

# Types of Machine Learning

Machine Learning (ML) techniques can be categorized into three main types based on how models learn from data:

1. Supervised Learning – Labeled data
2. Unsupervised Learning – No labels, finding hidden patterns
3. Reinforcement Learning – Learning through rewards and penalties

## 1. Supervised Learning

In Supervised Learning, the model is trained on a dataset with labeled inputs and outputs (i.e., the model learns from past examples). The goal is to map input data to correct outputs.

### 📌 Characteristics

- Uses labeled training data.
- The model learns from known inputs and outputs.
- Can be used for classification and regression tasks.

📌 Types		
Supervised Learning Type	Description	Example
<b>Classification</b>	Categorizes data into distinct labels (e.g., spam or not spam).	Email Spam Detection, Disease Prediction
<b>Regression</b>	Predicts continuous numerical values.	House Price Prediction, Stock Market Forecasting

### 📌 Examples

- Classification: Decision Trees, Random Forest, Support Vector Machines (SVM), Logistic Regression, Neural Networks.
- Regression: Linear Regression, Polynomial Regression, Support Vector Regression (SVR), Ridge Regression.

## 📌 Example Use Case

Predicting whether an email is spam or not (Classification). The model is trained with past emails labeled as spam or not spam.

## 2. Unsupervised Learning

In Unsupervised Learning, the model is given unlabeled data and must find patterns, relationships, or structures on its own.

### 📌 Characteristics

- Works with unlabeled data.
- The model finds patterns and structures without explicit training labels.
- Used for clustering, anomaly detection, and dimensionality reduction.

📌 Types		
Unsupervised Learning Type	Description	Example
Clustering	Groups similar data points together.	Customer Segmentation, Document Clustering
Dimensionality Reduction	Reduces the number of features while preserving information.	PCA (Principal Component Analysis) for Data Compression
Anomaly Detection	Identifies outliers in data.	Fraud Detection, Intrusion Detection

### 📌 Examples

- Clustering: K-Means, Hierarchical Clustering, DBSCAN.
- Dimensionality Reduction: PCA (Principal Component Analysis), t-SNE.
- Anomaly Detection: Isolation Forest, Autoencoders.

### 📌 Example Use Case

Segmenting customers into different groups based on purchasing behavior (Clustering). The model does not have predefined labels but learns similarities between customers.

### 3. Reinforcement Learning (RL)

In Reinforcement Learning, an agent learns by interacting with an environment, receiving rewards or penalties based on its actions.

#### 📌 Characteristics

- The model learns through trial and error.
- Uses a reward-based feedback system.
- Commonly used in robotics, gaming, and self-driving cars.

📌 Key Components	
Term	Description
<b>Agent</b>	The learning model that makes decisions (e.g., AI player in a game).
<b>Environment</b>	The surroundings in which the agent interacts (e.g., chessboard).
<b>Actions</b>	Choices the agent can make (e.g., moving left or right).
<b>Reward</b>	Feedback given after an action (positive or negative).

#### 📌 Examples

- Q-Learning (Value-based RL)
- Deep Q-Networks (DQN) (Deep RL)
- Proximal Policy Optimization (PPO) (Policy-based RL)

#### 📌 Example Use Case

Teaching an AI to play chess using RL. The model learns the best moves by playing multiple games and receiving rewards for winning and penalties for losing.

---

Key Differences Between		ML	Types
Feature	Supervised Learning	Unsupervised Learning	Reinforcement Learning
Data Type	Labeled	Unlabeled	No predefined labels, learns via rewards
Goal	Predict outcomes (Classification/Regression)	Discover patterns (Clustering, Anomaly Detection)	Learn optimal actions through trial and error
Examples	Spam Detection, Fraud Prediction	Customer Segmentation, Image Compression	Self-Driving Cars, Game AI
Dependency	Requires labeled training data	No labeled data required	Requires an interactive environment

## 5. Other Categories of Machine Learning

Aside from these three main types, ML can be categorized based on model functionality:

- 📌 1. Semi-Supervised Learning
  - Uses a mix of labeled and unlabeled data.
  - Example: Google Photos Face Recognition – Starts with labeled faces but learns to recognize new ones.
- 📌 2. Self-Supervised Learning (Advanced ML)
  - The model generates labels from input data itself.
  - Example: BERT (Natural Language Processing), GPT Models.
- 📌 3. Deep Learning (DL)

- Uses Neural Networks for complex data like images, speech, and text.
- Example: CNNs (Convolutional Neural Networks) for Image Recognition, RNNs for Time-Series.

