

# Simplified Machine Learning

---

*The essential building blocks for  
Machine Learning expertise*

---

**Dr. Pooja Sharma**



[www.bpbonline.com](http://www.bpbonline.com)

First Edition 2024

Copyright © BPB Publications, India

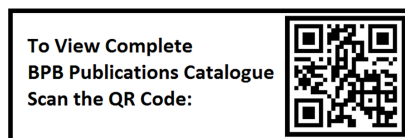
ISBN: 978-93-55516-145

*All Rights Reserved.* No part of this publication may be reproduced, distributed or transmitted in any form or by any means or stored in a database or retrieval system, without the prior written permission of the publisher with the exception to the program listings which may be entered, stored and executed in a computer system, but they can not be reproduced by the means of publication, photocopy, recording, or by any electronic and mechanical means.

### **LIMITS OF LIABILITY AND DISCLAIMER OF WARRANTY**

The information contained in this book is true to correct and the best of author's and publisher's knowledge. The author has made every effort to ensure the accuracy of these publications, but publisher cannot be held responsible for any loss or damage arising from any information in this book.

All trademarks referred to in the book are acknowledged as properties of their respective owners but BPB Publications cannot guarantee the accuracy of this information.



**Dedicated to**

*My students*

*and*

*My family*

## About the Author

**Dr. Pooja Sharma**, Assistant Professor, in Computer Science and Engineering has teaching and research experience of more than 17 years. She is a gold medalist in post-graduation and her other academic achievements include a fellowship for a regular PhD from UGC, New Delhi after qualifying UGC NET and JRF, several merit certificates, gold and silver medals in matric, higher secondary, undergraduate and postgraduate levels. She was awarded PhD in 2013 on Content-Based Image Retrieval under the supervision of Dr. Chandan Singh from Punjabi University, Patiala. She has several research publications in peer-reviewed International journals of Springer and Elsevier with significant Thomson Reuters impact factors. She is the author of various research book chapters and has published a book on “Programming in Python”. She is the reviewer of various International journals Elsevier, IET (IEEE Computer Society), and Scientific Research and Essays. She has participated in various conferences and workshops. Her areas of specialization include Data Analysis, Machine Learning, Content-Based Image Retrieval, Face Recognition, Pattern Recognition, and Digital Image Processing. She worked and was selected at various eminent Universities and Colleges including Central University. She had been the Head of Department at DAV University for 3 years. Currently, she holds a position as an Assistant Professor in the Department of Computer Science and Engineering, IKG Punjab Technical University, Main Campus, Kapurthala.

---

## Acknowledgement

I am always indebted and grateful to the Almighty for making me capable of writing this book. I extend my gratitude to my parents for their encouragement and support in every sphere of life. I would like to acknowledge my husband Rajkumar for his persistent encouragement and support throughout writing this book and my daughter Angel.

I convey my heartiest thanks to the authorities of IKG Punjab Technical University, Kapurthala and my colleagues in the Department of Computer Science and Engineering for their support and cooperation.

I would like to acknowledge BPB Publications for their direction and expertise in bringing this book to completion. It was a long journey of writing and revising this book, with valuable assistance and guidance from reviewers, technical experts, and editors.

Finally, I would like to thank all the readers who have taken an interest in my book and for their support in making it a reality. Your encouragement has been invaluable.

# Preface

Welcome to the world of Machine Learning! This book is designed to be your companion on a journey through the exciting and rapidly evolving field of Artificial Intelligence. Whether you are a student, researcher, or industry professional, this comprehensive guide aims to demystify complex concepts and equip you with practical skills to navigate the vast landscape of Machine Learning.

In recent years, Machine Learning has emerged as a transformative force, revolutionizing industries and shaping our everyday experiences. From personalized recommendations on streaming platforms to advanced medical diagnostics, the applications of Machine Learning are ubiquitous and ever-expanding. As you embark on this learning adventure, our goal is to provide you with a solid foundation in the fundamental principles, algorithms, and techniques that power these innovations.

This book is structured to cater to learners at all levels. We start by laying down the groundwork, explaining core concepts such as supervised and unsupervised learning, regression, classification, and clustering. We then dive into the intricacies of various Machine Learning algorithms, including decision trees, support vector machines, neural networks, and association rule mining. Each chapter is enriched with intuitive explanations, illustrative examples, and hands-on exercises to reinforce your understanding.

