

UNIT-1

Introduction to Data Science

Definition

Data Science is a **multidisciplinary field** that combines **statistics**, **computer science**, **domain knowledge**, and **machine learning** to extract useful insights, patterns, and knowledge from data (structured, semi-structured, and unstructured).

Data Science = Collecting + Processing + Analyzing + Interpreting data → Insights + Decisions.

Why Data Science?

- Today, organizations generate massive amounts of data (from apps, websites, IoT devices, social media, transactions, etc.).
- Traditional methods (like simple statistics or Excel) cannot handle this **volume, velocity, and variety** of data.
- Data Science helps to:
 - Make better decisions (data-driven decision-making).
 - Predict future trends.
 - Automate business processes.
 - Personalize customer experiences.

Components of Data Science

1. **Data Collection & Storage** – from databases, APIs, logs, sensors, etc.
2. **Data Cleaning & Preparation** – removing duplicates, handling missing values, standardization.
3. **Exploratory Data Analysis (EDA)** – understanding patterns, trends, and relationships.
4. **Data Modeling (ML/AI)** – applying algorithms for prediction, classification, clustering, etc.
5. **Interpretation & Visualization** – using dashboards, graphs, reports for decision-making.
6. **Deployment** – integrating the model into real-time systems (e.g., recommendation engines).

Benefits of Data Science

1. Better Decision-Making

- Converts raw data into actionable insights.
- Helps managers make data-driven decisions instead of relying only on intuition.

2. Predicting Future Trends

- Machine Learning models forecast demand, sales, risks, and opportunities.
- Businesses can plan resources and strategies in advance.

3. Improved Efficiency

- Automates repetitive tasks.
- Optimizes processes (e.g., supply chain, inventory management).

4. Personalization

- Provides tailored experiences for customers.
- Example: Netflix recommending movies, Amazon suggesting products.

5. Fraud & Risk Detection

- Detects unusual patterns in finance, banking, and cybersecurity.
- Reduces losses by identifying fraud before it happens.

6. Innovation

- Helps in developing new products, technologies, and services.
- Example: Self-driving cars, smart assistants (Alexa, Siri), drug discovery.

7. Competitive Advantage

- Businesses using Data Science gain an edge over those relying on guesswork.
- Example: Companies like Google, Amazon, and Facebook thrive by leveraging data.

Uses of Data Science

In Different Industries:

1. Healthcare

- Disease prediction and prevention.
- Analyzing medical images (X-rays, MRI).
- Personalized treatment recommendations.

2. Finance

- Fraud detection in transactions.
- Credit scoring for loans.
- Stock market prediction.

3. **Retail & E-commerce**

- Product recommendation systems.
- Customer segmentation.
- Demand forecasting.

4. **Transportation**

- Self-driving cars (computer vision).
- Route optimization.
- Predictive maintenance of vehicles.

5. **Marketing & Sales**

- Sentiment analysis of social media.
- Targeted advertisements.
- Customer lifetime value prediction.

6. **Manufacturing**

- Quality control using sensors.
- Predictive maintenance of machines.
- Optimizing supply chain.

7. **Government & Public Sector**

- Crime pattern analysis.
- Disaster management & weather forecasting.
- Smart city planning.

Facets of Data

Very large amount of data will generate in big data and data science. These data is various types and main categories of data are as follows:

- a) Structured
- b) Natural language
- c) Graph-based
- d) Streaming
- e) Unstructured
- f) Machine-generated
- g) Audio, video and images

a) Structured Data

- Structured data is arranged in rows and column format. It helps for application to retrieve and process data easily. Database management system is used for storing structured data.
- The term structured data refers to data that is identifiable because it is organized in a structure. The most common form of structured data or records is a database where specific information is stored based on a methodology of columns and rows.
- An Excel table is an example of structured data.

b) Unstructured Data

- Unstructured data is data that does not follow a specified format. Row and columns are not used for unstructured data. Therefore it is difficult to retrieve required information. Unstructured data has no identifiable structure.
- The unstructured data can be in the form of Text: (Documents, email messages, customer feedbacks), audio, video, images. Email is an example of unstructured data.

c) Natural Language

- Natural language is a special type of unstructured data.
- Natural language processing enables machines to recognize characters, words and sentences, then apply meaning and understanding to that information. This helps machines to understand language as humans do.
- For natural language processing to help machines understand human language, it must go through speech recognition, natural language understanding and machine translation. It is an iterative process comprised of several layers of text analysis.

d) Machine - Generated Data

- Machine-generated data is an information that is created without human interaction as a result of a computer process or application activity. This means that data entered manually by an end-user is not recognized to be machine-generated.
- Machine data contains a definitive record of all activity and behavior of our customers, users, transactions, applications, servers, networks, factory machinery and so on.

e) Graph-based or Network Data

- Graphs are data structures to describe relationships and interactions between entities in complex systems. In general, a graph contains a collection of entities called nodes and another collection of interactions between a pair of nodes called edges.
- Nodes represent entities, which can be of any object type that is relevant to our problem domain. By connecting nodes with edges, we will end up with a graph (network) of nodes.

f) Streaming Data

Streaming data is data that is generated continuously by thousands of data sources, which typically send in the data records simultaneously and in small sizes (order of Kilobytes).

- Streaming data includes a wide variety of data such as log files generated by customers using your mobile or web applications, ecommerce purchases, in-game player activity, information from social networks, financial trading floors or geospatial services and telemetry from connected devices or instrumentation in data centers.

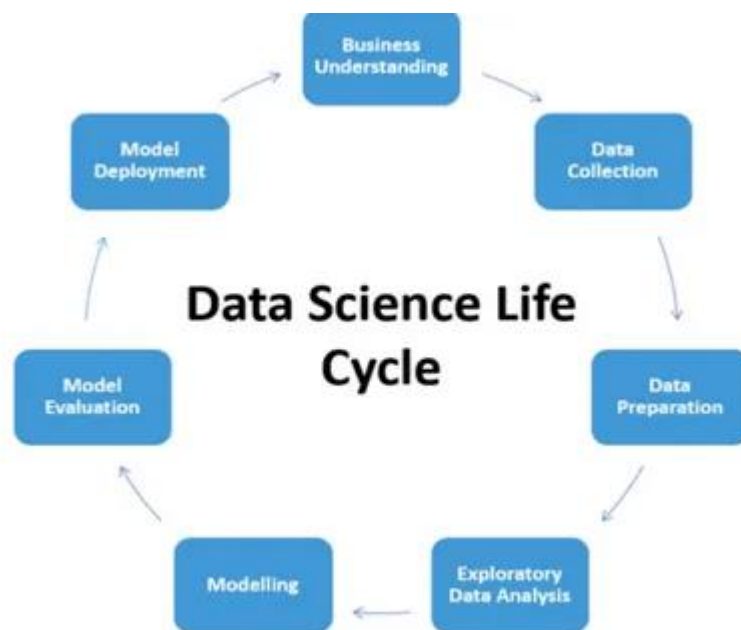
Data Science Process in Brief

Definition

Data Science is a **multidisciplinary field** that combines **statistics**, **computer science**, **domain knowledge**, and **machine learning** to extract useful insights, patterns, and knowledge from data (structured, semi-structured, and unstructured).

Data Science = **Collecting** + **Processing** + **Analyzing** + **Interpreting data** → **Insights** + **Decisions**.

Life Cycle:



1. Problem definition

- This is the **starting point**. We must understand **what needs to be solved** and **why it matters**.
- Without a clear problem, analysis has no direction.
- Example: A bank wants to reduce loan defaults → The problem is “*predict which customers are likely to default on loans.*”

2. Data Collection

- Once the problem is defined, the next step is to **collect the right data**.
- Data can come from multiple sources: databases, Excel/CSV files, APIs, sensors, logs, or even web scraping.
- Example: Collect customer's income, credit history, past loans, repayment behavior.

3. Data Cleaning

- Real-world data is **messy**. It often has missing values, duplicates, errors, or inconsistent formats.
- Cleaning means **fixing and preparing** the data so the model can understand it.
- Tasks include:
 - Removing duplicate rows
 - Handling missing values (fill, drop, or estimate)
 - Converting data into correct formats (dates, categories, numbers)
- Example: Replace missing salary values with the average salary.

4. Data Explore

- Also called **Exploratory Data Analysis (EDA)**.
- Here we **analyze data patterns, trends, and relationships** using summary statistics and visualizations.
- Helps us understand the dataset before modeling.
- Example:
 - Plot income vs. loan defaults to see if low income leads to higher defaults.
 - Check distributions, correlations, and outliers.

5. Model Building

- Now we apply **Machine Learning (ML) or statistical models** to solve the problem.
- The choice of model depends on the type of problem:
 - Prediction → Regression/Classification
 - Grouping → Clustering
 - Recommendation → Collaborative filtering
- Example: Use logistic regression or random forest to predict loan default.

6. Data Evaluate

- A model must be **tested** to check how well it performs.
- We split the dataset into **training and testing sets**.
- Common evaluation metrics:
 - Classification → Accuracy, Precision, Recall, F1-score
 - Regression → RMSE, MAE, R^2
- Example: If the loan default model has 90% accuracy but low recall, it may miss many defaulters → not good.

7. Data Deployment

- Once the model works well, it's put into **real-world use**.
- Deployment means integrating the model into applications, websites, or business systems.
- Example: Bank's loan approval system automatically uses the model to predict if a new applicant might default.

8. Data Monitor

- The process doesn't end after deployment.
- Models **degrade over time** because data changes (this is called **data drift**).
- Monitoring ensures the model stays accurate. If performance drops, we retrain with new data.
- Example: If the bank's loan model starts giving wrong predictions due to new customer behavior, it must be retrained.

Big Data Ecosystem and Data science

The **Big Data Ecosystem** refers to the wide variety of tools, technologies, frameworks, and processes that work together to handle, store, process, and analyze massive volumes of data that traditional systems cannot manage efficiently. It is designed to tackle the challenges posed by big data's characteristics: volume, velocity, variety, veracity, and value.

Volume: Massive data sizes

Velocity: Speed of data generation

Variety: Different formats (text, video, logs)

Veracity: Data quality and trustworthiness

Value: Business insights derived from data

Key Components of Big Data Ecosystem:

1. **Data Sources:**
 - Social media feeds (Twitter, Facebook)
 - Sensor data (IoT devices)
 - Transactional databases
 - Log files
 - Streaming data (e.g., clickstreams)
2. **Data Ingestion:**

Tools that collect and import data into the ecosystem.

- **Apache Kafka:** Distributed streaming platform.
 - **Apache Flume:** Collects log data from many servers.
 - **Sqoop:** Transfers data between Hadoop and relational databases.
3. **Data Storage:**
Systems to store vast amounts of data reliably and scalably.
- **HDFS (Hadoop Distributed File System):** Stores large files across multiple machines.
 - **NoSQL Databases:** Such as HBase, Cassandra, MongoDB — store unstructured or semi-structured data.
4. **Data Processing:**
Tools to process data in batch or real-time.
- **Apache Hadoop MapReduce:** Batch processing framework.
 - **Apache Spark:** Fast in-memory processing for batch and real-time.
 - **Apache Flink:** Stream processing framework.
5. **Data Warehousing and Query Engines:**
Tools to query large datasets using SQL-like languages.
- **Hive:** Data warehouse infrastructure built on Hadoop.
 - **Impala, Presto:** Fast query engines.
6. **Data Analytics and Visualization:**
Platforms to analyze and visualize data results.
- **Tableau, Power BI:** Visualization tools.
 - **Elasticsearch & Kibana:** Search and visualize data in real time.
7. **Machine Learning and AI:**
Tools to build predictive and intelligent models.
- **Apache Spark MLlib:** Machine learning library.
 - **TensorFlow, PyTorch:** Deep learning frameworks.
8. **Workflow Orchestration and Management:**
Manage data pipelines and automate workflows.
- **Apache Airflow:** Workflow scheduling and management.
 - **Apache NiFi:** Data flow automation.

Data Science

Data Science is the discipline that uses **scientific methods, algorithms,** and systems to extract knowledge and insights from **structured and unstructured data.** It leverages the

tools and frameworks from the big data ecosystem to process data and solve real-world problems.

