\

# ACT-UNIT 5

# Data Science and Google technologies

# Data-Information-Knowledge:

Data is raw, unprocessed information. It can be anything from numbers to text to images. Data is not meaningful on its own. It needs to be processed and analyzed in order to become information
Or
Data: Data refers to raw, unprocessed facts, figures, or observations collected from various sources. It can be in the form of numbers, text, images, audio, or any other structured or unstructured format. Data can be collected from sensors, databases, files, web scraping, surveys, or any other means. However, on its own, data lacks context and meaning.

Information is data that has been processed and analyzed. It is meaningful and can be used to answer questions, make decisions, or take action. Information can be stored in a variety of formats, such as databases, spreadsheets, and documents.
Or
Information: Information is the result of organizing, structuring, and analyzing data to derive meaning and context. It involves extracting relevant patterns, relationships, or insights (**a clear, deep**) from the data. Information provides answers to specific questions, reveals trends, and helps in making informed decisions. It transforms data into a more understandable and usable form.

Knowledge is information that has been internalized (**of an idea, opinion, belief, and feeling**) and applied. It is the ability to use information to solve problems, make decisions, and take action. Knowledge is often tacit, meaning that it is difficult to express in words. It is often gained through experience and practice.
Or
Knowledge: Knowledge is a higher-level understanding gained from the analysis and interpretation of information. It goes beyond individual data points and insights to capture broader concepts, principles, and connections. Knowledge represents the understanding of how different pieces of information fit together and can be used to solve problems or make predictions. It involves the synthesis and integration of various information sources and experiences.

## Uses or Purposes of Data-Information-Knowledge

The concepts of data, information, and knowledge find numerous applications in data science and Google technologies. Here are some examples of their uses:

### Data:

- Data collection: Gathering raw data from various sources such as databases, APIs, or sensors.

- Data preprocessing: Cleaning, filtering, and transforming data to ensure its quality and readiness for analysis.
- Data storage: Storing data in databases, data lakes, or cloud platforms like Google Cloud Storage.

Information:

- Data visualization: Creating visual representations (charts, graphs, dashboards) to present data in a more understandable and informative manner.
- Business intelligence: Analyzing data to generate reports, identify trends, and gain insights into business operations and performance.
- Data exploration: Conducting exploratory data analysis to discover patterns, correlations, or anomalies within the data.

Knowledge:

- Predictive analytics: Applying statistical, historical data and machine learning models to make predictions based on historical data and derive actionable insights.
- Recommendation systems: Using knowledge about user preferences and behavior to provide personalized recommendations, as seen in Google's recommendation algorithms for products, videos, or search results.
- Natural language processing: Extracting knowledge and understanding from unstructured text data, enabling tasks like sentiment analysis, text classification, or language translation.

## Applications of Data-Information-Knowledge

- Healthcare
- Finance
- Retail
- Manufacturing and Supply Chain
- Government
- E-commerce and Marketing
- Social Media and Content Recommendations

- Healthcare: Data from medical records can be used to identify trends in diseases, develop new treatments, and improve patient care. For example, researchers at the University of California, San Francisco, used data from electronic health records to identify a new risk factor for heart disease.

- Finance: Financial institutions use data to track market trends, manage risk, and make investment decisions. For example, banks use data from credit card transactions to identify potential fraud.

- Retail: Retailers use data to track customer behavior, optimize inventory, and target marketing campaigns. For example, Amazon uses data from past purchases to recommend products to customers.

- Manufacturing: Manufacturers use data to improve product design, optimize production processes, and reduce costs. For example, General Electric uses data from sensors on its jet engines to predict when they need to be repaired.

- Personalize experiences: By understanding a user's preferences, businesses can personalize their experience with a product or service. For example, Netflix uses data about what movies and TV shows users watch to recommend new content.

- Make predictions: By analyzing data, businesses can make predictions about future events. For example, weather forecasting companies use data about past weather patterns to predict future weather conditions.

- Improve decision-making: By having access to data and information, businesses can make better decisions about how to allocate resources and operate their businesses. For example, businesses can use data to identify which products or services are most popular with customers and then focus their marketing efforts on those products or services.