**6.**

## **Summary**

**Table**

<b>Machine Learning Type</b>	<b>Goal</b>	<b>Example Algorithms</b>	<b>Applications</b>
<b>Supervised Learning</b>	Predict outcomes	Linear Regression, SVM, Decision Trees	Spam Filtering, Fraud Detection
<b>Unsupervised Learning</b>	Find hidden patterns	K-Means, PCA, Hierarchical Clustering	Customer Segmentation, Anomaly Detection
<b>Reinforcement Learning</b>	Learn through trial & error	Q-Learning, Deep Q-Networks	Robotics, Self-Driving Cars

## **Characteristics of Decision Tree Algorithm**

The Decision Tree Algorithm is a supervised machine learning algorithm used for classification and regression tasks. It follows a tree-like structure to make decisions. The key characteristics are:

1. Hierarchical Structure: It divides data recursively into subsets based on feature values, forming a tree structure with nodes and branches.
2. Root Node: Represents the entire dataset and splits into child nodes based on the best feature selection.
3. Internal Nodes: Each node represents a decision based on a feature condition.
4. Leaf Nodes: These contain the final class labels (in classification) or numerical values (in regression).
5. Splitting Criteria: Uses statistical methods like Gini Index, Entropy (Information Gain), and Chi-square to determine the best split.
6. Prone to Overfitting: Complex trees can capture noise, leading to overfitting, which can be controlled using pruning techniques.
7. Handles Categorical & Numerical Data: Works well with both types of data.
8. Non-Parametric Algorithm: No assumption about data distribution is required.
9. Interpretable & Easy to Visualize: Unlike some black-box models, decision trees are intuitive and explainable.
10. Computationally Efficient: Training and prediction times are generally fast for small to medium datasets.

11. Pruning for Generalization: Techniques like pre-pruning (stopping early) and post-pruning (removing unnecessary branches) help reduce overfitting.
  12. Handles Missing Values: Some implementations of decision trees can handle missing data efficiently.
- 

## **Classifiers of Decision Tree Algorithm**

Decision Trees are primarily used for classification tasks. The most common classifiers include:

1. ID3 (Iterative Dichotomiser 3)
  - Uses Entropy and Information Gain to determine the best attribute for splitting.
  - Suitable for small datasets but can lead to overfitting.
2. C4.5 (Successor of ID3)
  - Handles both categorical and continuous data.
  - Uses Gain Ratio instead of pure Information Gain, which normalizes splits to avoid bias.
  - Can handle missing values by assigning probabilities based on available data.
3. CART (Classification and Regression Trees)
  - Can be used for both classification and regression tasks.
  - Uses Gini Index instead of Entropy to decide splits.
  - Produces binary trees (each node has two children).
4. CHAID (Chi-Square Automatic Interaction Detector)
  - Uses Chi-Square statistics to determine the best feature split.
  - Produces multi-way splits (unlike CART).
  - Best suited for categorical variables.
5. Random Forest (Ensemble of Decision Trees)
  - Combines multiple decision trees to improve accuracy and reduce overfitting.
  - Uses Bagging (Bootstrap Aggregating) for randomness in training data selection.
  - The final prediction is based on majority voting (classification) or averaging (regression).
6. Extra Trees (Extremely Randomized Trees)
  - Similar to Random Forest but introduces more randomness in split selection.
  - Reduces variance further than Random Forest but may increase bias.
7. Gradient Boosted Trees (GBDT, XGBoost, LightGBM, CatBoost)
  - Sequentially builds decision trees, correcting errors in the previous models.
  - More computationally expensive but provides higher accuracy.

