

Chapter 1

Basic AI Concepts and Terminology

THE AWS CERTIFIED AI PRACTITIONER EXAM OBJECTIVES COVERED IN THIS CHAPTER MAY INCLUDE:

✓ **Domain 1: Fundamentals of AI and ML**

- Task Statement 1.1: Explain basic AI concepts and terminology.





Artificial intelligence (AI) has become one of the most transformative technologies of our time, revolutionizing industries, reshaping our daily lives, and challenging our understanding of human cognition. But what exactly is AI, and how has our understanding of it evolved over the decades?

At its core, artificial intelligence refers to the development of computer systems capable of performing tasks that typically require human intelligence. These tasks include visual perception, speech recognition, decision-making, and language translation, among others. However, the definition of AI has been far from static, evolving significantly since the term was first coined in 1956 at the Dartmouth Conference. AI has come a long way since then—at the time of writing this book, the 2024 Nobel Prize in Physics has been awarded to John Hopfield and Geoffrey Hinton for their fundamental discoveries in Machine Learning, which have directly or indirectly lead to all the advanced AI that we see around us today.

The First Nobel Prize for Work in AI

The 2024 Nobel Prize in Physics shines a spotlight on the fascinating intersection of artificial intelligence and fundamental physics. John J. Hopfield and Geoffrey E. Hinton, two visionaries in the field, have been honored for their groundbreaking contributions that have reshaped our understanding of both machine and human intelligence.

Hopfield's ingenious creation, the Hopfield network, emerged in the 1980s as a revolutionary approach to mimicking the human brain's ability to store and recall information. This clever system uses binary patterns to represent memories, much like the on-off firing of neurons in our brains. What makes Hopfield's work particularly intriguing is its ability to retrieve complete memories from partial or noisy inputs—imagine recognizing a friend's face in a blurry photo, and you'll get the idea of what these networks can do.

Geoffrey Hinton, affectionately dubbed the "godfather of deep learning," has been instrumental in bringing neural networks to life. His co-invention of the backpropagation algorithm in the 1980s was nothing short of revolutionary. This mathematical technique allows neural networks to learn from their mistakes, gradually refining their performance through a process of trial and error. It's akin to how humans learn—by recognizing our errors and adjusting our approach.

The impact of Hopfield and Hinton’s work extends far beyond the realm of computer science. Their theories have paved the way for AI to tackle complex problems in fields as diverse as healthcare, autonomous driving, and even artistic creation. Imagine AI systems that can detect diseases from medical images with uncanny accuracy, or virtual assistants that can engage in natural conversations—these are direct descendants of the foundational work done by our Nobel laureates.

As we look to the future, the theories proposed by Hopfield and Hinton continue to fuel innovation in AI. Their work has opened doors to understanding not just artificial intelligence, but human intelligence as well. Who knows? The next breakthrough in understanding consciousness or cognition might very well stem from the seeds planted by these pioneering minds.

In its early days, AI was primarily focused on mimicking human reasoning through symbolic logic and rule-based systems. This approach, known as “Good Old-Fashioned AI” (GOFAI), aimed to create machines that could solve complex problems by manipulating symbols according to predefined rules. The goal was to build systems that could reason about the world in a way similar to humans.

As research progressed, it became clear that many aspects of human intelligence, such as perception and learning, were not easily reducible to symbolic manipulation. This realization led to a shift in focus toward machine learning approaches, where systems learn from data rather than being explicitly programmed. This paradigm shift marked a significant evolution in the definition of AI.

Today, the field of AI encompasses a broad spectrum of approaches and techniques. Modern AI systems can learn from vast amounts of data, recognize patterns, make predictions, and even generate human-like text or images. The definition of AI has expanded to include not just systems that mimic human reasoning, but also those that can perform tasks beyond human capabilities in specific domains.

One of the most significant developments in recent years has been the rise of deep learning, a subset of machine learning based on artificial neural networks. Deep learning has enabled breakthroughs in areas such as computer vision and natural language processing, pushing the boundaries of what we consider possible in AI.

As we delve deeper into the world of AI, it’s crucial to understand the various components and concepts that make up this complex field. The following sections explore key terms and ideas that form the foundation of modern AI. The chapter starts by defining basic AI terms, including *machine learning* (ML) and *deep learning* (DL), and explores how these approaches differ from traditional AI methods. The chapter then describes the similarities and differences among AI, ML, and deep learning and discusses the concept of models and algorithms, the core components that enable AI systems to learn and make predictions. This leads to how models are trained (supervised, unsupervised, or reinforcement learning). The chapter covers the various types of inferencing (batch,

real-time, and async inference) and describes various types of data in AI models such as tabular, time-series, image, text, and so on.

By the end of this chapter, you'll have a comprehensive understanding of the key concepts in AI. Whether you're a student, professional, or simply curious about this transformative technology, this exploration of AI will provide you with the knowledge to navigate and contribute to the AI-driven world we live in.

Let's get started!

A Brief History of AI

Artificial intelligence (AI), the field dedicated to creating machines or systems that can simulate human intelligence, has captured humanity's imagination since the mid-20th century. From its conceptual beginnings in 1956 at the Dartmouth Conference, AI has represented our quest to understand and replicate human cognitive abilities. While the field encompasses various approaches and methodologies, one of its most successful and rapidly evolving branches has been machine learning.

While AI and machine learning are often used interchangeably in popular discourse, they represent distinct but related concepts. AI is the broader field that aims to create systems capable of intelligent behavior, encompassing everything from rule-based expert systems to natural language processing. Machine learning, on the other hand, is a specific approach within AI that focuses on developing algorithms that can learn from and make predictions based on data. Think of AI as the overarching goal of creating intelligent machines, while machine learning is one of the primary tools we use to achieve that goal.

Machine learning (ML) has evolved significantly over the past several decades, with its roots firmly planted in mathematics and statistics. Beginning in the early 1940s, foundational concepts like the McCulloch-Pitts neuron model, which represented a simple mathematical abstraction of how neurons might function, set the stage for later developments. In the 1950s, Alan Turing proposed the Turing Test, a benchmark for machine intelligence, and Frank Rosenblatt's invention of the perceptron provided the first neural network model. These early ideas, shown on the timeline in Figure 1.1, were crucial in forming the language and concepts that would later define the field of machine learning.

Throughout the 1960s and 1970s, symbolic AI emerged as a dominant paradigm, with experts working on the development of rule-based systems. However, the limitations of early neural networks were highlighted in the 1969 *Perceptrons* book by Minsky and Papert, which led to a decline in funding and interest in neural network research, a period often referred to as the "AI Winter." Despite this, significant progress continued in the 1980s, with backpropagation algorithms allowing networks to be trained more effectively and convolutional neural networks (CNNs) being applied to pattern recognition. Key events in this period are shown in Figure 1.2.

The 2000s marked the resurgence of deep learning, with advancements like Geoffrey Hinton's deep belief networks and the breakthrough success of AlexNet in 2012, which used deep learning to win the ImageNet competition. The introduction of transformers

FIGURE 1.1 Highlights from 1940s to 50s.