Relationship with Big Data Ecosystem:

- The Big Data Ecosystem provides the infrastructure and tools necessary for **data scientists** to collect, store, and process large datasets.
- Data science focuses on **analyzing** this data through statistical methods, machine learning, and visualization to **derive insights**.
- Together, they enable businesses and organizations to make **data-driven decisions** by turning massive data into actionable knowledge.

Data Science Workflow (in Big Data Context):

1. **Data Acquisition:** Use ingestion tools to gather big data.
2. **Data Storage:** Store data on distributed file systems or NoSQL databases.
3. **Data Processing:** Use frameworks like Spark for cleaning and transforming data.
4. **Exploratory Analysis:** Analyze and visualize data to find patterns.
5. **Model Building:** Apply machine learning models to predict outcomes or classify data.
6. **Model Deployment:** Use models for real-time predictions or batch analysis.
7. **Monitoring & Updating:** Keep models effective by updating them with new data.

How They Interact

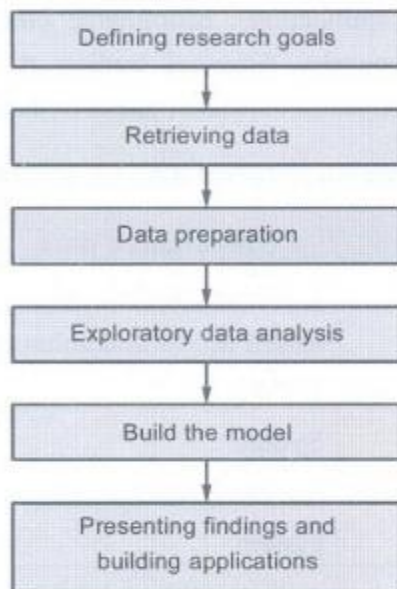
- **Big Data Ecosystem** provides the infrastructure and raw data.
- **Data Science** transforms that data into actionable insights.
- Together, they power innovations in AI, automation, and decision-making.

Data Science Process: OVERVIEW

Data science process consists of six stages :

1. Discovery or Setting the research goal
2. Retrieving data
3. Data preparation
4. Data exploration
5. Data modeling
6. Presentation and automation

- shows data science design process.



Step 1: Discovery or Defining research goal

This step involves acquiring data from all the identified internal and external sources, which helps to answer the business question.

• Step 2: Retrieving data

- It collection of data which required for project.
- This is the process of gaining a business understanding of the data user have and deciphering what each piece of data means.
- This could entail determining exactly what data is required and the best methods for obtaining it.
- This also entails determining what each of the data points means in terms of the company.
- If we have given a data set from a client, for example, we shall need to know what each column and row represents.

- **Step 3: Data preparation**

- Data can have many inconsistencies like missing values, blank columns, an incorrect data format, which needs to be cleaned.
- We need to process, explore and condition data before modeling. The cleandata, gives the better predictions.

- **Step 4: Data exploration**

- Data exploration is related to deeper understanding of data.
- Try to understand how variables interact with each other, the distribution of the data and whether there are outliers.
- To achieve this use descriptive statistics, visual techniques and simple modeling. This steps is also called as Exploratory Data Analysis.

- **Step 5: Data modeling**

- In this step, the actual model building process starts. Here, Data scientist distributes datasets for training and testing.
- Techniques like association, classification and clustering are applied to the training data set. The model, once prepared, is tested against the "testing" dataset.

- **Step 6: Presentation and automation**

- Deliver the final baselined model with reports, code and technical documents in this stage.
- Model is deployed into a real-time production environment after thorough testing. In this stage, the key findings are communicated to all stakeholders.
- This helps to decide if the project results are a success or a failure based on the inputs from the model.

Defining Goals and Creating a Project Charter in Data Science

This is the **first and most important step** in the Data Science process. It sets the **direction, scope, and success criteria** of the project.

1. What is a Goal in Data Science?

A **goal** is a **clear and measurable business objective** you want to achieve using data.

Example: "Reduce customer churn by 20% within 6 months using predictive modeling."

2. What is a Project Charter?

A **Project Charter** is a short document that outlines the **why, what, and how** of the data science project. It acts like a **blueprint** or **agreement** between stakeholders and the data team.

Components of a Good Project Charter

Component	Description	Example
Project Title	Name of the project	"Customer Churn Prediction"
Business Goal	What you want to achieve	Reduce customer churn
Problem Statement	The issue to solve	Customers are leaving after 3 months of joining
Success Criteria	How success will be measured	Churn reduced by 20%, Model accuracy \geq 85%
Data Sources	Where the data will come from	CRM database, Customer feedback
Team Members	Who is involved	Data scientist, domain expert, business analyst
Timeline	Start and end dates	Aug 1 – Sep 30
Deliverables	What the team will produce	Predictive model, dashboard, report
Constraints & Risks	Limitations	Limited labeled data, data privacy issues

Retrieving data

Retrieving data means **collecting or getting data from different sources** like databases, files, APIs, websites, sensors, or cloud platforms so it can be used for analysis or building models.

Sources of Data

1. Databases

- Databases store data in an organized way.
- Two main types:
 - **SQL (Relational)**: Data stored in tables (e.g., MySQL, PostgreSQL).
 - **NoSQL**: Data stored in flexible formats (e.g., MongoDB).
- We use **queries (SQL commands)** or **Python libraries** to get the data.
- **Example**: Getting customer details from a PostgreSQL database.

2. APIs

- APIs let us access data from other applications or services through the internet.
- We make requests like GET or POST to get data.
- **Example**: Using a weather API to get today's temperature for different cities.

3. Files

- Data often comes in files such as **CSV, Excel, JSON**.
- We read these files from a local computer or a shared drive.
- **Example**: Reading a CSV file of sales records using Excel or Python.

4. Websites

- Sometimes data is available on web pages, but not as a downloadable file.
- In such cases, we do **web scraping** to collect that information.
- **Example**: Collecting product prices from an e-commerce website.

5. Sensors

- Devices like IoT sensors or machines generate real-time data (temperature, speed, etc.).
- This data is usually sent to a system through **networks** or APIs.
- **Example**: A weather station sending temperature and humidity data every minute.

6. Cloud Platforms

- Companies store large amounts of data in **cloud storage** like AWS, Google Cloud, or Azure.
- We can download or connect to these platforms using **cloud tools or Python SDKs**.
- **Example:** Downloading a big dataset from AWS S3 for analysis.

Cleansing

Cleansing means cleaning the data so it is accurate, complete, and consistent before analysis.

When we collect raw data, it often has problems like missing values, duplicates, wrong formats, or spelling mistakes.

If we don't clean it, our results and models will be wrong.

Common Steps in Cleansing

1. **Remove Duplicates**
Delete repeated records that appear more than once.
Example: Same customer entry appearing twice in the database.
2. **Handle Missing Values**
Fill empty values with average/median or remove those rows.
Example: If age is missing, replace it with the average age.
3. **Correct Errors**
Fix spelling mistakes or wrong data entries.
Example: "Indai" → "India".
4. **Standardize Formats**
Make all data follow the same format.
Example: Dates changed from 12/31/23 to 2023-12-31.
5. **Fix Data Types**
Convert data into correct types like number, text, or date.
Example: Change "twenty" → 20.

Example:

- **Before Cleaning:**
Name: Jon, Age: twenty, Date: 12/31/23
- **After Cleaning:**
Name: John, Age: 20, Date: 2023-12-31

Integrating

Integrating means **combining data from different sources into a single dataset** so that it can be analyzed together.

When data comes from multiple sources like databases, files, or APIs, we need to merge them in a structured way.

Steps in Integration

- 1. Match Common Fields**
Combine data using common identifiers (like Customer ID).
Example: Customer details from a database + purchase history from a CSV file.
- 2. Merge Tables or Files**
Join multiple datasets into one unified table.
Example: Sales data from Excel + inventory data from another file.
- 3. Remove Duplicates**
Avoid repeated records after merging.

Transforming

Transforming means **changing data into the right format or structure for analysis or modeling**.

Raw data is often not in a usable format, so we transform it into a clean, structured, and consistent format.

Steps in Transformation

- 1. Normalize Values**
Make values consistent (e.g., same currency, same units).
Example: Converting all prices to USD.
- 2. Change Formats**
Standardize formats for fields like date or time.
Example: Change 12/31/23 to 2023-12-31.
- 3. Encode Categories**
Convert text values into numbers for models.
Example: Yes → 1, No → 0.
- 4. Scale Data**
Adjust numbers to a standard range (like 0 to 1) for machine learning.

Exploratory Data Analysis (EDA)

Exploratory Analysis means **examining and understanding data before building models**. It helps find patterns, trends, relationships, and problems in the data.

EDA is important because it gives insights that guide feature selection and model building.

What We Do in EDA

- 1. Summarize Data**
Check basic statistics like mean, median, min, max.
Example: Average age of customers is 32.
- 2. Check Data Distribution**
See how values are spread using histograms or boxplots.
Example: Most customers spend between \$50 and \$150.
- 3. Find Missing Values & Outliers**
Identify where data is incomplete or unusual.
Example: Income column has 10% missing values.
- 4. Understand Relationships**
Check correlation between variables.
Example: Higher income → higher spending.
- 5. Visualize Data**
Use charts to understand trends and patterns.
Example: Bar chart for sales by region, scatter plot for age vs. spending.

UNIT-2

Applications of Machine Learning in Data Science

1. Predictive Analytics

- **Meaning:** Uses past data patterns to predict future outcomes.
- **How it helps:** Data scientists use ML models like regression, decision trees, or neural networks to make forecasts.
- **Examples:**
 - Retail: Predicting product demand for inventory planning.
 - Healthcare: Forecasting the likelihood of a patient being readmitted to the hospital.
 - Finance: Predicting loan default risks.

2. Classification

- **Meaning:** Classifies data into predefined categories or labels.
- **How it helps:** Automates decision-making where outcomes are discrete.
- **Examples:**
 - Emails: Spam vs. non-spam.
 - Banking: Fraudulent vs. legitimate transactions.
 - Healthcare: Benign vs. malignant tumors.

3. Regression

- **Meaning:** Predicts a continuous numerical value from input features.
- **How it helps:** Useful when outcome is a number, not a category.
- **Examples:**
 - Predicting housing prices based on location, size, and features.
 - Forecasting revenue for the next quarter.
 - Estimating fuel consumption of a vehicle.

4. Clustering (Unsupervised Learning)

- **Meaning:** Groups similar data points when labels are unknown.
- **How it helps:** Data scientists discover hidden structures/patterns in data.
- **Examples:**
 - Marketing: Customer segmentation (grouping customers by behavior).

- Retail: Grouping similar products bought together (market basket analysis).
- Cybersecurity: Identifying unusual patterns of network activity.

5. Natural Language Processing (NLP)

- **Meaning:** Helps machines read, understand, and respond to human language.
- **How it helps:** Converts unstructured text into usable insights.
- **Examples:**
 - Business: Analyzing customer reviews to see satisfaction levels.
 - Support: Chatbots for automated customer service.
 - Media: Automatic captioning or translation.

6. Computer Vision

- **Meaning:** Enables computers to “see” and analyze images or videos.
- **How it helps:** Automates tasks that require visual understanding.
- **Examples:**
 - Healthcare: Detecting diseases from X-rays or MRIs.
 - Transportation: Self-driving cars recognizing pedestrians and traffic signals.
 - Security: Face recognition for identity verification.

7. Recommendation Systems

- **Meaning:** Suggests items based on user preferences and past behavior.
- **How it helps:** Increases personalization and customer engagement.
- **Examples:**
 - Netflix recommending movies based on what you’ve watched.
 - Amazon suggesting “Frequently Bought Together” products.
 - Spotify creating daily personalized playlists.

8. Anomaly & Fraud Detection

- **Meaning:** Identifies unusual or suspicious activity.
- **How it helps:** Protects businesses from financial or operational risks.
- **Examples:**
 - Banking: Detecting unusual spending behavior on credit cards.
 - Manufacturing: Spotting early signs of machine failure.
 - IT: Detecting abnormal traffic patterns in a network (possible cyber-attack).

9. Time Series Forecasting

- **Meaning:** Uses historical time-based data to predict future trends.
- **How it helps:** Crucial for planning and resource management.
- **Examples:**
 - Energy: Predicting daily electricity demand.
 - Supply Chain: Forecasting demand for seasonal products.
 - Environment: Predicting air pollution levels or rainfall.

