# DATA MINING (COCSC16)
## UNIT 1

## What is Data Mining?

➢ Data mining is the process of extracting knowledge or insights from large amounts of data using various statistical and computational techniques.

➢ The data can be structured, semi-structured or unstructured, and can be stored in various forms such as databases, data warehouses, and data lakes.

➢ The primary goal of data mining is to discover hidden patterns and relationships in the data that can be used to make informed decisions or predictions. This involves exploring the data using various techniques such as clustering, classification, regression analysis, association rule mining, and anomaly detection.

➢ Data mining has a wide range of applications across various industries, including marketing, finance, healthcare, and telecommunications. For example, in marketing, data mining can be used to identify customer segments and target marketing campaigns, while in healthcare, it can be used to identify risk factors for diseases and develop personalized treatment plans.

➢ Data mining is one of the most useful techniques that help entrepreneurs, researchers, and individuals to extract valuable information from huge sets of data.

➢ **Data mining is also called *Knowledge Discovery in Database (KDD)*. The knowledge discovery process includes Data cleaning, Data integration, Data selection, Data transformation, Data mining, Pattern evaluation, and Knowledge presentation.**

## Advantages of Data Mining

o The Data Mining technique enables organizations to obtain knowledge-based data.

o Data mining enables organizations to make lucrative modifications in operation and production.

o Compared with other statistical data applications, data mining is a cost-efficient.

o Data Mining helps the decision-making process of an organization.

o It Facilitates the automated discovery of hidden patterns as well as the prediction of trends and behaviors.

o It can be induced in the new system as well as the existing platforms.

o It is a quick process that makes it easy for new users to analyze enormous amounts of data in a short time.

## Disadvantages of Data Mining

o There is a probability that the organizations may sell useful data of customers to other organizations for money. As per the report, American Express has sold credit card purchases of their customers to other organizations.

- o Many data mining analytics software is difficult to operate and needs advance training to work on.
- o Different data mining instruments operate in distinct ways due to the different algorithms used in their design. Therefore, the selection of the right data mining tools is a very challenging task.
- o The data mining techniques are not precise, so that it may lead to severe consequences in certain conditions.

## Data Mining Applications

- o Data Mining is primarily used by organizations with intense consumer demands- Retail, Communication, Financial, marketing company, determine price, consumer preferences, product positioning, and impact on sales, customer satisfaction, and corporate profits.
- o Data mining enables a retailer to use point-of-sale records of customer purchases to develop products and promotions that help the organization to attract the customer.



These are the following areas where data mining is widely used:

1. **Data Mining in Healthcare:**
   - ➢ Data mining in healthcare has excellent potential to improve the health system. It uses data and analytics for better insights and to identify best practices that will enhance health care services and reduce costs.
   - ➢ Analysts use data mining approaches such as Machine learning, multi-dimensional database, Data visualization, soft computing, and statistics. Data Mining can be used to forecast patients in each category. The procedures ensure

that the patients get intensive care at the right place and at the right time. Data mining also enables healthcare insurers to recognize fraud and abuse.

2. **Data Mining in Market Basket Analysis:**
   - Market basket analysis is a modeling method based on a hypothesis. If you buy a specific group of products, then you are more likely to buy another group of products. This technique may enable the retailer to understand the purchase behavior of a buyer.
   - This data may assist the retailer in understanding the requirements of the buyer and altering the store's layout accordingly. Using a different analytical comparison of results between various stores, between customers in different demographic groups can be done.

3. **Data mining in Education:**
   - Education data mining is a newly emerging field, concerned with developing techniques that explore knowledge from the data generated from educational Environments.
   - EDM objectives are recognized as affirming student's future learning behavior, studying the impact of educational support, and promoting learning science.
   - An organization can use data mining to make precise decisions and also to predict the results of the student. With the results, the institution can concentrate on what to teach and how to teach.

4. **Data Mining in Manufacturing Engineering:**
   - Knowledge is the best asset possessed by a manufacturing company. Data mining tools can be beneficial to find patterns in a complex manufacturing process.
   - Data mining can be used in system-level designing to obtain the relationships between product architecture, product portfolio, and data needs of the customers.
   - It can also be used to forecast the product development period, cost, and expectations among the other tasks.