Beyond theory, this book emphasizes practical applications. We walk you through real-world case studies, providing insights into deploying Machine Learning models, interpreting results, and addressing ethical considerations. Whether you are interested in healthcare, finance, or natural language processing, you will find actionable insights and project-based learning exercises tailored to diverse domains.

As you embark on this journey into the world of Machine Learning, we invite you to embrace curiosity, embrace challenges, and embrace the possibilities that await. Let us embark on this adventure together and discover the transformative potential of Machine Learning. The Chapterization of the book is given as follows:

**Chapter 1: Introduction to Machine Learning** – This chapter covers Machine Learning and how it is related to Artificial Intelligence. Various types of Machine Learning along with their applications are also discussed in this chapter. Furthermore, the step-by-step guide is given to install Python Jupyter for implementing Machine Learning algorithms.

**Chapter 2: Data Pre-processing** – This chapter covers the pre-processing of data before applying any of the Machine Learning algorithms. This chapter covers various types of datasets, the need for data pre-processing, data cleaning, data transformation, data splitting, data normalization and scaling, data integration and aggregation, text processing and more.

**Chapter 3: Supervised Learning: Regression** - This chapter covers the detailed concept of supervised learning. Various types of regression techniques are discussed such as simple linear regression, multiple linear regression, ridge regression, lasso regression, polynomial regression, and applications of regression are discussed with appropriate programming codes and examples. Various model evaluation methods and errors are also covered in this chapter.

**Chapter 4: Supervised Learning: Classification** - This chapter covers the fundamental concepts related to the second major technique of supervised learning i.e., classification. This chapter includes logistic regression, K nearest neighbours, decision tree, random forest, naïve bayes classifier, and support vector machines.

**Chapter 5: Unsupervised Learning: Clustering** - This chapter covers unsupervised learning and the most important concept under it, i.e., clustering. The needs and applications of clustering are also discussed in detail. All types of clustering such as partition-based, hierarchical, and density-based clustering techniques are explored. Principal component analysis and anomaly detection are also covered in this chapter.

**Chapter 6: Dimensionality Reduction and Feature Selection** – This chapter covers feature engineering using dimensionality reduction, feature selection, and recursive feature elimination method and much more.

**Chapter 7: Association Rule Mining** - This chapter covers how association is important in Machine Learning. This chapter explains the apriori algorithm in detail with the key concept of association rule mining.

**Chapter 8: Artificial Neural Network** - This chapter covers the comprehensive concept of artificial neural networks (ANN) which leads to deep learning. All the prominent ANNs viz. perceptron, feedforward, backpropagation, convolutional, recurrent, long short-term memory, gated recurrent, and autoencoders are explored in detail with their programming codes.

**Chapter 9: Reinforcement Learning** - This chapter covers reinforcement learning with its components and applications. Various methods such as epsilon greedy, softmax, upper confidence bound, and Markov decision process are explored in this chapter.

**Chapter 10: Project** – This chapter covers real-world problems for a deep understanding of Machine Learning models. A comprehensive practical approach is followed to apply a Machine Learning model with graphical results and measurement of accuracy.

## Code Bundle and Coloured Images

Please follow the link to download the *Code Bundle* and the *Coloured Images* of the book:

**<https://rebrand.ly/7be7b5>**

The code bundle for the book is also hosted on GitHub at

**[https://github.com/bpbpublications/Simplified-Machine Learning](https://github.com/bpbpublications/Simplified-Machine-Learning)**.

In case there's an update to the code, it will be updated on the existing GitHub repository.

We have code bundles from our rich catalogue of books and videos available at **<https://github.com/bpbpublications>**. Check them out!

## Errata

We take immense pride in our work at BPB Publications and follow best practices to ensure the accuracy of our content to provide with an indulging reading experience to our subscribers. Our readers are our mirrors, and we use their inputs to reflect and improve upon human errors, if any, that may have occurred during the publishing processes involved. To let us maintain the quality and help us reach out to any readers who might be having difficulties due to any unforeseen errors, please write to us at :

**[errata@bpbonline.com](mailto:errata@bpbonline.com)**

Your support, suggestions and feedbacks are highly appreciated by the BPB Publications' Family.

Did you know that BPB offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at [www.bpbonline.com](http://www.bpbonline.com) and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at :

**[business@bpbonline.com](mailto:business@bpbonline.com)** for more details.

At **[www.bpbonline.com](http://www.bpbonline.com)**, you can also read a collection of free technical articles, sign up for a range of free newsletters, and receive exclusive discounts and offers on BPB books and eBooks.