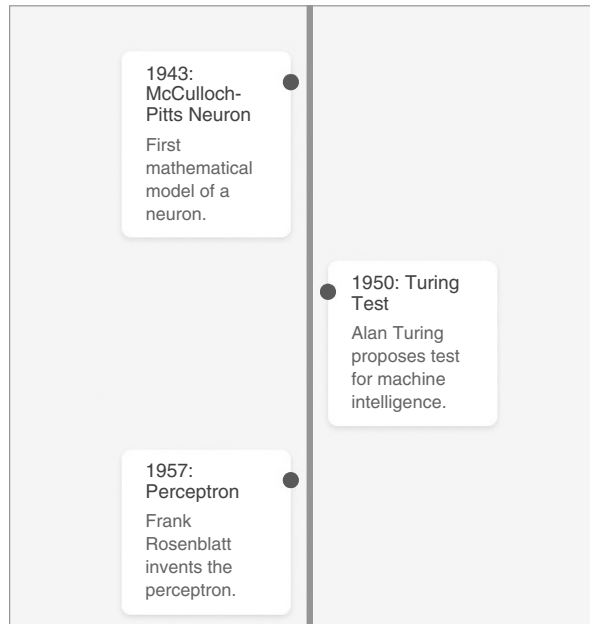
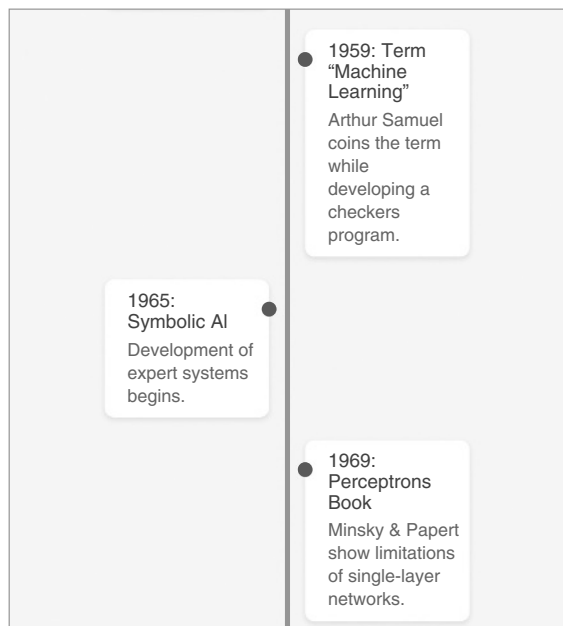


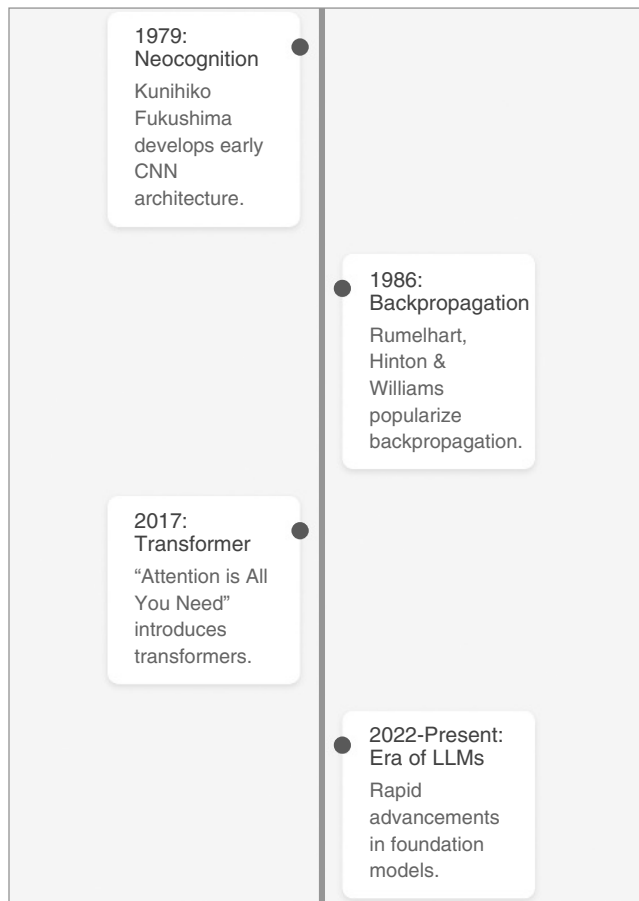
FIGURE 1.2 Highlights from 1960s to 1970s.



in 2017 revolutionized the natural language processing landscape. At the same time, the massive increase in computational power, facilitated by advances in cloud computing, GPUs, and distributed systems, played a pivotal role in scaling up machine learning models to the size we see today. With models like GPT-3 and the continued expansion of large language models (LLMs) and other types of foundation models (FMs), we are witnessing the rapid advancement of AI, fueled by the algorithms developed over decades and the unprecedented computational resources available now. A timeline of this revolution is shown in Figure 1.3.

Machine learning has truly come a long way, moving from theoretical models to practical applications that continue to impact nearly every industry. As ML continues to evolve, the combination of mathematical rigor, statistical models, and massive computational power will keep driving innovation forward. Although the history of AI and ML will not be tested

FIGURE 1.3 Highlights from 1980s to present.



in the exam, we believe this is useful and interesting context for you to know. Next, the chapter dives deeper into specific terms you should know in this domain, which set the right foundational understanding and make the rest of the book easier to read.

Diving Deeper into Terms You Should Know

As you saw from the timeline, the field of AI has evolved from early rule-based systems to today's sophisticated, large-scale neural networks, fundamentally changing our approach to problem-solving, automation, and human-computer interaction. For practitioners such as yourself, AI encompasses a broad spectrum of capabilities that you can bring to applications, from visual perception and speech recognition to complex decision-making and natural language understanding. What makes modern AI particularly powerful is its ability to handle uncertainty and adapt to new situations, moving beyond the rigid, predefined rules that characterized early computer programs.

The heart of modern AI lies in machine learning (ML), a revolutionary approach that transforms how people create intelligent systems. Rather than following explicit programming instructions, ML systems learn from experience, discovering patterns in data that would be impossible or impractical to encode manually. This paradigm shift has profound implications for system development and deployment. But how does this actually happen? When you're training a machine learning model, you're essentially creating a system that can recognize complex patterns in data and apply this knowledge to new, unseen situations (also called *generalization*). This approach has proven remarkably effective across diverse fields, from medical diagnoses and drug discovery to financial forecasting and climate modeling.

The Learning Paradigms: Supervised, Unsupervised, and Reinforcement Learning

Machine learning encompasses several distinct learning paradigms, each suited to different types of problems and data situations. *Supervised learning*, perhaps the most widely used approach, involves training models on labeled data—examples where you know the correct answer. These “correct answers” are often referred to as labels. This process mirrors how you might teach a child to recognize objects by showing them examples and providing the correct labels. In a supervised learning context, you present the algorithm with input-output pairs, such as images paired with their corresponding object labels, or housing features paired with their market prices. The algorithm learns to recognize patterns that map inputs to outputs, gradually refining its ability to make accurate predictions through sophisticated optimization techniques.

Consider a supervised learning system designed to detect fraudulent financial transactions. The system would be trained on a dataset of historical transactions, each labeled as either

legitimate or fraudulent. During training, the algorithm learns to identify patterns associated with fraudulent activity—perhaps unusual spending patterns, geographical anomalies, or suspicious timing. The power of supervised learning lies in its ability to discover subtle combinations of features that might not be obvious to human analysts.

The following table includes a small sample of transactional data. This represents a small excerpt of the massive data that might be needed to train your own fraud detection model. Look at the amounts, times, and merchant descriptions and finally the label that might be provided by an expert. Do you agree with the transactions that are labeled fraudulent? If you can spot the transactions that look unusual or out of place due to the amounts or locations, these patterns are the same signals a supervised learning algorithm might use to identify fraudulent transactions. However, note that for training real-world models, we expect to provide millions of rows of real transaction data, from which the models need to infer anomalous patterns without us specifying the rules of what makes each transaction anomalous. This is similar to teaching by example. With millions of rows, a single human may take years to go through the raw data and may end up not inferring any patterns from the data.