5. **Data Mining in CRM (Customer Relationship Management):**
   - Customer Relationship Management (CRM) is all about obtaining and holding Customers, also enhancing customer loyalty and implementing customer-oriented strategies.
   - To get a decent relationship with the customer, a business organization needs to collect data and analyze the data. With data mining technologies, the collected data can be used for analytics.

6. **Data Mining in Fraud detection:**
   - Billions of dollars are lost to the action of frauds. Traditional methods of fraud detection are a little bit time consuming and sophisticated.
   - Data mining provides meaningful patterns and turning data into information. An ideal fraud detection system should protect the data of all the users. Supervised

methods consist of a collection of sample records, and these records are classified as fraudulent or non-fraudulent.

➢ A model is constructed using this data, and the technique is made to identify whether the document is fraudulent or not.

# Data Mining Techniques

➢ Data mining includes the utilization of refined data analysis tools to find previously unknown, valid patterns and relationships in huge data sets.
➢ These tools can incorporate statistical models, machine learning techniques, and mathematical algorithms, such as neural networks or decision trees. Thus, data mining incorporates analysis and prediction.
➢ Depending on various methods and technologies from the intersection of machine learning, database management, and statistics, professionals in data mining have devoted their careers to better understanding how to process and make conclusions from the huge amount of data, but what are the methods they use to make it happen?

In recent data mining projects, various major data mining techniques have been developed and used, including association, classification, clustering, prediction, sequential patterns, and regression.

## 1. Classification:

This technique is used to obtain important and relevant information about data and metadata. This data mining technique helps to classify data in different classes.

Data mining techniques can be classified by different criteria, as follows:

i. **Classification of Data mining frameworks as per the type of data sources mined:** This classification is as per the type of data handled. For example, multimedia, spatial data, text data, time-series data, World Wide Web, and so on.

ii. **Classification of data mining frameworks as per the database involved:** This classification based on the data model involved. For example. Object-oriented database, transactional database, relational database, and so on.

iii. **Classification of data mining frameworks as per the kind of knowledge discovered:** This classification depends on the types of knowledge discovered or data mining functionalities. For example, discrimination, classification, clustering, characterization, etc. some frameworks tend to be extensive frameworks offering a few data mining functionalities together.

iv. **Classification of data mining frameworks according to data mining techniques used:** This classification is as per the data analysis approach utilized, such as neural networks, machine learning, genetic algorithms, visualization, statistics, data warehouse-oriented

or database-oriented, etc. The classification can also take into account, the level of user interaction involved in the data mining procedure, such as query-driven systems, autonomous systems, or interactive exploratory systems.

## 2. Clustering:

- ➢ Clustering is a division of information into groups of connected objects. Describing the data by a few clusters mainly loses certain confine details, but accomplishes improvement.
- ➢ It models data by its clusters. Data modeling puts clustering from a historical point of view rooted in statistics, mathematics, and numerical analysis.
- ➢ From a machine learning point of view, clusters relate to hidden patterns, the search for clusters is unsupervised learning, and the subsequent framework represents a data concept.
- ➢ From a practical point of view, clustering plays an extraordinary job in data mining applications. For example, scientific data exploration, text mining, information retrieval, spatial database applications, CRM, Web analysis, computational biology, medical diagnostics, and much more.

**In other words, we can say that Clustering analysis is a data mining technique to identify similar data. This technique helps to recognize the differences and similarities between the data. Clustering is very similar to the classification, but it involves grouping chunks of data together based on their similarities.**

## 3. Regression:

- ➢ Regression analysis is the data mining process used to identify and analyze the relationship between variables because of the presence of the other factor. It is used to define the probability of the specific variable.
- ➢ Regression, primarily a form of planning and modeling. For example, we might use it to project certain costs, depending on other factors such as availability, consumer demand, and competition. Primarily it gives the exact relationship between two or more variables in the given data set.

## 4. Association Rules:

- ➢ This data mining technique helps to discover a link between two or more items. It finds a hidden pattern in the data set.
- ➢ Association rules are if-then statements that support to show the probability of interactions between data items within large data sets in different types of databases. Association rule mining has several applications and is commonly used to help sales correlations in data or medical data sets.

### 5. Outer detection:

➢ This type of data mining technique relates to the observation of data items in the data set, which do not match an expected pattern or expected behavior.

