3	0
4	1
5	1

If the trained model

- Healthcare
- Banking & finance
- Email filtering
- Education

- $\beta_1 > 0$: As x increases, probability of class 1 increases
- $\beta_1 < 0$: As x increases, probability of class 1 decreases
- Exponentiating β_1 gives the ratio odds ratio:

$$\text{Odds Ratio} = e^{\beta_i}$$

$$\text{Ex- } \beta_1 = 1 \Rightarrow e^1 \approx 2.0 \rightarrow 18$$

\rightarrow extra hour of study doubles the odds of passing

Model fit statistics

Model fit statistics help us evaluate how well a logistic regression model represents the data i.e., how accurate, meaningful, and reliable it is.

1. Deviance

It is like the "error" in logistic regression

• Null Deviance: error of a model with only the intercept (no predictors)

• Residual Deviance: error after adding predictors (features)

Interpretation:

- lower residual deviance = better fit

- compare both using Likelihood Ratio Test (LRT)

2. Likelihood Ratio Test (LRT)

It compares the null model and the full model.

$$G^2 = -2 \times (\text{Log } L_{\text{null}} - \text{Log } L_{\text{full}})$$

• If G^2 is large, it means the model with predictors is significantly better.

• Compare G^2 to the chi-square table for significance.

3. Pseudo R^2

Since regular R^2 used in linear regression doesn't apply, we use pseudo R^2 to measure how

much variance is explained.

$$\text{McFadden's } R^2 = 1 - \frac{\log L_{\text{full}}}{\log L_{\text{null}}}$$

cor. Incl. similar to R^2 but not scaled to 1

4. Akaike Information Criterion

$$AIC = 2k - 2 \log L$$

• k = no. of model parameters

• L : likelihood of the model

• Lower AIC = better model

Used for comparing multiple models.

5. Confusion matrix

A table showing the actual vs predicted values

Actual	0	1
0	TP	FP
1	FN	TP

$$\text{Accuracy} = \frac{TP + TN}{(TP + FN) + (FP + TN)}$$

6. Precision and Recall

Precision: How many predicted positives are actually correct.

$$\text{Precision} = \frac{TP}{TP + FP}$$

• Recall (sensitivity): How many actual positives are correctly predicted.

$$\text{Recall} = \frac{TP}{TP + FN}$$

7. F1 score

It is the harmonic mean of precision and recall. It balances the two when you need a single performance measure.

$$F_1 \text{ score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

c. ROC curve (receiver operating characteristic)

Plots True positive rate (Recall) vs False positive rate.

$$\text{false positive rate} = \frac{FP}{(FP+TN)}$$

d. AUC (area under curve)

measures the model's ability to distinguish between classes.

1.0 - perfect model, 0.5 - No better than random

<0.5 - worse the random.

Model construction:-

Assumptions:-

- Binary dependent variable (0 or 1)
- Linearity the logit for continuous variables
- No multicollinearity
- Independence of observations.

Logistic function

$$P(Y=1 | X) = \frac{1}{1 + e^{-(B_0 + B_1 X_1 + B_2 X_2)}}$$

Steps in construction:-

1. Define the problem
2. Select feature (X) and target (Y)
3. Split data into training / testing.
4. Check multicollinearity
5. Train the logistic regression model.
6. Predict & evaluate using accuracy, AUC etc.

Detailed flow of logistic regression.

$$\text{Input features } (x_1, x_2, \dots, x_n) \rightarrow \text{linear combination} \rightarrow \text{sigmoid fn} \\ z = B_0 + B_1 x_1 + \dots + B_n x_n \rightarrow \sigma(z) = \frac{1}{1 + e^{-z}}$$

Output probability

$$\downarrow \\ P(Y=1 | X)$$

Thresholding

$$\downarrow \\ \text{If } P > 0.5 \rightarrow 1 \text{ else } 0$$

Final classification

$$\downarrow \\ 0 \text{ or } 1$$

Analytics - applications to various business

Domains:

Data - analytics is the process of examining large datasets to discover patterns, trends and insights.

It is widely used across different business

domains to improve decision-making, enhance customer experience, reduce costs, and predict future outcomes.

1. healthcare

Predictive analytics for patient health monitoring disease prediction and hospital resource management

e.g:- Using machine learning to predict patient

readmission rates and optimize hospital bed allocation.

2. Retail - e-commerce

Customer segmentation, demand forecasting and recommendation systems.

e.g:- Amazon uses customer behavior data to recommend personalized products and optimize inventory.

3. Banking and finance

Credit risk modeling, fraud detection, customer segmentation.

e.g:- Banks use analytics to detect suspicious transactions in real-time to prevent fraud.

4. Manufacturing

Predictive maintenance, quality control, and chain optimization.

e.g:- An automobile company uses sensor data from machines to predict equipment failures before they occur.

5. Telecommunications

Churn analysis, network optimization, and pricing strategies.

e.g:- Telecom companies analyze call data to detect customers likely to switch to competitors and provide offers to retain them.

6. Logistics and supply chain

Route optimization, demand forecasting, and warehouse management.

7. Marketing and sales

Campaign performance analysis, customer lifetime value prediction and ROI measurement.

8. Education

Student performance prediction, adaptive learning and resource allocation.

Analytics has become an essential tool in modern business operations. By converting raw data into actionable insights, organizations across domains can enhance decision-making, reduce cost

Unit - IV

Object Segmentation:

- Regression vs Segmentation
- supervised and unsupervised learning
- Tree Building - Regression
- Classification

Overfitting

- Pruning and Complexity
- Multiple Decision Trees etc.