Transaction ID	Amount (\$)	Location	Time	Merchant	Labeled Fraud?
1	500	New York, USA	2:00 PM	Electronics Store	Legitimate
2	7500	Tokyo, Japan	3:00 AM	Jewelry Shop	Fraudulent
3	15	New York, USA	12:00 PM	Coffee Shop	Legitimate
4	10	London, UK	4:00 AM	Electronics Store	Fraudulent
5	150	New York, USA	6:00 PM	Grocery Store	Legitimate

Unsupervised learning tackles a more challenging and open-ended problem: finding structure in data without explicit guidance about what patterns to look for. This approach proves particularly valuable when you don't know exactly what insights might exist in the data. Consider the same transactional dataset shown earlier, but now imagine that the labels indicating whether a transaction is fraudulent or legitimate are missing. An unsupervised learning algorithm might analyze the data to identify transactions that deviate significantly from the norm. For example, it might flag unusually large amounts, like \$7,500 for a jewelry purchase, or unexpected behaviors, such as a purchase at 4:00 AM in a different country. These anomalies could serve as potential indicators of fraud, even without explicit labels.

By relying on statistical patterns and clustering techniques, unsupervised learning allows you to detect suspicious activity in situations where labeled examples are unavailable. This flexibility makes it a powerful tool for exploratory analysis and discovering insights that might otherwise go unnoticed.

Consider another example in customer behavior analysis—an unsupervised learning algorithm might discover natural groupings of customers based on hundreds of interaction patterns, purchase histories, and demographic features. These discovered patterns might reveal market segments that traditional analysis would miss, leading to more targeted and effective business strategies.

More About Supervised and Unsupervised Learning Algorithms

Several supervised learning algorithms have been developed that are powerful tools when labeled data is available. Logistic regression offers a simple and interpretable baseline, while decision trees and random forests handle complex, nonlinear relationships effectively. Gradient boosting methods, like XGBoost and LightGBM, are highly accurate and identify important patterns in transactional data. Neural networks are particularly suited to large datasets with intricate patterns, although they require more computational resources, and support vector machines are effective for datasets with clear class separations.

Unsupervised learning algorithms excel when labels are unavailable, focusing on detecting patterns and anomalies in the data. k-means clustering and DBSCAN group data into clusters, flagging transactions that deviate significantly as anomalies. Isolation forests isolate potential fraud cases by identifying unique data points, while autoencoders reconstruct data to find irregularities. Principal Component Analysis (PCA) reduces data dimensions, identifying outliers based on deviations from typical patterns. These approaches are critical for exploratory analysis and uncovering hidden fraud signals.

Although specific algorithms and their implementation are not tested in the exam, you are encouraged to dive deeper into examples of algorithms in each category.

Moving on, *reinforcement learning* represents a fundamentally different approach to machine learning, one that more closely mirrors how humans and animals learn through interaction with their environment. In this paradigm, an AI agent learns by taking actions and receiving feedback in the form of rewards or penalties. The complexity lies in the temporal nature of the learning process—the agent must learn to make sequences of decisions where the consequences of actions might only become apparent much later.

Consider an AI system learning to play chess through reinforcement learning. The system isn't told explicitly what moves are good or bad; instead, it learns by playing many games and receiving a reward only when it wins. The challenge lies in connecting individual moves to the eventual outcome—a move that appears advantageous in the short term might lead to a losing position several turns later. This “credit assignment problem” is one of the fundamental challenges in reinforcement learning, requiring sophisticated algorithms that can evaluate the long-term consequences of actions.

Although this section has covered the main types of learning—supervised, unsupervised, and reinforcement learning—there are other newer types of learning relevant to today’s applications in AI, such as self-supervised learning.

Self-Supervised Learning

Self-supervised learning is a training approach that enables AI models to learn from data without requiring manually labeled examples. In self-supervised learning, the model generates its own labels by leveraging inherent patterns or structures within the data. This approach is particularly useful for working with vast amounts of unlabeled data, as it reduces the dependency on costly and time-consuming labeling processes.

Self-supervised learning often involves predicting one part of the data from another. For example, in natural language processing, models can be trained to predict missing words in a sentence (like BERT) or the next word in a sequence (like GPT). In computer vision, models might learn by predicting the rotation angle of an image or by reconstructing a partially obscured section. These tasks provide the model with pseudo-labels that it uses to learn meaningful representations of the data.

This method has proven highly effective in applications such as speech recognition, image processing, and language understanding, where large amounts of unlabeled data are available. Self-supervised learning also facilitates transfer learning, as the models can be pre-trained on large, general datasets and later fine-tuned for specific tasks, significantly enhancing their adaptability and performance.

When to Use the Different Types of Learning

This section goes through a simple mental model for navigating through a given problem or use case with the aim of identifying which type of algorithm to use. For instance, if an AI practitioner has a well-labeled dataset where each data point is matched with a desired outcome, supervised learning is likely the best fit. This approach is ideal for applications where the goal is to make accurate predictions or classifications based on clear historical examples. For example, in a medical setting, if a dataset contains labeled images of X-rays marked as “healthy” and “diseased,” a supervised learning model can be trained to classify new images based on these labels, assisting with diagnosis.

In scenarios where labeled data is not available, unsupervised learning is often more appropriate. This approach is particularly useful when the goal is to uncover underlying patterns or groupings in the data, such as in market segmentation or anomaly detection. An AI practitioner analyzing customer data, for example, may use clustering techniques to identify different customer segments based on purchasing behavior. Here, unsupervised learning helps to reveal relationships and structures that were not immediately obvious, offering valuable insights without predefined outcomes.

For tasks involving sequential decision-making or scenarios where the AI system interacts with its environment, reinforcement learning can be a powerful tool. This approach is well-suited for dynamic environments where outcomes are influenced by a series of actions, and immediate feedback is available. A common use case is in training robots to navigate complex physical environments, where reinforcement learning allows the robot to learn optimal strategies through trial and error. Practitioners designing autonomous vehicles, for example, often use reinforcement learning to help the vehicle make real-time decisions, balancing safety, efficiency, and speed.

When the data available is a mix of labeled and unlabeled information or when the problem involves understanding both structure and classification, practitioners may even opt to combine supervised and unsupervised learning techniques. This can be the case in semi-supervised learning scenarios, where a model is initially trained on a small, labeled dataset and later fine-tuned with additional unlabeled data, as seen in applications like natural language processing.

These scenarios illustrate the flexibility of AI training methods and the importance of choosing an approach aligned with the specific nature of the data and the goals of the project. By carefully evaluating these factors, AI practitioners can select the most effective method to train their models, enabling them to maximize the value of their data and achieve the desired outcomes. You can take a deeper look into the various use cases across domains in Chapter 3.

The Deep Learning Revolution

The breakthrough that has revolutionized AI in recent years is deep learning, built upon artificial neural networks inspired by our understanding of biological neural processing. While the analogy to biological brains is often overstated, the architectural principles of neural networks have proven remarkably effective for machine learning tasks. The “deep” in deep learning refers to the use of multiple layers of artificial neurons, each capable of learning increasingly abstract representations of the input data. Neural networks are also called universal function approximators that can learn any mapping. Going back to the example of fraudulent transactions from a table, you can easily construct a deep learning neural network that “maps” the input transaction to a probability score of whether the transaction looks fraudulent.