➢ This technique may be used in various domains like intrusion, detection, fraud detection, etc. **It is also known as Outlier Analysis or Outlier mining.**

➢ The outlier is a data point that diverges too much from the rest of the dataset. The majority of the real-world datasets have an outlier.

➢ Outlier detection plays a significant role in the data mining field. Outlier detection is valuable in numerous fields like network interruption identification, credit or debit card fraud detection, detecting outlying in wireless sensor network data, etc.
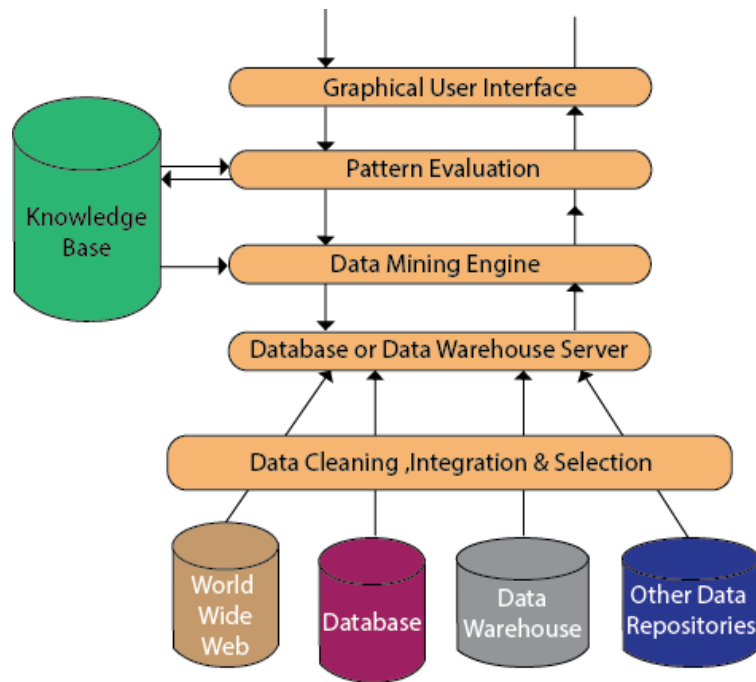
### 6. Sequential Patterns:

➢ The sequential pattern is a data mining technique specialized for **evaluating sequential data** to discover sequential patterns. It comprises of finding interesting subsequences in a set of sequences, where the stake of a sequence can be measured in terms of different criteria like length, occurrence frequency, etc.

➢ In other words, this technique of data mining helps to discover or recognize similar patterns in transaction data over some time.

### 7. Prediction:

➢ Prediction used a combination of other data mining techniques such as trends, clustering, classification, etc. It analyzes past events or instances in the right sequence to predict a future event.

## Data Mining Architecture

The significant components of data mining systems are a data source, data mining engine, data warehouse server, the pattern evaluation module, graphical user interface, and knowledge base.

## Data Source:

➢ The actual source of data is the Database, data warehouse, World Wide Web (WWW), text files, and other documents. You need a huge amount of historical data for data mining to be successful.

➢ Organizations typically store data in databases or data warehouses. Data warehouses may comprise one or more databases, text files spreadsheets, or other repositories of data.

➢ Sometimes, even plain text files or spreadsheets may contain information. Another primary source of data is the World Wide Web or the internet.

## Different processes:

➢ Before passing the data to the database or data warehouse server, the data must be cleaned, integrated, and selected. As the information comes from various sources and in different formats, it can't be used directly for the data mining procedure because the data may not be complete and accurate.

➢ So, the first data requires to be cleaned and unified. More information than needed will be collected from various data sources, and only the data of interest will have to be selected and passed to the server.

➢ These procedures are not as easy as we think. Several methods may be performed on the data as part of selection, integration, and cleaning.

## Database or Data Warehouse Server:

➢ The database or data warehouse server consists of the original data that is ready to be processed. Hence, the server is cause for retrieving the relevant data that is based on data mining as per user request.

### Data Mining Engine:

➢ The data mining engine is a major component of any data mining system. It contains several modules for operating data mining tasks, including association, characterization, classification, clustering, prediction, time-series analysis, etc.