Time series method:

- ARIMA
- Measures of forecast accuracy
- STL approach
- extract features from generated model
 - as Height, charge, energy etc
- analyze for prediction.

Regression

Aspect

Regression is a supervised learning technique used to predict a continuous numeric value based on input feature.

definition

An unsupervised learning technique used to group similar data points into clusters.

Supervised learning vs unsupervised learning.

Learning Type

To predict a dependent variable (e.g. sales, price, temperature) to group similar data points without predefined labels.

Output - A numeric value (e.g. house price = ₹ 50 lakhs) - A cluster or segment label (e.g., customer Group A, B, C)

examples predicting house prices, customer segmentation, student scores, sales forecasting.

linear regression, k-means clustering, polynomial regression, hierarchical clustering, decision tree regression

Alg 0

use case forecasting future trends in buying or demand

Identifying customer groups for targeted marketing.



Segmentation

Supervised and unsupervised learning

What is ML (machine learning)

It is a branch of artificial intelligence where machines learn from data without being explicitly programmed. ML is broadly classified into :-

1. Supervised learning
2. Unsupervised learning.

Supervised learning

It is a machine learning approach where the model is trained using labeled data (i.e., each input has a known output). The model learns the mapping functions from input to outputs.

Key concepts:-

- Labeled data - Data with input-output pairs ($x \& y$ known)
- Goal
 - predict output values for new inputs.

Output types - continuous (regression) or categorical (classification)

Types of supervised learning

1. Regression :-

- It is a type of supervised learning that is used to predict continuous values.

example:- Predicting house price, temperature etc.

Algorithms:-

- linear regression
- polynomial regression
- Lasso regression
- Ridge regression.

2. Classification

Classification is a type of supervised learning that is used to predict categorical values, such as whether a customer will churn or not, whether an email is spam or not, or whether a medical image shows a tumor or not.

- classification algo learn a function that maps from the input features to probability distribution over the output classes.

Algorithms:-

- logistic regression, support vector machine (SVM)
- decision tree, - random forest
- naive bayes.

Applications:-

- emails - spam detection
- healthcare - disease diagnosis
- e-commerce - product recommendation
- finance :- credit fraud detection
- NLP :- sentiment analysis, translation.

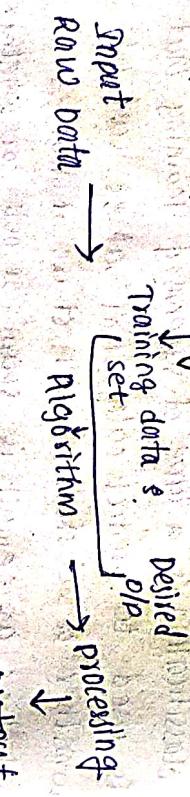
Advantages

- Builds on previous experience
- Gets better with more data
- solves real-world classification & prediction tasks
- offers clear, interpretable results.

Disadvantages:-

- Needs large labeled datasets (costly to create)
- can struggle with complex or unstructured data
- high computation and training time for large data.

~~Unsupervised learning~~ ~~supervisor~~



Unsupervised learning

It is a machine learning method where the model is trained on unlabeled data. The goal is to discover patterns, structures or relationships without any predefined output.

Key :-

- Data - unlabeled (no output provided)
- Goal - find clusters, associations or hidden structures
- No supervision - No teacher or known answers
- Output - Often groups, clusters or rules rather than specific predictions.

Input raw data → Interpretation → Algorithm → processing → output

Model training

Types of unsupervised learning :-

1. Clustering :-
Groups similar data points based on clusters based on distance / similarity.

clustering algorithms work by iteratively moving data points closer to their cluster centers and further away from data points in other clusters.

Algo :-
- k-means clustering
- hierarchical clustering

- DBSCAN

- Gaussian Mixture Models (GMM)

2. Association rule learning.

Finds interesting rules between variables
e.g:- If a customer buys Bread, they are likely to buy Butter.

Algo :-
- prior algorithm, eclat algorithm,
fb growth.

Applications:-
- e-commerce - customer segmentation for targeted marketing

- Recommendations - suggest movies / music / products.
- science - finding hidden patterns in research data.
- Anomaly detection - fraud or intrusion detection.

Advantages :-

- No need for labeled data
- useful for data exploration
- can reveal unknown patterns in complex datasets.
- handles large and high-dimensional data well.

Disadvantages :-

- No predefined OLP → Hard to measure accuracy
- Results may need manual validation
- sensitive to outliers, noise and missing values
- may produce less precise groupings than supervised learning.

TREE BUILDING - Regression, Classification

Tree based models are powerful tools for both regression and classification tasks. These models split data into smaller and smaller groups based on feature values, creating tree-like structure.

Tree building is the process of constructing a decision tree, which uses a series of decisions (nodes) to predict an output. It splits dataset into branches based on feature values until a final decision (leaf) is made.

CART (Classification And Regression Trees) is a variation of the decision tree algorithm. It can handle both classification and regression tasks.

CART is a predictive algorithm used in Machine learning and it explains how the target variables values can predicted based on other variables. It is a decision tree where each fork is split into a predictor variable and each node has a prediction for the target variable at the end.

The term CART serves as a generic term for the following categories of decision trees:

- Classification trees:- The tree is used to determine which "class" the target variable is most likely to fall into when it is continuous.
 - Regression trees:- These are used to predict a continuous variable value.

CART Algorithm:-
It is a decision tree algorithm that is used for both classification and regression tasks. It is a supervised learning algorithm that learns from labelled data to predict unseen data.