# Databases - Data warehouse - Big Data:

Databases are a collection of data that is organized in a way that makes it easy to store, retrieve, and update. They can be used to store all sorts of data, including customer information, product data, and financial data. In the context of Google technologies, Google Cloud provides several database options, including:

- Cloud SQL: A fully managed relational database service compatible with MySQL, PostgreSQL, and SQL Server.( SQL stands for Structured Query Language. It is a programming language designed for managing data in relational database management systems (RDBMS). SQL is used to perform tasks such as updating data on a database, retrieving data from a database, and creating and deleting databases.)( MySQL is an open-source relational database management system (RDBMS). It is one of the most popular database systems in the world, and is used by a wide range of applications, from small websites to large enterprise systems.)
- Cloud Spanner: A globally distributed, horizontally scalable, and strongly consistent relational database service.
- Firestore: A NoSQL document database that provides real-time data synchronization and offline support for web and mobile applications.( NoSQL is a term used to describe a wide variety of non-relational database management systems. NoSQL databases are designed to store and retrieve large amounts of data in a flexible and scalable way.)

Data warehouses are a type of database that is specifically designed for storing large amounts of data. They are often used to store data that is collected from a variety of sources, such as customer transactions, sensor data, and social media data. It is designed for reporting, analytics, and decision-making purposes. In the Google ecosystem, Google BigQuery is the primary data

warehouse offering. BigQuery is a fully managed, serverless, and highly scalable data warehouse that allows you to run fast and cost-effective SQL queries over vast amounts of data. It supports both structured and semi-structured data formats.

Big data is a term used to describe data sets that are so large or complex that they are difficult to process using traditional data processing methods. Big data can be used to gain insights into customer behavior, identify trends, and make predictions. In Google's ecosystem, several tools are available for processing and analyzing big data:

- Google Cloud Storage: A scalable object storage service that can store and retrieve large amounts of unstructured data, such as files and backups.
- Google Cloud Dataproc: A managed service for running Apache Hadoop and Apache Spark, It is a good choice for businesses which enables distributed processing of big data workloads.
- Google Cloud Dataflow: A fully managed service for building and executing data processing pipelines using popular frameworks like Apache Beam.
- Google Cloud Pub/Sub: A messaging service that allows you to ingest and process streaming data in real-time.

** Relational database service: it easy to set up, operate, and scale a relational database in the cloud. It provides cost-efficient and resizable capacity while automating time-consuming administration tasks such as hardware provisioning, database setup, patching and backups.

# Applications of Databases - Data warehouse - Big Data

- Customer relationship management (CRM): Databases are used to store customer data, such as names, addresses, and purchase history. This data can be used to track customer interactions, identify trends, and provide personalized service.
- Fraud detection: Data warehouses are used to store large amounts of data, such as financial transactions and customer information. This data can be used to identify fraudulent transactions and prevent fraud.
- Risk management: Big data is used to assess risk and make better decisions. For example, big data can be used to predict customer churn(Customer churn is the rate at which customers stop doing business with a company), identify potential problems with products, and assess the risk of natural disasters.
- Healthcare: Data warehouses are used to store patient data, such as medical records and test results. This data can be used to improve patient care, identify trends, and develop new treatments.
- Retail: Data warehouses are used to store sales data, such as product sales and customer demographics. This data can be used to improve inventory management, identify trends, and target marketing campaigns.
- Manufacturing: Data warehouses are used to store production data, such as machine performance and product quality. This data can be used to improve efficiency, identify problems, and improve product quality.
- Financial services: Data warehouses are used to store financial data, such as customer transactions and market data. This data can be used to improve risk management, identify fraud, and make better investment decisions.

- **Business Intelligence and Analytics:** Organizations use data science and Google technologies to gain insights into their operations, customer behavior, and market trends. By leveraging databases, data warehouses, and big data tools, they can analyze large datasets, perform advanced analytics, and generate reports and visualizations to make data-driven decisions.
- **Internet of Things (IoT) Analytics:** The proliferation of IoT devices generates massive amounts of data. Data science techniques, combined with big data platforms, can analyze this data to gain insights, optimize operations, and improve decision-making. For example, analyzing sensor data from manufacturing equipment can help identify patterns and optimize maintenance schedules.
- **Supply Chain Optimization:** By analyzing supply chain data using data science techniques, organizations can optimize inventory management, predict demand, improve logistics, and reduce costs. This helps streamline operations and enhance overall efficiency.

# Supervised learning, Unsupervised learning, Reinforcement learning with respect to Datamining - Machine Learning - Artificial Intelligence - Deep Learning:

1. **Supervised learning** is a type of machine learning where the model is trained on a dataset that includes both input and output data. The model learns to predict the output data given the input data. For example,1. a supervised learning model could be trained to predict the price of a house given its features, such as the number of bedrooms, the square footage, and the location.
   2.Chatbots,Etc.
   3.A bank could use supervised learning to build a model that predicts which customers are likely to default on their loans. The bank would train the model on a dataset that includes data on past customers, such as their income, their credit score, and their history of making payments on time. The model would learn to identify patterns in the data that are associated with customers who are likely to default on their loans.

