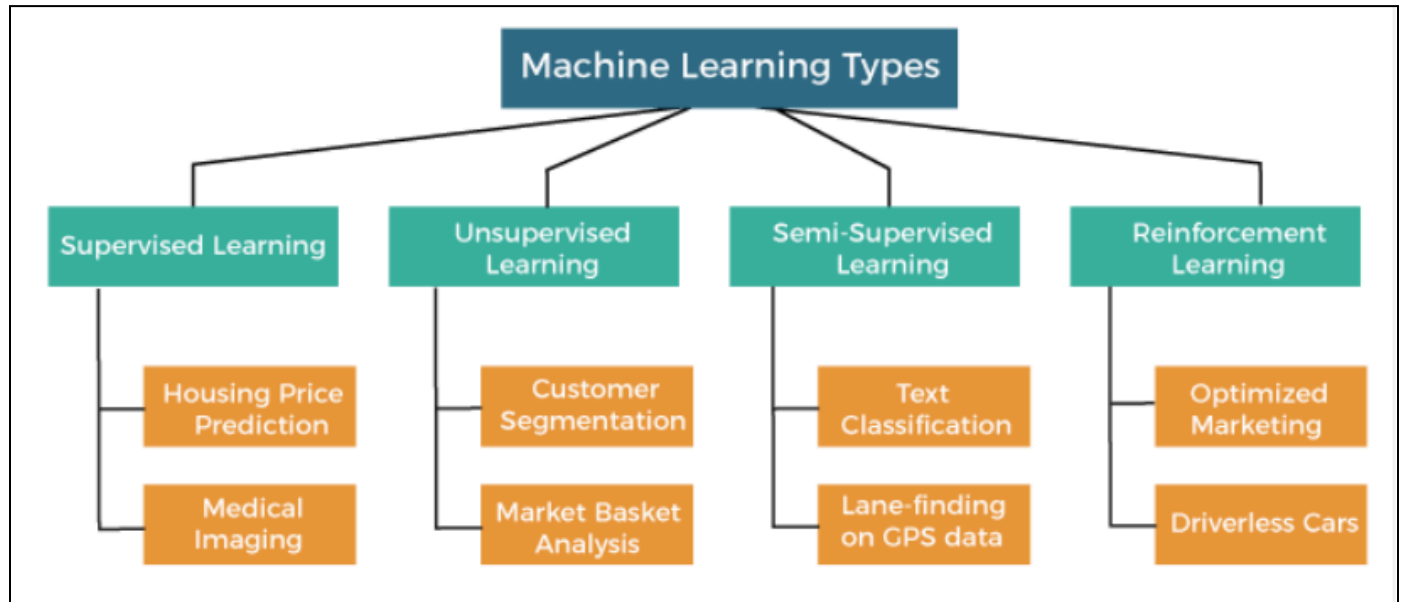


## UNIT - IV TYPES OF LEARNING

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.



## 4.1 TYPES OF LEARNING: SUPERVISED, UNSUPERVISED, SEMI-SUPERVISED LEARNING

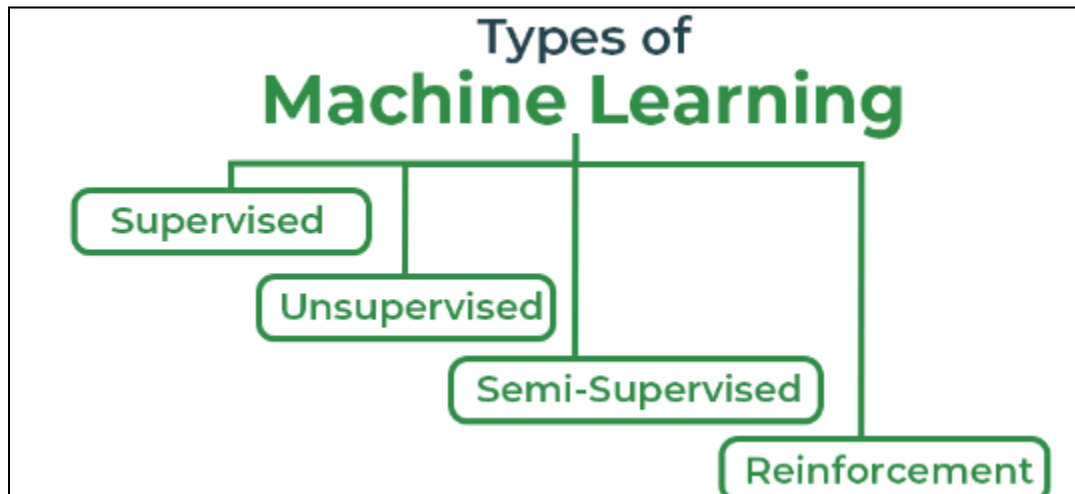
---

### 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 Machine Learning
2. Unsupervised Machine Learning
3. Reinforcement Learning

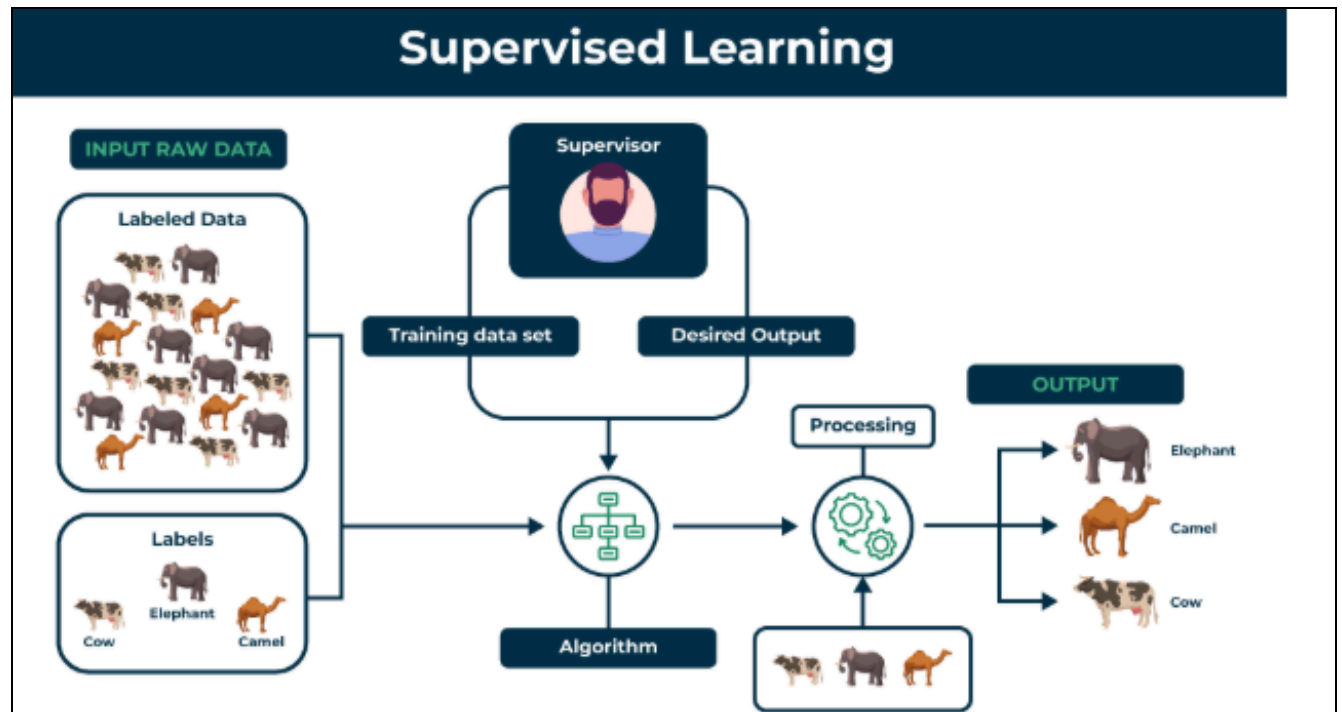
Additionally, there is a more specific category called semi-supervised learning, which combines elements of both supervised and unsupervised learning.



### Types of Machine Learning

#### 1. Supervised Machine 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.



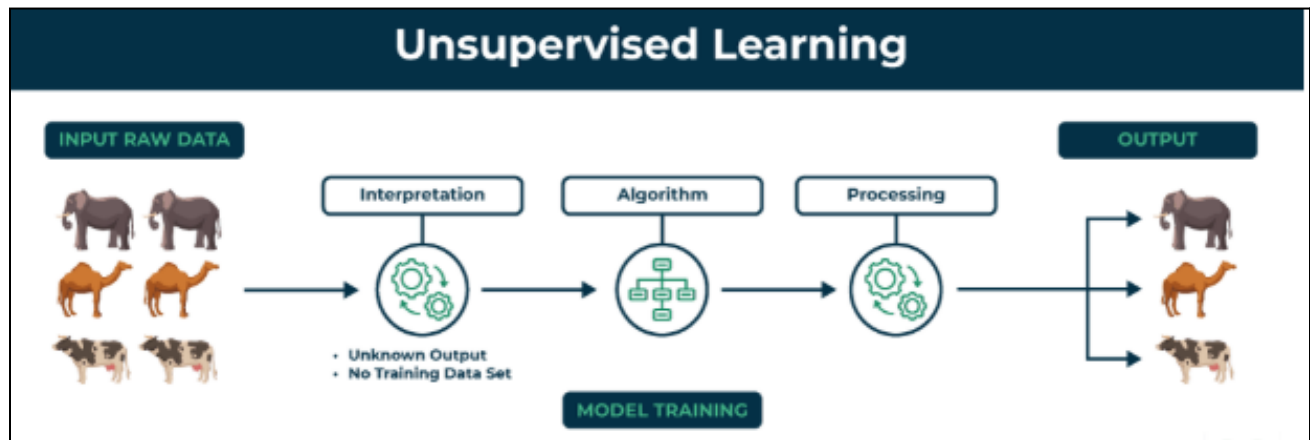
## *Supervised Learning*

Let's understand it with the help of an example.

Example: Consider a scenario where you have to build an image classifier to differentiate between cats and dogs. If you feed the datasets of dogs and cats labelled images to the algorithm, the machine will learn to classify between a dog or a cat from these labeled images. When we input new dog or cat images that it has never seen before, it will use the learned algorithms and predict whether it is a dog or a cat. This is how supervised learning works, and this is particularly an image classification.