10. Automation & Decision-Making

- **Meaning:** ML models enable systems to take actions automatically.
- **How it helps:** Reduces human effort, speeds up processes, and increases accuracy.
- **Examples:**
 - Airlines: Dynamic ticket pricing based on demand.
 - Retail: Automated chatbots handling customer queries.
 - Transportation: Autonomous vehicles making real-time driving decisions.

Role of Machine Learning in Data Science

1. Core Component of Data Science

- Data Science is about extracting insights and value from data.
- ML is the **engine** that allows data scientists to **learn patterns** from data automatically instead of relying only on manual rules.

2. From Data to Predictions

- ML helps **analyze large volumes of data** and **predict future outcomes**.
- Example: A data scientist studying sales trends can apply ML models to **forecast future sales**.

3. Automation of Insights

- Without ML: data scientists manually write rules to detect fraud, anomalies, etc.
- With ML: models **learn rules automatically** from data, reducing manual effort.

4. Handling Complexity

- Many datasets are **too large, unstructured, or complex** for traditional analytics.
- ML handles:

- **Text (NLP)**
- **Images (Computer Vision)**
- **Streaming Data (Real-time ML)**

5. Decision-Making Support

- ML provides **actionable insights** that support better business decisions.
- Example: In healthcare, ML models help doctors detect diseases early from X-rays or lab results.

6. Personalization

- Data Science projects often aim to **improve customer experience**.
 - ML enables personalization (recommendations, targeted ads, dynamic pricing).
-

7. Improving Accuracy

- Traditional statistics explains data.
- ML goes further to **predict with high accuracy**, improving outcomes in:
 - Finance (credit scoring, fraud detection)
 - Retail (inventory forecasting)
 - Manufacturing (predictive maintenance)

Python tools like sklearn

In Data Science, we use different Python libraries (tools) for different tasks. The main ones are:

1. **Scikit-learn (sklearn)**
 - Used for **machine learning models**.
 - Helps in classification, regression, clustering, and model evaluation.
 - Example: Predicting house prices or checking if an email is spam.
2. **TensorFlow & PyTorch**
 - Used for **deep learning** (neural networks).
 - These are powerful when working with images, speech, or natural language.
 - Example: Face recognition, self-driving car vision.
3. **Pandas & NumPy**
 - Used for **data handling and preprocessing**.

- NumPy works with numbers and arrays (matrices).
- Pandas helps in working with tables (rows & columns like Excel).
- Example: Cleaning a CSV file, removing missing values, filtering data.

4. Matplotlib & Seaborn

- Used for **data visualization** (drawing graphs and charts).
- Matplotlib is the base library, Seaborn makes graphs more attractive and easy.
- Example: Plotting sales trends or showing customer age distribution.

Modelling process for feature engineering

What is Feature Engineering?

Feature engineering means **changing raw data into useful information** that helps your machine learning model make better predictions.

Simple Steps in Feature Engineering

1. Understand Your Data

- Look at your data carefully.
- Ask: What does each column mean? Are there missing values or strange numbers?

2. Clean the Data

- Remove or fix missing values.
- Fix wrong or messy data (like typos or outliers).

3. Choose Useful Features

- Pick columns that help solve your problem.
- Remove columns that don't add value or confuse the model.

4. Create New Features

- Combine columns (e.g., "height" and "weight" → "BMI").
- Make new columns from dates (e.g., "date" → "day of week").

5. Convert Text to Numbers

- ML models need numbers, not words.
- Use techniques like:
 - **Label Encoding**: Turn categories into numbers (e.g., "Red" → 1, "Blue" → 2)

- **One-Hot Encoding:** Make separate columns for each category

6. Scale the Data

- Make sure all numbers are on a similar scale.
- Example: Convert prices from ₹1 to ₹1,000,000 into a smaller range so the model doesn't get confused.

7. Split the Data

- Divide your data into:
 - **Training set:** To teach the model
 - **Test set:** To check how well it learned

Machine Learning

- ✓ A Machine learning is the branch of [Artificial Intelligence](#) that focuses on developing models and algorithms that let computers learn from data and improve from previous experience without being explicitly programmed for every task.
- ✓ In simple words, ML teaches the systems to think and understand like humans by learning from the data.
- ✓ In this article, we will explore the various **types of [machine learning algorithms](#)** that are important for future requirements. **Machine learning** is generally a training system to learn from past experiences and improve performance over time. [Machine learning](#) helps to predict massive amounts of data. It helps to deliver fast and accurate results to get profitable opportunities.

Types of Machine Learning

There are several types of machine learning, each with special characteristics and applications. Some of the main types of machine learning algorithms are as follows:

1. Supervised Learning
2. Unsupervised Learning
3. Reinforcement Learning
4. Semi-supervised Learning

Additionally, there is a more specific category called semi-supervised learning, which combines elements of both supervised and unsupervised learning.

1. Supervised Learning

[Supervised learning](#) is defined as when a model gets trained on a "**Labelled Dataset**". Labelled datasets have both input and output parameters.

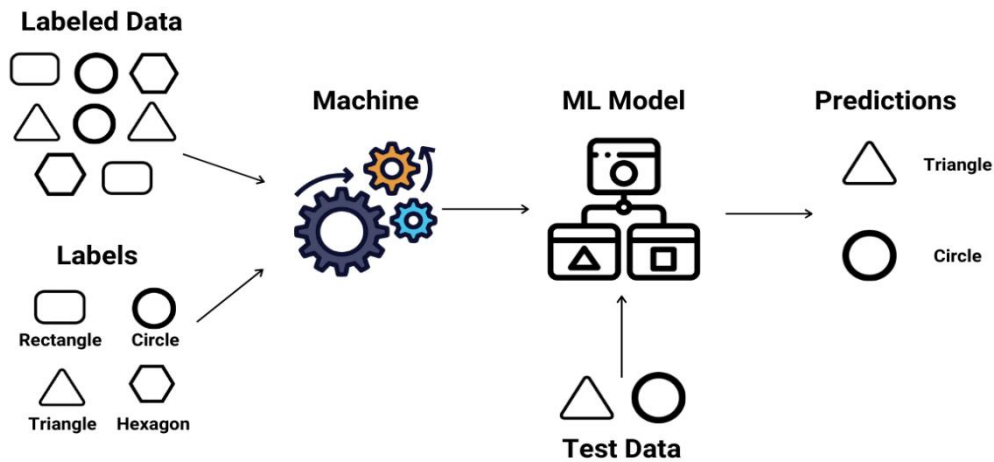
In **Supervised Learning** algorithms learn to map points between inputs and correct outputs. It has both training and validation datasets labelled.

✓ [Classification](#)

✓ [Regression](#)



Supervised Learning



Classification

- [Classification](#) deals with predicting **categorical** target variables, which represent discrete classes or labels.
- For instance, classifying emails as spam or not spam, or predicting whether a patient has a high risk of heart disease.
- Classification algorithms learn to map the input features to one of the predefined classes.

Here are some classification algorithms:

- [Logistic Regression](#)
- [Support Vector Machine](#)

Regression

- [Regression](#), on the other hand, deals with predicting **continuous** target variables, which represent numerical values.
- For example, predicting the price of a house based on its size, location, and amenities, or forecasting the sales of a product.
- Regression algorithms learn to map the input features to a continuous numerical value.

Here are some regression algorithms:

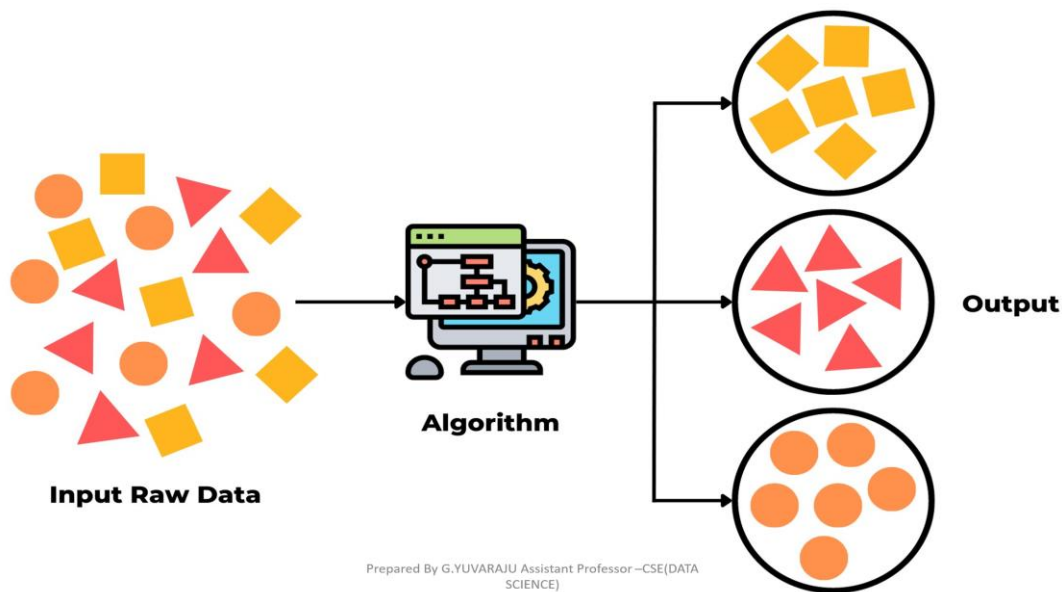
- [Linear Regression](#)
- [Polynomial Regression](#)

2. Unsupervised Learning

Unsupervised Learning Unsupervised learning is a type of machine learning technique in which an algorithm discovers patterns and relationships using unlabeled data.

Unlike supervised learning, unsupervised learning doesn't involve providing the algorithm with labeled target outputs.

The primary goal of Unsupervised learning is often to discover hidden patterns, similarities, or clusters within the data, which can then be used for various purposes, such as data exploration, visualization, dimensionality reduction, and more.



There are two main categories of unsupervised learning that are mentioned below:

- [Clustering](#)
- [Association](#)

Clustering

- [Clustering](#) is the process of grouping data points into clusters based on their similarity.
- This technique is useful for identifying patterns and relationships in data without the need for labeled examples.

Here are some clustering algorithms:

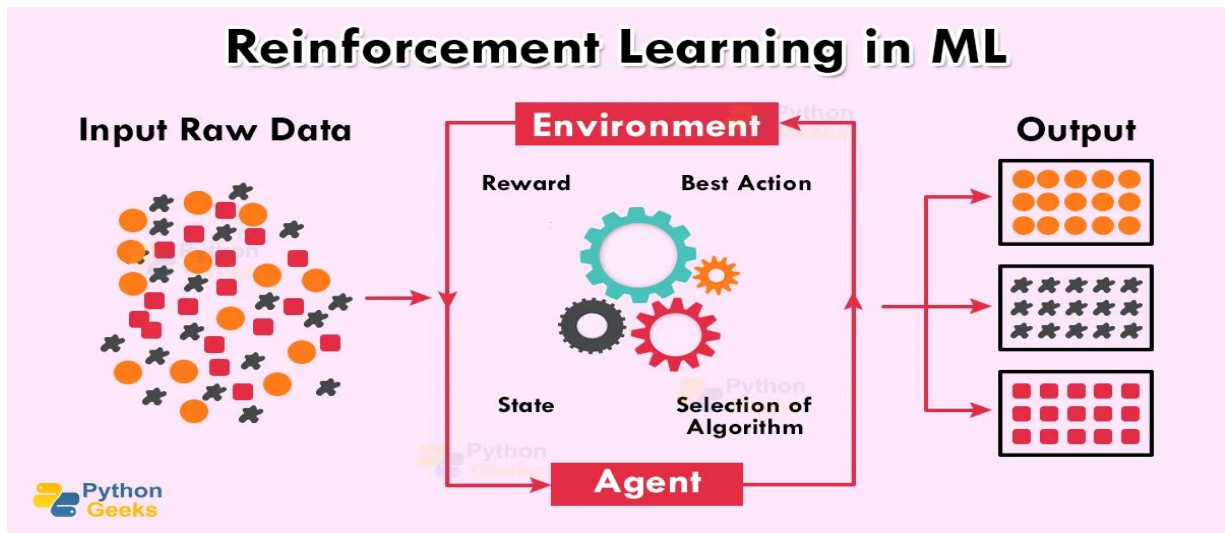
[K-Means Clustering algorithm](#)

[Mean-shift algorithm](#)

3. Reinforcement Machine Learning

- ✓ Reinforcement learning is like learning by playing a game .You try different actions and you get rewards for good actions and penalties for bad ones. you learn the best actions by practicing

- ✓ In this technique, the model keeps on increasing its performance using Reward Feedback to learn the behavior or pattern.
- ✓ These algorithms are specific to a particular problem e.g. Google Self Driving car, AlphaGo where a bot competes with humans and even itself to get better and better performers in Go Game.
- ✓ Each time we feed in data, they learn and add the data to their knowledge which is training data. So, the more it learns the better it gets trained and hence experienced.