Here are some of the benefits of using supervised learning:

- Accuracy: Supervised learning algorithms can be very accurate, especially when used with large datasets.
- Scalability: Supervised learning algorithms can be scaled to handle large datasets.
- Interpretability: Supervised learning algorithms can be interpreted, which can help to understand how the model works.

2. **Unsupervised learning** is type of machine learning (ML) that allows computers to learn from unlabeled data, is a type of machine learning where the model is trained on a dataset that only includes input data. The model learns to identify patterns in the data without any guidance from labeled output data.

For example, 1.an unsupervised learning model could be used to cluster customers into different groups based on their purchasing behavior.

2.A retailer could use unsupervised learning to cluster customers into different groups based on their purchasing behavior. The retailer would train the model on a dataset that includes data on past customer purchases. The model would learn to identify patterns in the data that are associated with different customer groups. For example, the model might identify a group of customers who are interested in buying organic food, a group of customers who are interested in buying electronics, and a group of customers who are interested in buying home goods.

Here are some of the benefits of using unsupervised learning:

- It can be used to find patterns in data that would be difficult or impossible to find with supervised learning.

- It can be used to identify relationships between different data points.

- It can be used to reduce the number of features in a dataset without losing too much information.

3. Reinforcement learning is a type of machine learning where the model learns to take actions in an environment in order to maximize a reward. The model is not explicitly trained on any data. Instead, it learns by trial and error.

   For example, 1. a reinforcement learning model could be used to train a robot to walk. The robot would start by randomly moving its legs. Over time, it would learn to walk more efficiently by trial and error.

   2.A self-driving car could use reinforcement learning to learn how to navigate traffic. The car would start by randomly driving around. Over time, it would learn to avoid obstacles and drive safely by trial and error.

Reinforcement learning (RL) is a type of machine learning where an agent learns to take actions in an environment in order to maximize a reward. RL is a trial-and-error process, where the agent learns by trying different actions and observing the consequences. The agent is not explicitly told what actions to take, but instead learns to associate actions with rewards.

RL has many benefits, including:

- Scalability: RL algorithms can be scaled to handle large and complex environments. This makes them well-suited for real-world problems, such as autonomous driving and robotics.

- Robustness: RL algorithms are robust to noise and disturbances in the environment. This makes them well-suited for real-world problems, where the environment is never perfectly known.

- Generalization: RL algorithms can generalize to new environments. This means that they can be trained on a small number of environments and then used in new environments without retraining.

- Adaptability: RL algorithms can adapt to changes in the environment. This means that they can continue to learn and improve even after they have been deployed in the real world.

| Feature | Supervised Learning | Unsupervised Learning | Reinforceme |
|---|---|---|---|
| Data | Labeled data | Unlabeled data | Trial and err |
| Goal | Predict the output for new data points | Identify patterns and relationships in the data | Take actions |
| Examples | Image classification, spam filtering, fraud detection | Clustering, dimensionality reduction, market basket analysis | Playing gam autonomous |

4. Data mining is the process of extracting knowledge from data. It is a process that involves finding patterns and trends in data that would otherwise be overlooked. Data mining can be used for a variety of purposes, such as:

- Predictive analytics: Data mining can be used to predict future events, such as customer behavior or product sales.
- Customer segmentation: Data mining can be used to segment customers into groups with similar characteristics, such as age, gender, or interests.
- Fraud detection: Data mining can be used to detect fraudulent activities, such as credit card fraud or insurance fraud.

**Examples:**

A retailer could use data mining to predict which customers are likely to churn. The retailer would train a model on a dataset that includes data on past customer purchases, such as the products they have purchased, the amount they have spent, and the frequency of their purchases. The model would learn to identify patterns in the data that are associated with customers who are likely to churn. For example, the model might identify customers who have not made a purchase in a long time, customers who have been spending less money, or customers who have been buying different products. The retailer could then use this information to target these customers with special offers or promotions in an effort to keep them as customers.

5. Machine learning is a type of artificial intelligence (AI) that allows software applications to become more accurate in predicting outcomes without being explicitly programmed to do so. Machine learning algorithms use historical data as input to predict new output values.

**Examples:**

A bank could use machine learning to build a model that predicts which customers are likely to default on their loans. The bank would train the model on a dataset that includes data on past customers, such as their income, their credit score, and their history of making payments on time. The model would learn to identify patterns in the data that are associated with customers

who are likely to default on their loans. The bank could then use this information to make more informed decisions about who to lend money to.