- Tree structure:- CART builds a tree like structure consisting of nodes and branches. The nodes represent different decision points, and the branches represent the possible outcomes of those decisions.

- Splitting criteria:-

CART uses a greedy approach to split the data at each node. For classification tasks, CART uses Gini impurity as splitting criterion. For regression tasks, CART uses residual reduction as splitting criterion.

- Pruning:- To prevent overfitting of the data, pruning is a technique used to remove the nodes that contribute little to the model's accuracy.

How does CART work

- Try all possible splits for each input variable
- Pick the best split (based on Gini Impurity or error reduction)
- Split the data into two parts using the split best split.
- Repeat for each new part (sub-tree) until:
 - Tree gets too deep
 - Or there's no better split left.

What is Gini Index / Gini impurity?

The Gini Index tells us how impure or mixed a group (subset) is.

- It is used in classification problems to decide where to split the data in a decision tree

Gini impurity checks: "what is the chance of making a wrong guess if we randomly assign a label from a group?"

$$\text{formula: } \text{Gini} = 1 - \sum_{i=1}^n p_i^2 \quad p_i = \text{probability}$$

Gini Index value:

• 0 → perfectly pure (only one class present)

, close to 1 → highly impure (many mixed classes)

• 0.5 → Two classes with equal 50-50 probability.

Ex:- Imagine you have a node with

- 4 apples
- 6 oranges

$$\text{papple} = 0.4 \quad \text{porange} = 0.6$$

$$\text{Gini} = 1 - (0.4^2 + 0.6^2) = 1 - (0.16 + 0.36) = 0.48$$

This node is not pure because it has a mix of 2 classes.

Why CART uses Gini?

- CART looks for the split that minimizes the Gini impurity.
- At each node, it tries all possible splits, calculates the Gini for each, & pick the best (lowest Gini).
- Lower Gini = more pure = better split.

CART for classification

In classification, CART is used when the output (target variable) is categorical - for example Yes/No, Pass/Fail, spam/not spam.

How it works

- Start with the full dataset (root node)
- Check all features
- For each feature, try different split values.
- Calculate Gini impurity for each possible split
- Choose the split that gives the lowest Gini impurity (most pure groups)

- Nodes become pure
- Tree reaches a maximum depth
- There's not enough data to continue split

CART for Regression

CART for regression is a decision tree learning method that creates a tree like structure to predict continuous target variables.

How does it work?

It works by splitting the training data recursively into smaller subsets on specific criteria. The objective is to split the data in a way that minimizes the residual reduction in each subset.

- **Residual Reduction:-**
It is a measure of how much the average squared difference b/w the predicted value and the actual values of the target variable is reduced by splitting the subsets.

- **Splitting Criterion:-**
It evaluates every possible split at each node and selects the one that results in the greatest reduction of residual error in the resulting subsets.

- **CART model representation:-**
It builds decision trees step by step based on data. Here is how the CART model is formed.

1. Greedy Algorithm
2. Stopping criterion.

3. Tree Pruning
4. Data preparation for CART.

Overfitting

It happens when a machine learning model learns the training data too well. Instead of just learning the general pattern, it also learns noise, outliers, or random fluctuations.

* Think of it like a student memorizing answers instead of understanding concepts they score high in practice tests but fail in actual exam.
causes:- model is too complex
• Not enough training data
• Too many parameters compared to the data size.
• Training for too many iterations.

Effects of Overfitting:-

- on training data
 - on test data
- very accurate (low error) poor performance (higher)
 - looks perfect fails to generalize.

example:- If you train a decision tree to keep splitting

- until each node has just 1 sample, it will:
- Give 100% accuracy on training data.
- But fail on new, unseen examples because it memorized, not learned.

How to detect overfitting

- Training accuracy is very high, but test accuracy is low.
- Performance changes a lot when switching between datasets (i.e., high variance).
- Validation error starts increasing while training error keeps decreasing.

Simple tricks to avoid:-

- Prune the decision tree
- Use cross-validation to monitor performance during training.
- Limit model complexity
- more training data.

Pruning

It is the process of removing unnecessary parts of a decision tree to make it simpler and less likely to overfit.

Why use pruning:-

- Reduces overfitting
- Improves model accuracy on test data
- speeds up predictions
- makes the tree easier to understand

Types of pruning

- i. Pre-pruning (Early stopping)
- & stop growing the tree before it becomes too complex.

- common stopping conditions

- Tree reaches a maximum depth.
- A node has too few samples
- The gain (impurity reduction) from a split is too small.

*Saves time and prevents the tree from getting too large in the first place.

2. Post-pruning (Simplification after growth)

1. First build the full tree, then cut back branches that are not useful.
 - *Keeps only the splits that improve accuracy or reduce complexity.
- This is often done using a validation set.

Ex:-

Imagine a tree branch that classifies just 2 out of 1000 training samples.

If those 2 samples don't really change the model's performance, we can prune that branch and avoid learning noise.

Pruning = "Trimming the decision tree to keep only the useful parts".

Complexity :-

Complexity refers to how large or detailed a decision tree becomes. This includes:

- Tree Depth - How many levels it has
- Number of nodes
- Splits - How many times the data is divided.

Why does complexity matter:-

High complexity tree:-

- May overfit:- memorizes training data, including noise.
- learns detailed patterns and can fit complex patterns

Low complexity tree :-

- Too simple to learn actual parameters
- May underfit:- performs poorly even on training data

Complexity control parameters:-

- Tree depth limit
- Minimum samples per leaf or split
- Minimum impurity decreases required for split.

