

1 Machine Learning: From Concepts to Practice

1.1 From AI to ML, DL, and NLP: Understanding the Landscape

1.1.1 Artificial Intelligence: The Big Picture

Artificial Intelligence (AI) is a broad field that aims to build computer systems capable of performing tasks that typically require human intelligence—such as perception, reasoning, learning, decision-making, and language understanding. As illustrated in Figure 1, AI contains many subfields, including machine learning, deep learning, natural language processing, and computer vision.

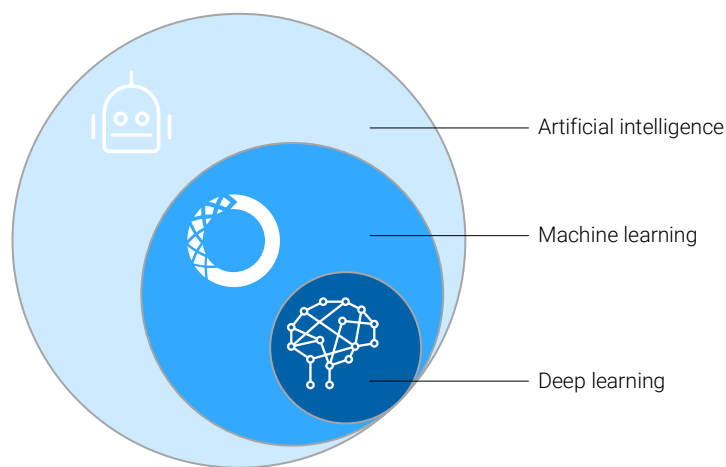


Figure 1. The relationship between AI, ML, and DL

1.1.2 Machine Learning: Learning from Data

Machine Learning is a core subfield of AI. Instead of writing explicit rules, machine learning enables computers to learn patterns from data and use what they learn to make predictions or decisions on new, unseen samples. Machine learning draws on multiple disciplines, including probability, statistics, linear algebra, optimization, and algorithms.

Machine learning is particularly valuable in situations where traditional rule-based or equation-driven methods fall short. It excels when data volumes are large, when the underlying system is too complex to be expressed with clear mathematical equations, and when the relationships within the data are nonlinear, uncertain, or influenced by many interacting factors.

1.1.3 Deep Learning: Hierarchical Feature Learning

Deep Learning is a specialized branch of machine learning that uses multi-layer neural networks to learn complex patterns automatically. Compared with traditional machine learning, deep learning can learn high-level features directly from raw data (such as images, text, or audio), reducing the need for manual feature engineering. Common deep learning frameworks in Python include TensorFlow, PyTorch, and Keras (which are outside the scope of this book).

1.1.4 Natural Language Processing: Teaching Machines to Understand Language

Natural Language Processing (NLP) focuses on enabling computers to understand, interpret, and generate human language, supporting a wide range of applications including text classification, sentiment analysis, information extraction, machine translation, and question-answering systems.

A recent and powerful development in this field is Large Language Models (LLMs), such as GPT. LLMs are trained on massive amounts of text data and can generate coherent, human-like text, answer questions, summarize documents, and even perform basic reasoning.

Unlike traditional NLP methods, which often rely on hand-crafted rules or task-specific models, LLMs can generalize across many language tasks, making them highly versatile tools for applications ranging from chatbots to content generation and coding assistance. These models highlight the growing capabilities of AI to interact with language in ways that feel increasingly natural and intelligent.

1.2 Labeled vs. Unlabeled Data: The Foundation of Learning

1.2.1 What Are Labels and Why They Matter

As shown in Figure 2, datasets may either contain labels, in which case supervised learning is used, or lack labels, in which case unsupervised learning is applied. Supervised learning relies on input–output pairs, using the features and their corresponding correct answers to train a model, whereas unsupervised learning only uses the input features to uncover hidden patterns and structure within the data.