➢ In other words, we can say data mining is the root of our data mining architecture. It comprises instruments and software used to obtain insights and knowledge from data collected from various data sources and stored within the data warehouse.

### Pattern Evaluation Module:

➢ The Pattern evaluation module is primarily responsible for the measure of investigation of the pattern by using a threshold value. It collaborates with the data mining engine to focus the search on exciting patterns.

➢ This segment commonly employs stake measures that cooperate with the data mining modules to focus the search towards fascinating patterns.

➢ It might utilize a stake threshold to filter out discovered patterns. On the other hand, the pattern evaluation module might be coordinated with the mining module, depending on the implementation of the data mining techniques used.

### Graphical User Interface:

➢ The graphical user interface (GUI) module communicates between the data mining system and the user.

➢ This module helps the user to easily and efficiently use the system without knowing the complexity of the process.

➢ This module cooperates with the data mining system when the user specifies a query or a task and displays the results.

### Knowledge Base:

➢ The knowledge base is helpful in the entire process of data mining. It might be helpful to guide the search or evaluate the stake of the result patterns.

➢ The knowledge base may even contain user views and data from user experiences that might be helpful in the data mining process.

➢ The data mining engine may receive inputs from the knowledge base to make the result more accurate and reliable. The pattern assessment module regularly interacts with the knowledge base to get inputs, and also update it.
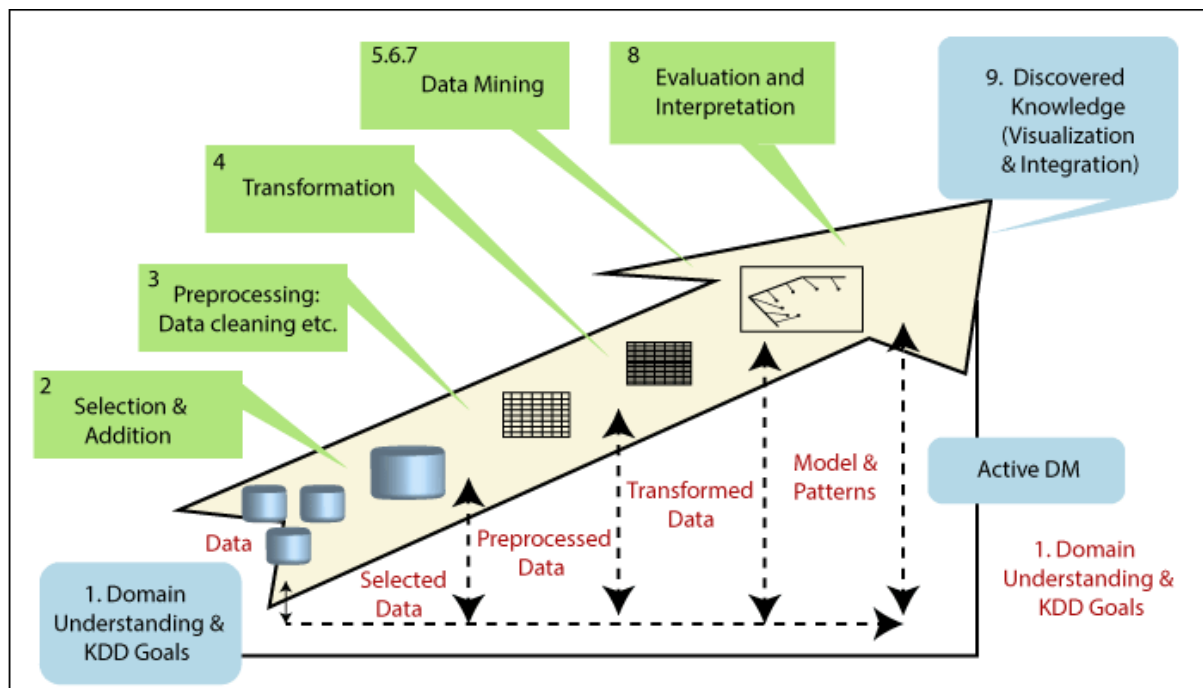
## KDD- Knowledge Discovery in Databases

➢ The term KDD stands for Knowledge Discovery in Databases. It refers to the broad procedure of discovering knowledge in data and emphasizes the high-level applications of specific Data Mining techniques.