Goal:-

To balance complexity so the model is not simple or too complex - this helps in achieving good accuracy and better generalization to unseen data.

Multiple Decision Trees

Instead of using just one decision tree, we use many decision trees together.

This helps:-

- Reduce overfitting
- Increase accuracy
- Improve generalization

This concept is called ensemble learning.

What use multiple trees

A single tree can:

- memorize training data (overfit)
- Be sensitive to noise.

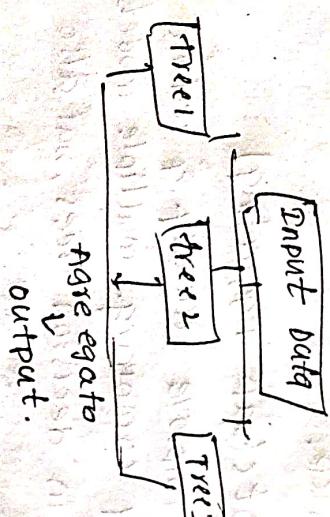
multiple trees:-

- work as a team (like a voting committee)
- give more stable and accurate predictions.

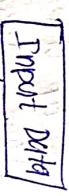
Type of ensemble methods:-

- Random forest
- Builds many trees on random subsets of data and feature. final output is majority vote (classification) or average (regression)

- uses bagging + random feature selection
- trees are trained independently
- works well for both classification & regression
- less overfitting than single tree.



are equal
outputs.



Boosting = many weak models

One strong model

- Build models one after another sequentially

- first model is trained on normal training data

- next model focuses on mistakes made by

the previous model.

- assign weights to data

lower weight - correctly predicted
higher weight - misclassified

- New model focuses more on hard classes

- This repeat

, until the model performs well

- A set number of model is built

- final prediction is a combination of all models

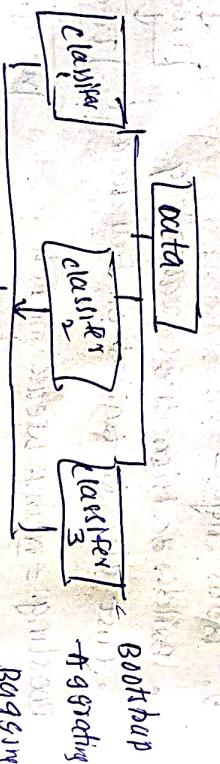
(like weighted voting).

* To make a strong model by fixing errors
step by step.

3. Bagging (Bootstrap Aggregating)

It used to improve accuracy & reduce overfitting by combining multiple models.

mostly used in decision trees, but also with other models.



ensemble classifier.

Advantages :-

- better accuracy, less overfitting

- more robust to noise

Disadvantages:-

- more computation, harder to interpret
- slower prediction.

Using multiple decision trees.

Improves accuracy

* Reduces overfitting
- gives stronger, more stable models.

Time Series :-

Time series analysis and forecast are crucial for predicting future trends, behaviours, and behaviours based on historical data.

* A time series is a sequence of data points collected, recorded or measured at successive, evenly-spaced time intervals.
e.g. Intervals may be monthly, yearly, daily etc.

How works:-

- create multiple datasets

- train multiple models

- train in parallel

- combine results.

Time series forecasting is the process of using a statistical model to predict future values of a time series based on past results.

Components of Time series data

There are four main components.

1. Trend (T)

A long-term upward or downward movement in the data.

- Shows overall direction of the data over time
- Can be increasing, decreasing or flat (no trend)

e.g.: - Company revenue increasing steadily over 5 years.

2. Seasonality (S)

A repeating pattern in data at regular intervals (days, months, quarters)

- Caused by seasonal effects such as holidays, weather etc.
- Pattern is predictable.

e.g.: - Ice cream sales increasing every summer.

3. Cyclical (C)

Similar to seasonality but not regular or fixed. Longer-term patterns often tied to economic cycles.

- Not always predictable like seasonality.
- Usually spans years (e.g. boom & recession).
- e.g.: - Real estate prices over decades.

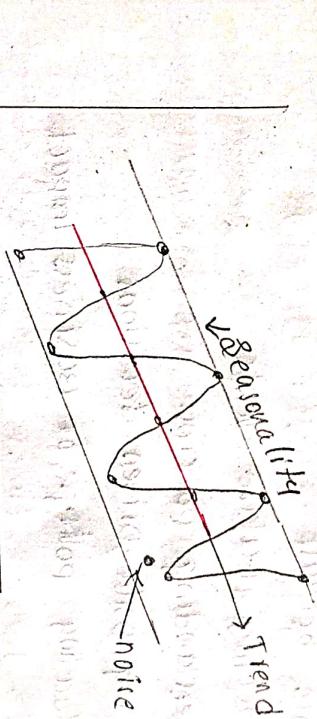
4. Irregular or Residual component (E)

Represents random noise or fluctuation in the data.

- Caused by unexpected events, like strike or natural disaster.

- Cannot be predicted.

e.g.: - Sudden drop in online traffic due to a server outage.



Characteristics:-

- Chronological order, - sequential order
- Temporal components - constant frequency

- Dynamic nature.

Mathematical representation / Decomposition

A time series $y(t)$ can be expressed as:

1. Additive model

$$y(t) = T(t) + S(t) + C(t) + E(t)$$

- Used when the components are independent & linear.
- Suitable when seasonal variation stays constant over time.