Understanding how neural networks function requires grasping several key concepts. Each artificial neuron performs a mathematical operation, combining its inputs with learned weights and passing the result through an activation function. Let’s look at the simplest version of a neural network.

The most basic neural network is a single-layer network. It consists of input features, weights that are learned during training, a simple summation operation, and an output prediction. Here’s how it works in the simple terms:

```
# Input: features [x1, x2, ..., xn]
# Weights to be learned: [w1, w2, ..., wn]
prediction = w1*x1 + w2*x2 + ... + wn*xn
error = (true_value - prediction)^2 # squared error
```

When referring to *model training* or “learning” the parameters of a neural network, you try to minimize the prediction error (called Mean Squared Error or MSE) by adjusting the weights.

This foundation is crucial because it reveals several basic concepts on which more complex neural networks can be built:

1. The weights w_1 , w_2 , ... are parameters. You may have heard of recent multi-billion parameter foundation models and LLMs. These represent learnable or adjustable weights similar to the ones you see here for this simple network.
2. Features [x_1 , x_2 , ...] are being multiplied by weights [w_1 , w_2 , ...]. This set of weights represents a single transformation layer in the network. Adding multiple layers of neurons instead of one leads to “deep” neural networks.
3. Introducing nonlinear activation functions, other types of neurons, and connecting these neurons in different ways gives rise to new, state-of-the-art architectures like transformers.
4. Using more sophisticated optimization techniques lets the training or learning process complete faster since it is not practical to explore every combination of weights that may reduce the predicted error or loss (imagine doing this for unique combinations of float numbers for billions of parameters!).
5. This example simply multiplies the weights and features and does not use a nonlinear function on top of this—once again, this is just linear regression. Applying nonlinear functions allows the network to be more expressive and can lead to better accuracy. These are often called activation functions.

Think of activation functions as transformations that help neural networks learn more complex patterns. In this simple example, you can only learn straight-line relationships between inputs and outputs. But real-world problems often involve curved or more complex relationships. This is where activation functions come in.

Imagine you’re trying to model how plants grow based on sunlight. A simple linear relationship might predict that more sunlight always means more growth. But in reality, there’s usually an optimal amount of sunlight—too little or too much can hurt growth. To capture this kind of relationship, you need nonlinear activation functions.

One popular activation function is called ReLU (Rectified Linear Unit). It’s like a switch that lets positive values pass through unchanged but turns negative values to zero. This simple behavior helps neural networks learn complex patterns while still being easy to compute. Other activation functions like ReLU (such as the sigmoid, or hyperbolic tangent [tanh]) introduce non-linearity into the network, allowing it to learn complex patterns that couldn’t be captured by simpler linear models.

Training: Learning Model Parameters Values

At the core of every machine learning system is its model—the mathematical framework used to represent patterns in data. The concept of a model in machine learning is broader and more dynamic than traditional statistical models. While a statistical model might

describe a fixed relationship between variables, a machine learning model is typically more flexible, with its exact form determined through the learning process. This flexibility comes from parameters—the adjustable values that the model learns during training.

Understanding the role of parameters and hyperparameters is crucial for effective machine learning. Parameters are the values that the model learns directly from the data, such as the weights and biases in a neural network or the coefficients in a regression model. Hyperparameters, on the other hand, are the configuration settings that control how the model learns. These include learning rate (how quickly the model updates its parameters), batch size (how many training examples are processed at once), and architectural choices like the number of layers in a neural network. The choice of hyperparameters can dramatically affect both training efficiency and final model performance, making this a crucial skill for machine learning practitioners.

The concept of model “fit”—or how well a model captures the underlying patterns in data—is central to machine learning success. Overfitting occurs when a model learns the training data too well, capturing noise and random fluctuations rather than genuine patterns. This results in poor generalization—the model performs well on training data but fails on new examples. Underfitting, conversely, happens when a model is too simple to capture important patterns in the data, resulting in poor performance on both training and test data.

The training process for neural networks relies heavily on *backpropagation*, a foundational algorithm for computing how each parameter (recall the weights $[w_1, w_2 \dots]$) in the network should be adjusted to improve performance. Backpropagation works by propagating error signals backward through the network, using the chain rule of calculus to determine how each parameter contributed to the network’s mistakes. This process is combined with optimization techniques like gradient descent, which iteratively refines the network’s parameters to minimize prediction errors.

Modern neural network training involves numerous sophisticated techniques for improving learning efficiency and preventing overfitting. Batch normalization helps maintain stable input distributions throughout the network, making training deeper networks more practical. Dropout randomly deactivates neurons during training, forcing the network to learn redundant representations and preventing over-reliance on any particular feature. These and other techniques have made it possible to train networks with millions or even billions of parameters effectively.

A fundamental algorithm in machine learning is called “gradient descent,” which minimizes a model’s error or loss by iteratively adjusting parameters to reduce errors. Variations such as stochastic gradient descent and mini-batch gradient descent have been developed to handle large datasets and improve efficiency. Ensemble methods like random forests and boosting algorithms combine multiple models to create more robust and accurate systems, leveraging diversity to improve overall performance and reduce overfitting.

Backpropagation and Training Neural Networks

Backpropagation is a fundamental algorithm for training neural networks, enabling them to learn by adjusting their internal parameters based on the errors in their predictions. This process is essential for optimizing the network's weights, which determine how inputs are transformed through the network to produce outputs. By minimizing the difference between predicted and actual outcomes, backpropagation allows the network to improve its accuracy over time.

Backpropagation operates by calculating the gradient of the loss function with respect to each weight in the network. This is achieved through a series of steps:

- 1. Forward Pass:** The input data is fed through the network, layer by layer, to generate an output. This output is compared with the actual target, and the difference between them is calculated using a loss function. Common loss functions include Mean Squared Error for regression tasks and Cross-Entropy Loss for classification tasks.
- 2. Error Calculation:** Once the loss is computed, the algorithm measures how far off the output is from the expected result. This error signal will guide the adjustments needed in the network's weights.
- 3. Backward Pass (Backpropagation):** Here, the error is propagated back through the network to update the weights. The algorithm uses the chain rule of calculus to calculate the gradient of the loss function with respect to each weight. This allows the network to understand which weights contributed most to the error and by how much they should be adjusted.
- 4. Weight Update:** Using the calculated gradients, the weights are adjusted according to the chosen optimization algorithm, most commonly gradient descent. The adjustments are made in the direction that minimizes the loss function. The learning rate, a hyperparameter that defines the step size during weight updates, influences how quickly or slowly the model learns. If the learning rate is too high, the model might converge too quickly to a suboptimal solution or even diverge. If it is too low, the model may take too long to converge.
- 5. Iteration:** The entire process of forward pass, error calculation, backpropagation, and weight update is repeated over many iterations, or epochs, until the model achieves satisfactory accuracy or the error rate stops improving significantly.

Feature Engineering and Data Preprocessing

The success of any machine learning system heavily depends on how data is prepared and represented. Feature extraction and engineering represent crucial steps in the machine learning pipeline, often making the difference between success and failure. In the simple example with a single layer neuron from earlier, the values $[x_1, x_2, \dots]$ were called features. These can be raw data coming from a real-world system—for example, sensors from a factory or a system that records financial transactions. Raw data typically contains noise, irrelevant information, and suboptimal representations that can confuse learning algorithms. The art of feature engineering involves identifying and extracting the most relevant aspects of the data for the task at hand.

Feature engineering takes many forms, from simple transformations like normalization and scaling to complex operations that combine multiple features in meaningful ways. Domain knowledge plays a crucial role here—understanding the problem space can suggest useful feature transformations that might not be obvious from the data alone. For example, in a financial prediction system, raw price data might be less useful than carefully engineered features like moving averages, volatility measures, or price ratios.