4. Reinforcement Learning

- **What it means:**
The model learns by **trial and error**, getting rewards for good actions and penalties for bad actions.
- **Goal:** Learn the best sequence of actions to maximize reward.
- **Examples:**
 - Self-driving cars learning to navigate roads.
 - Robots learning to walk.
 - Game AI learning to play chess or Go.

Problems and general techniques for handling large data

Problems When Handling Large Data

1. Memory Issues

- Big data may not fit in your computer's memory (RAM).
- If you try to load everything at once, your computer can slow down or crash.
- *Example:* Loading 50 GB of sales data on a laptop with 8 GB RAM.

2. Slow Processing

- Doing calculations or analysis on millions of records takes a long time.
- *Example:* Finding the total sales per month from 10 million rows.

3. Data Quality Problems

- Large data can have missing values, duplicates, or wrong information.
- Cleaning this data is harder because there is so much of it.
- *Example:* Some customers may have missing email addresses or wrong phone numbers.

4. I/O Bottlenecks (Read/Write Problems)

- Reading and saving big files takes a long time.
- *Example:* Saving a 100 GB CSV file to your computer repeatedly can be very slow.

5. Scalability Problems

- Methods or programs that work on small data may not work well on big data.
- *Example:* A loop in Python might work for 1,000 rows but fail for 10 million rows.

General Techniques for Handling Large Data

1. Sampling

- Use a small part of the data to test or analyze first.
- *Example:* Take 10% of your data to check the results before using the full dataset.

2. Batch Processing / Chunking

- Work on the data in small pieces instead of all at once.
- *Example:* Read 100,0

- x00 rows at a time using Python or pandas.

3. Distributed Computing

- Split data across many computers to process it faster.
- *Tools:* Hadoop, Spark, Dask.

4. Data Compression

- Save data in smaller files using formats like **Parquet** or **Avro**.
- This saves space and makes reading/writing faster.

5. Indexing & Partitioning

- Organize data so you can find and use it faster.
- *Example:* Divide a table by year, region, or category.

6. Streaming Processing

- For data that comes continuously (like social media or sensors), process it immediately instead of saving it all first.
- *Tools:* Kafka, Spark Streaming.

Programming tips for dealing large data

1. Use Efficient Libraries

- Instead of normal Python lists or loops, use **pandas**, **NumPy**, **Dask**, or **PySpark**.
- These libraries are optimized for big data and are much faster.
- *Example:* Calculating totals or averages for millions of records is quicker using pandas.

2. Load Data in Chunks

- Don't try to load the entire large file at once. Read it in small pieces.
- *Example in Python:*

```
import pandas as pd

for chunk in pd.read_csv('big_file.csv', chunksize=100000):
    process(chunk)
```

- This saves memory and prevents crashes.

3. Avoid Loops; Use Vectorized Operations

- Loops are slow for large datasets. Use vectorized operations instead.
- *Example:*

```
df['total'] = df['price'] + df['tax'] # Faster than looping row by row
```

4. Use Generators

- Generators load data **one item at a time**, saving memory.

```
def read_file(file):  
    for line in open(file):  
        yield line
```

- Useful for reading huge files line by line.

5. Use Sparse Data Structures

- If your dataset has a lot of zeros, use **sparse matrices** to save memory.

6. Save Intermediate Results

- Don't repeat heavy calculations on large data. Save partial results to disk.

Example: Save processed chunks as CSV or Parquet files.

7. Profile and Optimize Your Code

- Check which part of your code is slow using **time** or **cProfile**.
- Optimize only the slowest parts to save time.

Case studies on DS projects for predicting malicious URLs

Predicting Malicious URLs

- **Goal:** To classify a website link (URL) as safe or malicious.
- **Why:** Hackers use fake links (phishing URLs) to steal data, so detecting them is important for security.
- **Data:**
 - We look at details of the URL such as:
 - URL length (long links are often suspicious).
 - Number of dots (.) in the link.
 - Special characters like @, ?, =.
 - Whether the link uses HTTPS or not.
- **Model:**
 - Machine Learning algorithms like:
 - **Logistic Regression** → Simple and interpretable.
 - **Random Forest** → Good for complex patterns.
- **Process:**
 1. Collect a dataset of URLs (safe and malicious).

2. Extract features from each URL (length, dots, HTTPS).
 3. Train the model to learn the difference.
 4. Test the model on new URLs.
- **Outcome:**
 - The system can **predict if a URL is safe or malicious**.
 - Example: Email providers warning about phishing links.
-

Building Recommender Systems

- **Goal: Suggest items to users** such as movies, products, or songs.
- **Why:** Helps users find what they like and improves engagement (used by Netflix, Amazon, Spotify).
- **How It Works:**
 - There are two main methods:
 1. **Collaborative Filtering:**
 - Based on what users like or interact with.
 - Example: If you and another person both like Movie A, and they liked Movie B, you might also like Movie B.
 2. **Content-Based Filtering:**
 - Based on item details (genre, category).
 - Example: If you like a sci-fi movie, it suggests other sci-fi movies.
- **Tools:**
 - **Scikit-learn:** For machine learning algorithms.
 - **Surprise library:** Specially designed for building recommendation models.
- **Outcome:**
 - Shows **personalized recommendations** to each user.
 - Example: “Recommended for you” section on Netflix or Amazon.

UNIT-3

NoSQL movement for handling Bigdata:

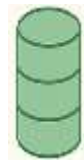
Distributing data storage and processing with Hadoop framework, case study on risk assessment for loan sanctioning, ACID principle of relational databases, CAP theorem, base principle of NoSQL databases, types of NoSQL databases, case study on disease diagnosis and profiling

Why NoSQL?



- Old databases (RDBMS) could not handle **huge, fast, and unstructured data**.
- Big Data needs **speed, flexibility, and scalability**.

What is NoSQL?



- **Not Only SQL** → databases that can store different types of data.
- Works well with **structured, semi-structured, and unstructured data**.
- Can **scale horizontally** (add more servers easily).

Types of NoSQL Databases



Key-Value → like a dictionary (e.g., Redis, DynamoDB).

- Every piece of data is stored as a **key → value pair**.
- **Key = unique identifier, Value = the data** (could be string, JSON, blob, etc.).

Example :

Key: "user101"

Value: {"name": "Sai", "age": 25, "email": "sai@email.com"}



Document → stores JSON/XML documents (e.g., MongoDB, CouchDB).

- Great for **semi-structured** data and dynamic schemas.

Example: {

```
  "username": 101,  
  "name": "Sai",  
  "age": 25,  
  "skills": ["Python", "SQL", "NoSQL"]  
}
```

 **Column-based** → stores data in columns (e.g., Cassandra, HBase).

- Optimized for queries over **large datasets** (good for analytics)
- Useful when you need to read/write a subset of columns frequently.

Example:

UserID | Name | Age | City

101 | Sai | 25 | Hyderabad

102 | Kiran | 30 | Bangalore



Graph → stores relationships (e.g., Neo4j).

- Focuses on storing **entities (nodes)** and their **relationships (edges)**.

Example (social network):

Node: User A (Sai)

Node: User B (Kiran)

Edge: Sai → Follows → Kiran

Summary

- **Key-Value** → Fast lookups, like a dictionary.
- **Document** → JSON/XML documents, flexible schema.
- **Column-based** → Good for analytics, stores by columns.
- **Graph** → Best for relationships and networks.

Benefits for Big Data

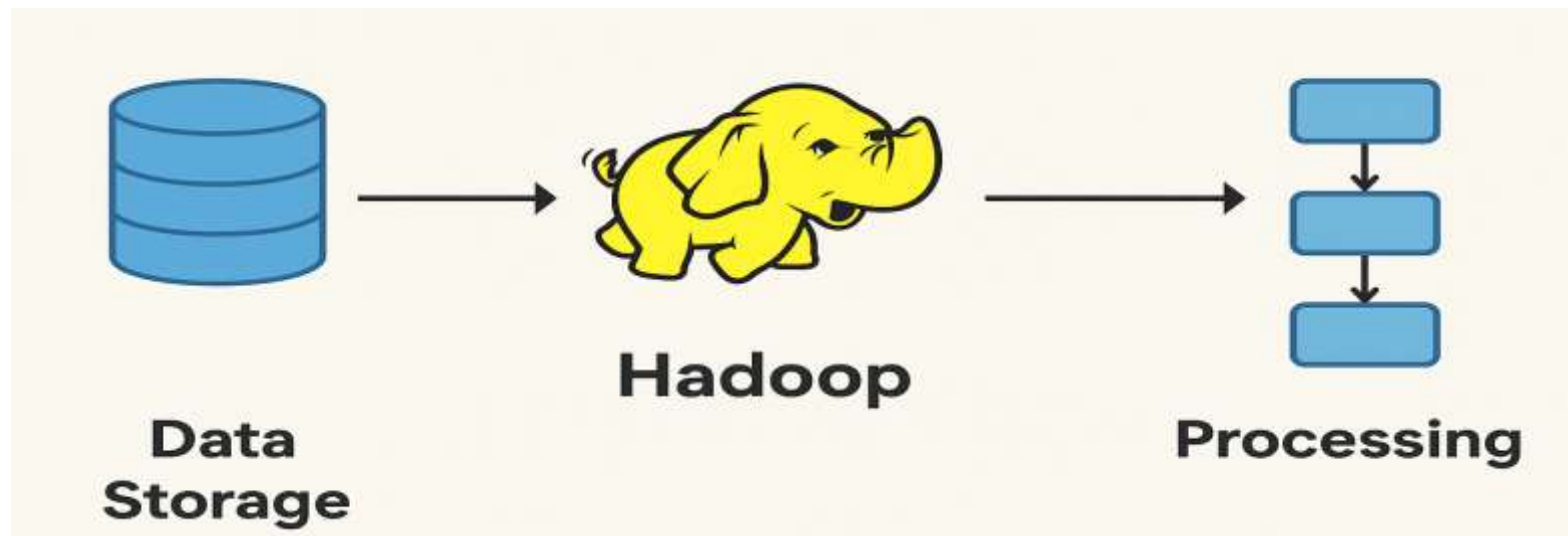


- Handles **large volume** of data.
- Works with **different formats**.
- **Fast read/write** performance.
- Supports **real-time analytics**.

Distributing Data Storage and Processing with Hadoop Framework

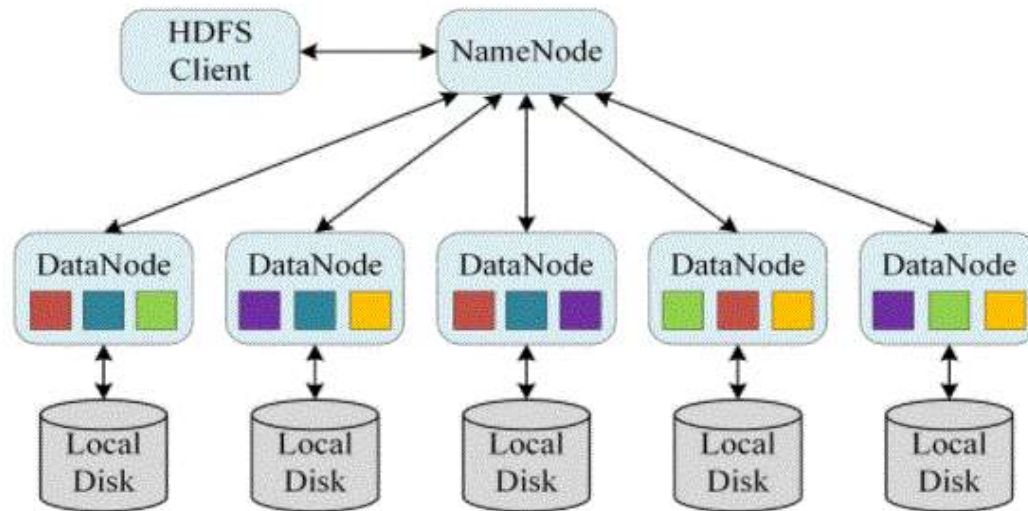
What is Hadoop?