## 2. Unsupervised Machine 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.



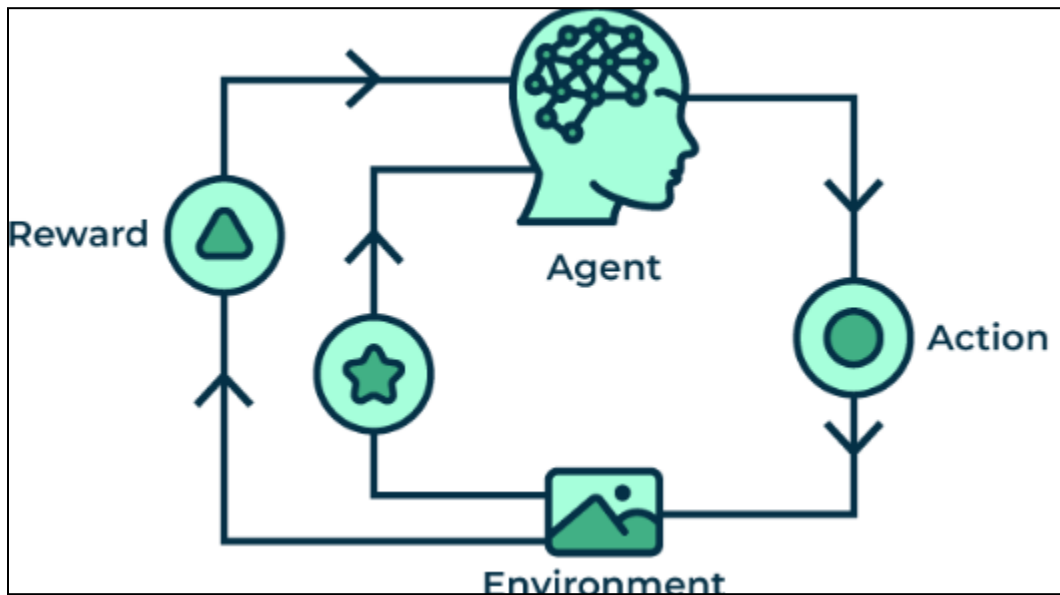
*Unsupervised Learning*

Let's understand it with the help of an example.

**Example:** Consider that you have a dataset that contains information about the purchases you made from the shop. Through clustering, the algorithm can group the same purchasing behavior among you and other customers, which reveals potential customers without predefined labels. This type of information can help businesses get target customers as well as identify outliers.

### 3. Reinforcement Machine Learning

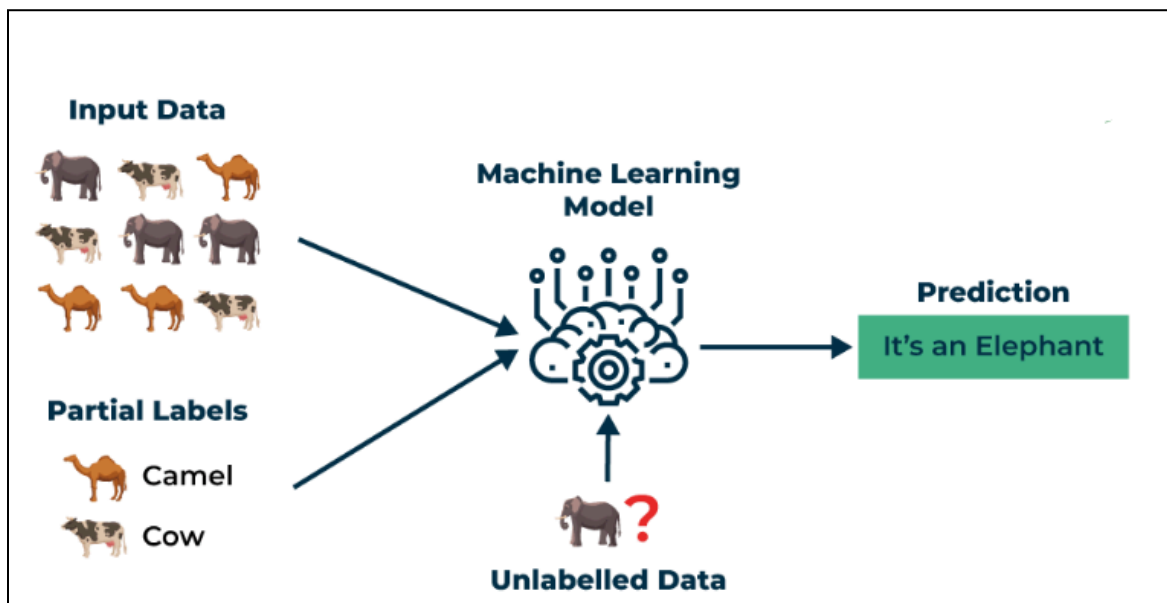
Reinforcement machine learning algorithm is a learning method that interacts with the environment by producing actions and discovering errors. **Trial, error, and delay** are the most relevant characteristics of reinforcement learning. In this technique, the model keeps on increasing its performance 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.



Example: Consider that you are training an AI agent to play a game like chess. The agent explores different moves and receives positive or negative feedback based on the outcome. Reinforcement Learning also finds applications in which they learn to perform tasks by interacting with their surroundings.

### 3. Semi-Supervised Learning: Supervised + Unsupervised Learning

Semi-Supervised learning is a machine learning algorithm that works between the supervised and unsupervised learning so it uses both **labelled** and **unlabelled** data. It's particularly useful when obtaining labeled data is costly, time-consuming, or resource-intensive. This approach is useful when the dataset is expensive and time-consuming. Semi-supervised learning is chosen when labeled data requires skills and relevant resources in order to train or learn from it.



**Example:** Consider that we are building a language translation model, having labeled translations for every sentence pair can be resource intensive. It allows the models to learn from labeled and unlabeled sentence pairs, making them more accurate. This technique has led to significant improvements in the quality of machine translation services.

## 4.2.1 SUPERVISED LEARNING: LEARNING A CLASS FROM EXAMPLES

---

Step-by-Step Process of Supervised Learning

### **1. Collect Labeled Data:**

**Input:** A dataset of input items (e.g., images of vegetables).

**Label:** Each item is associated with a label (e.g., "Carrot", "Tomato", "Bell Pepper").

In the image:

- \* The dataset includes various vegetables.
- \* Labels are provided for each image.

### **2. Train the Model:**

- \* The labeled data is fed into a machine learning model.
- \* The model learns the patterns in the data and how to associate features with labels.
- \* This stage is shown in the image as “**Model Training**”

### **3. Use Test Data:**

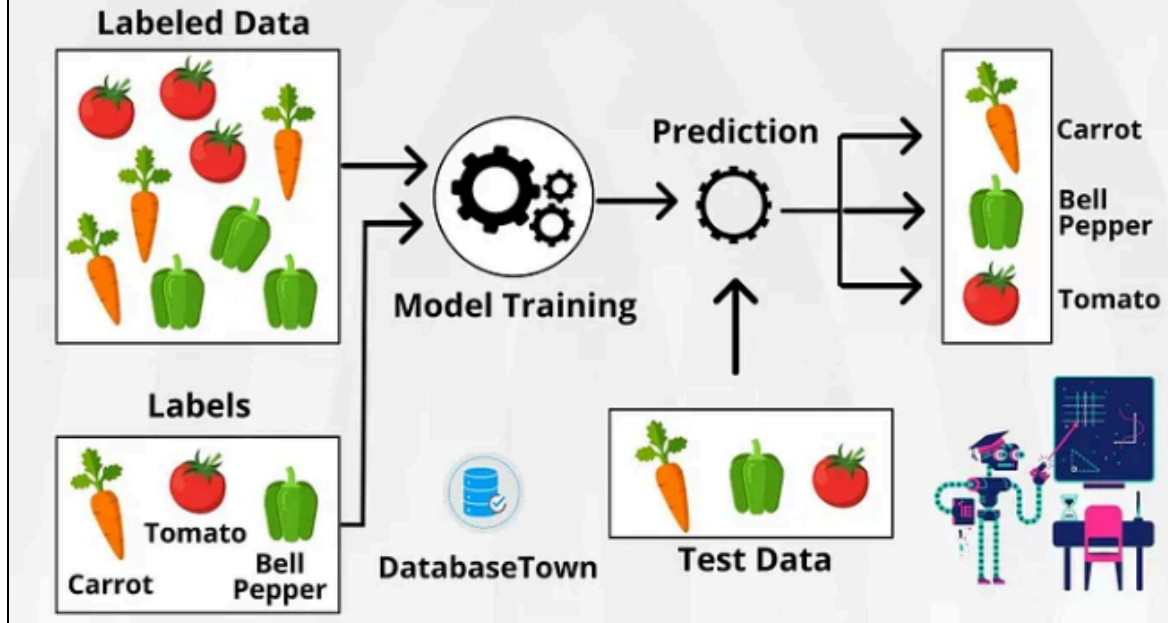
- \* Separate unlabeled or unseen data (test data) is used to evaluate the model's performance.
- \* This test data is not used during training.
- \* In the image: The model is given new vegetable images without labels.

### **4. Make Predictions:**

- \* The trained model predicts labels for the new test data.
- \* The goal is to see how well the model can generalize to new, unseen data.
- \* In the image: The model identifies the test vegetables as "Carrot", "Bell Pepper", and "Tomato".

# SUPERVISED LEARNING

Supervised machine learning is a branch of artificial intelligence that focuses on training models to make predictions or decisions based on labeled training data.



## 5. Evaluate the Results:

- \* Compare the model's predictions with the actual labels (if available).
- \* Assess the model's accuracy and performance.
- \* In the image, this is illustrated by the robot analyzing the results.



Step	Description	Image Reference
1	Collect Labeled Data	Labeled vegetables with names
2	Train the Model	“Model Training” gear icon
3	Test the Model	New vegetables (Test Data)
4	Predict Labels	Prediction arrows and output labels
5	Evaluate Model Performance	Robot reviewing results

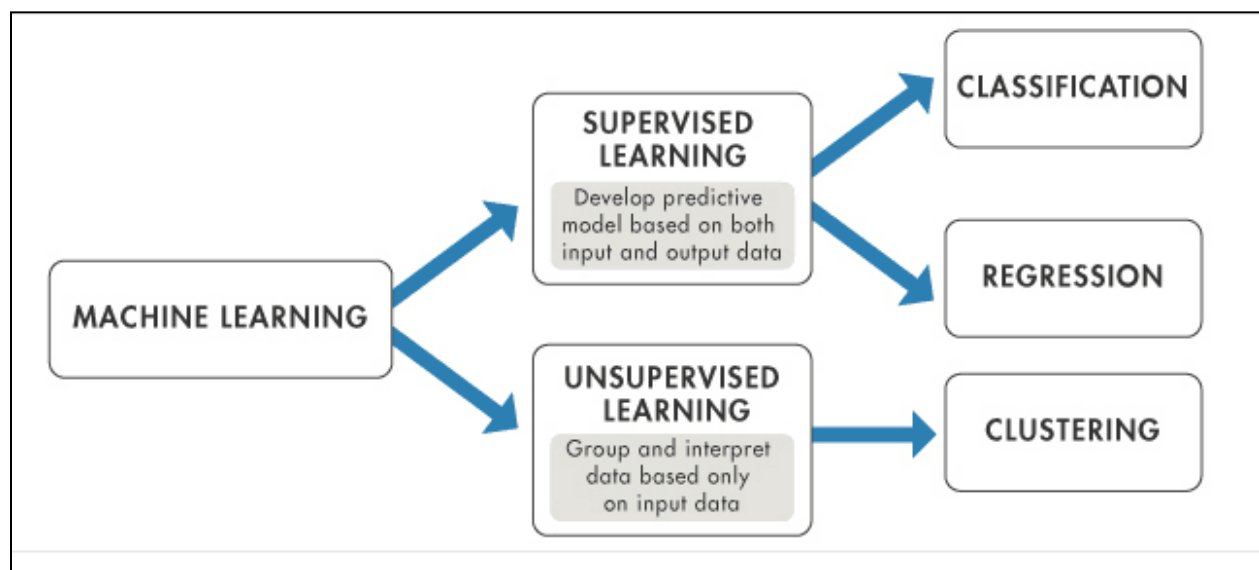
## 4.2.2 INTRODUCTION OF DIFFERENT TYPES OF SUPERVISED MACHINE LEARNING ALGORITHMS:

---

Supervised machine learning learns patterns and relationships between input and output data. It is defined by its use of labeled data. A labeled data is a dataset that contains a lot of examples of Features and Target. Supervised learning uses algorithms that learn the relationship of Features and Target from the dataset. This process is referred to as Training or Fitting.