- It is a field of interest to researchers in various fields, including artificial intelligence, machine learning, pattern recognition, databases, statistics, knowledge acquisition for expert systems, and data visualization.
- The main objective of the KDD process is to extract information from data in the context of large databases. It does this by using Data Mining algorithms to identify what is deemed knowledge.
- The Knowledge Discovery in Databases is considered as a programmed, exploratory analysis and modeling of vast data repositories.
- KDD is the organized procedure of recognizing valid, useful, and understandable patterns from huge and complex data sets.
- Data Mining is the root of the KDD procedure, including the inferring of algorithms that investigate the data, develop the model, and find previously unknown patterns.
- The model is used for extracting the knowledge from the data, analyze the data, and predict the data.

## The KDD Process

- The process begins with determining the KDD objectives and ends with the implementation of the discovered knowledge. At that point, the loop is closed, and the Active Data Mining starts.
- Subsequently, changes would need to be made in the application domain. For example, offering various features to cell phone users in order to reduce churn. This closes the loop, and the impacts are then measured on the new data repositories, and the KDD process again.
- Following is a concise description of the nine-step KDD process, Beginning with a managerial step:

**1. Building up an understanding of the application domain**

- This is the initial preliminary step. It develops the scene for understanding what should be done with the various decisions like transformation, algorithms, representation, etc.
- The individuals who are in charge of a KDD venture need to understand and characterize the objectives of the end-user and the environment in which the knowledge discovery process will occur (involves relevant prior knowledge).

**2. Choosing and creating a data set on which discovery will be performed**

- Once defined the objectives, the data that will be utilized for the knowledge discovery process should be determined.
- This incorporates discovering what data is accessible, obtaining important data, and afterward integrating all the data for knowledge discovery onto one set involves the qualities that will be considered for the process.
- This process is important because of Data Mining learns and discovers from the accessible data. This is the evidence base for building the models.
- If some significant attributes are missing, at that point, then the entire study may be unsuccessful from this respect, the more attributes are considered.
- On the other hand, to organize, collect, and operate advanced data repositories is expensive, and there is an arrangement with the opportunity for best understanding the phenomena.
- This arrangement refers to an aspect where the interactive and iterative aspect of the KDD is taking place. This begins with the best available data sets and later expands and observes the impact in terms of knowledge discovery and modeling.

**3. Preprocessing and cleansing**

- In this step, data reliability is improved. It incorporates data clearing, for example, Handling the missing quantities and removal of noise or outliers.
- It might include complex statistical techniques or use a Data Mining algorithm in this context. For example, when one suspects that a specific attribute of lacking reliability or has many missing data, at this point, this attribute could turn into the objective of the Data Mining supervised algorithm.
- A prediction model for these attributes will be created, and after that, missing data can be predicted.
- The expansion to which one pays attention to this level relies upon numerous factors. Regardless, studying the aspects is significant and regularly revealing by itself, to enterprise data frameworks.

**4. Data Transformation**

- In this stage, the creation of appropriate data for Data Mining is prepared and developed. Techniques here incorporate dimension reduction (for example, feature

selection and extraction and record sampling), also attribute transformation (for example, discretization of numerical attributes and functional transformation).

- This step can be essential for the success of the entire KDD project, and it is typically very project-specific. For example, in medical assessments, the quotient of attributes may often be the most significant factor and not each one by itself.

- In business, we may need to think about impacts beyond our control as well as efforts and transient issues. For example, studying the impact of advertising accumulation.

- However, if we do not utilize the right transformation at the starting, then we may acquire an amazing effect that insights to us about the transformation required in the next iteration.

- Thus, the KDD process follows upon itself and prompts an understanding of the transformation required.

## 5. Prediction and description

- We are now prepared to decide on which kind of Data Mining to use, for example, classification, regression, clustering, etc.

- This mainly relies on the KDD objectives, and also on the previous steps. There are two significant objectives in Data Mining, the first one is a prediction, and the second one is the description.

- Prediction is usually referred to as supervised Data Mining, while descriptive Data Mining incorporates the unsupervised and visualization aspects of Data Mining.

- Most Data Mining techniques depend on inductive learning, where a model is built explicitly or implicitly by generalizing from an adequate number of preparing models.