Q. Multiplicative Model

$$y(t) = T(s) \times S(t) \times C(t) \times I(t)$$

- used when components interact and amplify each other
- suitable when seasonal variations increases/decreases with trend.

Importance of time series analysis:-

Time series analysis is crucial both business and science because

- forecasting future events
- understanding patterns
- making data-driven decisions
- evaluating policy or business impact.

Advantages :-

1. Captures temporal structure
2. Helps discover hidden insights
3. enables accurate forecasting
4. useful in automation.

Types of time series methods

1. Decomposition methods like STL
2. Smoothing techniques moving average, exponential smoothing
3. Machine learning models LSTM, XGBoost

Forecasting models of time series:-

(AR, MA, ARIMA, ARIMAX)

1. AR (Autoregressive model) — (P)

uses past values to predict future
order of AR model is denoted by (P)
 $P = \text{no. of lagged observations}$

$$AR(P) \quad y_t = c + \varphi_1 y_{t-1} + \varphi_2 y_{t-2} + \dots + \varphi_p y_{t-p} + \varepsilon$$

2. MA (Moving Average model) — (q)

predicts current value using past forecast errors

$$q = \text{no. of lagged error terms.}$$

$$y_t = \varepsilon_t + \varphi_1 \varepsilon_{t-1} + \varphi_2 \varepsilon_{t-2} + \dots + \varphi_q \varepsilon_{t-q}$$

3. ARMA

AutoRegressive Moving Average
→ combines AR & MA model for forecasting

- captures both temporal dependencies & error terms

$$\rightarrow \text{Order } (P, q)$$

$$ARMA(P, q)$$

$$\rightarrow y_t = \varphi_1 y_{t-1} + \varphi_2 y_{t-2} + \dots + \varphi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}$$

It only works if time series is stationary.

4. ARIMA

These is our main topic ARIMA

AR - Auto Regressive

I - Integrated

MA - Moving Average

Components:-

$$(\Delta y_t) = \phi_1 (\Delta y_{t-1}) + \phi_2 (\Delta y_{t-2}) + \dots + \phi_p (\Delta y_{t-p}) + Q_1 e_{t-1} + \dots + Q_q e_{t-q}$$

predicts future values using past values
 $p = \text{no. of past values (lags) used.}$

- i. I (Integrated) $\rightarrow d$
- makes the data stationary (i.e. removes trend)

uses differencing, like

$$\Delta y_t = y_t - y_{t-1}$$

- $d = \text{no. of times differencing is applied.}$

3. MA - (q)

uses past forecast errors to predict future values

- $q = \text{no. of past errors used.}$

Notation

$$\text{ARIMA}(p, d, q)$$

p - no. of AR terms

d - differencing degree

q - no. of MA terms.

Importance:- One of the most popular & reliable

models for time series forecasting

- Used in finance, sales, energy, weather prediction
- helps make decision-driven approach based on trends.

When to use? -

use ARIMA when:-

- the data shows a trend
- The data is not stationary (initially becomes stationary after differencing)
- you want to forecast univariate time series.

Advantages :-

- accurate for short term and medium-term forecasting
- flexible - fits many types of time series
- can handle data with trends & noise.

Measures of forecast Accuracy

It tell us how close our predicted values (forecast) are to the actual observed values. They help us evaluate the performance of a time series model.

1. Mean Absolute Error (MAE)

$$\text{MAE} = \frac{1}{n} \sum_{t=1}^n |A_t - f_t|$$

A_t - actual value, f_t - forecast value

n - total no. of observations.

- measures average absolute error b/w predicted and actual values
- simple and easy to interpret.

2. Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{t=1}^n (A_t - f_t)^2$$

- squares the errors \rightarrow larger error have more impact
- helps detect larger deviations.

8. Root mean squared error (RMSE)

$$RMSE = \sqrt{MSE}$$

- Square root of MSE
- Penalizes big errors more than MAE

4. Mean Absolute Percentage Error (MAPE)

$$MAPE = \frac{100\%}{n} \sum_{t=1}^n \left| \frac{A_t - f_t}{A_t} \right|$$

- Expresses error as a percentage
- Easy to understand "the forecast was off by % on average"

Limitations:-

Not desirable when actual values are very close to zero.

Why accuracy matter measure matters :-

- Help to choose best model for forecasting
- Helps to detect overfitting or underfitting.

Advantages :-

- This loop improves accuracy until results stabilize.

STL approach :-

Seasonal - Trend decomposition using LOESS (Locally estimated scatterplot smoothing)

It is a powerful method to break down a time series into 3 clear parts:

1. Trend → long term direction of data
2. seasonality - regular, repeating patterns
3. Remainder / Residual → leftover part (noise)

components :- If you have a time series y_t , STL splits it into

$$y_t = T_t + S_t + R_t$$

- T_t : trend component, S_t : seasonal component

R_t : remainder (random noise)

How it works :-

1. DESS Smoothing is used:

DESS is a flexible smoothing technique that fits local regressions to small parts of data.

- It helps STL capture non-linear trends and complex seasonal effects.

2. STL works in iterations.

First estimates trend → removes it → estimates

seasonality → refines remainder.

This loop improves accuracy until results stabilize.

Extract features from generated model

as Height, Average, Energy etc.

When we build a time series model like STL, ARIMA, we can analyze the output components (trend, seasonal, residual) to extract useful numerical characteristics called features.

These features help us:

- compare different time series
- detect patterns
- Improve forecasting
- Use them in machine learning models.

Common Features extracted :-

Feature Name

Description

Peak (Height) — maximum value in the time series.