There are two types of supervised learning algorithms:

1. Classification
2. Regression



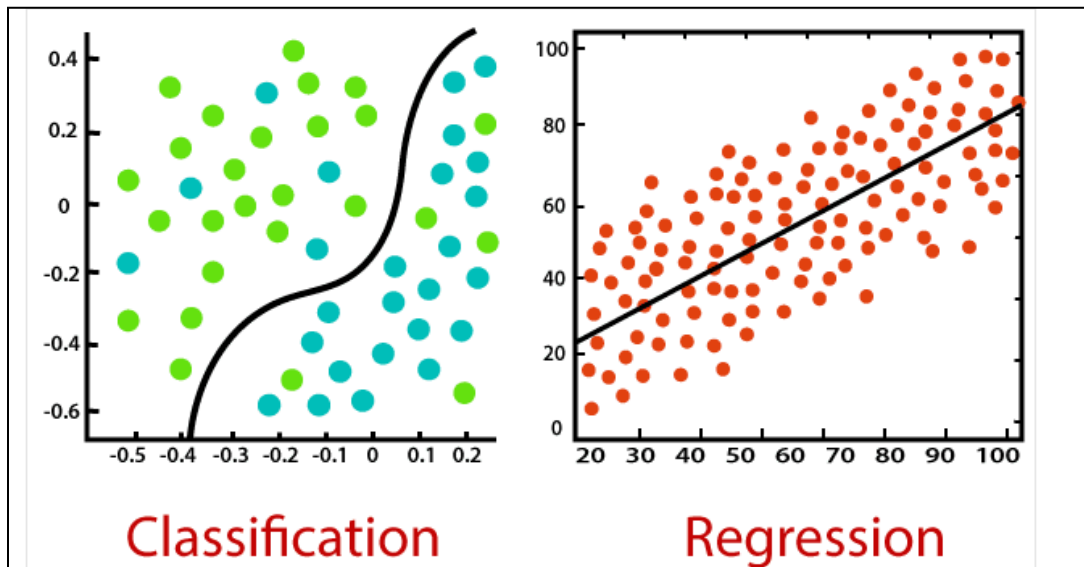
### 1. CLASSIFICATION:

**Classification** is a type of supervised machine learning where algorithms learn from the data to predict an outcome or event in the future. For example:

A bank may have a customer dataset containing credit history, loans, investment details, etc. and they may want to know if any customer will default. In the historical data, we will have Features and Target.

- Features will be attributes of a customer such as credit history, loans, investments, etc.
- Target will represent whether a particular customer has defaulted in the past (normally represented by 1 or 0 / True or False / Yes or No).

Classification algorithms are used for predicting discrete outcomes, if the outcome can take two possible values such as True or False, Default or No Default, Yes or No, it is known as Binary Classification. When the outcome contains more than two possible values, it is known as Multiclass Classification.



## 2. REGRESSION:

Regression is a type of supervised machine learning where algorithms learn from the data to predict continuous values such as sales, salary, weight, or temperature. For example:

A dataset containing features of the house such as lot size, number of bedrooms, number of baths, neighborhood, etc. and the price of the house, a Regression algorithm can be trained to learn the relationship between the features and the price of the house.

Category	Algorithms	Description
Regression	Linear Regression	Predicts a continuous numerical value based on input variables.
	Decision Tree (Regression)	Uses tree-based splits to predict continuous values.
	Random Forest (Regression)	An ensemble of multiple regression trees to improve accuracy and reduce overfitting.
Classification	Logistic Regression	Predicts categorical outcomes (binary or multi-class) based on probability.
	Decision Tree (Classification)	Uses tree-based splits to classify data into categories.
	K-Nearest Neighbors (K-NN)	Classifies based on the closest data points in feature space.
	Random Forest (Classification)	An ensemble of multiple classification trees for better accuracy.

### 4.2.3 LINEAR REGRESSION, LOGISTIC REGRESSION, DECISION TREE, K-NEAREST NEIGHBORS, RANDOM FOREST

---

#### 1. LINEAR REGRESSION

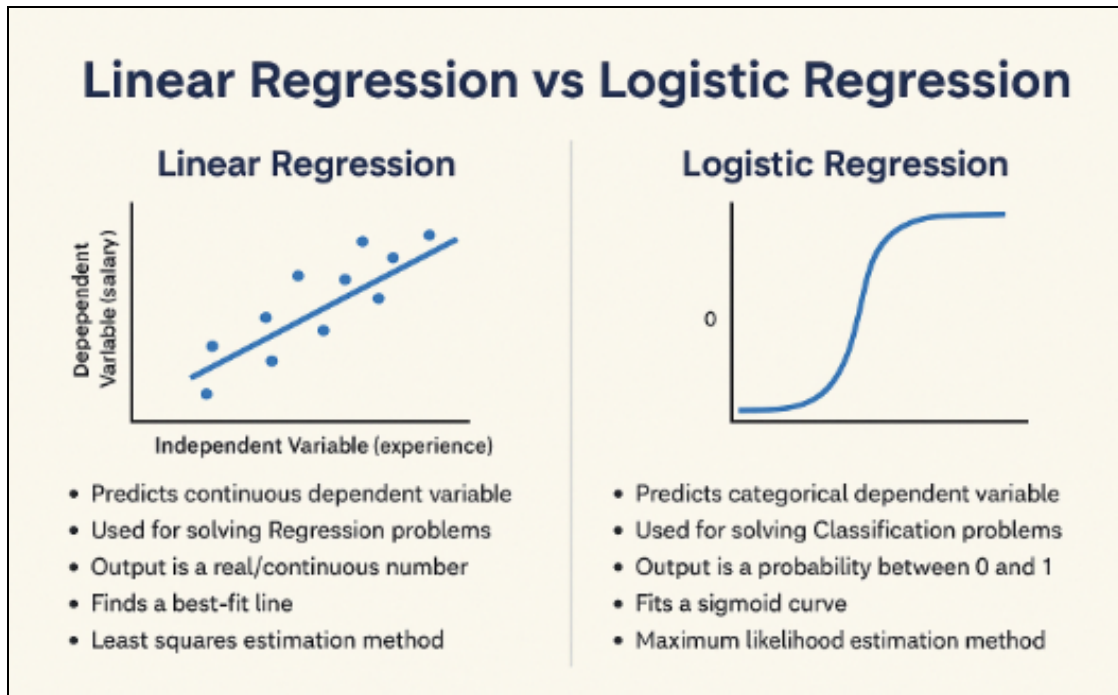
Linear regression is a supervised learning algorithm used to predict continuous values by modeling the relationship between a dependent (target) variable and one or more independent (predictor) variables. The objective of linear regression is to determine a linear relationship between the variables, which can then be used to predict outcomes for unseen data.

#### **Dependent and Independent Variables:**

The dependent variable is the value the model aims to predict, such as house prices or customer spending. On the other hand, independent variables are the features used to make the prediction, such as the size of a house, the number of rooms, or marketing spend. The model identifies how these independent variables influence the dependent variable.

#### **Example: Real-Life Applications of Linear Regression**

1. **Predicting House Prices:** A real estate company can use linear regression to predict the price of a house based on its size, number of bedrooms, and location.
2. **Sales Forecasting:** Companies can forecast future sales based on factors like advertising spend and seasonal trends, helping them plan inventory and campaigns.
3. **Medical Research:** In healthcare, linear regression helps predict patient outcomes (like blood pressure) based on variables such as age, weight, and lifestyle factors.



### What is Logistic Regression?

Logistic regression is a supervised learning algorithm used to estimate the probability that a given instance belongs to a particular class. It works by modeling the relationship between a categorical dependent variable and one or more independent variables. Logistic regression is particularly useful for binary outcomes, such as classifying whether an email is spam or not.

Unlike linear regression, which predicts continuous numerical values, logistic regression outputs probabilities between 0 and 1. It achieves this by passing the linear combination of input features through a sigmoid function, ensuring that predictions remain within a valid probability range.

### Real-life Applications of Logistic Regression in Machine Learning:

1. **Spam Detection:** Logistic regression classifies emails as spam or non-spam based on features like subject line and sender.
2. **Medical Diagnosis:** It predicts the presence or absence of a disease based on patient data.

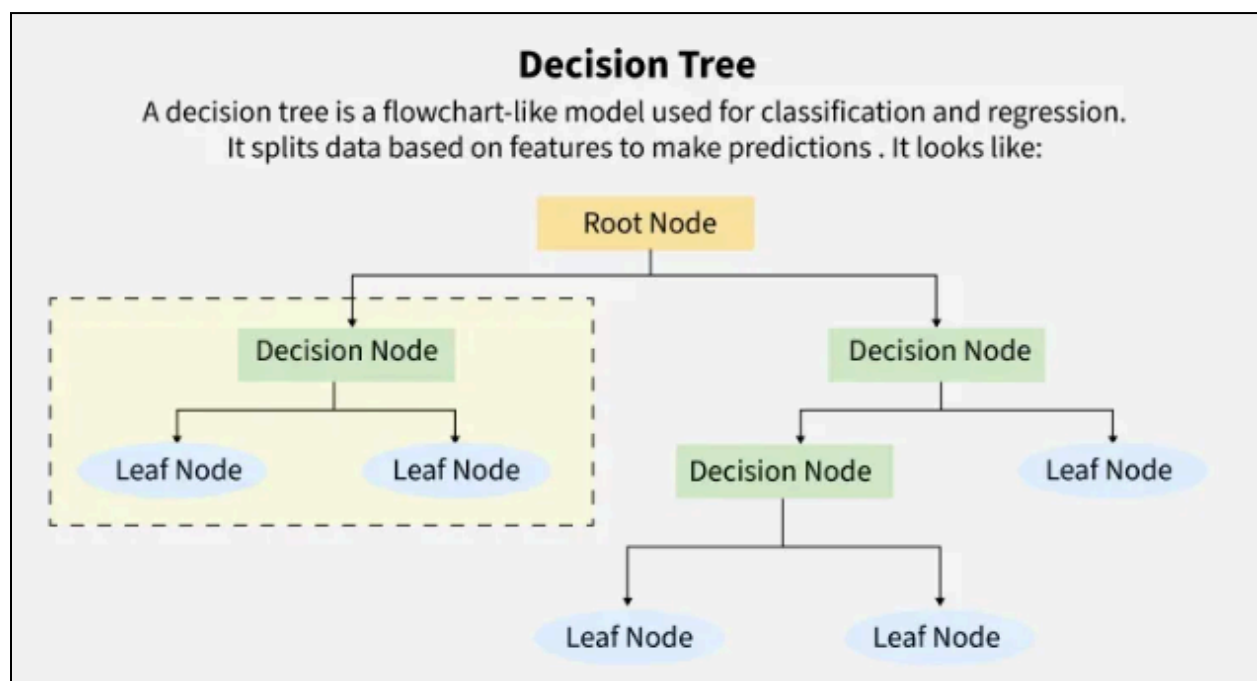
3. **Churn Prediction:** Logistic regression helps businesses identify customers likely to cancel their subscriptions.
4. **Loan Default Prediction:** Determines whether a loan applicant is likely to default.