- Hadoop is an **open-source framework** that helps in storing and processing **Big Data** using **clusters of computers** instead of a single machine.



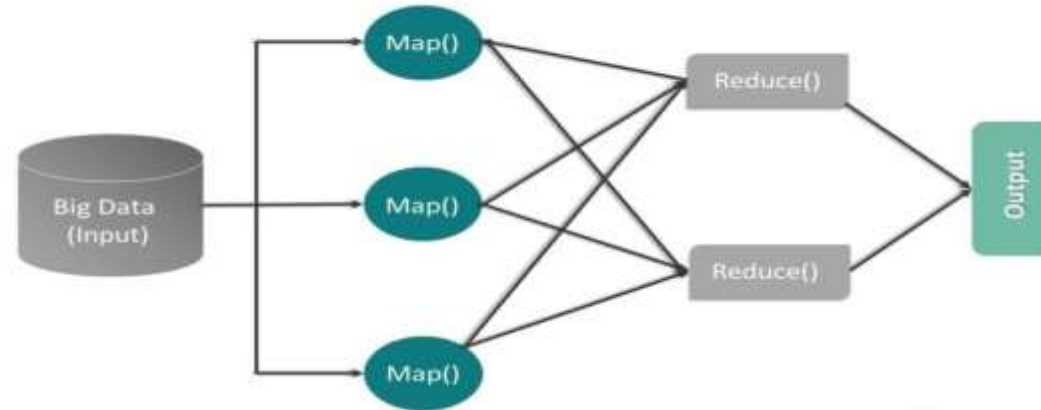
Main Components of Hadoop

1.HDFS (Hadoop Distributed File System) → for storage



- Splits large files into **blocks** (e.g., 128 MB each).
- Stores blocks across many machines (nodes).
- Keeps **replicas** of data for fault tolerance.

2. MapReduce → for processing



- **Map step** → divides the task into smaller subtasks and runs them on different machines.
- **Reduce step** → collects and combines the results.

3. YARN (Yet Another Resource Negotiator) → for resource management

- Assigns tasks and manages cluster resources.
- YARN is the resource manager of Hadoop that decides who gets how much CPU and memory, and helps run jobs smoothly.

How Hadoop Distributes Data & Work

- A **big file** is broken into blocks.
- These blocks are **stored across multiple nodes** in HDFS.
- When processing:
 - Hadoop sends the **program to the data location** (not data to program).
 - Each node processes its part of data using **Map tasks**.
 - Results are collected and combined by **Reduce tasks**.

Benefits

- **Scalable** → can add more machines easily.
- **Fault tolerant** → data is replicated, so no loss if one machine fails.
- **Cost effective** → uses low-cost hardware.
- **Efficient** → moves computation to where data is stored.

NoSQL + Hadoop Together

- Many NoSQL databases (like **HBase**, **Cassandra**) work **on top of Hadoop**.
- Example:
 - Store petabytes of data in **HDFS**.
 - Use **HBase (a NoSQL DB)** to provide fast random access to this distributed data.
 - Use **MapReduce/Spark** for large-scale processing.

Example

Imagine you have a **1 TB file** split into 10 blocks across 10 nodes.

- Instead of bringing 1 TB to one computer,
- Hadoop sends the **processing code** to all 10 nodes,
- Each node processes only its block (100 GB).
- Results are gathered → much **faster and efficient**.

case study on risk assessment for loan sanctioning

1. Objective

The main aim is to **predict the probability of loan default** before sanctioning a loan.

- Banks want to lend money but also **minimize losses** from borrowers who might fail to repay.
- A good risk assessment framework ensures only reliable customers are approved.

2.Data Collected

To assess risk, the bank collects both financial and personal data:

- **Income** → Determines repayment ability.
- **Employment** → Stable jobs = lower risk, unstable jobs = higher risk.
- **Credit Score** → (like CIBIL score in India) shows repayment history across banks.
- **Loan Amount Requested** → Higher amount = higher risk.
- **Repayment History** → Past defaults or delays raise risk

3.Method

The bank uses different **risk modeling techniques**:

- **Credit Scoring Models**

- Rule-based: e.g., “Reject if CIBIL < 650.”

- **Statistical Models**

- **Logistic Regression** → Predicts probability of default (0 = safe, 1 = risky).

- **Machine Learning Models**

- Random Forest, Gradient Boosting → Learn patterns from past borrower data.
- Example: If borrowers with high DTI (Debt-to-Income ratio) often defaulted, the model flags new applicants with high DTI as risky.

4. Decision Making

Applicants are classified into categories:

- **Low Risk** → Eligible, normal interest rate.
- **Medium Risk** → Eligible but with conditions (collateral, higher interest).
- **High Risk** → Rejected to avoid future losses.

Example:

- Applicant A (Income ₹80,000, stable job, CIBIL 720) → **Low Risk → Loan Approved.**
- Applicant B (Income ₹20,000, unstable job, CIBIL 580) → **High Risk → Loan Rejected.**

Outcome

- Reduces **Non-Performing Assets (NPAs)** for banks.
- Ensures **financial stability**.
- Helps in **faster loan processing** with automated decision models.
- Balances **profitability** (lending more) with **safety** (avoiding bad loans).

ACID principle of relational databases

Exploring "ACID"

This presentation explores ACID properties. ACID stands for **Atomicity, Consistency, Isolation, and Durability**. These properties are crucial in database transactions. They ensure reliable and accurate data management. We will cover the importance of ACID compliance in the up-coming slides.

Atomicity: All or Nothing

Explanation

- Atomicity means a transaction is indivisible. It's an "all or nothing" principle.
- Either all parts of the transaction complete, or none do.

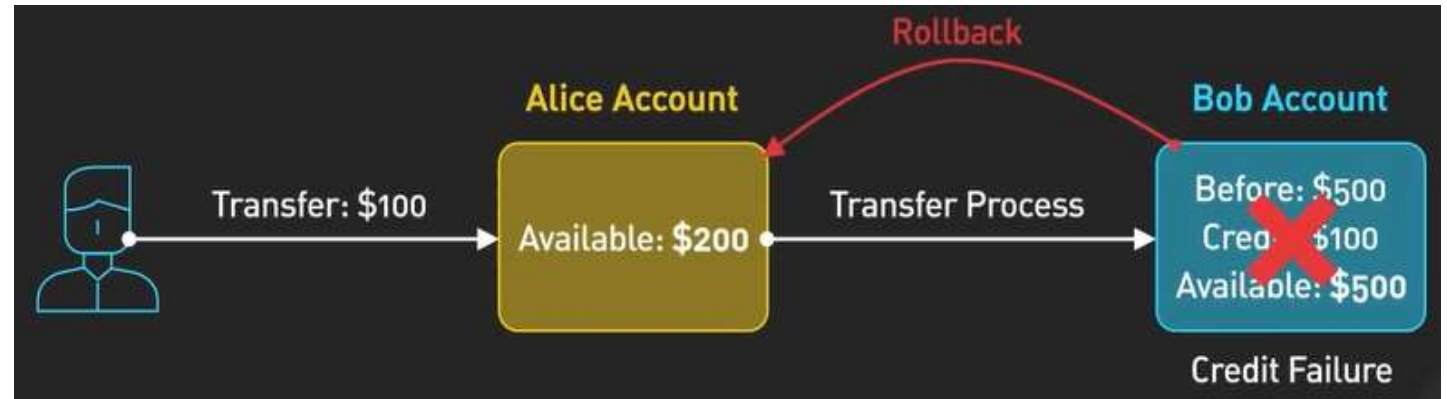
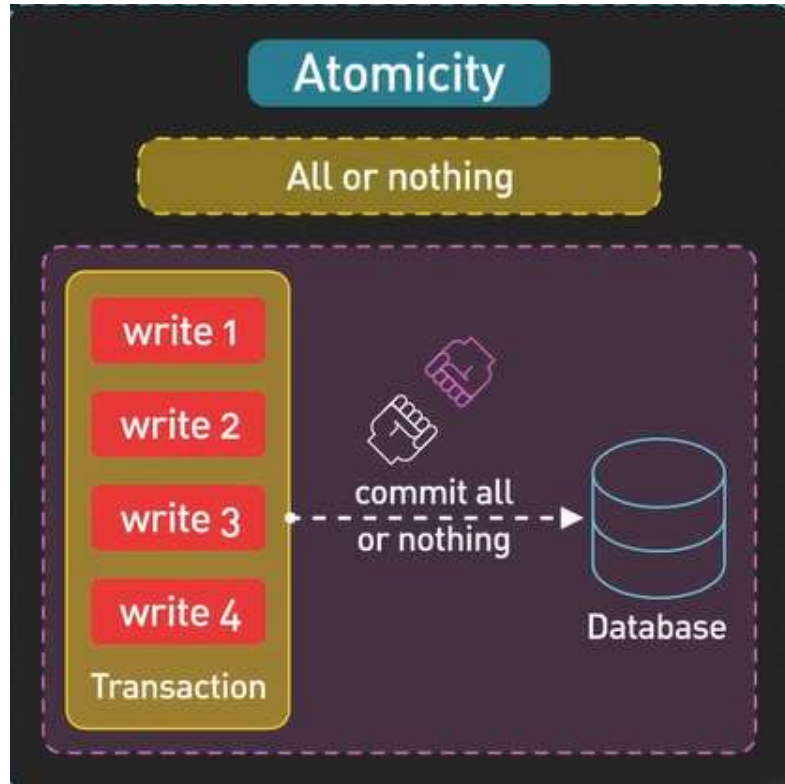
Bank transaction

- Either debit and credit happen, or neither.
- If one fails, the entire transaction rolls back.

E-commerce order

- The whole order is placed, or completely rolled back.
- This ensures a consistent state.

A - ATOMICITY



Consistency: Maintaining Data Integrity

Consistency Defined

Consistency ensures the database remains in a valid state. Data follows rules and constraints.

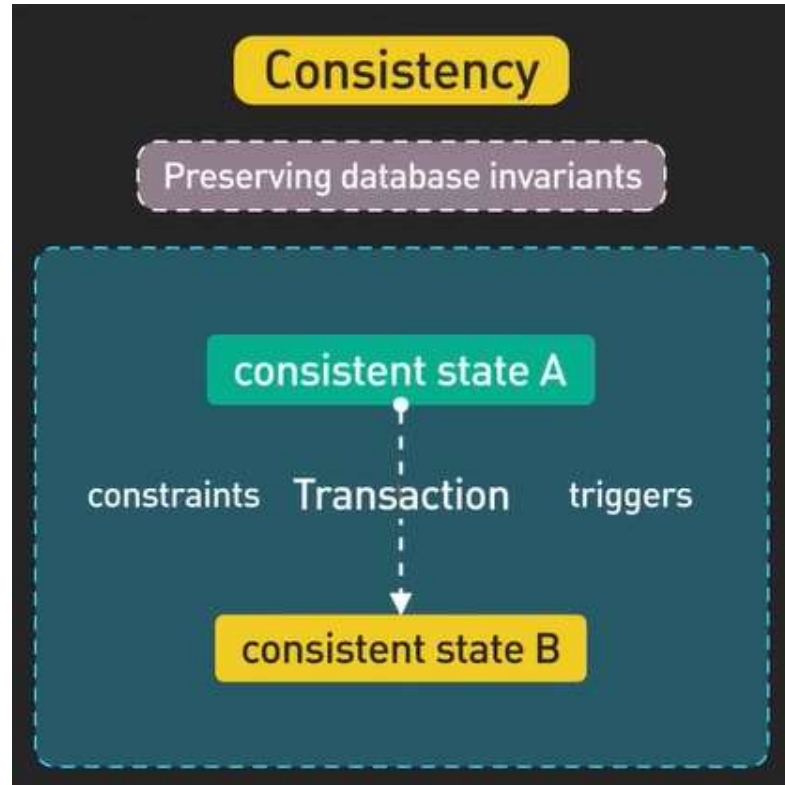
Account Balance

Account balance must be greater than or equal to the minimum balance after a transaction. Violations cause a roll back.