6. <mark>Artificial intelligence</mark> (AI) is a branch of computer science that deals with the creation of intelligent agents, which are systems that can reason, learn, and act autonomously. AI research has been highly successful in developing effective techniques for solving a wide range of problems, from game playing to medical diagnosis.

**Examples:**

A self-driving car uses AI to navigate the road without human input. The car uses a variety of sensors, such as cameras, radar, and lidar, to perceive its surroundings. It then uses AI algorithms to make decisions about how to move safely through traffic.

7. <mark>Deep learning</mark> is a type of machine learning that uses artificial neural networks to learn from data. Neural networks are inspired by the human brain and are able to learn complex patterns from data. Deep learning has been used to achieve state-of-the-art results in a variety of tasks, such as image classification, natural language processing, and speech recognition.

**Examples:**

A company could use deep learning to build a system that can automatically detect fraud in financial transactions. The system would be trained on a dataset of fraudulent and non-fraudulent transactions. The system would then use deep learning algorithms to learn to identify patterns in the data that are associated with fraudulent transactions. The system could then be used to flag suspicious transactions for further review.

# Real time applications on AI/ML/DL/DS:

- **Self-driving cars:** Self-driving cars use AI and machine learning to navigate the road without human input.
- **Virtual assistants:** Virtual assistants like Siri and Alexa use AI to understand natural language and respond to user requests.
- **Risk assessment:** AI and machine learning are used to assess risk in real time, such as the risk of a customer defaulting on a loan.
- **Healthcare:** AI and machine learning are used to diagnose diseases, recommend treatments, and monitor patients in real time.
- **Manufacturing:** AI and machine learning are used to optimize production processes, improve product quality, and prevent defects in real time.
- **Retail:** AI and machine learning are used to personalize shopping experiences, recommend products, and prevent fraud in real time.
- **Financial services:** AI and machine learning are used to manage risk, make investment decisions, and prevent fraud in real time.
- **Natural Language Processing (NLP) and Chatbots:** AI and ML techniques are used to develop chatbots and virtual assistants that can understand and respond to natural language queries in real time. These systems enable businesses to provide instant customer support, automate repetitive tasks, and enhance user experiences.

- **Image and Video Recognition:** AI and DL algorithms are used to analyze and recognize objects, faces, and scenes in real-time images and videos. Applications include facial recognition for security purposes, object detection in autonomous vehicles, and real-time video analytics for surveillance systems.
- **Fraud Detection and Anomaly Detection:** ML and DS techniques are employed to detect fraudulent activities in real time, such as credit card fraud, insurance fraud, or network intrusion. By analyzing patterns and deviations from normal behavior, these systems can quickly flag and respond to potential fraud.
- **Predictive Maintenance:** ML models can be used for real-time monitoring and prediction of equipment failures and maintenance needs. By analyzing sensor data and historical patterns, organizations can identify potential issues before they occur, optimize maintenance schedules, and reduce downtime.
- **Autonomous Vehicles:** AI and ML technologies are at the core of autonomous vehicles, enabling real-time perception, decision-making, and control. These systems analyze sensor data, such as cameras and LiDAR, to detect objects, navigate roads, and make driving decisions in real time.
- **Financial Trading and Algorithmic Trading:** AI and ML algorithms are utilized in real-time financial trading to analyze market data, identify patterns, and make trading decisions. High-frequency trading relies on real-time data analysis and rapid decision-making to execute trades at optimal prices.
- **Recommendation Systems:** Real-time recommendation systems, powered by ML and DL techniques, analyze user behavior and preferences to provide personalized recommendations in various domains such as e-commerce, streaming services, and content platforms.
- **Sentiment Analysis and Social Media Monitoring:** ML and NLP techniques are used to analyze social media streams and real-time customer feedback to understand sentiment, monitor brand reputation, and detect emerging trends or issues.

**Here are some specific examples of real-time applications of AI/ML/DL/DS in different industries:**

- In the healthcare industry, AI is being used to develop new treatments for diseases, diagnose diseases more accurately, and monitor patients remotely. For example, Google AI has developed a new AI system that can diagnose skin cancer as accurately as a dermatologist.

- In the retail industry, AI is being used to personalize shopping experiences, recommend products, and prevent fraud. For example, Amazon uses AI to recommend products to customers based on their past purchases.

- In the financial services industry, AI is being used to manage risk, make investment decisions, and prevent fraud. For example, BlackRock uses AI to manage risk across its $6.3 trillion asset portfolio.