## 2. DECISION TREE:

A decision tree is a supervised learning algorithm used for both classification and regression tasks. It has a hierarchical tree structure which consists of a root node, branches, internal nodes and leaf nodes. It works like a flowchart help to make decisions step by step where:

- Internal nodes represent attribute tests
- Branches represent attribute values
- Leaf nodes represent final decisions or predictions.

Decision trees are widely used due to their interpretability, flexibility and low preprocessing needs.



## How Does a Decision Tree Work?

A decision tree splits the dataset based on feature values to create pure subsets, ideally all items in a group belong to the same class. Each leaf node of the tree corresponds to a class label and the internal nodes are feature-based decision points. Let's understand this with an example.

Let's consider a decision tree for predicting whether a customer will buy a product based on age, income and previous purchases: Here's how the decision tree works:

### 1. Root Node (Income)

**First Question: "Is the person's income greater than \$50,000?"**

- If Yes, proceed to the next question.
- If No, predict "No Purchase" (leaf node).

### 2. Internal Node (Age):

**If the person's income is greater than \$50,000, ask: "Is the person's age above 30?"**

- If Yes, proceed to the next question.
- If No, predict "No Purchase" (leaf node).

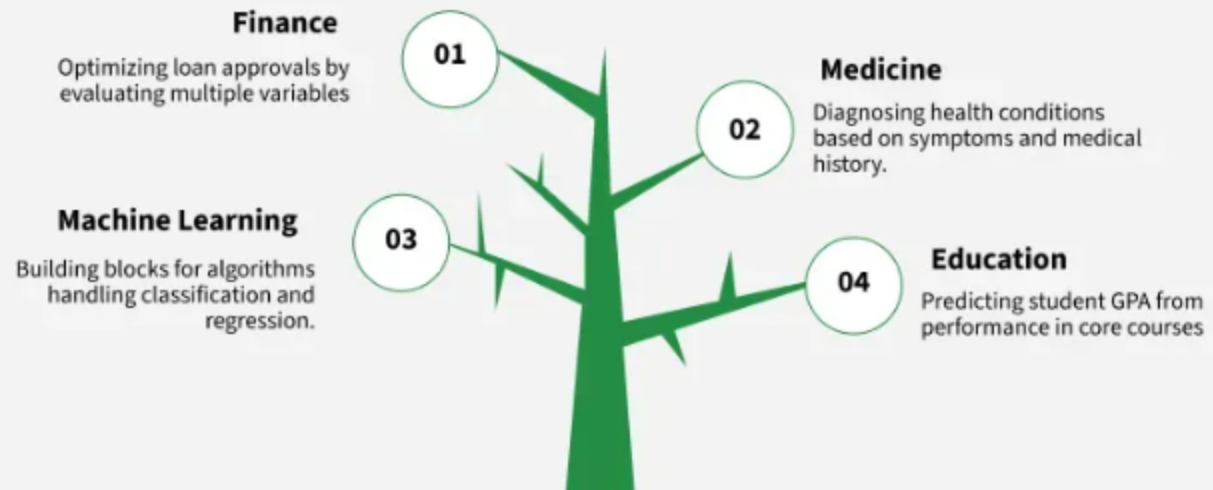
### 3. Internal Node (Previous Purchases):

- If the person is above 30 and has made previous purchases, predict "Purchase" (leaf node).
- If the person is above 30 and has not made previous purchases, predict "No Purchase" (leaf node).



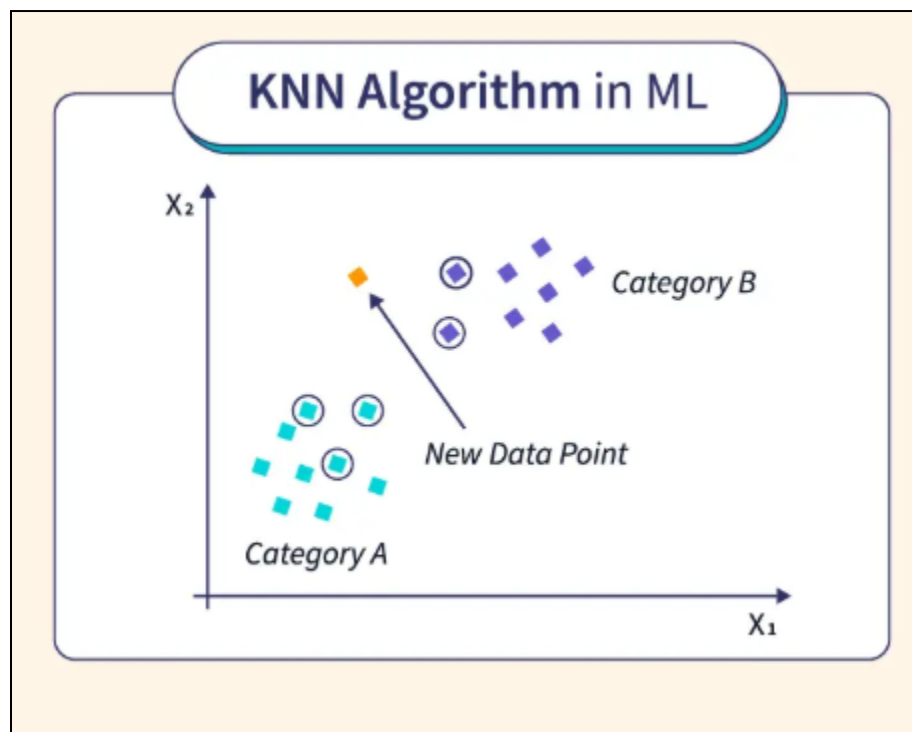
## Applications of Decision Trees

A decision tree is a flowchart-like model used for classification and regression. It splits data based on features to make predictions . It looks like:



### 3. K-Nearest Neighbors (KNN):

K-Nearest Neighbors (KNN) is a supervised machine learning algorithm generally used for classification but can also be used for regression tasks. It works by finding the "k" closest data points (neighbors) to a given input and makes predictions based on the majority class (for classification) or the average value (for regression). Since KNN makes no assumptions about the underlying data distribution it makes it a non-parametric and instance-based learning method.



In **classification**, KNN assigns a class label to a new data point based on the majority class of its nearest neighbors. For instance, if a data point has five nearest neighbors, and three of them belong to class A while two belong to class B, the algorithm will classify the point as class A.

In **regression**, KNN predicts continuous values by averaging the values of the k-nearest neighbors. For example, if you're predicting house prices, KNN will use the average prices of the k-nearest neighbors to estimate the price of a new house.

### Step 1: Determine the Number of Nearest Neighbors (k)

The first step is to select the number of neighbors (**k**) to consider. The value of **k** determines how many neighboring points will influence the classification or prediction of a new data point.

### Step 2: Calculate the Distance Between the Query Point and Dataset Points

For each data point in the dataset, the algorithm calculates the distance between the query point (the new point to be classified or predicted) and every other point. Various distance metrics can be used, such as **Euclidean distance**, **Manhattan distance**, or **Minkowski distance**.

$$d(P, Q) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

### Step 3: Sort and Select the k-Nearest Neighbors

After calculating the distances, the algorithm sorts all data points in ascending order of distance. It then selects the **k-nearest neighbors**—the data points that are closest to the query point.

### Step 4: Make a Prediction

- **For classification:** The algorithm assigns the query point to the class label that is most frequent among the k-nearest neighbors (majority voting).
- **For regression:** The algorithm predicts the value by averaging the values of the k-nearest neighbors.

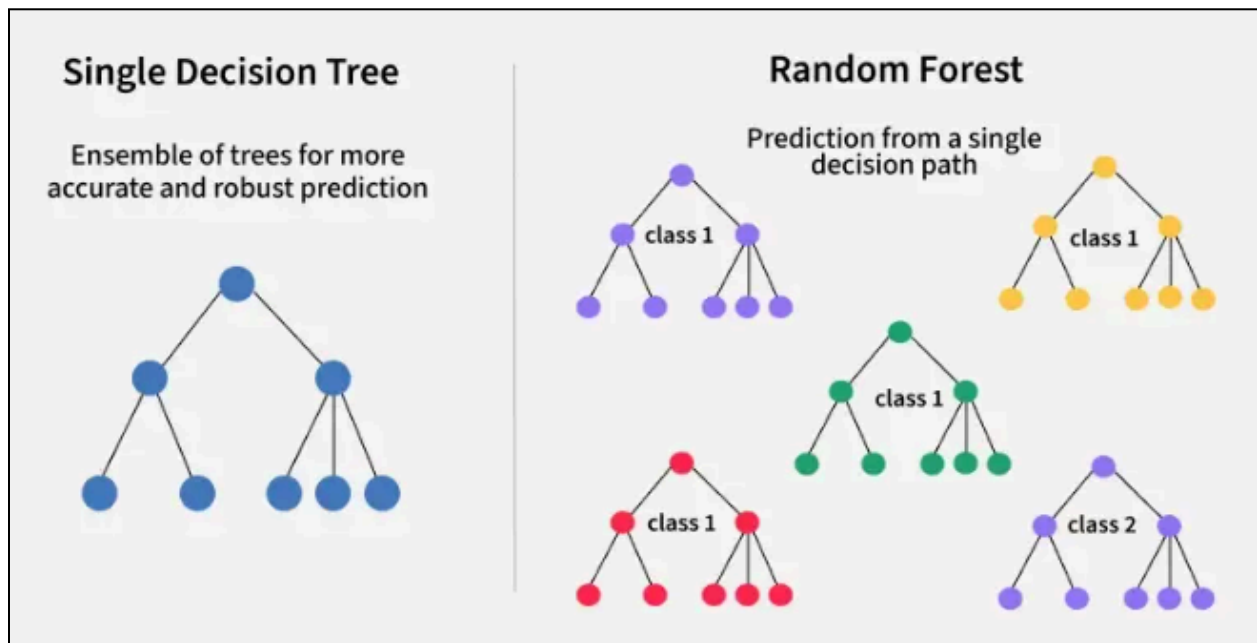
## APPLICATIONS OF KNN ALGORITHM IN MACHINE LEARNING

1. **Text Classification** – Classifies documents (e.g., spam detection) by comparing keyword patterns with labeled data.
2. **Image Recognition** – Identifies objects or faces by analyzing pixel similarity.

3. **Recommendation Systems** – Suggests items based on preferences of similar users.
4. **Medical Diagnosis** – Predicts diseases or treatments by comparing patient data with similar cases.