Data Types

Validates data types. For example, age should always be a number and never text.

C - CONSISTENCY



Isolation: Concurrent Transactions

Transactions Independent

- Transactions operate independently from each other.
- One transaction should not interfere with another.

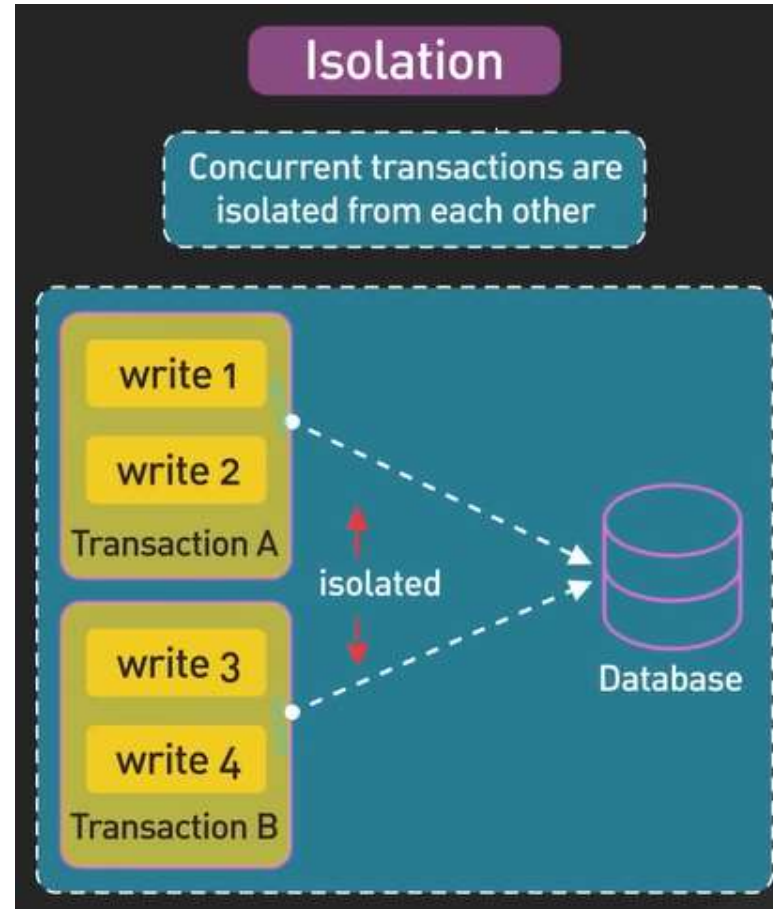
Isolation Levels

- Common levels include: Read Committed, Repeatable Read, and Serializable.
- These levels control the degree of isolation.

Booking Tickets

- Two users attempting to book the last train ticket concurrently.
- Isolation prevents overbooking.

I - ISOLATION



Durability: Guaranteed Data Persistence



- **Committed Transactions**

- Committed transactions are permanent and cannot be lost. They survive system failures.



- **System Failures**

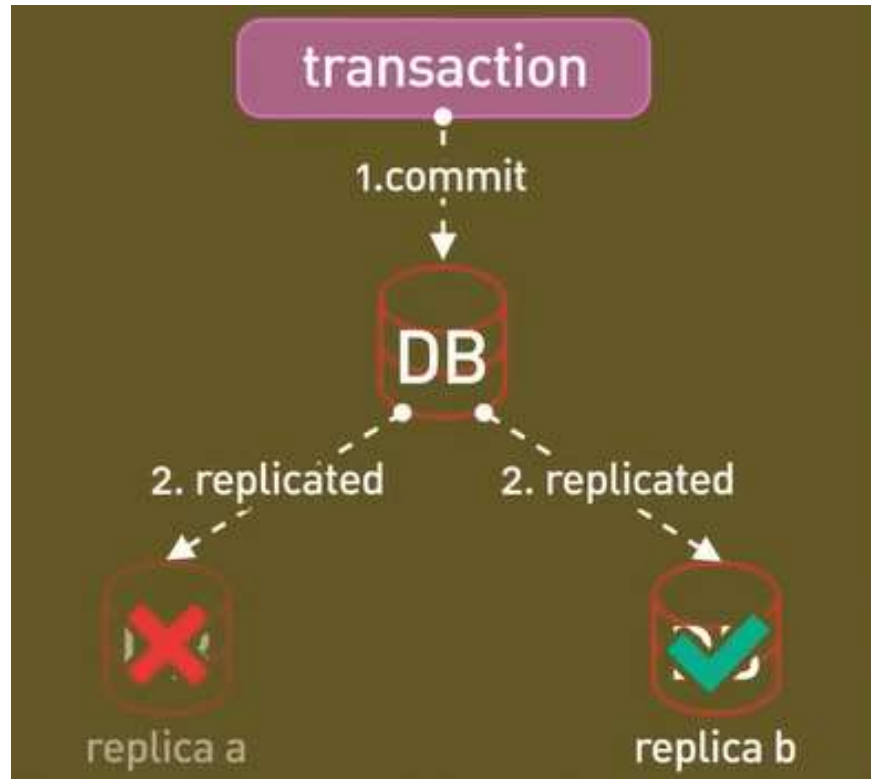
- Surviving power outages and crashes. Data is reliably stored.



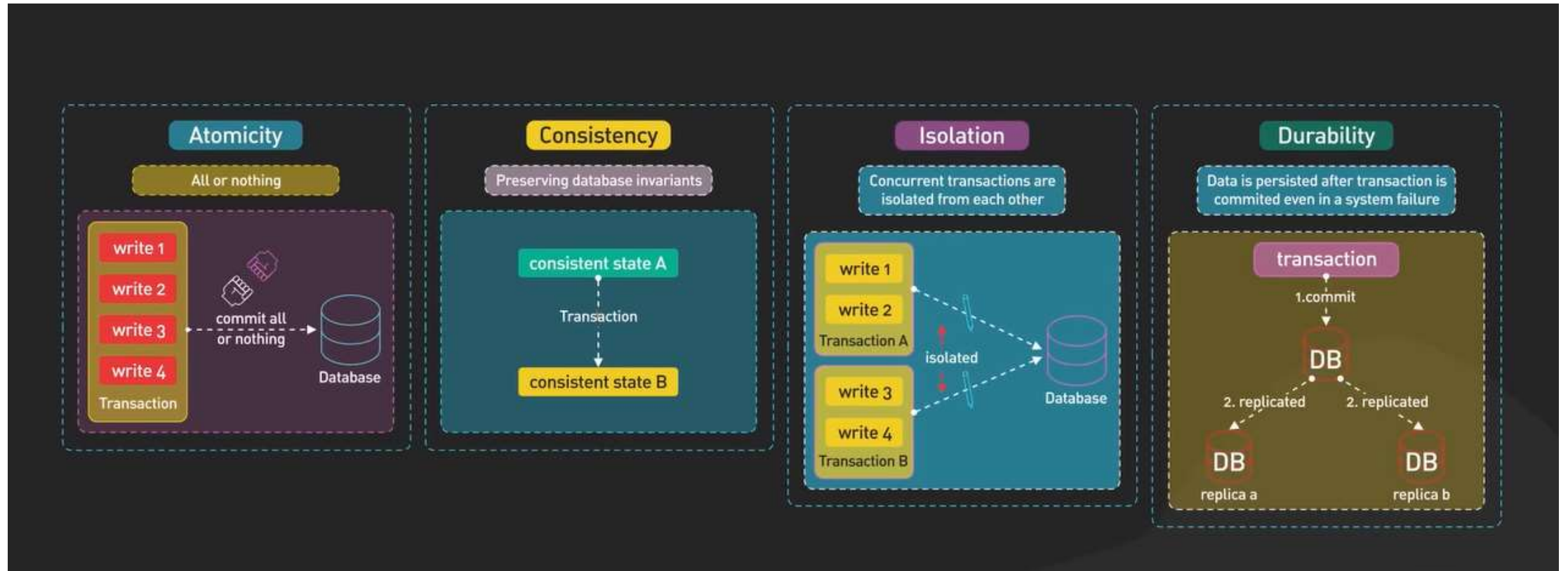
- **+Transaction Logs**

+Transaction logs ensure durability. Backup and recovery mechanisms support this.

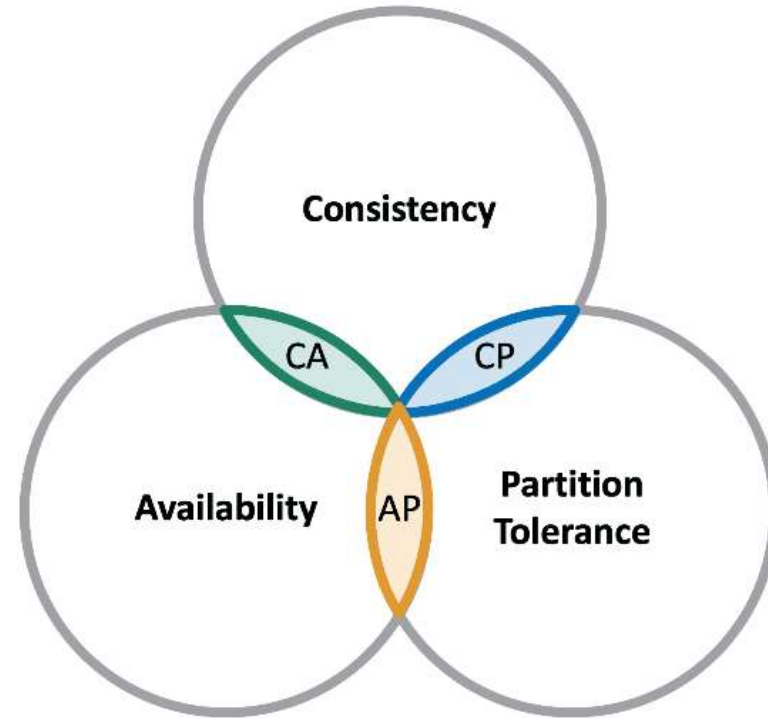
D - DURABILITY



SUMMARY



CAP Theorem



Introduction

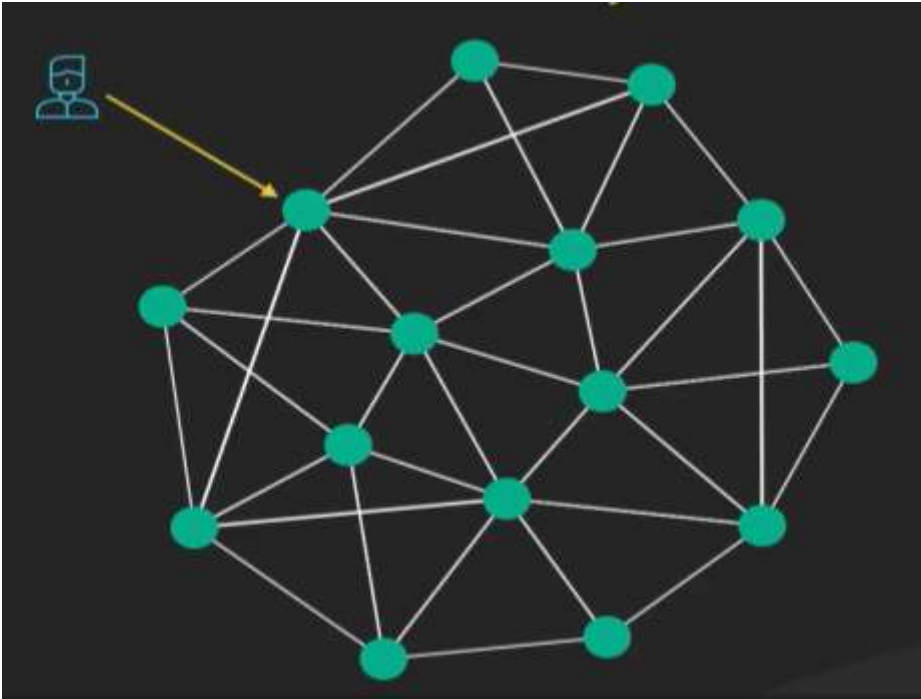
CAP theorem is a fundamental principle of distributed systems