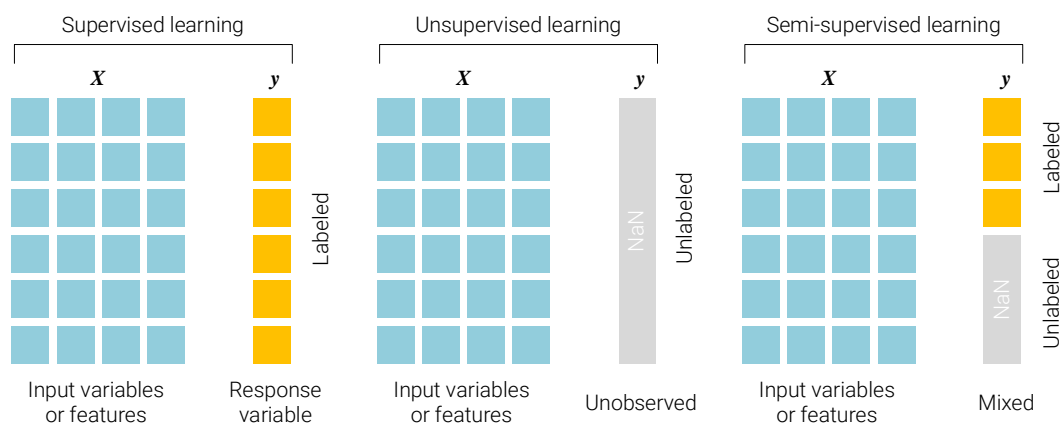


Figure 2. Data categorized by the presence of labels

1.2.2 Four Major Types of Machine Learning Tasks: Picking the Right Tool

Machine learning tasks can be broadly categorized based on the type of learning and the nature of the labels in the data. These categories help determine which algorithms and approaches are most appropriate for a given problem (see Figure 3).

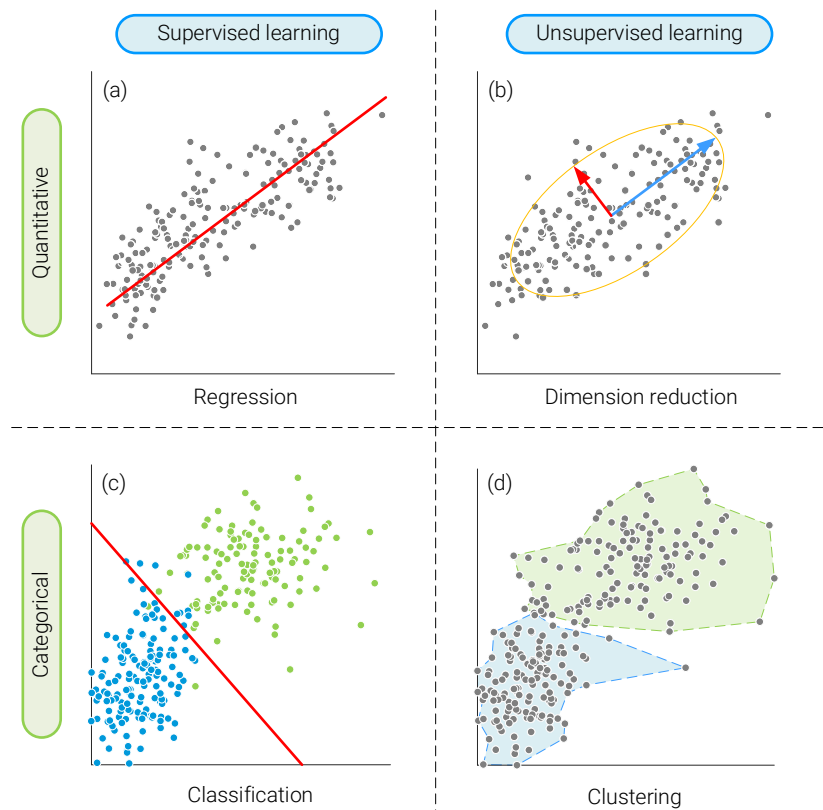


Figure 3. Overview of the four major types of machine learning tasks

In supervised learning, the data includes labels that the model tries to predict. When the labels are continuous values, the task is called regression, where the goal is to predict numerical outcomes such as prices, temperatures, or probabilities (Figure 4 (a)).

When the labels are categorical, the task is classification, which involves assigning samples to discrete categories like spam versus not-spam or different types of animals (Figure 4 (c)).

In unsupervised learning, the data does not include labels, so the model aims to discover patterns or structure within the input data. One common task is dimensionality reduction, where the goal is to simplify complex data while preserving its essential structure, often to make visualization or further analysis easier (Figure 4 (b)).

Another task is clustering, where the model groups samples based on similarity, identifying natural clusters within the data without any prior labels (Figure 4 (d)).

In summary, classification deals with discrete outputs, regression predicts continuous values, dimensionality reduction simplifies and summarizes the data, and clustering organizes unlabeled data into meaningful groups. Each type addresses a different type of problem, providing a framework for selecting appropriate machine learning methods.