#### 4. RANDOM FOREST:

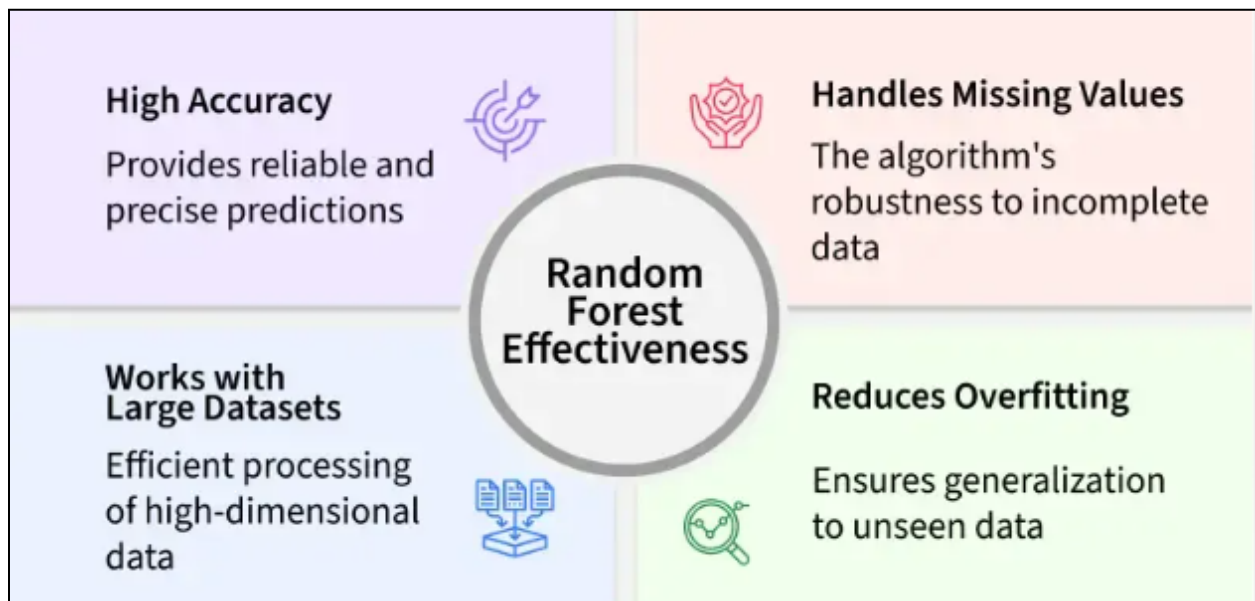
Random Forest is a machine learning algorithm that uses many decision trees to make better predictions. Each tree looks at different random parts of the data and their results are combined by voting for classification or averaging for regression. This helps in improving accuracy and reducing errors.



##### Working of Random Forest Algorithm:

- **Create Many Decision Trees:** The algorithm makes many decision trees each using a random part of the data. So every tree is a bit different.
- **Pick Random Features:** When building each tree it doesn't look at all the features (columns) at once. It picks a few at random to decide how to split the data. This helps the trees stay different from each other.

- Each Tree Makes a Prediction: Every tree gives its own answer or prediction based on what it learned from its part of the data.
- Combine the Predictions:
  - For classification we choose a category as the final answer is the one that most trees agree on i.e majority voting.
  - For regression we predict a number as the final answer is the average of all the tree predictions.
- Why It Works Well: Using random data and features for each tree helps avoid overfitting and makes the overall prediction more accurate and trustworthy.



### Applications of Random Forest

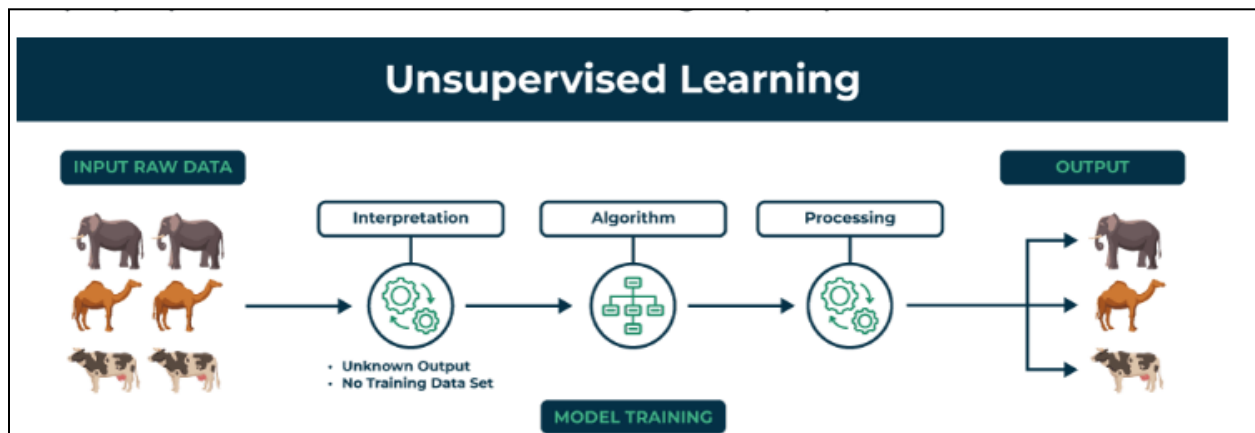
1. **Credit Card Fraud Detection:** Random Forest has been used to detect fraudulent credit card transactions. It can identify outliers and detect fraudulent activities by analyzing a variety of features such as amount, type of transaction, location etc.
2. **Stock Market Prediction:** Random Forest has been used to predict the stock market prices. It can identify patterns in the data and make predictions about future prices.

3. **Image Recognition:** Random Forest has been used in image recognition. It can be used to classify objects in images and identify patterns in the data.
4. **Customer Churn Prediction:** Random Forest has been used to predict customer churn. It can analyze customer data and predict the likelihood of a customer to leave the company.
5. **Diabetes Prediction:** Random Forest has been used to predict the onset of diabetes. It can analyze a variety of features such as age, lifestyle, diet and medical history to predict the risk of developing diabetes.

### 4.3. UNSUPERVISED LEARNING:

---

Unsupervised learning is a branch of machine learning that deals with unlabeled data. Unlike supervised learning, where the data is labeled with a specific category or outcome, unsupervised learning algorithms are tasked with finding patterns and relationships within the data without any prior knowledge of the data's meaning. Unsupervised machine learning algorithms find hidden patterns and data without any human intervention, i.e., we don't give output to our model. The training model has only input parameter values and discovers the groups or patterns on its own.



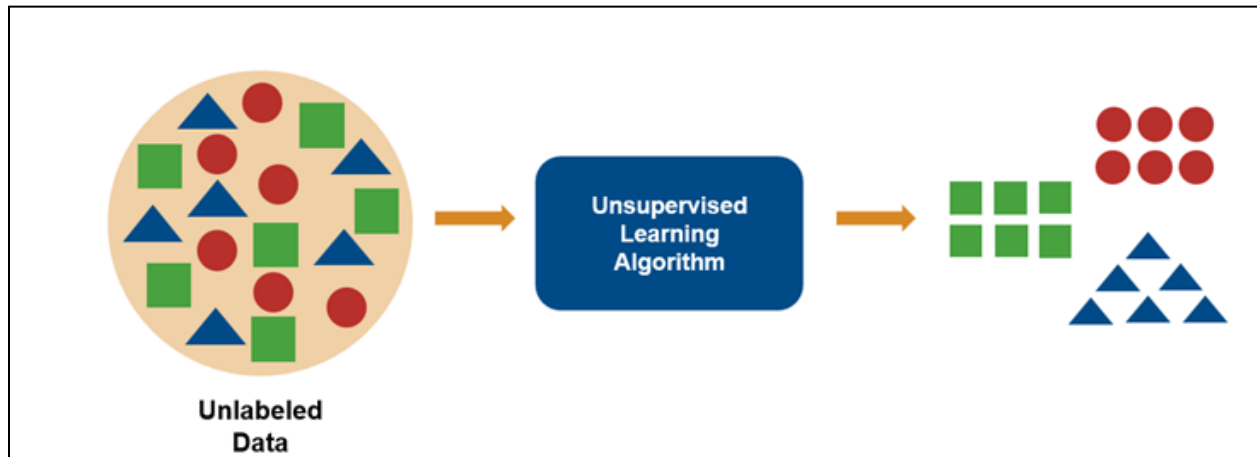
**The image shows a set of animals:** elephants, camels, and cows that represents raw data that the unsupervised learning algorithm will process.

- The "Interpretation" stage signifies that the algorithm doesn't have predefined labels or categories for the data. It needs to figure out how to group or organize the data based on inherent patterns.
- **Algorithms** represent the core of an unsupervised learning process using techniques like clustering, dimensionality reduction, or anomaly detection to identify patterns and structures in the data.
- **The processing** stage shows the algorithm working on the data.

The output shows the results of the unsupervised learning process. In this case, the algorithm might have grouped the animals into clusters based on their species (elephants, camels, cows).

## How does unsupervised learning work?

- Unsupervised learning algorithms discover hidden patterns, structures, and groupings within data, without any prior knowledge of the outcomes. These algorithms rely on unlabeled data, data that has no predefined labels.



- A typical unsupervised learning process involves data preparation, applying the right unsupervised learning algorithm to it, and, finally, interpreting and evaluating the results. This approach is particularly useful for tasks such as clustering, where the goal is to group similar data points together, and dimensionality reduction, which simplifies data by reducing the number of features (dimensions). By analyzing the inherent structure of the data, unsupervised learning can provide a better understanding of your data sets.
- Unsupervised learning can also be applied before supervised learning to identify features in exploratory data analysis and establish classes based on groupings. This is part of feature engineering, a process for transforming raw data into features suitable for supervised machine learning.



### 4.3.1 TYPES OF UNSUPERVISED LEARNING:

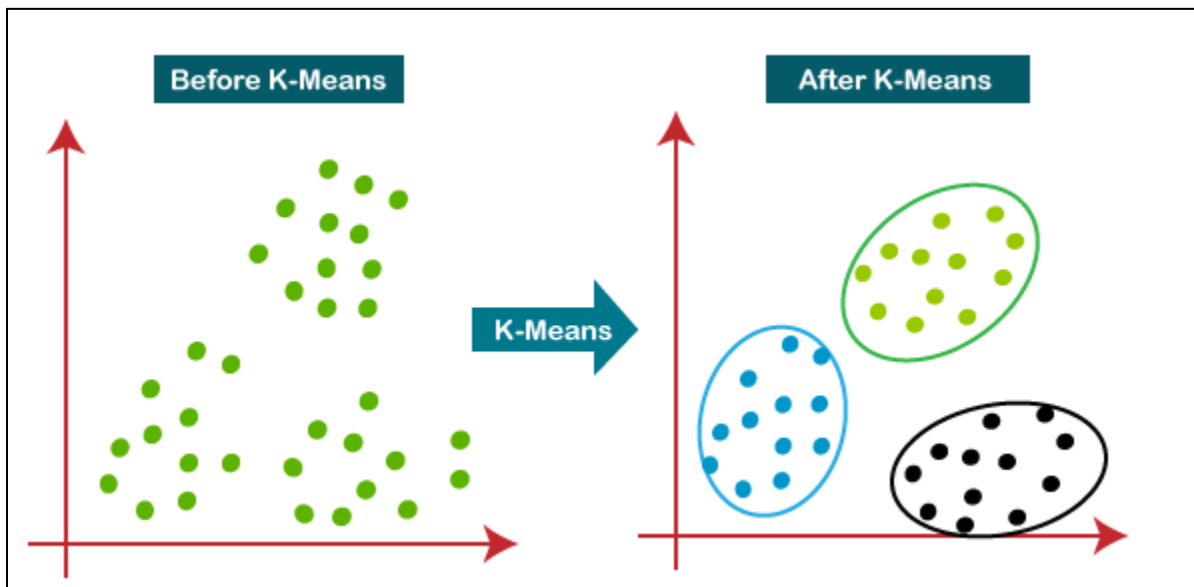
---

Clustering is one of the most popular techniques in unsupervised learning. It involves grouping data points with similar characteristics into clusters, allowing for easy categorization and analysis.