Through (Depth) — minimum value in time series

Trend slope — rate of change of the trend component.

Average level — mean value of the series or its trend.

Average energy — sum of squares of the values used to capture intensity or power.

Seasonal strength — how dominant the seasonal component is compared to the noise.

Volatility — standard deviation of the residuals

Cycle length — average time between peaks or recurring patterns.

Autocorrelation — measures how correlated a series is with a lagged version of itself. Outlier count — No. of unusual spikes or dips in the residuals.

Why this :-

- Height tells us about the peak demand usage
- Energy tells how strong or intense seasonal effects are
- Volatility helps in risk assessments
- Slope gives growth trend.
- Used in e.g., forecasting future sales
- Identifying unusual behavior.

Analyze for Prediction

Time series involves understanding past data to accurately forecast future values.

Once features like trend, seasonality, energy, & volatility are extracted, we analyze these components to make informed predictions.

Steps to analyze :-

1. Decompose the time series into component:

2. Break the series into component: Trend, seasonality, Residuals.

2. Feature Analysis

Extract useful features

- Height, energy, volatility, trend slope

3. Model fitting

Fit a model to the time series:

- ARIMA / ARMA

- Exponential smoothing (ets)

4. Evaluate Forecast Accuracy

use these metrics

- MAE, MSE, RMSE etc.

lower error = better prediction

5. Generate forecast

- Predict future time points using the fitted model

Adjust based on seasonal & trend behaviour.

The better you analyse the components, the

smarter your predictions. A lack of analysis of

various components will lead to bad

check feature importance

choose model that best match your data behaviour.

Advantages:-

- Improved accuracy

Applications:-

- Finance

- Healthcare

- Simplify complex data.

- Energy

- Faster model training.

- Refer!

- Model - Agnostic

- Manufacturing

- Robust to noise

(Q) Once features are extracted, we use them to

- Train machine learning model

- Cluster time series

- Detect anomalies.

Unit - V

Data Visualization

- Pixel - Oriented visualization Techniques
- Geometric projection visualization Techniques
- Icon Based visualization Techniques
- Hierarchical visualization Techniques
- visualizing complex Data & relations.

Data visualization

It is the graphical representation of information and data. By using visual elements like charts, graphs and maps, data visualization tools provide an accessible way to see and understand trends, outliers and patterns in data.

Why is it important?

- human process visuals faster than raw data
- helps in identifying trends, patterns, & outliers quickly
- enables better decision-making through insights.
- makes complex data easier to understand.
- useful in storytelling with data, especially in business and analytics.

Objectives :-

- make data clear and intuitive
- show relationships b/w variables
- support data driven decision making
- Detect anomalies & trends.

Common Tools :-

- .Tableau .power BI .excel .Python

- R • D3.js

Advantages :-

- speeds up data analysis
- makes reports engaging & easy to interpret
- useful for both technical and non-technical users
- helps in real time monitoring through dashboards
- Improves communication across teams & stakeholders.

Applications :- Business Intelligence, Healthcare,

Finance, marketing, Government.

usually by focusing on one key attribute
- this method allows large datasets to be shown
compactly on screen.

Categorization of visualization methods

- pixel oriented
 - dimension based
 - geometric projection
 - visualizing complex data and relations.

Pixel-Oriented Visualization Techniques

Pixel-oriented visualization is a method that uses color-coded pixels to represent individual data values from large and high dimensional datasets.

It lets us ideal what you want to see

many dimensions and records on a single screen.

→ It is a way to show large amounts of data using tiny colored squares (pixels)

each pixel represents one data value and the color of the pixel tells you how high or low that value is.

How it works :-

- each data value is shown as a colored pixel
- for a dataset with m dimensions, create m windows (one per dimension)
- all records in values are shown as m colored pixels, one in each window

- All windows follow the same global order

e.g. • customer data 4 dimensions: income, transaction volume, age etc

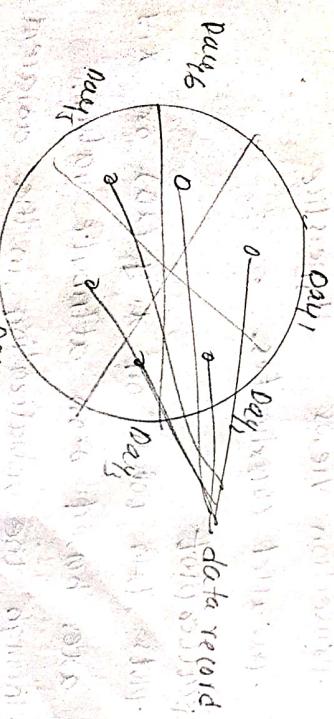
- data is sorted by income
- visualization shows how other attribute relate to income

- helps in spotting patterns & correlations.



laying out pixels in circle segments.

To save space and show the connection among multiple dimensions, space filling is often done in circle segment



Advantages

Drawbacks

- compact display
- detailed view
- multiple dimensions
- fast comparison
- color interpretation
- requires sorting
- overload risk
- limited interaction.

Geometric projection visualization techniques.

To visualize multi-dimensional data by projecting it into 2D or 3D space, making it easier to understand patterns, correlations and outliers.

key:-

- Helps in understanding data distribution, outliers and attribute relationships.

- often requires clustering before visualization to reduce cluster in large datasets.

Methods:-

- Direct visualization, scatterplot and scatterplot matrices

- Landscapes - projection pursuit technique

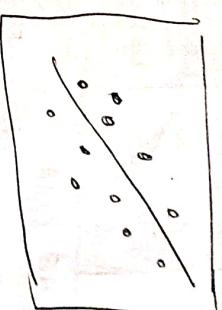
- prospection views - hyperslice

- parallel coordinates.