1.3 Supervised Learning: Mapping Inputs to Outputs

1.3.1 Regression: Predicting Continuous Values

Regression aims to predict a continuous numerical value, such as a house price, temperature, or stock return. The simplest and most widely used regression method is linear regression, which assumes a straight-line relationship between the input features and the target variable. When this relationship is not linear, we use

nonlinear regression, where polynomial regression is a common example—it models curved relationships by adding powers of input features as new predictors.

Regularization methods, such as Ridge regression (L2 penalty), Lasso regression (L1 penalty), and Elastic Net (a combination of both), help prevent overfitting by constraining the model's coefficients. Bayesian regression introduces prior beliefs about model parameters and updates them with data, providing a probabilistic view of prediction.

Some classification algorithms, such as k -nearest neighbors (kNN) and support vector machines (SVM), can be adapted for regression tasks as well. For instance, support vector regression (SVR) applies the principles of SVM to continuous outcomes, seeking an optimal function that fits the data within a certain margin of tolerance.

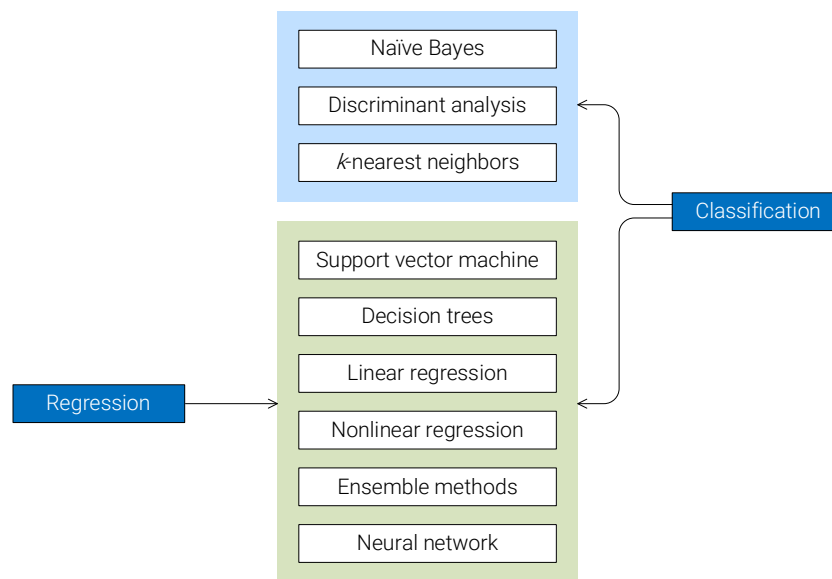


Figure 4. Commonly used algorithms in supervised learning

1.3.2 Classification: Predicting Categories

Classification focuses on predicting a discrete label, such as determining whether an email is “spam” or “not spam.” Various algorithms approach this task in different ways.

The k -Nearest Neighbors (kNN) method classifies a new sample by examining the labels of its closest neighbors in feature space, assigning the most common label among them.

Naïve Bayes applies Bayes’ theorem while assuming that features are conditionally independent, making it a simple yet effective choice for text and document classification.

Discriminant Analysis models the probability distribution of each class and uses these distributions to identify the most likely class for a new observation.

Support Vector Machines (SVM) aim to find the optimal boundary, or hyperplane, that separates classes with the maximum margin and can handle nonlinear boundaries using kernel functions.

Decision Trees take a different approach by recursively splitting data based on feature values to create a tree-like structure, which can be easily visualized and interpreted. Together, these methods provide a diverse toolkit for handling classification tasks in machine learning.

1.4 Unsupervised Learning: Exploring Patterns Without Labels

1.4.1 Dimensionality Reduction: Simplifying High-Dimensional Data

Dimensionality reduction is the process of transforming high-dimensional data into a lower-dimensional space while preserving its most important features. This not only reduces noise and redundancy but also lowers computational costs, accelerates model training, and often improves model interpretability.