Modern deep learning has somewhat reduced the need for manual feature engineering in some domains, as deep neural networks can learn useful features automatically from raw data. However, feature engineering still remains crucial in many applications, particularly when working with structured data or when domain expertise can suggest powerful hand-crafted features.

Evolution of Specialized Architectures for Complex Data

From years of research, it has become apparent that different types of data require different architectural approaches. For example, Convolutional Neural Networks (CNNs) have revolutionized computer vision by incorporating our understanding of how visual processing works. The key insight behind CNNs is the use of convolutional layers, which apply the same learned patterns across different regions of an image. This architectural choice dramatically reduces the number of parameters needed compared to fully connected networks while maintaining the ability to detect features regardless of their position in the image.

CNNs typically consist of multiple convolutional layers, each learning increasingly complex features. Early layers might detect simple edges and textures, while deeper layers combine these features to recognize complex objects, faces, or scenes. The hierarchical nature of this processing mirrors the organization of the human visual system, which may partly explain the remarkable success of CNNs in computer vision tasks.

For sequential data like text, speech, or time-series, Recurrent Neural Networks (RNNs) and their modern variants have proven particularly effective. The key innovation of RNNs is their ability to maintain an internal state or “memory” that gets updated as they process each element in a sequence. This allows them to capture temporal dependencies and patterns that extend over time.

However, traditional RNNs suffer from the vanishing gradient problem, making it difficult for them to learn long-range dependencies. This led to the development of more sophisticated architectures like Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRUs). These architectures include specialized mechanisms for controlling how information flows through the network, allowing them to maintain relevant information over longer sequences.

The Transformer Revolution

The introduction of the transformer architecture in 2017 marked a pivotal moment in deep learning, particularly for natural language processing. Unlike RNNs, which process sequences element by element, transformers use a mechanism called self-attention to process entire sequences in parallel. This allows them to model relationships between any elements in a sequence, regardless of their distance from each other. Transformers are also being used with various kinds of data apart from text data, including images and videos, tables, audio, and time-series data.

The key innovation of transformers is their attention mechanism, which allows the model to dynamically focus on different parts of the input when producing each element of the output. This has proven particularly powerful for tasks like machine translation, where the relationship between input and output words can be complex and non-local.

The success of transformers has led to the development of increasingly large and sophisticated language models. Large Language Models (LLMs) like GPT (Generative Pre-trained Transformer) represent a quantum leap in natural language processing capabilities. These models are trained on massive datasets of text and can perform a wide range of language tasks, from translation and summarization to question-answering and creative writing.

The scale of modern LLMs is staggering—models with hundreds of billions of parameters, trained on trillions of words of text. This scale brings both remarkable capabilities and significant challenges. While these models can generate impressively coherent and contextually appropriate text, they also raise important questions about bias, reliability, and the nature of language understanding. This has led the way to Generative AI and applications of foundation models in various data modalities.

Generative AI

Generative AI represents one of the most exciting frontiers in machine learning. These systems learn to create new content—whether text, images, music, or other media—that exhibits characteristics similar to their training data. Generative models come in several varieties, each with its own strengths and challenges.

Generative Adversarial Networks (GANs) work by setting up a competition between two neural networks: a generator that creates content and a discriminator that tries to distinguish generated content from real examples. This adversarial training process leads to increasingly sophisticated and realistic outputs. GANs have achieved remarkable success in image generation, able to create photorealistic images of faces, scenes, and objects that never existed.

Variational autoencoders (VAEs) take a different approach, learning to compress data into a compact representation and then reconstruct it. This creates a continuous “latent space”

of possible outputs, allowing for controlled generation of new content by manipulating the latent representation. While GANs often produce sharper results, VAEs offer more stable training and better control over the generation process.

The emergence of *diffusion models* marked another significant advance in generative AI. Unlike GANs or VAEs, diffusion models work by gradually adding noise to training data and then learning to reverse this process. This approach, while computationally intensive, has proven remarkably effective at generating high-quality images. Models like DALL-E, Stable Diffusion, and Midjourney demonstrate the power of this architecture, creating stunning images from text descriptions with unprecedented quality and control.

Foundation models represent a paradigm shift in machine learning, extending beyond language to encompass multiple modalities. These foundation models are typically trained on vast amounts of data using self-supervised learning techniques, creating versatile representations that can be adapted to numerous downstream tasks. This general pre-training approach represents a dramatic shift from traditional task-specific models, as it allows the model to develop a rich understanding of language patterns, relationships, and world knowledge that can be leveraged across multiple applications.

For instance, traditional named entity recognition (NER) systems required carefully labeled training data and hand-crafted features or specialized architectures trained specifically for entity extraction. In a sentence like this:

"Jeff Bezos founded Amazon in Seattle in 1994"

traditional NER systems need explicit training to identify “Jeff Bezos” as a PERSON, “Amazon” as an ORGANIZATION, “Seattle” as a LOCATION, and “1994” as a DATE entity.

In contrast, modern LLMs can perform entity extraction through simple prompting or fine-tuning on a small amount of task-specific data, leveraging their pre-trained understanding of entities and their relationships in natural language. Similar patterns emerge across other traditional NLP tasks like sentiment analysis, text classification, and relationship extraction—tasks that previously required separate specialized models can now be handled by a single pre-trained model with minimal task-specific adaptation.

The Impact of LLMs

LLMs have fundamentally transformed both research and industry landscapes. Academic institutions and research labs continue to push the boundaries of these models, exploring areas like instruction-tuning, constitutional AI, and chain-of-thought prompting, while also investigating crucial questions about model interpretability, alignment, and knowledge representation. In the commercial sphere, companies like Anthropic, OpenAI, and Google have developed increasingly sophisticated models that power a wide range of applications. These models serve as the backbone for modern enterprise solutions, enabling everything

from automated customer service and content generation to code assistance and data analysis. Startups have particularly embraced LLMs, building specialized applications on top of these models through techniques like retrieval-augmented generation (RAG), fine-tuning, and prompt engineering. This has led to a proliferation of domain-specific tools in areas like legal document analysis, medical research assistance, and financial analysis. The ability to leverage these powerful models through APIs has democratized access to advanced AI capabilities, spawning an entire ecosystem of AI-first applications and services. You explore this more in future chapters.

While Large Language Models like GPT and Claude represent one class of foundation models, similar approaches have been successfully applied to vision, audio, and multimodal tasks. Similar to LLMs, which can be considered textual foundation models, vision foundation models, such as ViT (Vision Transformer) and Swin Transformer, have demonstrated that architectures originally developed for language can be effectively adapted to image understanding. These models process images as sequences of patches, applying self-attention mechanisms to capture complex visual relationships. This approach has led to state-of-the-art performance across various computer vision tasks, from object detection to image segmentation.

Multimodal foundation models represent perhaps the most ambitious development in this space. Models like GPT-4V and Claude 3 can process and reason about both text and images, while systems like PaLM-E extend this capability to robotic control. These models demonstrate remarkable cross-modal understanding. They are able to answer questions about images, generate image descriptions, and even follow visual instructions. The success of these models suggests that the principles underlying foundation models—massive scale, self-supervised learning, and flexible architectures—can create systems with increasingly general capabilities across different types of data and tasks.

The intersection of foundation models with scientific domains has opened new frontiers in fields like protein structure prediction (AlphaFold), molecular design (ESM), and scientific discovery. These specialized foundation models demonstrate how the architectural insights from language and vision can be adapted to complex scientific data, potentially accelerating research and discovery across multiple disciplines.