1. **K-means Clustering:** K-means clustering is one in every of the only and most common unsupervised machine learning algorithms.

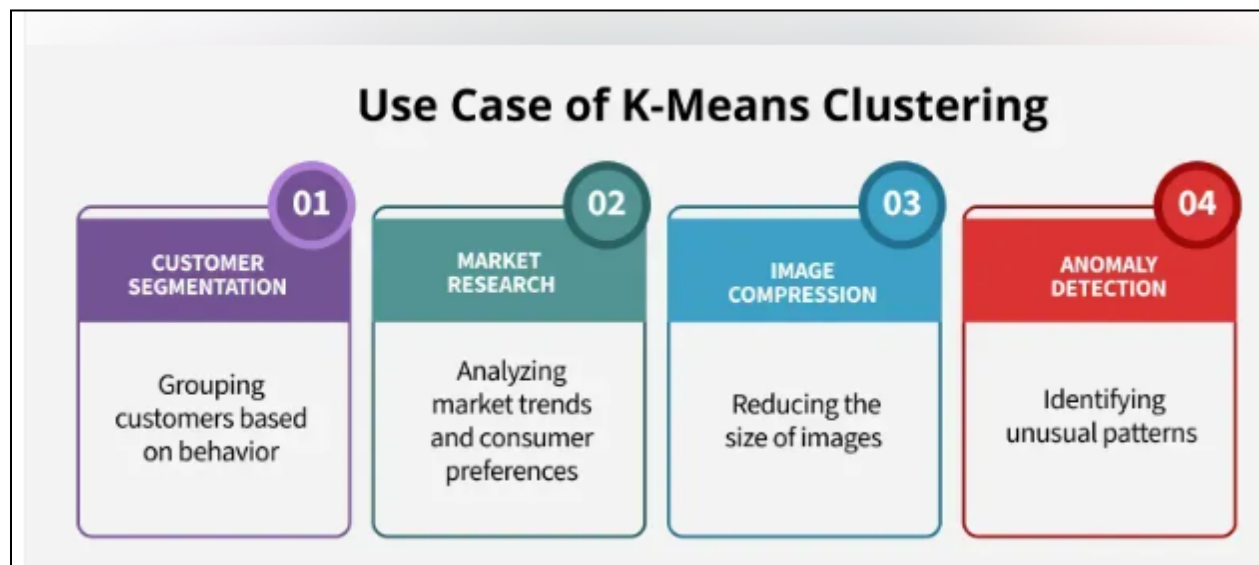
Typically, unsupervised algorithms create inferences from datasets mistreatment solely input vectors while not bearing on famed, or tagged, outcomes. A cluster refers to a collection of data points aggregated together because of certain similarities.

In this we will define the value of  $k$ , where  $k$  is the number of clusters we want to retain with the dataset. So, on defining the value for  $k$  we are defining the number of clusters we want to make.



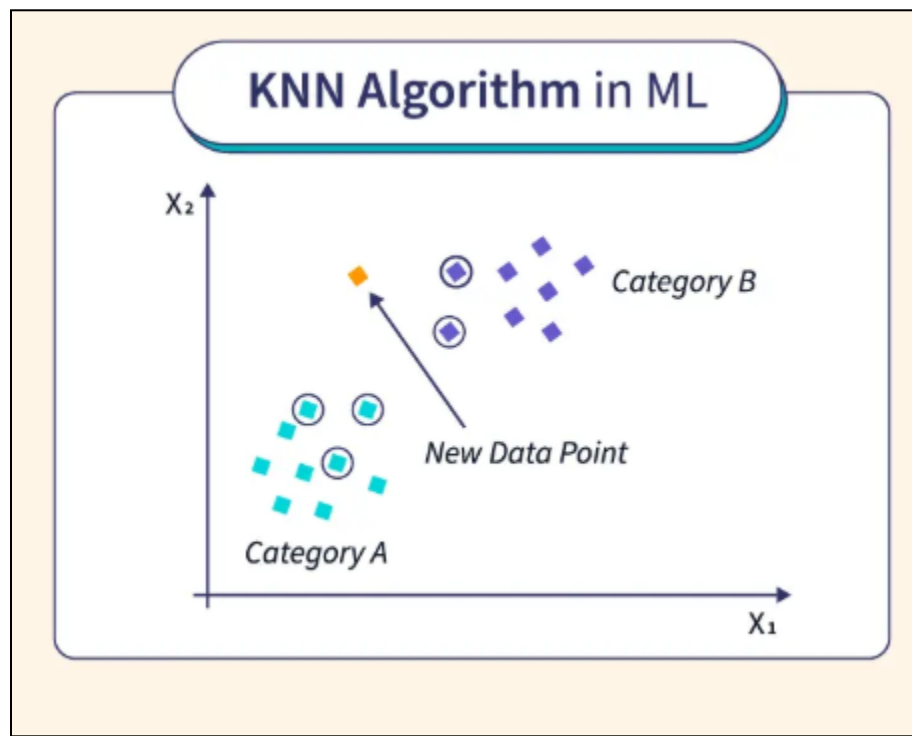
## How does it work?

1. First, we have to sort out the value of  $k$  in order to get the number of clusters to our data points on the graph
2. And secondly, we will assign some random points in our dataset as the centroids and these centroids are known as cluster centroids
3. Then assign the points to the respective cluster to which the cluster centroid is closest, among all the cluster centroids, and let's assign those data points to that cluster for instance
4. Then comes the moving of the centroid, where the assigned points to a particular cluster are taken and then points to the respective cluster are averaged, and then that averaged coordinate is assigned as the new cluster centroid
5. Then we follow up the cluster assignment and moving of centroid iteratively till we get our optimum state of clustering
6. After the iterative process, we are then left with the  $K$  clusters as we desired to form.



## 2. KNN (k-Nearest Neighbors),:

K-Nearest Neighbors (KNN) is a supervised machine learning algorithm generally used for classification but can also be used for regression tasks. It works by finding the "k" closest data points (neighbors) to a given input and makes predictions based on the majority class (for classification) or the average value (for regression). Since KNN makes no assumptions about the underlying data distribution it makes it a non-parametric and instance-based learning method.



In **classification**, KNN assigns a class label to a new data point based on the majority class of its nearest neighbors. For instance, if a data point has five nearest neighbors, and three of them belong to class A while two belong to class B, the algorithm will classify the point as class A.

In **regression**, KNN predicts continuous values by averaging the values of the k-nearest neighbors. For example, if you're predicting house prices, KNN will use the average prices of the k-nearest neighbors to estimate the price of a new house.

### Step 1: Determine the Number of Nearest Neighbors (k)

The first step is to select the number of neighbors (**k**) to consider. The value of **k** determines how many neighboring points will influence the classification or prediction of a new data point.

### Step 2: Calculate the Distance Between the Query Point and Dataset Points

For each data point in the dataset, the algorithm calculates the distance between the query point (the new point to be classified or predicted) and every other point. Various distance metrics can be used, such as **Euclidean distance**, **Manhattan distance**, or **Minkowski distance**.

$$d(P, Q) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

### Step 3: Sort and Select the k-Nearest Neighbors

After calculating the distances, the algorithm sorts all data points in ascending order of distance. It then selects the **k-nearest neighbors**—the data points that are closest to the query point.

### Step 4: Make a Prediction

- **For classification:** The algorithm assigns the query point to the class label that is most frequent among the k-nearest neighbors (majority voting).
- **For regression:** The algorithm predicts the value by averaging the values of the k-nearest neighbors.

## APPLICATIONS OF KNN ALGORITHM IN MACHINE LEARNING

5. **Text Classification** – Classifies documents (e.g., spam detection) by comparing keyword patterns with labeled data.
6. **Image Recognition** – Identifies objects or faces by analyzing pixel similarity.

7. **Recommendation Systems** – Suggests items based on preferences of similar users.
8. **Medical Diagnosis** – Predicts diseases or treatments by comparing patient data with similar cases.

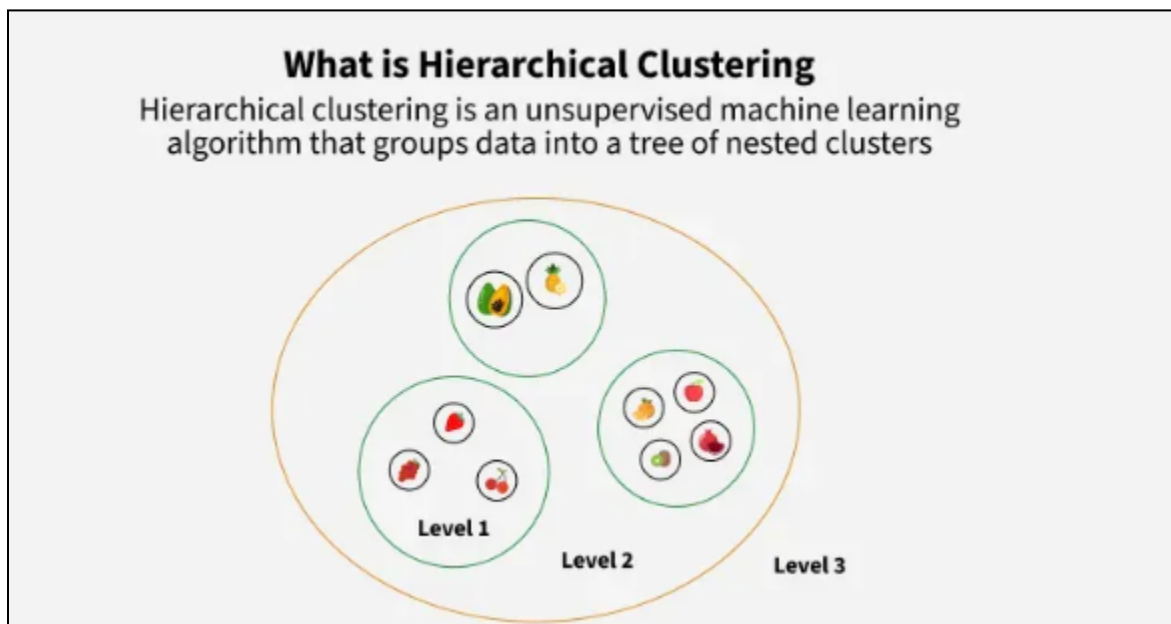
### 3. HIERARCHICAL CLUSTERING:

**Hierarchical clustering** is used to group similar data points together based on their similarity creating a **hierarchy or tree-like structure**. The key idea is to begin with each data point as its own separate cluster and then progressively merge or split them based on their similarity. Lets understand this with the help of an example

Imagine you have four fruits with different weights: an **apple (100g)**, a **banana (120g)**, a **cherry (50g)** and a **grape (30g)**. Hierarchical clustering starts by treating each **fruit as its own group**.