## Continuous Attributes in Machine Learning

Continuous attributes are numerical values that can take an infinite number of values within a given range. They are widely used in machine learning and require special handling in various models.

### Characteristics of Continuous Attributes

- Can take any value within a range (e.g., temperature, salary, weight).
- Require binning or discretization for some algorithms that prefer categorical inputs.
- Handled differently in different machine learning models:
  - Decision Trees → Threshold-based splitting
  - Linear Models → Directly used in equations
  - Neural Networks → Normalized for better performance

### Handling Continuous Attributes in ML

1. Normalization & Scaling
  - Standardization (Z-score normalization):  $(X - \mu) / \sigma$
  - Min-Max Scaling:  $X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$
2. Discretization (Binning)
  - Equal-width binning: Dividing into equal range intervals
  - Equal-frequency binning: Dividing based on percentiles (e.g., quartiles)
3. Feature Engineering
  - Polynomial transformations (e.g.,  $Age^2$ ,  $Age^3$ )
  - Log transformations to reduce skewness
4. Outlier Handling
  - Using IQR (Interquartile Range) or Z-score methods to detect and remove outliers.

## Approaches for Splitting Based on Continuous Attributes in Decision Trees

When dealing with continuous attributes (e.g., Age, Salary, Temperature), decision trees must determine the best way to split the data. Here are the primary approaches:

### 1. Binary Splitting Using a Threshold

- The continuous attribute is split into two groups based on a threshold value.
- The threshold is determined by evaluating different possible splits and choosing the one that minimizes impurity (e.g., Gini Index or Information Gain).

- Example:
  - If Age values are: {18, 22, 30, 40, 50}, possible thresholds could be 20, 26, 35, 45.
  - The best split could be  $\text{Age} \leq 30$  (Group 1:  $\text{Age} \leq 30$ , Group 2:  $\text{Age} > 30$ ).

## 2. Finding the Best Threshold (Sorted Splitting)

- The algorithm sorts the continuous values in ascending order.
- It evaluates each midpoint between consecutive values as a potential split.
- The best threshold is chosen based on impurity reduction.
- Example:
  - Age values: {18, 22, 30, 40, 50}
  - Midpoints:  $(18+22)/2 = 20$ ,  $(22+30)/2 = 26$ ,  $(30+40)/2 = 35$ , etc.
  - The threshold with the highest Information Gain or lowest Gini Index is selected.

## 3. Multi-way Splitting (Discretization)

- Instead of a single threshold, the continuous attribute is divided into multiple discrete intervals (bins).
- Can be done using equal-width binning (fixed range intervals) or equal-frequency binning (splitting based on percentiles).
- Example:
  - Age categories:
    - Young ( $\leq 25$ )
    - Middle-aged (26-45)
    - Senior ( $> 45$ )

## 4. Mean or Median-based Splitting

- A simple way to handle continuous attributes is by splitting based on the mean or median of the values.
- Example:
  - If the median Age in the dataset is 35, the split could be:
    - $\text{Age} \leq 35$
    - $\text{Age} > 35$

## 5. Entropy/Gini-based Splitting

- Decision trees like ID3, C4.5, and CART use Entropy (Information Gain) or Gini Index to determine the best split point.
- The algorithm selects the split that provides the greatest reduction in impurity.

# **Bias-Variance Tradeoff, Underfitting & Overfitting in Machine Learning**

The bias-variance tradeoff is a key concept in machine learning that explains the balance between model complexity and its ability to generalize to new data. It helps in understanding underfitting and overfitting.

## 1. Understanding Bias and Variance

Bias and variance contribute to a model's total error. The goal is to minimize both to achieve a model that generalizes well.