The Relationship Among AI, ML, and Deep Learning

Artificial intelligence, machine learning, and deep learning are interconnected fields within computer science that focus on creating intelligent systems. While they share common goals and methodologies, they also have distinct characteristics and applications. You

should know that one method does not supersede the others in this list—all categories of methods and concepts still exist in production applications today. This section explores their similarities and differences, highlighting their relationships, unique features, and roles in the broader landscape of intelligent systems.

Hierarchical Relationship

The relationship among AI, ML, and deep learning can be visualized as a series of nested circles, as shown in Figure 1.4, with AI encompassing ML, and ML encompassing deep learning.

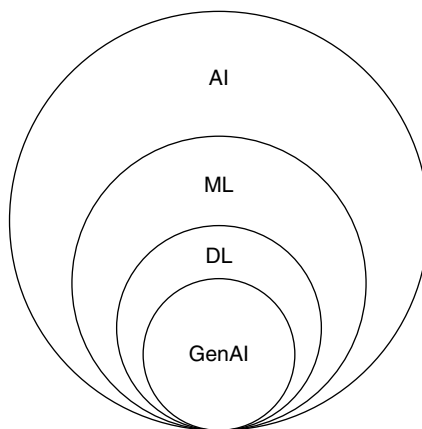
- Artificial intelligence is the broadest concept, representing the overall field of creating intelligent machines that can simulate human-like cognitive functions.
- Machine learning is a subset of AI that focuses on algorithms that can learn from and make predictions or decisions based on data.
- *Deep learning* (DL) is a specialized subset of ML that uses multilayered neural networks to learn complex patterns in data.
- *Generative AI* (GenAI) uses foundational models to generate new content in applications such as text, image, video and audio generation.

This hierarchical relationship is crucial to understanding how these fields relate to and build upon each other. The next section dives a little deeper into each of these fields, including their key characteristics, similarities, and differences.

Artificial Intelligence

AI is the overarching discipline that aims to create intelligent machines capable of performing tasks that typically require human intelligence. These tasks include

FIGURE 1.4 Hierarchical relationship among AI, ML, DL, and Gen AI.



problem-solving, learning, perception, language understanding, and reasoning. Key characteristics of AI include:

- **Goal-oriented:** AI systems are designed to achieve specific objectives or solve particular problems.
- **Adaptability:** AI can adjust its behavior based on new information or changing environments.
- **Reasoning:** AI systems can make inferences and draw conclusions from available data.
- **Human-like interaction:** Many AI systems aim to interact with humans in natural ways, such as through language or visual interfaces.

AI encompasses a wide range of approaches, including rule-based systems, expert systems, and machine learning. It's important to note that not all AI systems use machine learning; some rely on predefined rules and logic to make decisions.

Machine Learning

Machine learning (ML) is a subset of AI that focuses on algorithms and statistical models that enable computer systems to improve their performance on a specific task through experience. Unlike traditional programming, where rules are explicitly coded, ML algorithms learn patterns from data. Key characteristics of ML include:

- **Data-driven:** ML models rely on large amounts of data to learn patterns and make predictions.
- **Iterative improvement:** ML algorithms refine their performance over time as they process more data.
- **Generalization:** ML models aim to perform well on new, unseen data, not just the data they were trained on.
- **Feature extraction:** Many ML techniques involve identifying relevant features in the data that are predictive of the desired outcome.

Three main types of ML have been developed so far:

- **Supervised learning:** The algorithm learns from labeled data, where both input features and target outputs are provided.
- **Unsupervised learning:** The algorithm identifies patterns in unlabeled data without predefined target outputs.
- **Reinforcement learning:** The algorithm learns through interaction with an environment, receiving feedback in the form of rewards or penalties.

ML encompasses a wide range of algorithms and techniques, including decision trees, support vector machines, k-means clustering, and neural networks.

Deep Learning

Deep learning is a specialized subset of machine learning that uses artificial neural networks with multiple layers (hence “deep”) to learn hierarchical representations of data. Inspired by the structure and function of the human brain, deep learning models can automatically learn to extract features from raw data. Key characteristics of DL include:

- **Hierarchical feature learning:** Deep neural networks learn to represent data at multiple levels of abstraction.
- **Large-scale data processing:** Deep learning models excel at handling vast amounts of high-dimensional data.
- **End-to-end learning:** Deep learning often eliminates the need for manual feature engineering, learning directly from raw data to final output.
- **Representation learning:** Deep networks can discover intricate structures in high-dimensional data without human guidance.

You may have already heard of these common deep learning architectures:

- **Convolutional Neural Networks (CNNs):** Particularly effective for image and video processing tasks.
- **Recurrent Neural Networks (RNNs):** Designed for sequential data, such as time-series or natural language.
- **Transformer Architecture:** Excels at natural language processing tasks and has led to breakthroughs like BERT and GPT models.

Now that you know the basics around AI, ML, and DL, the following sections look at some similarities and differences.

Similarities Among AI, ML, and DL

AI, ML, and DL share these similarities:

- **Data-driven approach:** All three fields rely heavily on data to inform decision-making and improve performance.
- **Goal of creating intelligent systems:** AI, ML, and deep learning all aim to develop systems that can perform tasks that traditionally required human intelligence.
- **Iterative improvement:** Systems in all three fields typically improve their performance over time through exposure to more data or experiences.
- **Interdisciplinary nature:** All three fields draw from various disciplines, including computer science, statistics, neuroscience, and cognitive psychology.
- **Application to complex problems:** AI, ML, and deep learning are all applied to solve complex, real-world problems across various domains.

Differences Among AI, ML, and DL

While AI, ML, and DL are closely related and often used interchangeably in casual conversation, understanding their distinctions is crucial for grasping the current state and future potential of intelligent systems. AI provides the overarching vision and goal of creating intelligent machines. ML offers a data-driven approach to achieving this goal, allowing systems to learn from experience. Deep learning pushes the boundaries further by enabling machines to learn complex representations from raw data, often achieving human-like or superhuman performance on specific tasks. The following table provides differences among AI, ML, and DL.

Characteristic	Artificial Intelligence (AI)	Machine Learning (ML)	Deep Learning (DL)
Scope and Breadth	Broadest field encompassing any approach to creating intelligent machines	Focuses specifically on algorithms that can learn from data	Narrows to multilayered neural networks
Level of Human Intervention	Requires significant human input in designing rules and knowledge bases	Reduces human intervention but may need feature engineering	Less human intervention in feature design; more focus on architecture design and hyperparameter tuning
Data Requirements	Can work with various data types and quantities, depending on approach	Requires substantial amounts of structured or semi-structured data	Requires very large datasets, particularly for complex tasks
Interpretability	Highly interpretable; decision-making process can be traced in rule-based systems	Offers a balance between performance and interpretability	Often considered a “black box” due to complexity; ongoing efforts to improve interpretability
Computational Resources	Generally requires less computational power	Computational needs vary by algorithm	Requires significant computational resources, especially for training
Handling of Unstructured Data	Often struggles with unstructured data	Can handle semi-structured data with feature engineering	Excels at processing unstructured data (images, audio, text)
Flexibility and Generalization	Can be brittle when facing unprogrammed scenarios	Better generalization within trained domain	Strong generalization across related tasks and domains; supports transfer learning

Understanding Data Types in AI Models

Data is (still) the cornerstone of developing AI models. Even with the latest developments around foundation models and GenAI, we repeatedly see that high-quality, curated data leads to better models, all other factors remaining constant. The data type an AI system works with dictates not only the design of the model and the choice of algorithm but also influences training techniques, hyperparameter tuning, and even deployment strategies. The diversity in data types—ranging from tabular data to complex image and audio inputs—means that each type has unique characteristics and challenges, which in turn shape the way AI models are developed and utilized. This section explores the different types of data that AI models encounter and how they align with various training and deployment options.