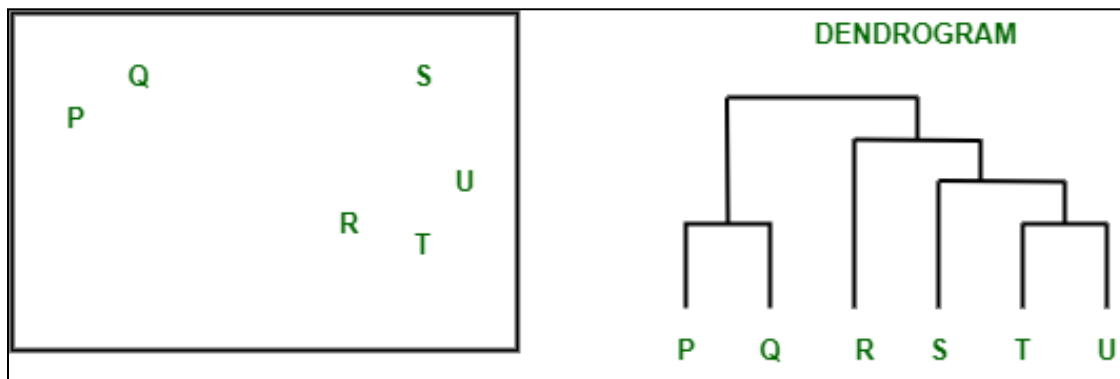
- It then merges the closest groups based on their weights.
- First the cherry and grape are grouped together because they are the lightest.
- Next the apple and banana are grouped together.

Finally all the fruits are merged into one large group, showing how hierarchical clustering progressively combines the most similar data points.



## Dendrogram:

A **dendrogram** is like a family tree for clusters. It shows how individual data points or groups of data merge together. The bottom shows each data point as its own group, and as you move up, similar groups are combined. The lower the merge point, the more similar the groups are. It helps you see how things are grouped step by step. The working of the dendrogram can be explained using the below diagram:



*Dendrogram*

In the above image on the left side there are five points labeled P, Q, R, S and T. These represent individual data points that are being clustered. On the right side there's a **dendrogram** which show how these points are grouped together step by step.

- At the bottom of the dendrogram the points P, Q, R, S and T are all separate.
- As you move up, the closest points are merged into a single group.
- The lines connecting the points show how they are progressively merged based on similarity.
- The height at which they are connected shows how similar the points are to each other; the shorter the line the more similar they are

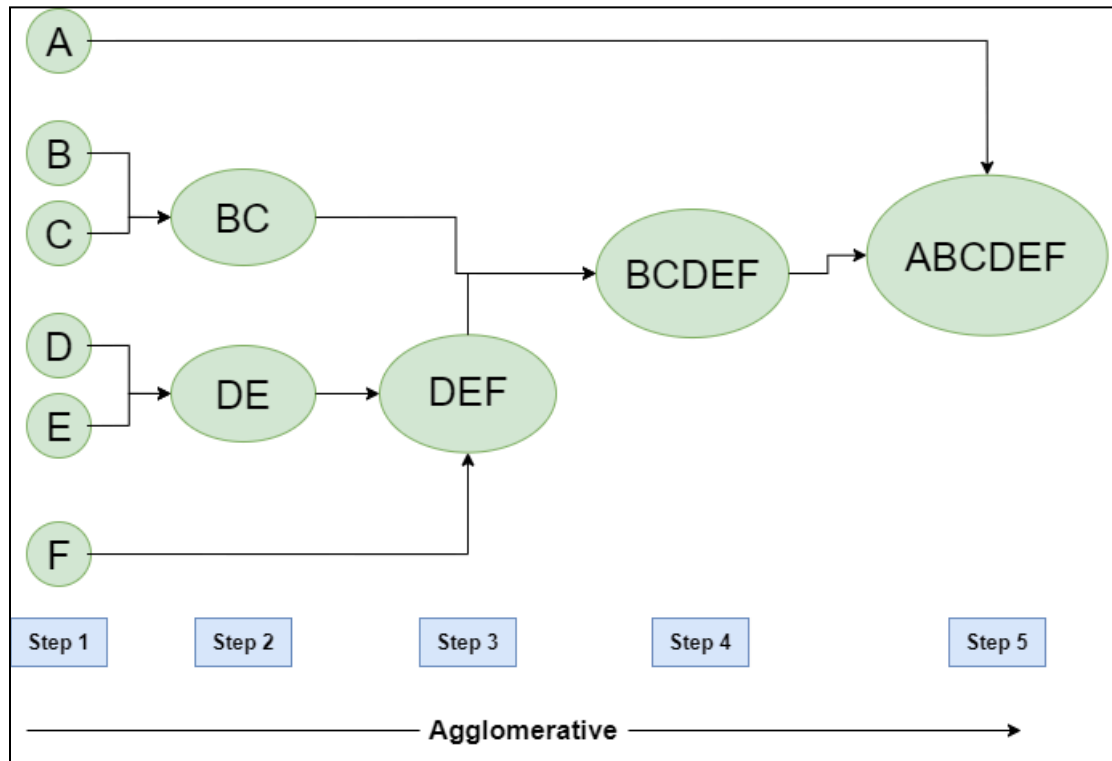
## Types of Hierarchical Clustering

Now we understand the basics of hierarchical clustering. There are two main types of hierarchical clustering.

1. Agglomerative Clustering
2. Divisive clustering

## Hierarchical Agglomerative Clustering

It is also known as the **bottom-up approach** or **hierarchical agglomerative clustering (HAC)**. Unlike flat clustering hierarchical clustering provides a structured way to group data. This clustering algorithm does not require us to prespecify the number of clusters. Bottom-up algorithms treat each data as a singleton cluster at the outset and then successively agglomerate pairs of clusters until all clusters have been merged into a single cluster that contains all data.



*Hierarchical Agglomerative Clustering*

### Workflow for Hierarchical Agglomerative clustering

1. **Start with individual points:** Each data point is its own cluster. For example if you have 5 data points you start with 5 clusters each containing just one data point.
2. **Calculate distances between clusters:** Calculate the distance between every pair of clusters. Initially since each cluster has one point this is the distance between the two data points.

3. **Merge the closest clusters:** Identify the two clusters with the smallest distance and merge them into a single cluster.
4. **Update distance matrix:** After merging you now have one less cluster. Recalculate the distances between the new cluster and the remaining clusters.
5. **Repeat steps 3 and 4:** Keep merging the closest clusters and updating the distance matrix until you have only one cluster left.
6. **Create a dendrogram:** As the process continues you can visualize the merging of clusters using a tree-like diagram called a **dendrogram**. It shows the hierarchy of how clusters are merged.

### **Hierarchical Divisive clustering**

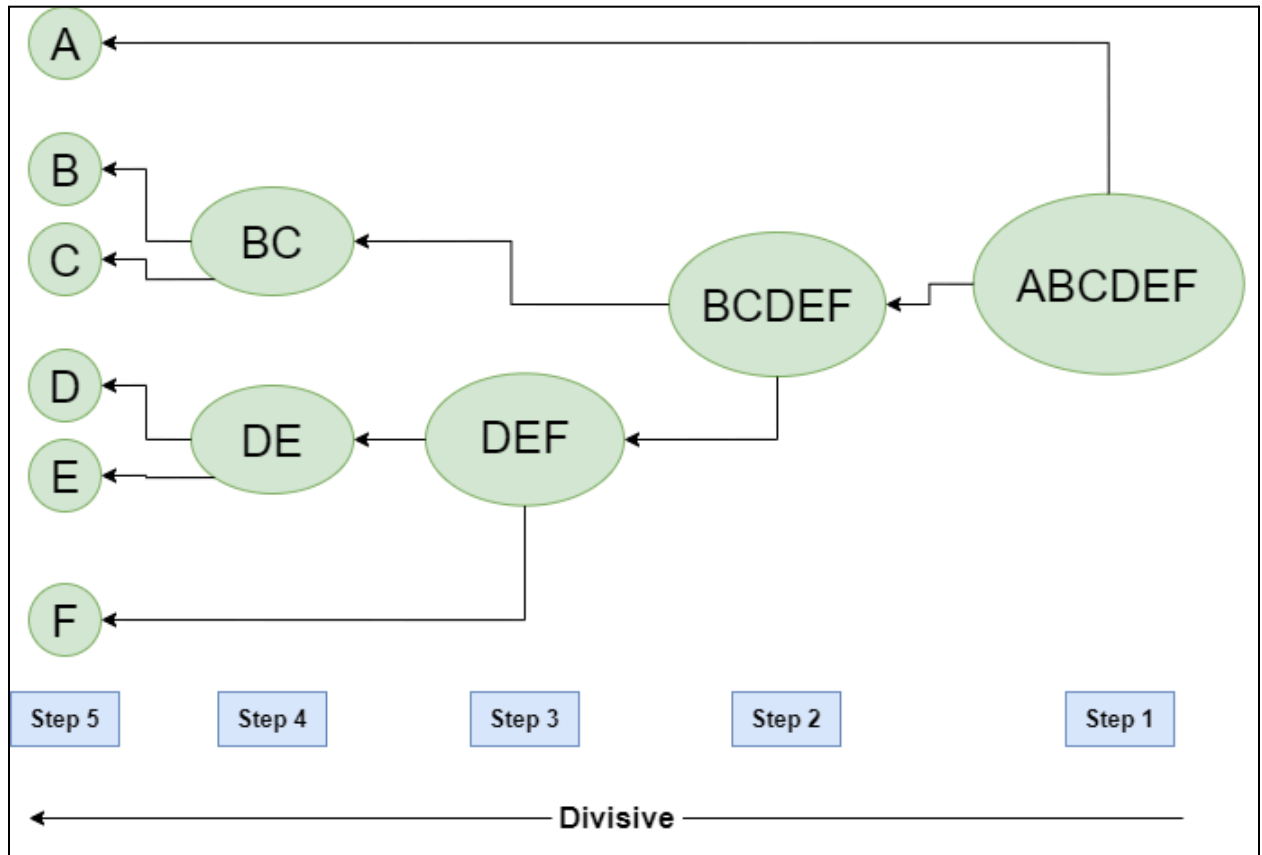
It is also known as a **top-down approach**. This algorithm also does not require to prespecify the number of clusters. Top-down clustering requires a method for splitting a cluster that contains the whole data and proceeds by splitting clusters recursively until individual data have been split into singleton clusters.