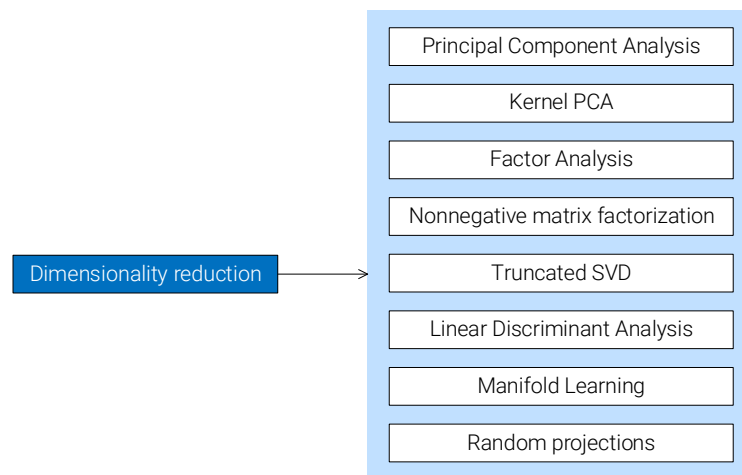


Figure 5. Commonly used dimensionality reduction techniques

Common techniques for dimensionality reduction aim to simplify high-dimensional data while preserving its most important structure.

Principal Component Analysis (PCA) is a linear method that identifies a new set of axes, called principal components, along which the data varies the most, and it can be computed using eigenvalue decomposition or singular value decomposition (SVD).

Kernel PCA (KPCA) extends this approach to handle nonlinear relationships by first mapping the data into a higher-dimensional space using a kernel function and then performing PCA in that space.

Truncated SVD is a variant of SVD designed to efficiently reduce dimensionality in sparse datasets, making it particularly useful for large-scale applications.

Linear Discriminant Analysis (LDA), although commonly applied in supervised learning, can also serve as a dimensionality reduction technique by finding projections that maximize the separation between classes, thereby simplifying the data while retaining discriminative information.

1.4.2 Clustering: Grouping Data Without Labels

Clustering is a technique that groups samples in a dataset based on their similarity, allowing meaningful patterns to emerge from unlabeled data. Algorithms for clustering approach this task in different ways.

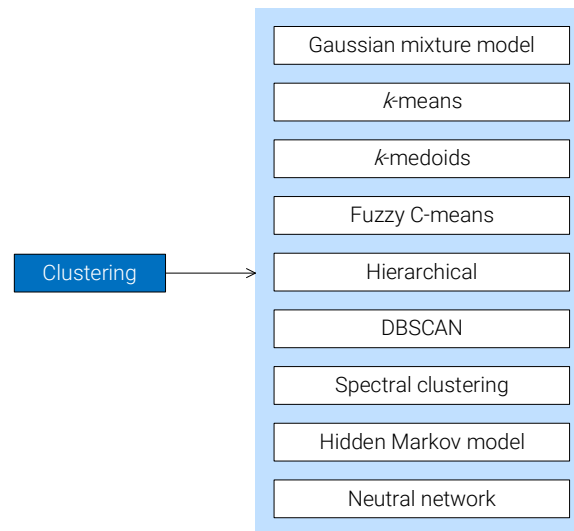


Figure 6. Commonly used clustering algorithms

k -Means partitions the data into a predefined number of clusters, each represented by the centroid of its points, while Gaussian Mixture Models (GMM) assume the data originates from a mixture of several Gaussian distributions and use the Expectation-Maximization algorithm to estimate the parameters of each cluster.

k -Medoids is similar to k -Means but selects actual data points as cluster centers, making it more robust to outliers. Hierarchical clustering constructs a tree of clusters either by merging smaller clusters in an agglomerative manner or by splitting larger clusters divisively.

DBSCAN, or Density-Based Spatial Clustering of Applications with Noise, forms clusters based on the density of points, allowing it to identify clusters of arbitrary shape and detect noise.

Spectral clustering, on the other hand, leverages the eigenvalues of a similarity matrix, often a Laplacian, to perform clustering in a transformed space, capturing more complex relationships between data points.

By understanding the differences between these methods and selecting the right one for your data, you can uncover meaningful structures and patterns that guide further analysis or predictive modeling.

1.5 Machine Learning Process: From Data to Deployment

The general flow of a machine learning project follows a series of structured steps, starting from raw data and ending with models deployed in real-world applications.