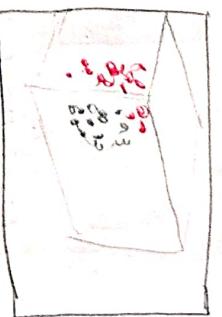
Scatterplot

- plots data points on x, y (2D) or x, y, z (3D) axes based on attribute value.

Through the visualization, in the adjacent figure we can see three types of types '+' , "x"



2D scatter plot



3D scatter plot

Scatterplot matrix (SPLOM)

A scatterplot matrix is a grid of scatter plots used to visualize relationships between multiple attributes (dimensions) of a dataset.

- For a dataset with n attributes, SPLOM shows $n \times n$ scatter plots
- Each cell shows a plot of one attribute vs another.
- Diagonal cells (e.g. A vs A, B vs B) are usually

Hyperslice

Shows 2D slices of a multi-dimensional function as a matrix

- Each slice shows how the function behaves with 2 variables while others are fixed.

- Advantage: easy interaction; users can explore dimensions directly.

use:- Great for scalar function mathematical data analysis.

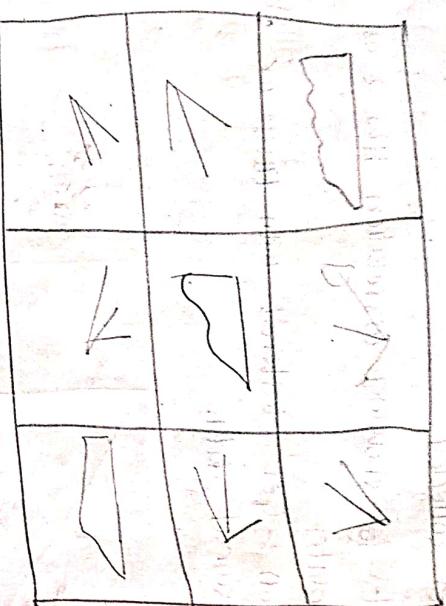
and to be collected.

use:- Best for showing correlation b/w 2 or 3 attributes

limitation :- Get clustered with huge or high dimensional data.

drawback:-

- cluster with large data set
- can be hard to interpret without presence of filtering

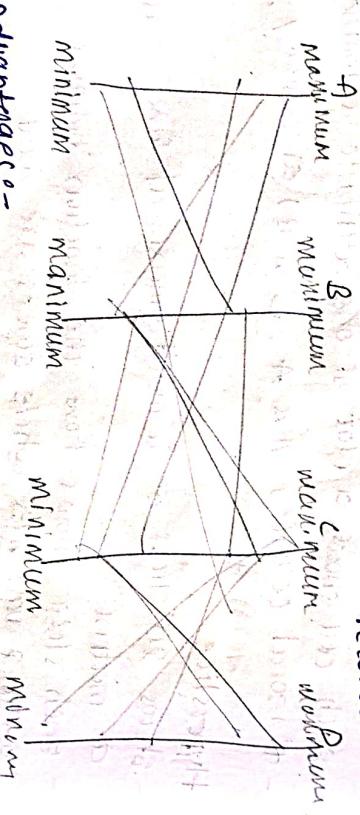


parallel coordinates

How it works: Draws parallel axes for each dimension.

- each data point is a line crossing the axes at its attribute values.
- use: useful for identifying patterns and clusters across many dimensions.

Limitations: Becomes hard to read with too many records.



Icon-based visualization techniques

These techniques use small icons to represent multidimensional data. Each icon visually encodes several variables using shapes, angles, sizes or colors.

- Instead of showing data in traditional graphs or charts, each data record is converted into a small icon
- Each part of the icon reflects a different data value
- Icons are arranged in a grid or layout for comparison.

common techniques are:-

- Reveals hidden patterns, clusters and correlations
- suitable for multivariate data
- Interactive exploration.

* Color Icons

* Chevrons

* Stack figures

* Star glyphs

General techniques :- shape coding, color icons, title bars.

- good for visualizing motion or variation b/w records.

Chenoff faces

Each data attribute controls a part of cartoon face (e.g. mouth, face, eyes & human eyes easily detect patterns or anomalies in social expressions).

- use cartoon faces to represent data

- each facial feature (like eyes, mouth, nose, face represent a variable)

- can display up to 16 dimensions

- easy for humans to notice patterns or anomalies using facial expression differences.

Ex:- each face = one record, eye size \rightarrow age, mouth curve \rightarrow income etc.



- viewing large data can be tedious
- Chenoff faces make the data easier for users to digest.

Octree figure

- represent data as a 5-part stick figure (4 limbs + body)

- variables are mapped to x on 4 axes (position)
- Other dimensions are mapped to the lengths of angles of the arms and legs.



Advantages :-

- multidimensional display
- human pattern recognition
- compact representation
- intuitive interpretation.

Limitations / Drawbacks

- limited precision
- difficult scaling
- visual clutter
- interpretation bias
- applications
- pattern discovery & customer profiling
- medical data analysis - psychology & social science
- education & research.