Labeled vs. Unlabeled Data

The earlier section on supervised vs. unsupervised learning briefly covered the concept of labels. One of the first considerations in working with AI data is whether it is labeled or unlabeled, as this determines the training methodology—supervised or unsupervised learning. Labeled data includes input-output pairs, where each data point is annotated with a corresponding label. For example, a set of medical images labeled with diagnoses allows a model to learn specific patterns linked to certain conditions. This labeled dataset is essential for supervised learning, where the model's goal is to predict the output for new, unseen data. Models trained with labeled data often benefit from extensive hyperparameter tuning to optimize performance, given the clear feedback during training.

On the other hand, unlabeled data lacks such annotations, making it suitable for unsupervised learning. In this scenario, the model seeks to find hidden patterns or relationships within the data, such as grouping similar items together or identifying anomalies. An e-commerce site might use unsupervised learning on customer browsing data to segment users into distinct categories for targeted marketing. This method is valuable when obtaining labels is difficult or expensive, which is often the case with large datasets like video or audio recordings.

Structured Data

Structured data is highly organized and often resides in tabular forms, such as rows and columns in a spreadsheet or database. Each piece of data is linked to predefined fields, making it easier for algorithms to process and interpret. For example, a bank might use structured data, including account balances, transaction history, and customer demographics, to develop a fraud detection model. Such datasets are often straightforward to prepare and analyze, as they follow a consistent schema and contain categorical and numerical features, allowing for direct manipulation using statistical methods or machine learning algorithms. Models that process structured data are often deployed on cloud platforms, where the data can be updated in real time, enabling batch or real-time inference depending on operational needs.

Tabular Data

Tabular data remains one of the most common types encountered in AI projects, especially in domains like finance, healthcare, and marketing. This structured data format is well-suited for machine learning algorithms, as it allows easy access to individual features and straightforward calculation of statistical measures. A retail company might use a tabular dataset containing sales figures, product categories, and customer demographics to forecast future sales or identify products likely to sell out. When working with tabular data, hyperparameter tuning is often essential to optimize model accuracy, especially when dealing with imbalanced datasets where certain outcomes are rare.

Relational database data is another form of structured, tabular data commonly used in business operations. It consists of information organized into tables with predefined relationships between them. For example, a CRM system might store customer data, purchase history, and support tickets in related tables, enabling efficient querying and reporting. Structured Query Language (SQL) is widely used to manage and retrieve this type of data.

Time-Series Data

Time-series data captures observations collected at sequential time intervals and is often used to track trends, detect seasonality, or forecast future events. This data type is prevalent in industries like finance, where stock prices are monitored, and in IoT, where sensor data is collected over time. Time-series data often exhibits autocorrelation, meaning that past values influence future values, which affects the choice of algorithm and model architecture. Models trained on time-series data might be deployed using asynchronous inference, where they process incoming data as it arrives, making them ideal for monitoring applications like anomaly detection in network traffic or predictive maintenance in industrial equipment.

Log Data

Log data, generated by systems, applications, or devices, is another example of structured data. It typically includes timestamps, error codes, and other predefined fields. For instance, server logs in IT operations might help identify performance issues or detect security breaches through anomaly analysis.

Unstructured Data

Unstructured data lacks a predefined format and includes types such as text, images, and videos. These data types require additional processing before they can be used in AI models, as they often contain complex information that doesn't fit neatly into rows and columns. Natural language data, for example, requires tokenization and embedding techniques to convert words into a numerical format suitable for a model. Models trained on unstructured data are increasingly deployed at the edge, especially in applications like object detection in images or video surveillance, where quick processing and immediate

responses are essential. Edge deployment allows AI systems to perform real-time inference directly on devices like smartphones or security cameras, reducing the need for constant cloud connectivity.

Text Data

Text data is fundamental in natural language processing (NLP) tasks, such as sentiment analysis, language translation, and information retrieval. Handling text data requires converting it into numerical format through techniques like word embeddings or transformers, which capture semantic relationships between words. A customer service chatbot trained on text data can provide real-time responses by deploying models in an online inference setting. This allows the model to handle queries immediately, making it suitable for customer-facing applications where quick interaction is essential.

Text Data and NLP Techniques

Natural Language Processing (NLP) has evolved significantly from its early reliance on rule-based systems and statistical methods like bag-of-words (BoW) and TF-IDF. These traditional approaches, which treated words as individual tokens without context, gave way to more sophisticated word embedding techniques such as Word2Vec and GloVe. These embeddings represented words as continuous vectors, enabling models to capture semantic relationships and contextual similarities more effectively. A major breakthrough came with the introduction of transformer architectures in 2017, which revolutionized NLP through their self-attention mechanism. This innovation allowed models to process relationships between words in a text simultaneously, regardless of their position. Transformers became the foundation for Large Language Models (LLMs) like GPT and BERT, which are trained on vast datasets and can generate coherent, contextually relevant text for various applications, from chatbots to content generation. The practical implementation of NLP systems requires converting text data into numerical formats through these advanced techniques. Modern applications, such as customer service chatbots, can now provide real-time responses by deploying these models in online inference settings. As foundation models continue to evolve, they can be fine-tuned for specialized applications across different domains, from legal document analysis to medical record interpretation, demonstrating the versatility and growing capabilities of NLP technology.

Image Data

Image data presents unique challenges and opportunities for AI due to its high dimensionality and complexity. In tasks like facial recognition or medical imaging, the model must learn to interpret intricate visual details within the pixel arrays. Convolutional neural networks (CNNs) are a popular model architecture for image data, as they are designed to capture spatial hierarchies within images. A CNN trained to detect defects in a manufacturing pipeline can be deployed at the edge, allowing real-time quality control

without the latency of cloud processing. Model deployment for image data often involves optimization techniques to reduce computational load, enabling efficient, on-device inference in environments where quick decisions are crucial.

Video Data

Video data, an extension of image data, consists of sequences of frames combined with temporal information. Video analytics is used in areas like autonomous driving, where models analyze footage to detect obstacles, or in retail, where surveillance systems monitor customer behavior. Processing video data requires advanced models capable of capturing both spatial and temporal features, often involving convolutional and recurrent neural networks.

Audio Data

Audio data is another form of unstructured data widely used in applications like speech recognition, music classification, and voice-based assistants. Converting audio data into usable formats often involves feature extraction techniques such as Mel-frequency cepstral coefficients (MFCCs) or spectrogram analysis. Models trained on audio data are typically deployed for real-time inference in applications like voice-controlled IoT devices or automated transcription services.

Understanding the various types of data that AI models can handle is essential for making informed decisions on model selection, training methods, and deployment strategies. Each data type presents unique challenges that influence the overall design and implementation of an AI system, from selecting the right hyperparameters to choosing a suitable inference type. As AI systems continue to evolve, the ability to effectively process diverse data types enables more sophisticated applications across industries, from real-time image processing at the edge to batch predictions on structured datasets in the cloud. By aligning data types with appropriate AI methods and deployment options, practitioners can leverage the full potential of their models to drive meaningful results in the real world.