## Piracy

If you come across any illegal copies of our works in any form on the internet, we would be grateful if you would provide us with the location address or website name. Please contact us at **business@bpbonline.com** with a link to the material.

## If you are interested in becoming an author

If there is a topic that you have expertise in, and you are interested in either writing or contributing to a book, please visit **www.bpbonline.com**. We have worked with thousands of developers and tech professionals, just like you, to help them share their insights with the global tech community. You can make a general application, apply for a specific hot topic that we are recruiting an author for, or submit your own idea.

## Reviews

Please leave a review. Once you have read and used this book, why not leave a review on the site that you purchased it from? Potential readers can then see and use your unbiased opinion to make purchase decisions. We at BPB can understand what you think about our products, and our authors can see your feedback on their book. Thank you!

For more information about BPB, please visit **www.bpbonline.com**.

## Join our book's Discord space

Join the book's Discord Workspace for Latest updates, Offers, Tech happenings around the world, New Release and Sessions with the Authors:

<https://discord.bpbonline.com>



# Table of Contents

<b>1. Introduction to Machine Learning.....</b>	<b>1</b>
Introduction.....	1
Structure.....	2
Objectives .....	2
Need for Machine Learning.....	2
Relation between Artificial Intelligence and Machine Learning .....	4
<i>Automation</i> .....	4
<i>Adaptability</i> .....	4
<i>Natural language processing</i> .....	4
<i>Computer vision</i> .....	5
Types of Machine Learning.....	5
<i>Supervised learning</i> .....	5
<i>Unsupervised learning</i> .....	5
<i>Semi-supervised learning</i> .....	5
<i>Reinforcement learning</i> .....	6
<i>Self-supervised learning</i> .....	6
<i>Online learning</i> .....	6
<i>Meta-learning</i> .....	6
<i>Ensemble learning</i> .....	6
Applications of Machine Learning .....	6
Lifecycle of Machine Learning .....	8
Steps to install Anaconda and Python.....	9
Conclusion.....	14
Questions.....	14
<b>2. Data Pre-processing .....</b>	<b>15</b>
Introduction.....	15
Structure.....	15
Objectives .....	16
Datasets.....	16
<i>CSV file</i> .....	17
Sources to obtain datasets .....	18

---

Need of data pre-processing.....	19
Data pre-processing .....	21
Data cleaning.....	21
Data transformation.....	22
<i>Data integration and aggregation</i> .....	22
Data splitting.....	22
<i>Handling imbalanced data</i> .....	22
Data normalization and scaling .....	23
<i>Text pre-processing for NLP</i> .....	23
<i>Handling missing data</i> .....	23
Data cleaning in detail .....	24
<i>Handling duplicates</i> .....	25
<i>Handling inconsistent data</i> .....	26
<i>Handling outliers</i> .....	28
<i>Noise reduction</i> .....	29
Data transformation.....	31
<i>Encoding categorical data</i> .....	31
<i>Feature engineering</i> .....	34
<i>Data reduction</i> .....	35
<i>Datetime transformation</i> .....	36
<i>Log transformation</i> .....	38
Data integration.....	39
<i>Aggregation</i> .....	40
<i>Data splitting</i> .....	42
<i>Normalization</i> .....	43
<i>Standardization</i> .....	44
<i>Standardizing data formats</i> .....	45
<i>Text processing using NLP</i> .....	46
Binning and discretization.....	48
Conclusion.....	50
Programming exercises .....	50
Questions.....	51
<b>3. Supervised Learning: Regression .....</b>	<b>53</b>
Introduction.....	53
Structure.....	53

Objectives .....	54
Key characteristics of supervised learning .....	54
Regression.....	55
Need of regression.....	55
<i>Applications of regression</i> .....	56
Simple linear regression .....	59
<i>Key points about simple linear regression</i> .....	60
Multiple linear regression .....	62
Polynomial regression .....	66
Ridge regression .....	69
Lasso regression.....	72
<i>Performance evaluation of regression models</i> .....	76
Conclusion.....	79
Programming exercises .....	79
Questions.....	79
<b>4. Supervised Learning: Classification .....</b>	<b>81</b>
Introduction.....	81
Structure.....	81
Objectives .....	82
Key elements of classification.....	82
Need for classification .....	82
<i>Applications of classification</i> .....	83
Logistic regression.....	83
<i>Key characteristics of logistic regression</i> .....	84
K nearest neighbors.....	87
Decision tree.....	91
<i>Key characteristics of decision trees</i> .....	91
<i>How decision trees work</i> .....	91
<i>Advantages of decision trees</i> .....	92
<i>Disadvantages of decision trees</i> .....	92
<i>Tree induction algorithm</i> .....	94
<i>Classification and regression trees</i> .....	95
<i>Split algorithm based on information theory</i> .....	97
<i>ID3 algorithm</i> .....	97