Proposed by Eric Brewer in 2000

Explains the trade-offs among three key properties:

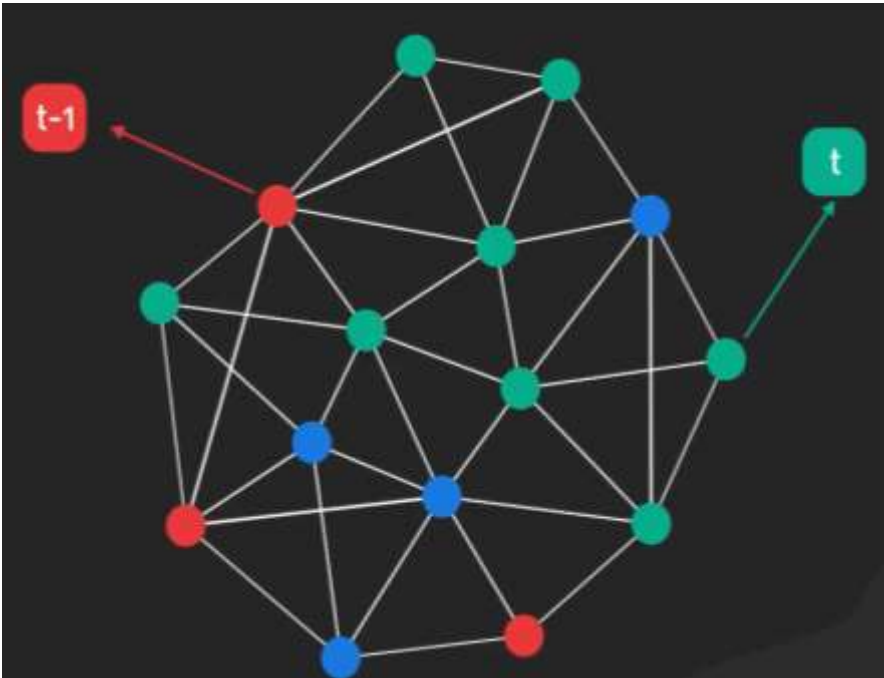
- Consistency (C)
- Availability (A)
- Partition Tolerance (P)

Consistency (C)



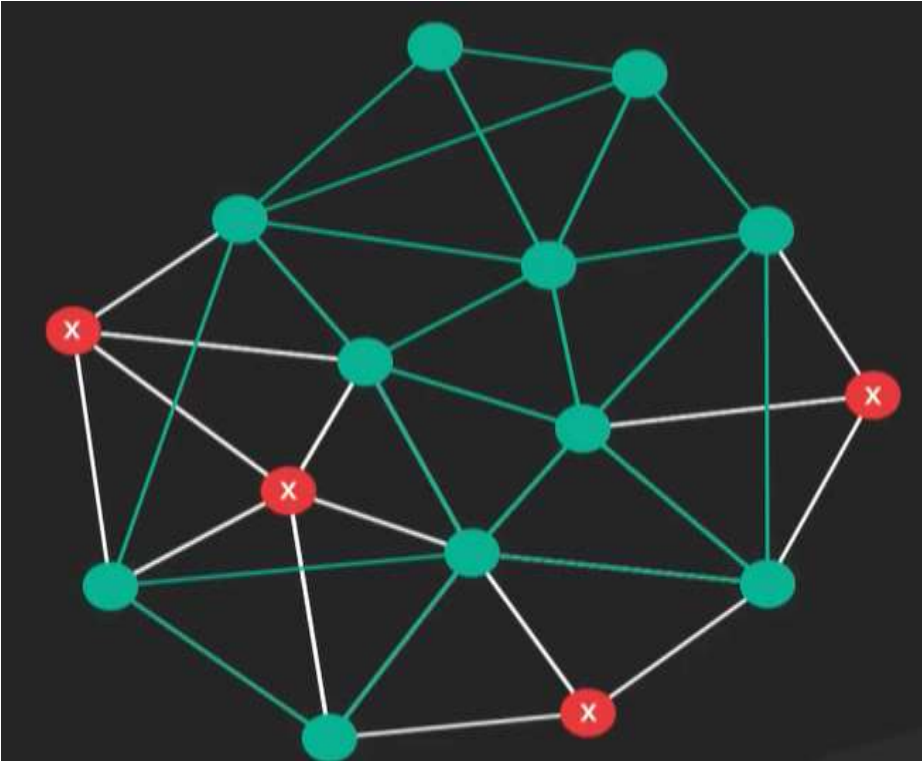
- + Every read receives the most recent write or an error
- + All nodes show the same data at the same time
- + Example: Banking systems require strict consistency

Availability (A)

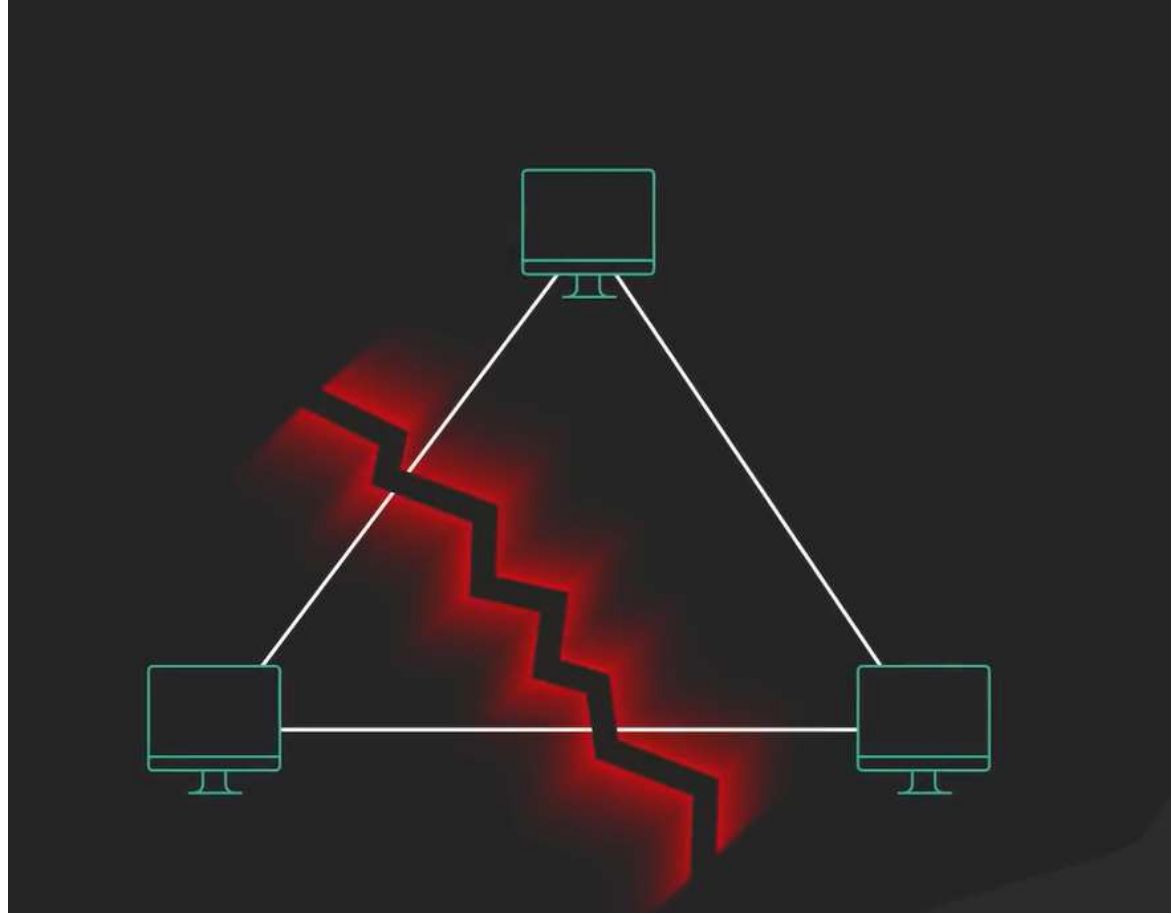


- + Every request gets a response (even if it's not the most recent data)
- + System remains operational at all times
- + Example: E-commerce websites prefer availability during high load

Partition Tolerance (P)



- + System continues to function despite network failures or message loss between nodes
- + Essential for any distributed system
- + Example: Microservices spread across datacenters



CAP Theorem Summary

- In a distributed system, it's impossible to guarantee **Consistency, Availability, and Partition Tolerance** all at the same time.
During a network partition, you must choose **either Consistency or Availability**, but not both.

Base Principles of NoSQL Databases

1. Schema-less (Flexible Data Model)

- Unlike RDBMS (which requires a fixed schema), NoSQL allows **dynamic and flexible schemas**.
- Example: In MongoDB, one document may have fields name, age, while another has name, email.

2. Horizontal Scalability

- Instead of scaling up with bigger servers (vertical scaling), NoSQL databases **scale out** by adding more machines (nodes).
- This helps handle **big data and high traffic** efficiently.

3. Distributed Architecture

- Data is distributed across multiple nodes (servers).
- Provides **fault tolerance, high availability, and faster performance**.

4. BASE Properties (instead of ACID)

- RDBMS uses **ACID** (Atomicity, Consistency, Isolation, Durability).
- NoSQL often uses **BASE** for flexibility:
 - **Basically Available** → system is available most of the time.
 - **Soft State** → data may change over time without input (due to replication delays).
 - **Eventually Consistent** → consistency is guaranteed, but not immediately.

5. Optimized for Specific Data Models

- Different types of NoSQL databases are optimized for different use cases:
 - **Key-Value Stores** (Redis, DynamoDB) → fast lookups
 - **Document Stores** (MongoDB, CouchDB) → semi-structured JSON/XML data
 - **Column Stores** (Cassandra, HBase) → analytical workloads, wide tables
 - **Graph Stores** (Neo4j) → relationships, networks

6. High Performance and Low Latency

- Designed to handle **large-scale read/write operations** with low response time.
- Often used in **real-time applications** (e.g., chat apps, recommender systems).

Summary

- NoSQL databases are built on **flexibility, scalability, and performance**, trading strict consistency (ACID) for availability and speed, making them suitable for **Big Data and real-time applications**.

Types of NoSQL Databases

NoSQL databases are generally divided into **four main types** (sometimes 5, if multi-model is included):

1. Key–Value Stores

- **Data Structure:** Key (unique identifier) → Value (data)
- **Use Case:** Very fast lookups, caching, session management, user preferences
- **Examples:** Redis, Amazon DynamoDB, Riak

2. Document Stores

- **Data Structure:** Documents (usually JSON, BSON, or XML format)
- Each document can have a different structure (schema-less).
- **Use Case:** Content management systems, e-commerce, catalogs, blogs, apps with flexible data.
- **Examples:** MongoDB, CouchDB

3. Column-Oriented Databases (Wide-Column Stores)

- **Data Structure:** Data stored in **columns instead of rows** (like RDBMS).
- Designed for **analytical workloads** and handling large volumes of structured data.
- **Use Case:** Big Data analytics, time-series data, IoT, logs.
- **Examples:** Apache Cassandra, HBase, ScyllaDB

4. Graph Databases

- **Data Structure:** Nodes (entities) + Edges (relationships between nodes).
- Ideal for storing and querying **relationships and networks**.
- **Use Case:** Social networks, fraud detection, recommendation engines.
- **Examples:** Neo4j, Amazon Neptune

5. Multi-Model Databases (Hybrid)

- Some databases support **multiple NoSQL models in one system** (e.g., key-value + document).
- **Examples:** ArangoDB, OrientDB, Cosmos DB

Case Study: Disease Diagnosis and Profiling

1. Title

- **Disease Diagnosis and Profiling using Data Analysis**

2. Objective

- To use data to **identify diseases early** and **understand patient health profiles** for better treatment.

3. Problem

- Doctors handle large amounts of data (test results, reports, patient details).

It's **hard and time-consuming** to check everything manually.

A computer system can help by **analyzing data automatically** and **finding patterns** that show diseases early.