- The fundamental assumption of the inductive approach is that the prepared model applies to future cases.

## 6. Selecting the Data Mining algorithm

- Having the technique, we now decide on the strategies. This stage incorporates choosing a particular technique to be used for searching patterns that include multiple inducers.

- For example, considering precision versus understandability, the previous is better with neural networks, while the latter is better with decision trees.

- For each system of meta-learning, there are several possibilities of how it can be succeeded. Meta-learning focuses on clarifying what causes a Data Mining algorithm to be fruitful or not in a specific issue.

- Thus, this methodology attempts to understand the situation under which a Data Mining algorithm is most suitable.

- Each algorithm has parameters and strategies of leaning, such as ten folds cross-validation or another division for training and testing.

**7. Utilizing the Data Mining algorithm**

- At last, the implementation of the Data Mining algorithm is reached. In this stage, we may need to utilize the algorithm several times until a satisfying outcome is obtained. For example, by turning the algorithms control parameters, such as the minimum number of instances in a single leaf of a decision tree.

**8. Evaluation**

- In this step, we assess and interpret the mined patterns, rules, and reliability to the objective characterized in the first step.
- Here we consider the preprocessing steps as for their impact on the Data Mining algorithm results. For example, including a feature in step 4, and repeat from there.
- This step focuses on the comprehensibility and utility of the induced model. In this step, the identified knowledge is also recorded for further use. The last step is the use, and overall feedback and discovery results acquire by Data Mining.

**9. Using the discovered knowledge**

- Now, we are prepared to include the knowledge into another system for further activity. The knowledge becomes effective in the sense that we may make changes to the system and measure the impacts.
- The accomplishment of this step decides the effectiveness of the whole KDD process. There are numerous challenges in this step, such as losing the "laboratory conditions" under which we have worked.
- For example, the knowledge was discovered from a certain static depiction, it is usually a set of data, but now the data becomes dynamic.
- Data structures may change certain quantities that become unavailable, and the data domain might be modified, such as an attribute that may have a value that was not expected previously.

# Data Mining tools

- ➢ Data Mining is the set of techniques that utilize specific algorithms, statical analysis, artificial intelligence, and database systems to analyze data from different dimensions and perspectives.
- ➢ Data Mining tools have the objective of discovering patterns/trends/groupings among large sets of data and transforming data into more refined information.
- ➢ It is a framework, such as Rstudio or Tableau that allows you to perform different types of data mining analysis.
- ➢ We can perform various algorithms such as clustering or classification on data set and visualize the results itself. It is a framework that provides us better insights for our data and the phenomenon that data represent. Such a framework is called a data mining tool.
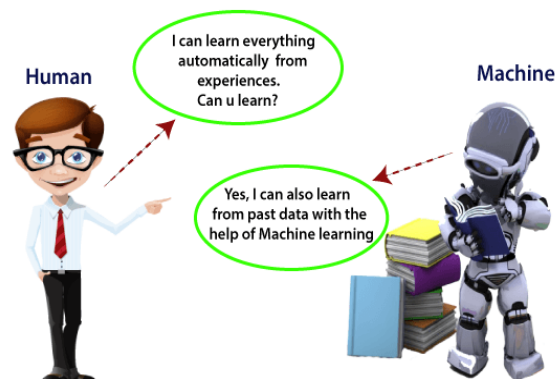
The Market for Data Mining tool is shining: as per the latest report from ReortLinker noted that the market would top **$1 billion** in sales by **2023**, up from **$ 591** million in **2018**

These are the most popular data mining tools:



## Machine Learning

- ➢ Machine learning is a growing technology which enables computers to learn automatically from past data.
- ➢ Machine learning uses various algorithms for **building mathematical models and making predictions using historical data or information**.
- ➢ Currently, it is being used for various tasks such as **image recognition**, **speech recognition**, **email filtering**, **Facebook auto-tagging**, **recommender system**, and many more.
- ➢ In the real world, we are surrounded by humans who can learn everything from their experiences with their learning capability, and we have computers or machines which work on our instructions. But can a machine also learn from experiences or past data like a human does? So here comes the role of **Machine Learning**.