### **Bias (Underfitting)**

- Bias is the error introduced by approximating a real-world problem with a simplistic model.
- High bias models make strong assumptions and fail to capture complex relationships in the data.
- Results in: Underfitting → Poor performance on both training and test data.

 Low Bias: Model captures important trends.

 High Bias: Model is too simple and cannot learn patterns.

- Example: Linear regression on a non-linear dataset.

### **Variance (Overfitting)**

- Variance measures how much the model's predictions change when trained on different data.
- High variance models are too sensitive to training data and learn noise instead of patterns.
- Results in: Overfitting → High accuracy on training data but poor performance on test data.

 Low Variance: Model generalizes well.

 High Variance: Model captures noise and fails on unseen data.

- Example: A deep decision tree that perfectly classifies training data but fails on new data.

## 2. Bias-Variance Tradeoff

A model must balance bias and variance to achieve good performance.

Scenario	Bias	Variance	Model Complexity	Performance
Underfitting	High	Low	Too Simple	Poor (on both training & test)
Optimal Model	Low	Low	Balanced	Best Generalization

Overfitting      Low      High      Too Complex      Poor (good on training, bad on test)

### 3. Underfitting vs. Overfitting

Factor	Underfitting	Overfitting
Model Complexity	Too Simple	Too Complex
Bias	High	Low
Variance	Low	High
Training Error	High	Low
Test Error	High	High
Performance	Fails on both training & test data	Works well on training, fails on test data

### 4. Techniques to Avoid Overfitting & Underfitting

#### → Techniques to Avoid Overfitting (High Variance)

- ✓ Cross-Validation – Ensures the model is tested on different data (e.g., K-Fold Cross-Validation).
- ✓ Regularization – Adds penalties for complexity (L1/L2 regularization in regression, Dropout in neural networks).
- ✓ Pruning (for Decision Trees) – Removes unnecessary branches.
- ✓ Reduce Model Complexity – Use shallower neural networks or fewer parameters.
- ✓ Dropout & Early Stopping – Stops training before the model memorizes noise.
- ✓ Ensemble Methods – Bagging (Random Forest), Boosting (XGBoost).
- ✓ More Data & Data Augmentation – Helps generalization (e.g., flipping images in deep learning).

#### → Techniques to Avoid Underfitting (High Bias)

- ✓ Use More Complex Models – If linear regression underfits, try polynomial regression or neural networks.
- ✓ Feature Engineering – Add relevant features to capture important patterns.

- ✓ Reduce Regularization – Too much regularization can prevent learning.
- ✓ Train Longer – Increase epochs for better learning in deep models.
- ✓ Hyperparameter Tuning – Adjust learning rate, tree depth, number of layers, etc.

## **Model Overfitting: Causes, Solutions, and Examples**

### What is Overfitting?

Overfitting occurs when a machine learning model learns the training data too well, capturing noise and random fluctuations instead of the underlying pattern. As a result, the model performs exceptionally well on training data but fails to generalize to new, unseen data.

### **Causes of Overfitting**

#### 1. Too Complex Models (High Variance)

- Models with too many parameters (e.g., deep decision trees, high-degree polynomial regression) can memorize the training data instead of learning general patterns.

#### 2. Insufficient Training Data

- When there is too little data, the model learns patterns that are not representative of the real-world distribution.

#### 3. Too Many Features (Curse of Dimensionality)

- If the dataset has many features but not enough samples, the model can find spurious correlations and overfit.

#### 4. Training for Too Many Epochs (Neural Networks)

- In deep learning, training a model for too long without early stopping can cause it to learn noise rather than meaningful patterns.

#### 5. Noisy or Unclean Data

- If the training data contains errors, duplicates, or irrelevant details, the model will capture noise as if it were a pattern.

#### 6. Improper Model Selection

- Using a complex model when a simpler one would suffice can lead to overfitting (e.g., using a deep neural network when linear regression would work).

## **Solutions to Overfitting**

### 1. Train with More Data

- Collecting more data helps the model learn more generalizable patterns and reduces the effect of noise.