4.Data Description

The dataset consists of patient health information:

Feature Type	Example Attributes	Description
Demographic Data	Age, Gender, Region	Basic patient details
Clinical Measurements	Blood pressure, Sugar level, Heart rate	Vital signs
Laboratory Results	Cholesterol, Hemoglobin, WBC count	Diagnostic test data
Lifestyle Information	Smoking, Alcohol, Exercise	Risk factor profiling
Target Variable	Disease label (e.g., Diabetes: Yes/No)	Outcome to predict

5. Steps Followed

- **Collect Data** – from hospital records or online health datasets.
- **Clean Data** – fill missing values, remove duplicates.
- **Analyze Data** – find which health factors affect diseases.
- **Build Model** – use machine learning (like Random Forest) to predict disease.
- **Test Model** – check how accurate the model is.

6. Results

- The computer model could **predict diseases with 90% accuracy.**
- Found important patterns like:
 - High sugar → **Diabetes risk**
 - High BP + Cholesterol → **Heart disease risk**
 - Low Hemoglobin → **Anemia**

7. Tools Used

- **Python** (pandas, sklearn, matplotlib)
- **Excel / Power BI** for graphs
- **Databases** (MySQL, PostgreSQL)

8. Outcome

- Faster diagnosis
- Personalized health reports for patients
- Early detection of diseases

9. Conclusion

- Data analysis and machine learning can help doctors:
- Find diseases early
- Save time
- Give better treatment to patients

Case Study: Classifying Reddit Posts using Text Mining and Machine Learning

1. Title

Automatic Classification of Reddit Posts using Text Mining Techniques

2. Objective

To build a system that can **automatically identify the topic** of Reddit posts — such as *sports, technology, or health* — using **Python, text mining, and machine learning**.

3. Problem

Reddit is a popular discussion platform with millions of posts shared daily.

Each post belongs to a different topic, but sorting and organizing them manually takes too much time.

We need a computer program that can read the text of each post, understand its meaning, and classify it into the correct category.

4. Tools and Libraries Used

Purpose	Tools / Libraries
Programming	Python
Text Processing	NLTK (Natural Language Toolkit)
Data Storage	SQLite
Data Handling	pandas, NumPy
Machine Learning	scikit-learn
Visualization (optional)	Matplotlib

5. Steps (Methodology)

Step 1: Data Collection

Posts are collected from Reddit using the Reddit API (like **PRAW**) or a sample dataset.

Each post contains:

- Title
- Text content

- Subreddit name (category)

Example data:

Post Title	Category
“India wins T20 World Cup 2025!”	Sports
“New AI beats humans in chess”	Technology
“10 Tips for a Healthy Heart”	Health

Step 2: Text Cleaning (Using NLTK)

Before analysis, the text is cleaned:

- Convert to lowercase
- Remove punctuation and stopwords (like “the”, “and”)
- Tokenize text into words
- Apply stemming or lemmatization

Example:

Original: "India wins the cricket match!"

Cleaned: ['india', 'win', 'cricket', 'match']

Step 3: Feature Extraction

Convert cleaned text into numeric format using **TF-IDF** (Term Frequency–Inverse Document Frequency).

This helps the computer understand how important each word is in the text.

Step 4: Model Training

Use machine learning models like:

- **Naive Bayes**
- **Logistic Regression**
- **Support Vector Machine (SVM)**

The model learns patterns in text and their corresponding categories.

Example:

Input: "Virat hits a century"

Predicted Category: Sports

Step 5: Model Evaluation

Test the model on new posts to check accuracy.

Example Result:

Accuracy: 88%

Post	Predicted Category
“Apple releases new iPhone”	Technology
“India wins T20 match”	Sports
“Best exercises for heart health”	Health

Step 6: Store Results

Store post title, predicted category, and confidence score in a **SQLite** database for future analysis.

6. Results

- The model correctly classifies most Reddit posts.
- Accuracy achieved around **85–90%**.
- Helps in organizing, filtering, and analyzing large text data automatically.

7. Conclusion

Using **text mining** (with NLTK) and **machine learning**, we can efficiently classify Reddit posts into meaningful categories.

The combination of **NLTK** for text cleaning and **SQLite** for data storage makes the system simple and practical for real-world use.

8. Future Work

- Add **sentiment analysis** (positive or negative posts).
- Use **deep learning models (BERT, LSTM)** for higher accuracy.
- Deploy the model as a web app for real-time Reddit post classification.

Data Visualization Options

Data Visualization means showing data in a picture or graphical form so that it is easy to understand.

Instead of looking at long tables or numbers, we can use charts and graphs to quickly see patterns, comparisons, and trends in the data.

Different visualization options are used depending on what kind of data we have and what we want to show.

1. Bar Chart

- Used to compare different categories.
- Data is shown using rectangular bars (vertical or horizontal).
- Example: Showing sales for each month.

Example:

Month Sales

Jan 100

Feb 150

Mar 200

A **bar chart** can show that sales increased every month.

2. Line Chart

- Used to show **trends over time**.
- Data points are connected by lines.
- Example: Showing how temperature changes every day for a week.

You can easily see if the temperature is rising or falling over time.

3. Pie Chart

- Used to show **parts of a whole**.
- Each slice represents a percentage of the total.
- Example: Showing market share of different companies.

If Company A has 50%, Company B has 30%, and Company C has 20%, the pie chart shows this clearly.

4. Histogram

- Used to show the **distribution of numerical data**.
- Similar to a bar chart but used for continuous data (like ages or marks).
- Example: Showing how many students scored within certain marks ranges.

You can quickly see if most students scored high or low.

5. Scatter Plot

- Used to show the **relationship between two variables**.
- Each point represents one data observation.
- Example: Comparing hours studied vs exam marks.

If the points form an upward line, it means students who studied more scored higher.

6. Area Chart

- Similar to a line chart, but the area below the line is filled with color.
- Used to show changes over time and total value.
- Example: Showing how website traffic increased month by month.

7. Heatmap

- Uses **color** to represent data values.
- Example: Showing how sales vary across different cities and months.
- Darker or brighter colors can mean higher or lower values.

8. Map Visualization

- Used for **geographical data**.
- Example: Showing COVID-19 cases by country or rainfall by region.

Helps see where values are high or low on a map.

9. Box Plot

- Shows **data spread and outliers** (unusual values).
- Example: Showing marks of students to see average and extreme scores.

10. Dashboard (Combination of Charts)

- Combines many charts together for a complete view.
- Example: A **Sales Dashboard** might include:
 - Bar chart (sales by region)
 - Line chart (sales trend)
 - Pie chart (sales by product type)
 - Map (sales by city)

Crossfilter

Crossfilter is a **JavaScript library** used for **fast, interactive data filtering** on the web. It helps you explore large datasets **right inside your browser** — without needing a big database or server.

It is mainly used to build **interactive dashboards** where multiple charts update automatically when you click or filter something.

What Crossfilter Does

Crossfilter allows you to:

1. **Filter data instantly** (like show data only for 2024 or only for “South Region”).
2. **Group data** by a field (like sum of sales per region).
3. **Link multiple charts together** — if you select a value in one chart, all others update automatically.

Why Crossfilter is Important

Normally, when you filter data in a web dashboard, it can be slow because:

- The browser has to reload data.
- Each chart works separately.

But with **Crossfilter**, filtering and grouping happen in memory (in the browser), so everything updates **instantly**.

It’s **fast, lightweight**, and perfect for **interactive dashboards**.

How Crossfilter Works (Simple Example)

Imagine you have sales data like this:

Region Year Sales

North 2023 100

South 2023 150

North 2024 200

South 2024 250

Now, using Crossfilter:

- You can **filter by year = 2024**, and instantly see:
 - Total sales by region.
 - Average sales for 2024.
- If you select **South region**, all charts will automatically show data only for the South region.

Simple Working Steps

1. Load data

Crossfilter loads your dataset (usually in JSON or CSV form).

2. Create dimensions

Dimensions are fields you can filter by (like “Year” or “Region”).

3. Create groups

Groups calculate totals or counts for each value in a dimension (like sum of sales per region).

4. Connect to charts

Use a library like **dc.js** or **D3.js** to display charts based on those dimensions and groups.

Advantages

- ✓ Very fast filtering (even with thousands of records)
 - ✓ Interactive and dynamic dashboards
 - ✓ Works completely in the browser (no server needed)
 - ✓ Easy integration with dc.js
-

Example Use Cases

- **Sales dashboards** – filter by year, region, or product.
- **Population data** – see age, gender, and city distributions interactively.
- **Website analytics** – see user data by country and device type.

JavaScript MapReduce Library

MapReduce is a concept originally developed for processing big data by dividing tasks into smaller parts. In JavaScript, some libraries implement a similar idea to process data quickly in the browser.

The **MapReduce model** works in two main steps:

1. **Map Step:** Each element in the dataset is processed or transformed individually.
Example: Extract the year from each sales record.
2. **Reduce Step:** All processed results are combined or summarized.
Example: Add up total sales for each year.

In JavaScript, this can be done easily using array functions like `.map()` and `.reduce()`.

For example:

```
let sales = [100, 200, 300];  
let total = sales.reduce((a, b) => a + b);  
console.log(total); // 600
```

MapReduce is useful in dashboards or web applications where large datasets must be summarized quickly without needing a back-end database. Libraries like **crossfilter.js** or **d3.map()** internally use MapReduce principles to make data aggregation fast.

Creating an Interactive Dashboard with dc.js

dc.js is a JavaScript library used to build **interactive, dynamic dashboards** in web browsers. It is built on top of **D3.js** and **Crossfilter**, which means it can visualize data efficiently and respond instantly to user actions.

With **dc.js**, you can create different types of charts — such as bar charts, pie charts, line charts, and data tables — and link them together.

When a user clicks on a specific part of one chart (for example, a region in a pie chart), all the other charts automatically update to show related data only for that region.

For example, in a **Sales Dashboard**, you can show:

- A **bar chart** for monthly sales
- A **pie chart** for sales by region
- A **line chart** for yearly sales growth
- A **data table** for viewing detailed transactions

All these charts are connected using Crossfilter, and the interactivity is handled by dc.js. This makes dc.js very useful for building real-time data analysis dashboards on websites without needing heavy software like Tableau or Power BI.

Dashboard Development Tools

Dashboard development tools are used to create **visual interfaces** that display data, key metrics, and trends. These tools help analysts and business users interact with data and make decisions faster.

Some popular tools include:

- **Power BI:** A Microsoft tool for creating interactive dashboards with drag-and-drop features. It connects easily to Excel, SQL databases, and cloud services.
- **Tableau:** A powerful data visualization tool for building complex charts, dashboards, and data stories with little or no coding.
- **Google Data Studio:** A free online tool that connects with Google Sheets, Analytics, and other Google services to create live dashboards.
- **Plotly Dash:** A Python-based framework that allows developers to build analytical web applications using code.
- **Streamlit:** A simple Python library used for quick data science app development and prototyping.

- **dc.js + Crossfilter:** A JavaScript-based combination for building web dashboards directly in the browser.

These tools make it easy to design, test, and deploy dashboards for different use cases like sales analysis, business performance tracking, and predictive analytics.

Applying the Data Science Process for Real-World Problem Solving (Case Study)

Let's take an example to understand how the **data science process** can be applied to a real-world problem using visualization and prototype development.

Case Study Example: Predicting Student Performance

Step 1: Problem Definition

A school wants to analyze student data to find the factors that affect exam performance and predict which students may need extra help.

Step 2: Data Collection

Data is collected from school records — including attendance, test scores, homework completion, and participation.

Step 3: Data Cleaning

Missing or incorrect data is corrected or removed. For example, missing marks are filled with the average, and spelling errors in subjects are fixed.

Step 4: Data Analysis

The relationship between attendance, study hours, and marks is analyzed using Python and visualization tools like Matplotlib or Power BI.

Step 5: Modeling

A machine learning model (like Linear Regression) is built to predict the final exam score based on the input factors.

Step 6: Visualization and Prototype

The results are shown in an **interactive dashboard** using Power BI or dc.js. For example:

- A **bar chart** shows average marks per subject.
- A **line chart** shows progress over time.
- A **pie chart** shows the percentage of students at risk.

Step 7: Decision Making

Teachers can use the dashboard to identify weak students early and provide additional help.