Split algorithm based on Gini index.....	100
Random forest.....	103
<i>Key concepts</i> .....	103
<i>Working principle of random forest</i> .....	103
<i>Advantages of random forest</i> .....	104
<i>Limitations and considerations</i> .....	104
Naïve Bayes classifier .....	106
<i>Working principle of Naïve Bayes</i> .....	106
<i>Types of Naïve Bayes classifiers</i> .....	106
Support vector machine .....	108
<i>Working principle of SVM</i> .....	108
Evaluation metrics for classifiers .....	110
Conclusion.....	113
Programming exercises .....	113
Questions.....	114
<b>5. Unsupervised Learning: Clustering .....</b>	<b>117</b>
Introduction.....	117
Structure.....	117
Objectives .....	118
Need for clustering.....	118
Applications of clustering .....	118
Partition-based clustering methods.....	119
<i>K-means clustering</i> .....	120
<i>Example</i> .....	120
Hierarchical clustering.....	123
<i>Key characteristics of hierarchical clustering methods</i> .....	123
<i>Applications of hierarchical clustering methods</i> .....	123
Density-based clustering .....	125
<i>Key characteristics density-based clustering methods</i> .....	126
<i>Applications of density-based clustering methods</i> .....	126
Principal component analysis.....	129
<i>Anomaly detection</i> .....	131
<i>Types of anomalies</i> .....	131
<i>Applications of anomaly detection</i> .....	132

---

<i>Methods for anomaly detection</i> .....	132
Conclusion.....	134
Programming exercises .....	134
Questions.....	135
<b>6. Dimensionality Reduction and Feature Selection.....</b>	<b>137</b>
Introduction.....	137
Structure.....	137
Objectives .....	138
Need for dimensionality reduction .....	138
<i>Feature importance</i> .....	139
Recursive feature elimination.....	140
<i>Working of RFE</i> .....	140
<i>Advantages of RFE</i> .....	140
Feature selection .....	141
<i>Variance threshold</i> .....	142
<i>SelectKBest</i> .....	142
<i>Recursive feature elimination</i> .....	144
Trade-offs in dimensionality reduction.....	144
Conclusion.....	146
Programming exercises .....	146
Questions.....	147
<b>7. Association Rule Mining.....</b>	<b>149</b>
Introduction.....	149
Structure.....	149
Objectives .....	150
Brief working principle of association rule mining.....	150
Need for association rule mining.....	150
Real-time example of association rule mining.....	151
<i>Real-time example 1: Retail market</i> .....	151
<i>Association rule mining</i> .....	151
<i>Real-time example 2: Online streaming service</i> .....	152
<i>Association rule mining</i> .....	152
<i>Applications of association rule mining</i> .....	153
Algorithms for association rule mining .....	154

Apriori algorithm .....	155
<i>Apriori algorithm steps</i> .....	156
Conclusion.....	158
Programming exercises .....	158
Questions.....	159
<b>8. Artificial Neural Network.....</b>	<b>161</b>
Introduction.....	161
Structure.....	162
Objectives .....	162
Components of an artificial neural network .....	162
<i>Working on an artificial neural network</i> .....	163
Types of artificial neural networks .....	164
<i>Need for artificial neural networks</i> .....	165
<i>Applications of artificial neural networks</i> .....	166
Activation functions.....	167
<i>Significance of activation functions</i> .....	171
Perceptron neural network .....	171
Feedforward neural network.....	173
Backpropagation network.....	175
<i>Key concepts of backpropagation</i> .....	176
Backpropagation process .....	176
Convolutional neural network.....	179
<i>Key components of a convolutional neural network</i> .....	180
Recurrent neural network.....	182
<i>Key components of a recurrent neural network</i> .....	182
Long short-term memory network.....	184
<i>Long short-term memory network structure</i> .....	184
<i>Learning in long short-term memory networks</i> .....	185
Gated recurrent unit network.....	187
<i>Key components of a gated recurrent unit</i> .....	187
<i>Learning in gated recurrent unit networks</i> .....	188
Autoencoders.....	189
<i>Autoencoder architecture</i> .....	190
<i>Learning in autoencoders</i> .....	191