#### **Workflow for Hierarchical Divisive clustering :**

1. **Start with all data points in one cluster:** Treat the entire dataset as a single large cluster.
2. **Split the cluster:** Divide the cluster into two smaller clusters. The division is typically done by finding the two most dissimilar points in the cluster and using them to separate the data into two parts.
3. **Repeat the process:** For each of the new clusters, repeat the splitting process:
  1. Choose the cluster with the most dissimilar points.
  2. Split it again into two smaller clusters.



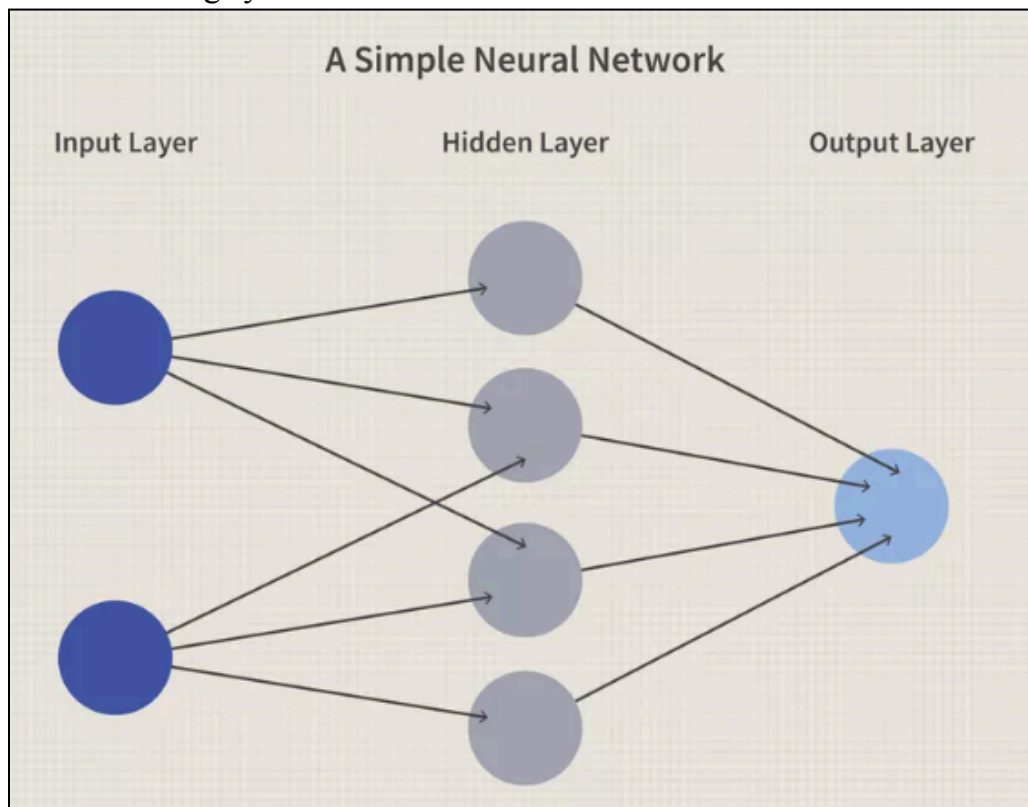
4. **Stop when each data point is in its own cluster:** Continue this process until every data point is its own cluster, or the stopping condition (such as a predefined number of clusters) is met.



#### 4. NEURAL NETWORKS:

A neural network is a series of algorithms that endeavors to recognize underlying relationships in a set of data through a process that mimics the way the human brain operates. In this sense, neural networks refer to systems of neurons, either organic or artificial in nature.

Neural networks can adapt to changing input; so the network generates the best possible result without needing to redesign the output criteria. The concept of neural networks, which has its roots in artificial intelligence, is swiftly gaining popularity in the development of trading systems.



Neural networks (NNs) are a core component of machine learning, especially in the advanced field of deep learning. They are computational models inspired by the structure and function of biological neural networks, like the human brain.

## **What are neural networks?**

Neural networks consist of interconnected nodes, or artificial neurons, organized in layers: input, hidden, and output. Each connection between neurons has an associated weight and bias, which determine the influence a neuron has on the next.

## **How do neural networks learn?**

Neural networks learn through a process of training, often utilizing labeled datasets (supervised learning). The network processes input data and generates an output. During training, a loss function measures the error between the predicted and actual output. The network then uses backpropagation to adjust the weights and biases of the connections between neurons, minimizing the error. This iterative process, often facilitated by optimization algorithms like gradient descent, allows the network to learn complex patterns and improve its accuracy over time.

## **Types of neural networks**

Neural networks come in various architectures, including:

- Feedforward Neural Networks (FNNs): Data flows from input to output without loops.
- Multilayer Perceptrons (MLPs): FNNs with one or more hidden layers for complex relationships.
- Convolutional Neural Networks (CNNs): Suited for image and video tasks by extracting features with convolutional layers.
- Recurrent Neural Networks (RNNs): Process sequential data using loops to retain information from previous inputs.
- Long Short-Term Memory (LSTM) Networks: A type of RNN that can retain information over longer sequences.
- Generative Adversarial Networks (GANs): Use competing networks to create realistic data.
- Transformer Networks: Use a self-attention mechanism for sequential data processing.

## **Advantages of neural networks:**

Neural networks offer several advantages:

- **Adaptability:** They can learn from new data.
- **Pattern Recognition:** They excel at identifying complex patterns.
- **Non-Linearity:** They can model non-linear data relationships.
- **Parallel Processing:** They can handle multiple tasks concurrently.
- **Generalization:** They can apply learned knowledge to new data.

## **Limitations of neural networks:**

Despite their strengths, neural networks have limitations:

- **High Computational Requirements:** Training often requires significant computing power.
- **Large Data Dependency:** They generally need extensive labeled data.
- **"Black Box" Nature:** Their decision processes can be difficult to understand.
- **Risk of Overfitting:** They might perform poorly on new data if they only memorize the training data.

## **Applications of neural networks:**

Neural networks are used across various fields:

- **Computer Vision:** Image and object recognition.
- **Natural Language Processing (NLP):** Machine translation and chatbots.
- **Healthcare:** Disease diagnosis and medical analysis.
- **Finance:** Fraud detection and forecasting.
- **Autonomous Systems:** Self-driving cars and robotics.

## 4.4 MODE EVALUATION:

---

Model evaluation is a process that uses some metrics which help us to analyze the performance of the model. Think of training a model like teaching a student. Model evaluation is like giving them a test to see if they *truly* learned the subject—or just memorized answers. It helps us answer:

- Did the model learn patterns?
- Will it fail on new questions?

Model development is a multi-step process and we need to keep a check on how well the model does future predictions and analyze a model's weaknesses. There are many metrics for that. Cross Validation is one technique that is followed during the training phase and it is a model evaluation technique.

### **Cross-Validation: The Ultimate Practice Test:**

Cross-validation is a technique used in machine learning to evaluate the performance of a model. It provides a better estimate of how well the model will generalize to unseen data by splitting the data into multiple subsets.

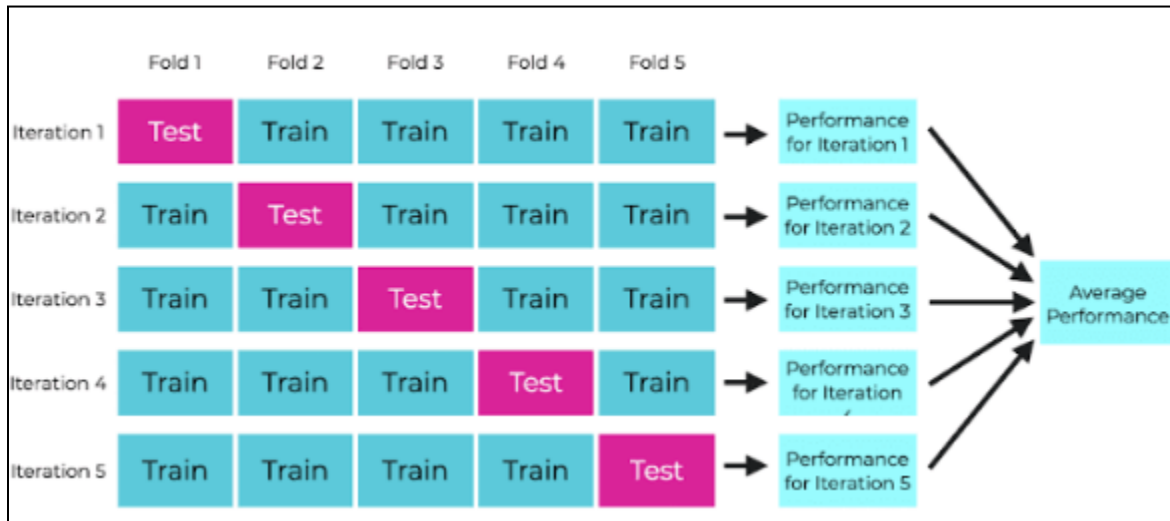
It helps in assessing the model's ability to handle different scenarios and identify potential issues such as overfitting or underfitting.

Cross-validation plays a crucial role in assessing the performance and reliability of machine learning models. It provides a more accurate evaluation of the model's performance by using multiple partitions of the data.

This helps in making informed decisions about model selection, hyperparameter tuning, and assessing the robustness of the chosen model.

## Why is Cross-Validation Used?

Cross-validation is a widely used technique in machine learning and statistical analysis to evaluate the performance of a model.



## 2. Benefits of cross-validation:

- **More accurate performance estimation:** Cross-validation provides a more robust and reliable estimate of model performance compared to a single train-test split by evaluating the model on different subsets of the data. This helps reduce the bias and variance associated with relying on a single split.
- **Reduced overfitting:** By training and testing on multiple, distinct folds, cross-validation helps prevent the model from memorizing the training data and improving its ability to generalize to unseen data.
- **Efficient data utilization:** Cross-validation maximizes the use of the available data, especially beneficial for smaller datasets, as every data point is used for both training and testing at some point.
- **Improved model selection:** When comparing different machine learning algorithms, cross-validation provides a fair and consistent way to evaluate their performance, aiding in the selection of the best model for a given task.
- **Hyperparameter tuning:** Cross-validation is crucial for optimizing the hyperparameters of a model by testing various configurations across folds to identify those that yield the best generalization performance.

### 3. Cross-validation with positive and negative classes:

In classification problems, especially binary classification, it is important to consider the concepts of positive and negative classes.

The **positive class** represents the outcome or characteristic the model is primarily designed to detect (e.g., presence of a disease), while the **negative class** represents the absence of that characteristic. When dealing with imbalanced datasets, where one class is significantly more prevalent than the other (e.g., few positive examples compared to many negative ones), using standard cross-validation can lead to biased model evaluations.

This is where Stratified K-Fold Cross-Validation becomes essential. Stratified K-Fold CV ensures that each fold maintains the same proportion of samples from each class as the original dataset. This is particularly important for imbalanced datasets, as it helps to:

- **Prevent biased evaluation:** By ensuring that the validation set has a representative distribution of both positive and negative classes, it allows for a more accurate estimation of the model's ability to handle imbalanced data.
- **Improve model performance on the minority class:** This method ensures the model is trained and evaluated on minority class examples in each fold, preventing it from being overly biased towards the majority class and improving its ability to predict the minority class accurately.