Hierarchical visualization Techniques

Hierarchical visualization is used to display data that is organized in a tree-like structure (parent-child relationships). These techniques help us easily to understand nested relationships like:

- organization charts

- file / folder systems
- customer segmentation
- cluttered data

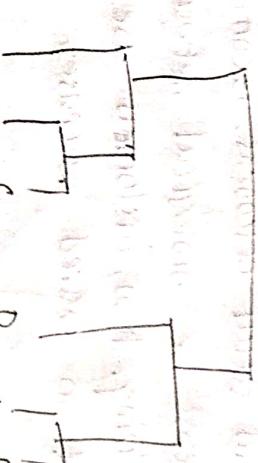
Common Techniques

- Tree Map:

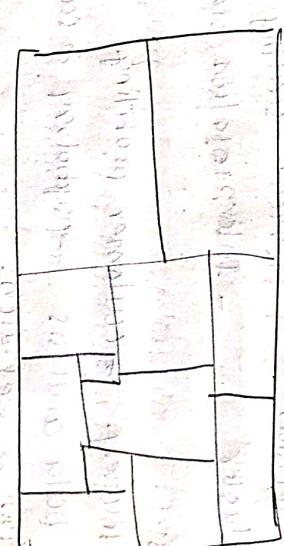
- A rectangular layout where each rectangle represents a data point or category

Working:-

- Nested boxes show hierarchy
- Size and color of the boxes represent data values like sales, sizes etc.



use case:-
Visualizing resource usage, financial data, or folder structure.



2. Dendograms:

- A tree like diagram that shows how clusters are formed
- used in **Hierarchical clustering**.

Working:- Starts with individual data points

- merges them step by step into clusters based on similarity
- continues until one large cluster remains.

Two main clustering approaches:
- Agglomerative (bottom-up) -

starts with single data point and merge clusters

- **Divisive (top-down)**: start with all data in one cluster and split into smaller ones

Applications:-

- customer segmentation
- organizational charts
- gene expression data (biostatistics)
- Document clustering.

Advantages:-
Shows nested structure & relationships clearly

- Helps with cluster analysis

Limitations:-

- may become cluttered with too many levels
- not ideal for flat or non-hierarchical data
- requires proper scaling for readability.

Visualizing Complex Data and Relationships:

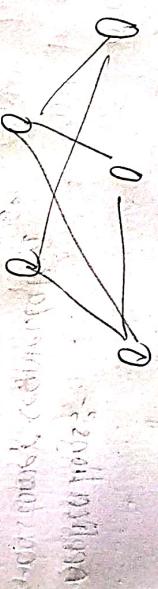
when dealing with complex data - such as large, multi-dimensional, or interconnected datasets, we use advanced visualization techniques to reveal hidden patterns, relationships and insights.

Techniques to visualize complex data

1. Network graphs

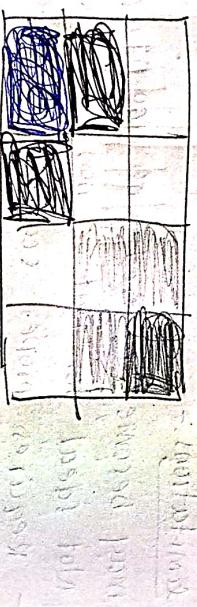
- show relationships between items (nodes & edges)
- common in social networks, recommendation systems, or biological networks.

e.g:- visualizing how customers are connected based on shared purchases.



2. Heatmaps

- display values as colors (p. 94)
- used for correlation matrices, activity logs, gene generation, easy to detect clusters, anomalies or high/low values



3. Dimensionality reduction techniques

used to reduce complex high-dimensional data into 2D or 3D for visualization!

- PCA
- t-SNE (t-distributed stochastic neighbour embedding)
- UMAP

4. 3D scatter plots & Interactive dashboard

- used to explore multiple variables with zooming, rotation & filtering tools like Power BI, Tableau and Plotly.

5. Line plot

A line plot shows data points connected by straight lines

used for showing trends over time

e.g:- monthly revenue for 1 year

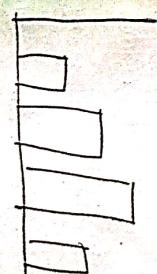


6. Bar plot / Bar Graphs

Represents categorical data with rectangular bars. The height or length of the bar shows the value.

Types:-
Vertical bars : Bar plot
Horizontal : Bar graph

used for comparing categories.



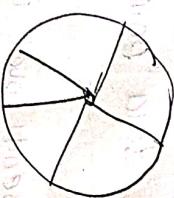
Bar plot



Bar graph

7. pie chart

A circular chart divided into slices to show proportions. Used for displaying part to whole relationship.



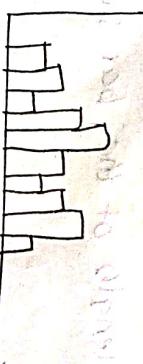
8. histogram

A special type of bar plot for continuous numerical data, grouped into ranges (bins).

Used for : showing frequency distributions

Best for : identifying data spread, skewness,

outliers & normality.



Advantages:-

- Improved Understanding
- Quick Decision making
- pattern Recognition
- Enhanced communication

- Data exploration

- Time saving.

- error detection.

ETL (extract, transform, load)

3 step process used in data warehousing and analytics to move and prepare data for analysis

1. extract
collect data from different sources

2. Transform
 - clean & convert the data into a standard usable format (removing duplicates, formats)

3. Load
 - store the transformed data into a target database or datawarehouse for reporting and analysis

Popular ETL tools

- Informatica Powercenter
- IBM InfoSphere Data Stage
- Apache spark
- pentaho (kettle).

Word Cloud :-

A word cloud is a visual display of words from document or dataset.

- Bigger words appear more frequently in text
- Helps highlight keywords, themes, trends in textual data
- common in text mining & social media analysis.