### 2. Use Cross-Validation

- K-Fold Cross-Validation ensures the model is tested on different subsets of the data, preventing reliance on a single training set.

### 3. Feature Selection & Dimensionality Reduction

- Removing irrelevant features or using PCA (Principal Component Analysis) can reduce complexity and improve generalization.

### 4. Regularization (L1 & L2)

- L1 (Lasso) Regularization: Shrinks some feature weights to zero, effectively selecting only important features.
- L2 (Ridge) Regularization: Penalizes large weights, preventing the model from becoming too complex.

### 5. Pruning (for Decision Trees)

- Pre-pruning: Stops tree growth before it becomes too complex.
- Post-pruning: Removes unnecessary branches after training.

### 6. Dropout (for Neural Networks)

- Randomly deactivates neurons during training to prevent the network from relying on specific features.

### 7. Early Stopping

- Monitors validation loss and stops training when performance starts degrading on validation data.

### 8. Data Augmentation

- Artificially increasing the dataset size by adding slight variations (e.g., flipping, rotating images in deep learning).

## **Examples of Overfitting**

## Example 1: Overfitting in Polynomial Regression

- Consider a dataset where we predict house prices based on size.
- A linear model ( $y = ax + b$ ) may underfit the data.
- A high-degree polynomial (e.g., 10th-degree polynomial) may fit every data point exactly but fail on new data.

● Solution: Use a lower-degree polynomial and apply L2 regularization.

## Example 2: Overfitting in Decision Trees

- A deep decision tree may perfectly classify training data but fail on unseen data.
- If the tree has too many nodes, it memorizes noise instead of general rules.

● Solution: Apply pruning to reduce tree depth and use ensemble methods like Random Forest.

## Example 3: Overfitting in Neural Networks

- A deep neural network trained for too many epochs memorizes specific training examples.
- Training accuracy is 99%, but test accuracy is 60% (poor generalization).

● Solution: Use dropout layers, early stopping, and data augmentation.

	Predicted No	Predicted Yes
Actual No	45	5
Actual Yes	5	95

- Accuracy
- Precision
- Recall
- F1-Score

- Accuracy: Overall, how often is the classifier correct?

$$\begin{aligned} \text{Accuracy} &= \frac{TN+TP}{TN+FP+FN+TP} \\ &= \frac{45 + 95}{150} = \underline{\underline{93.33\%}} \end{aligned}$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95

- Misclassification Rate: Overall, how often is it wrong?

$$\begin{aligned} \text{Missclassification Rate} &= \frac{FN+FP}{TN+FP+FN+TP} \\ &= \frac{5 + 5}{150} = \underline{\underline{6.67\%}} \end{aligned}$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95

- **True Positive Rate:** When it's actually yes, how often does it predict yes?
- also known as "Sensitivity" or "Recall"
- $True\ Positive\ rate = \frac{TP}{Actual\ Yes}$

$$= \frac{95}{100} = 95\%$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95

- **False Positive Rate:** When it's actually no, how often does it predict yes?

$$• False\ Positive\ rate = \frac{FP}{Actual\ No}$$

$$= \frac{5}{50} = 10\%$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95

- **True Negative Rate:** When it's actually no, how often does it predict no?

- also known as "Specificity"

$$• True\ Negative\ rate = \frac{TN}{Actual\ No}$$

$$= \frac{45}{50} = 90\%$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95

- **Precision:** When it predicts yes, how often is it correct?

$$\bullet \text{ Precision} = \frac{\text{TP}}{\text{Predicted Yes}}$$

$$= \frac{95}{100} = 95\%$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95

- **Prevalence:** How often does the yes condition actually occur in our sample?

$$\bullet \text{ Prevalence} = \frac{\text{Actual Yes}}{\text{Total}}$$

$$= \frac{100}{150} = \underline{66.67\%}$$

	Predicted No	Predicted Yes
Actual No	TN = 45	FP = 5
Actual Yes	FN = 5	TP = 95