---

Conclusion.....	193
Programming exercises .....	193
Questions.....	194
<b>9. Reinforcement Learning .....</b>	<b>195</b>
Introduction.....	195
Structure.....	196
Objectives .....	196
Example of reinforcement learning .....	196
<i>Environment</i> .....	196
<i>Agent</i> .....	196
<i>Key components of RL</i> .....	197
<i>Reinforcement learning process</i> .....	197
Applications of reinforcement learning .....	198
<i>Major algorithms of reinforcement learning</i> .....	199
Q-values, V-values and Alpha-values .....	200
Exploration vs. exploitation.....	201
<i>Q-learning algorithm</i> .....	202
<i>Epsilon greedy algorithm</i> .....	205
Softmax action selection .....	206
Upper confidence bound.....	208
Markov decision process.....	210
Conclusion.....	213
Programming exercises .....	213
Questions.....	214
<b>10. Project.....</b>	<b>215</b>
Introduction.....	215
Structure.....	216
Objectives .....	216
Detailed Machine Learning project.....	216
Conclusion.....	232
<b>Appendix .....</b>	<b>233</b>
<b>Bibliography .....</b>	<b>243</b>
<b>Index.....</b>	<b>245-250</b>



# CHAPTER 1

# Introduction to Machine Learning

## Introduction

**Machine Learning (ML)** is a transformative field within **Artificial Intelligence (AI)** that empowers computers to learn and make predictions or decisions without explicit programming. At its core, it simulates the human learning process by using data and mathematical algorithms to identify patterns, make inferences, and improve performance over time. In essence, ML allows computers to generalize from data. It starts with a dataset containing examples and corresponding outcomes, and the ML model learns to recognize underlying patterns or relationships within this data. These patterns can range from recognizing handwritten characters, predicting stock prices and diagnosing diseases from medical scans, to recommending products based on user behavior.

Supervised learning, one of the core branches of ML, involves training a model on labelled data, where the correct outcomes are known. In contrast, unsupervised learning deals with unlabeled data, aiming to discover hidden structures or groupings. Reinforcement learning focuses on training agents to make decisions by interacting with their environment and receiving feedback. Based on this, ML has a broad array of applications across industries, from healthcare and finance to e-commerce and self-driving cars. Its success is driven by advances in computing power, the availability of massive datasets, and improvements in algorithms. Popular ML libraries and frameworks like TensorFlow and sci-kit-learn have democratized the field, enabling researchers and developers to build and deploy powerful models.

Therefore, we can say that ML is the science of teaching computers to learn from data, paving the way for intelligent systems that can automate tasks, make predictions, and continually improve their performance. As it continues to evolve, ML holds immense potential to revolutionize various aspects of our lives and industries.

## Structure

The chapter includes the following topics:

- Need for Machine Learning
- Relation between Artificial Intelligence and Machine Learning
- Types of Machine Learning
- Applications of Machine Learning
- Lifecycle of Machine Learning
- Steps to install Anaconda and Python

## Objectives

At the end of this chapter, you will be able to understand about basic concept of ML and why it is needed. You will go through all the types of ML and its major applications. Apart from that you will learn about ML lifecycle and steps to install Anaconda for implementing ML algorithms using Python libraries.

## Need for Machine Learning

ML is a powerful and versatile field that offers numerous benefits and opportunities, making it a compelling choice for various applications and industries. Here are some key reasons why ML is widely adopted:

- **Data-driven insights:** ML excels at extracting valuable insights and patterns from vast amounts of data. It can uncover trends and relationships that may not be apparent through traditional statistical analysis.
- **Automation:** ML algorithms can automate repetitive and labor-intensive tasks, freeing up human resources for more creative and strategic work. This is particularly valuable in industries like manufacturing, finance, and customer service.
- **Personalization:** ML enables businesses to provide personalized experiences to customers. This includes tailored product recommendations, content suggestions, and targeted marketing campaigns, which can improve customer satisfaction and retention.