- Machine Learning is said as a subset of **artificial intelligence** that is mainly concerned with the development of algorithms which allow a computer to learn from the data and past experiences on their own.
- The term machine learning was first introduced by **Arthur Samuel** in **1959**. We can define it in a summarized way as:
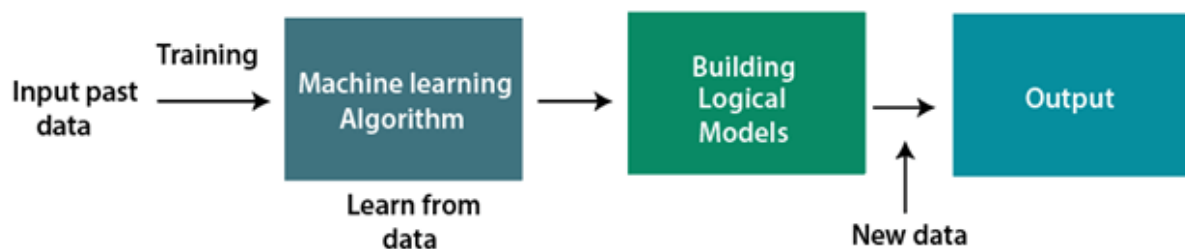
Machine learning enables a machine to automatically learn from data, improve performance from experiences, and predict things without being explicitly programmed.

- With the help of sample historical data, which is known as **training data**, machine learning algorithms build a **mathematical model** that helps in making predictions or decisions without being explicitly programmed.
- Machine learning brings computer science and statistics together for creating predictive models. Machine learning constructs or uses the algorithms that learn from historical data.

**A machine has the ability to learn if it can improve its performance by gaining more data.**

# How does Machine Learning work

- A Machine Learning system **learns from historical data, builds the prediction models, and whenever it receives new data, predicts the output for it**. The accuracy of predicted output depends upon the amount of data, as the huge amount of data helps to build a better model which predicts the output more accurately.
- Suppose we have a complex problem, where we need to perform some predictions, so instead of writing a code for it, we just need to feed the data to generic algorithms, and with the help of these algorithms, machine builds the logic as per the data and predict the output.
- Machine learning has changed our way of thinking about the problem. The below block diagram explains the working of Machine Learning algorithm:



**Features of Machine Learning:**

- Machine learning uses data to detect various patterns in a given dataset.
- It can learn from past data and improve automatically.
- It is a data-driven technology.

o Machine learning is much similar to data mining as it also deals with the huge amount of the data.
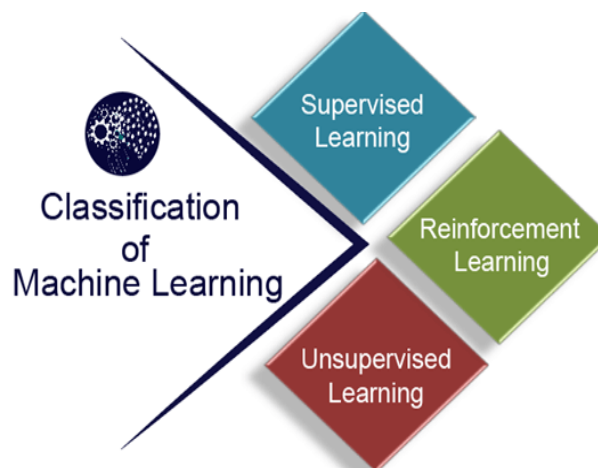
### Need for Machine Learning

➢ The need for machine learning is increasing day by day. The reason behind the need for machine learning is that it is capable of doing tasks that are too complex for a person to implement directly.

➢ As a human, we have some limitations as we cannot access the huge amount of data manually, so for this, we need some computer systems and here comes the machine learning to make things easy for us.

➢ We can train machine learning algorithms by providing them the huge amount of data and let them explore the data, construct the models, and predict the required output automatically.

➢ The performance of the machine learning algorithm depends on the amount of data, and it can be determined by the cost function. With the help of machine learning, we can save both time and money.

➢ The importance of machine learning can be easily understood by its uses cases, Currently, machine learning is used in **self-driving cars**, **cyber fraud detection**, **face recognition**, and **friend suggestion by Facebook**, etc.

➢ Various top companies such as Netflix and Amazon have build machine learning models that are using a vast amount of data to analyze the user interest and recommend product accordingly.