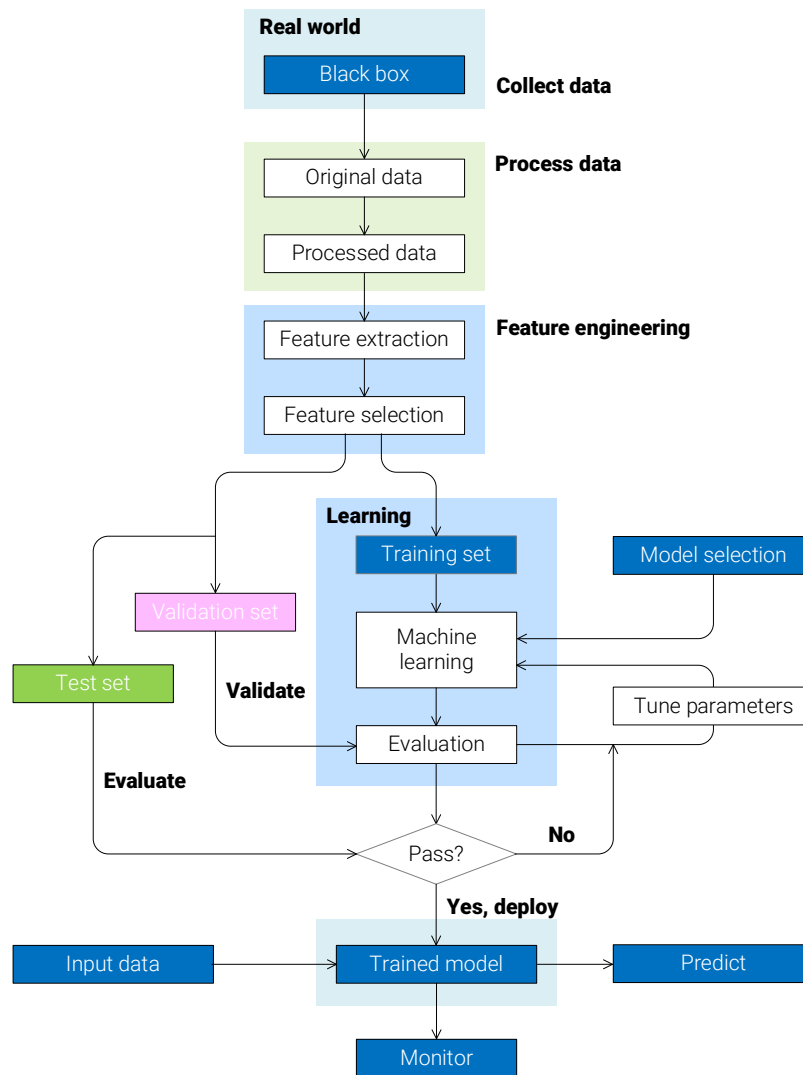


Figure 7. General flow of a machine learning project

First, data collection gathers datasets from various sources, which often requires cleaning to remove invalid entries, correct errors, and handle missing values.

Next, feature engineering transforms raw data into a suitable format for modeling, extracting informative attributes, selecting relevant features, scaling or normalizing values, and sometimes creating entirely new features. After the data is prepared, it is divided into training, validation, and test sets. The training set is used to build the model, the validation set is used to tune parameters and select the best model, and the test set evaluates the final model's performance.

Once the data is ready, the next step is to choose a model appropriate for the problem, such as linear regression, decision trees, or neural networks. The selected model is then trained using the training set, and techniques like cross-validation help assess and optimize its performance.

After training, the model is tested on unseen data to evaluate its accuracy, robustness, and generalization ability. If necessary, the model can be refined by adjusting features, parameters, or even the model type itself.

Finally, the model is applied to real-world tasks, such as prediction, classification, or recommendation, and its performance is continuously monitored to ensure it remains effective as new data or conditions change.

1.6 Conclusion

This chapter provides a comprehensive introduction to machine learning, explaining its place within the broader field of artificial intelligence and its connections to deep learning and natural language processing. Machine learning enables computers to learn patterns from data and make predictions or decisions without explicit programming.

The chapter covers the distinction between labeled and unlabeled data, introducing supervised learning, which predicts outputs from known labels, and unsupervised learning, which discovers hidden structures in data. Key tasks such as regression, classification, dimensionality reduction, and clustering are described along with common algorithms used for each.

The chapter also outlines the machine learning process, including data collection, cleaning, feature engineering, model training, validation, testing, and deployment. Emphasis is placed on understanding different methods, selecting the appropriate approach, and ensuring data quality, highlighting that effective machine learning combines algorithmic knowledge, careful data preparation, and practical application.

Now, let us officially begin our journey through this book. We will explore common machine learning algorithms using clear visual illustrations and practical coding examples, which help reveal how these methods work and why they are effective. Along the way, we will also uncover the mathematical tools that underpin these algorithms, such as probability, statistics, linear algebra, and optimization, showing how theory translates into practical application.