- **Efficiency:** ML can optimize processes and resource allocation, leading to cost savings and improved operational efficiency. For example, predictive maintenance in manufacturing can reduce downtime and maintenance costs.
- **Scalability:** ML algorithms can handle large-scale data analysis and decision-making, making them suitable for applications ranging from e-commerce to healthcare.
- **Improved decision-making:** ML models can make data-driven decisions in real-time, which can be invaluable in fields like finance for algorithmic trading or in healthcare for treatment recommendations.
- **Problem solving:** ML can tackle complex problems that may have no straightforward algorithmic solution. This includes tasks like image recognition, language translation, and game playing.
- **Adaptability:** ML models can adapt to changing conditions and new data, allowing systems to remain relevant and effective over time.
- **Innovation:** ML has led to breakthroughs in various domains, including autonomous vehicles, natural language processing, and medical diagnostics, driving innovation across industries.
- **Competitive advantage:** Organizations that harness the power of ML can gain a competitive edge by offering better products, services, and customer experiences.
- **Scientific discovery:** In fields like genomics and materials science, ML accelerates research by analyzing complex datasets and predicting discoveries.
- **Accessibility:** With the availability of open-source ML libraries and cloud-based ML platforms, businesses and researchers have easy access to powerful tools and resources.
- **Sustainability:** ML can be used to optimize resource usage, reduce waste, and support sustainability efforts in areas such as agriculture, energy management, and transportation.
- **Healthcare advancements:** In healthcare, ML assists in disease diagnosis, drug discovery, and personalized treatment plans, potentially saving lives and improving patient outcomes.
- **Cybersecurity:** ML helps organizations detect and respond to cybersecurity threats by identifying anomalies and patterns in network traffic and user behavior.

ML offers the potential to solve complex problems, improve efficiency, and drive innovation across a wide range of fields. Its ability to learn from data and make data-driven decisions makes it a valuable tool for businesses, researchers, and industries looking to harness the power of data and automation to achieve their goals.

# Relation between Artificial Intelligence and Machine Learning

AI and ML are closely related fields, with ML being a subset of AI. Here is how they are connected:

- **AI:** It refers to the broader concept of creating intelligent machines that can mimic human-like cognitive functions, such as reasoning, problem-solving, learning, perception, and language understanding. It encompasses a wide range of techniques, including rule-based systems, expert systems, knowledge representation, and more.
- **ML:** It is a subset of AI that focuses on the development of algorithms and statistical models that enable computers to learn from and make predictions or decisions based on data. It is a specific approach within AI that deals with the learning aspect.
- **Learning from data:** ML is a fundamental component of many AI systems. It provides the ability for AI systems to learn patterns, behaviors, and insights from large datasets.

AI systems may use ML techniques to improve their performance or adapt to changing conditions.

## Automation

AI and ML often go hand in hand in automating tasks and decision-making. AI systems can use ML models to make informed decisions based on data.

For example, in autonomous vehicles, AI algorithms use ML models to process sensor data and make real-time driving decisions.

## Adaptability

ML enables AI systems to adapt and improve their performance over time. AI systems can learn from new data and adjust their behavior accordingly. This adaptability is crucial for AI systems to handle complex and dynamic environments.

## Natural language processing

**Natural Language Processing (NLP)** is a subset of AI that deals with human language understanding and generation. ML plays a significant role in NLP, as it is used to build models for tasks like language translation, sentiment analysis, and chatbots.

## Computer vision

Computer vision is another AI subfield that focuses on enabling computers to interpret and understand visual information from the world. ML, particularly deep learning, has revolutionized computer vision, allowing AI systems to recognize objects, faces, and scenes in images and videos.

ML is a core component of many AI systems, providing them with the ability to learn, adapt, and make data-driven decisions. While AI encompasses a broader set of goals and techniques, ML is a crucial tool within the AI toolkit, enabling AI systems to perform tasks that involve learning from data and improving their performance over time.

## Types of Machine Learning

ML can be categorized into several types, each with its approach and characteristics. The primary types of ML are as follows.

### Supervised learning

In supervised learning, the algorithm is trained on a labelled dataset, where each input example is paired with its corresponding output or target. The goal is to learn a mapping from inputs to outputs, allowing the model to make predictions or classifications on new, unseen data. Common algorithms include linear regression, logistic regression, decision trees, support vector machines, and neural networks.

### Unsupervised learning

Unsupervised learning deals with unlabeled data, where the algorithm seeks to discover patterns, structures, or relationships within the data without explicit guidance. It is often used for tasks like clustering (grouping similar data points) and dimensionality reduction (reducing the number of features while preserving important information). Common algorithms include k-means clustering, hierarchical clustering, **Principal Component Analysis (PCA)**, and autoencoders.

### Semi-supervised learning

Semi-supervised learning is a combination of supervised and unsupervised learning. It uses both labelled and unlabeled data to improve model performance. This is especially useful when obtaining labelled data is expensive or time-consuming. Techniques may include using a small amount of labelled data to guide the model's learning on the larger unlabeled dataset.