# Classification of Machine Learning

At a broad level, machine learning can be classified into three types:

1. **Supervised learning**
2. **Unsupervised learning**
3. **Reinforcement learning**

## 1) Supervised Learning

➢ Supervised learning is a type of machine learning method in which we provide sample labeled data to the machine learning system in order to train it, and on that basis, it predicts the output.

➢ The system creates a model using labeled data to understand the datasets and learn about each data, once the training and processing are done then we test the model by providing a sample data to check whether it is predicting the exact output or not.

➢ The goal of supervised learning is to map input data with the output data. The supervised learning is based on supervision, and it is the same as when a student learns things in the supervision of the teacher.

➢ The example of supervised learning is **spam filtering**.

Supervised learning can be grouped further in two categories of algorithms:

- o **Classification**
- o **Regression**

## 2) Unsupervised Learning

➢ Unsupervised learning is a learning method in which a machine learns without any supervision.

➢ The training is provided to the machine with the set of data that has not been labeled, classified, or categorized, and the algorithm needs to act on that data without any supervision.

➢ The goal of unsupervised learning is to restructure the input data into new features or a group of objects with similar patterns.

➢ In unsupervised learning, we don't have a predetermined result. The machine tries to find useful insights from the huge amount of data.

It can be further classifieds into two categories of algorithms:

- o **Clustering**
- o **Association**

## 3) Reinforcement Learning

➢ Reinforcement learning is a feedback-based learning method, in which a learning agent gets a reward for each right action and gets a penalty for each wrong action.

➢ The agent learns automatically with these feedbacks and improves its performance. In reinforcement learning, the agent interacts with the environment and explores it.

➢ The goal of an agent is to get the most reward points, and hence, it improves its performance.

# Difference between supervised and unsupervised learning

| Supervised Learning | Unsupervised Learning |
|---|---|
| Supervised learning algorithms are trained using labeled data. | Unsupervised learning algorithms are trained using unlabeled data. |
| Supervised learning model takes direct feedback to check if it is predicting correct output or not. | Unsupervised learning model does not take any feedback. |
| Supervised learning model predicts the output. | Unsupervised learning model finds the hidden patterns in data. |
| In supervised learning, input data is provided to the model along with the output. | In unsupervised learning, only input data is provided to the model. |
| The goal of supervised learning is to train the model so that it can predict the output when it is given new data. | The goal of unsupervised learning is to find the hidden patterns and useful insights from the unknown dataset. |
| Supervised learning needs supervision to train the model. | Unsupervised learning does not need any supervision to train the model. |
| Supervised learning can be categorized in **Classification** and **Regression** problems. | Unsupervised Learning can be classified in **Clustering** and **Associations** problems. |
| Supervised learning can be used for those cases where we know the input as well as corresponding outputs. | Unsupervised learning can be used for those cases where we have only input data and no corresponding output data. |
| Supervised learning model produces an accurate result. | Unsupervised learning model may give less accurate result as compared to supervised learning. |
| Supervised learning is not close to true Artificial intelligence as in this, we first train the model for each data, and then only it can predict the correct output. | Unsupervised learning is more close to the true Artificial Intelligence as it learns similarly as a child learns daily routine things by his experiences. |
| It includes various algorithms such as Linear Regression, Logistic Regression, Support Vector Machine, Multi-class Classification, Decision tree, Bayesian Logic, etc. | It includes various algorithms such as Clustering, KNN, and Apriori algorithm. |

# Difference between Classification and Clustering

| Classification | Clustering |
| --- | --- |
| Classification is a supervised learning approach where a specific label is provided to the machine to classify new observations. Here the machine needs proper testing and training for the label verification. | Clustering is an unsupervised learning approach where grouping is done on similarities basis. |
| Supervised learning approach. | Unsupervised learning approach. |
| It uses a training dataset. | It does not use a training dataset. |
| It uses algorithms to categorize the new data as per the observations of the training set. | It uses statistical concepts in which the data set is divided into subsets with the same features. |
| In classification, there are labels for training data. | In clustering, there are no labels for training data. |
| Its objective is to find which class a new object belongs to form the set of predefined classes. | Its objective is to group a set of objects to find whether there is any relationship between them. |
| It is more complex as compared to clustering. | It is less complex as compared to clustering. |