Making Predictions Using Trained Models

Once an AI model has been trained using the appropriate training method, it produces a model artifact—a file or set of files that contain all the learned parameters and configurations needed to make predictions on new data. This artifact, which can be stored in various formats such as ONNX, TensorFlow SavedModel, or PyTorch model file, essentially represents the model's knowledge and serves as the basis for deployment and inference.

The trained model artifact is typically saved and stored in a model repository, where it can be versioned, shared, and accessed as needed. Websites like Hugging Face have revolutionized

model sharing by offering an extensive repository of pre-trained models for a wide range of applications, from natural language processing to computer vision. These repositories make it easy for practitioners to access state-of-the-art models and either use them directly or fine-tune them for specific tasks, accelerating the model deployment process significantly.

Once the model artifact is ready, it can be deployed for inference—the process of using the model to make predictions on new data. There are several types of inference that cater to different deployment needs, each with unique advantages and applications. The following sections explain those inference types.

Batch Inference

Batch inference involves processing large volumes of data at once, rather than responding to individual requests in real-time. This approach is suitable for scenarios where predictions can be generated periodically, such as in overnight data processing or analyzing data at regular intervals. For example, a retail company might use batch inference to generate product recommendations for all users at the end of each day based on their latest browsing data. Batch inference is highly efficient for processing large datasets but is not suitable for applications requiring immediate responses. Figure 1.5 illustrates batch inference, where input data is processed in bulk by a batch orchestrator, which sends the data to the model at scheduled intervals. The model generates output data and stores it in a designated location.

Real-Time Inference

Real-time inference, as the name suggests, involves generating predictions on the fly, typically in response to user requests. This method is ideal for applications where quick response times are essential, such as fraud detection in banking, chatbots, or personalized recommendations. For instance, an e-commerce website might use real-time inference to provide personalized product recommendations as a user browses, ensuring that the suggestions are contextually relevant and timely. Real-time inference typically requires robust infrastructure to handle incoming requests with low latency, often utilizing dedicated servers or cloud-based services optimized for quick response times. Figure 1.6 shows real-time inference, where a client sends a request directly to the model (or specifically the model inference stack), and the model responds immediately with an output, facilitating quick, synchronous interaction.

FIGURE 1.5 Batch inference.

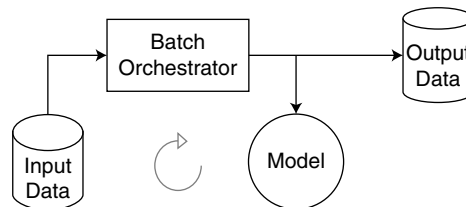
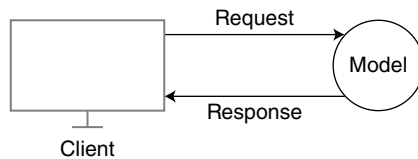
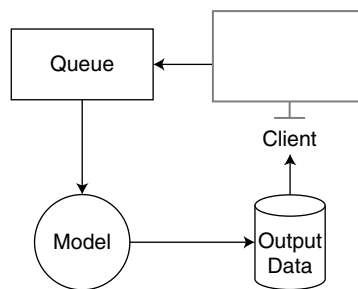


FIGURE 1.6 Real-time inference.**FIGURE 1.7** Asynchronous inference.

Asynchronous Inference

Asynchronous inference is a hybrid approach that accommodates tasks where predictions are needed quickly but do not require an immediate response. In this setup, prediction requests are queued, and the model processes them as resources become available. This is useful for applications where a delay is acceptable, but the predictions need to be made within a certain timeframe. This is especially useful when the payload size is large, and processing a single payload will take time. An example might be in video processing, where individual frames are sent to a model for analysis, and the results are returned when processing is complete. Asynchronous inference helps balance resource allocation by allowing systems to manage workloads flexibly without the stringent requirements of real-time processing. Figure 1.7 depicts asynchronous inference, where the client sends data to a queue, which the model processes when resources are available. The output data is then made accessible to the client.

In choosing among batch, real-time, and asynchronous inference, practitioners must consider factors like the desired response time, the available infrastructure, and the specific needs of the application. Each method provides a way to utilize the model artifact effectively, enabling AI systems to deliver predictions in a manner that aligns with their operational requirements. By selecting the appropriate inference type, organizations can maximize the efficiency and utility of their AI models, delivering meaningful results to end users in the most suitable way.

Summary

This chapter established a solid foundation in artificial intelligence by defining essential terms and concepts, including AI, ML, DL, neural networks, computer vision, NLP models, algorithms, training, inferencing, bias, and large language models (LLMs). It highlighted the interconnections among AI, ML, and deep learning, clarifying their similarities and differences.

The chapter covered the three core learning paradigms—supervised learning, unsupervised learning, and reinforcement learning—emphasizing their respective roles and suitability for different data types and problem domains.

The chapter also explored various inferencing methods, such as batch, real-time, and asynchronous inferencing, demonstrating their significance in model deployment and application. It examined different data types utilized in AI, including labeled and unlabeled data, structured data like tabular and time-series formats, and unstructured data encompassing images, text, and videos.

Exam Essentials

Understand key AI terminology. Familiarize yourself with fundamental terms such as artificial intelligence (AI), machine learning (ML), deep learning, neural networks, computer vision, natural language processing (NLP), models, algorithms, training, inferencing, and large language models (LLMs). While the chapter explored several deeper concepts, focus additionally on terms that represent techniques that are still used today, for example back propagation, feature engineering, and inference techniques.

Differentiate among AI, ML, and deep learning. Know the similarities and distinctions among these three fields, particularly how they relate to each other and their roles in developing intelligent systems.

Identify types of inferencing methods. Be able to explain batch processing, real-time, and asynchronous inferencing, along with their respective use cases and implications for AI model deployment.

Recognize various data types in AI models. Understand the differences between labeled and unlabeled data, and the distinctions among structured data (tabular, time-series) and unstructured data (images, text, video).

Differentiate learning paradigms. Be familiar with supervised learning, unsupervised learning, and reinforcement learning, including their applications and how they relate to the data used.

Apply knowledge to real-world scenarios. Be prepared to discuss how the concepts of AI, data types, and learning methods can be applied in practical situations, such as developing applications for image recognition or natural language processing.

Review Questions

1. Which statements correctly describe the hierarchical relationship among AI, ML, and deep learning? (Choose two.)
 - A. Machine learning is the broadest field, encompassing both AI and deep learning.
 - B. Deep learning is a specialized subset of machine learning that uses multilayered neural networks.
 - C. Traditional AI methods always require less computational power than deep learning.
 - D. Artificial intelligence is the broadest field, with machine learning being a subset of AI.
2. Which of the following is NOT a valid, key difference between traditional AI and deep learning approaches?
 - A. Traditional AI systems require extensive human intervention in rule design, while deep learning requires less feature engineering.
 - B. Deep learning typically requires larger computational resources than traditional AI methods.
 - C. Traditional AI methods excel at processing unstructured data compared to deep learning.
 - D. Deep learning models are generally less interpretable than traditional AI systems.
3. What is a key difference between labeled and unlabeled data in AI models? (Choose two.)
 - A. Labeled data is used in unsupervised learning.
 - B. Unlabeled data is suitable for unsupervised learning.
 - C. Labeled data includes input-output pairs.
 - D. Unlabeled data includes annotated labels.
4. Which of the following statements about inference types is true? (Choose two.)
 - A. Batch inference is ideal for applications needing immediate response times.
 - B. Real-time inference generates predictions in response to user requests instantly.
 - C. Asynchronous inference is suitable for processing large datasets all at once.
 - D. Batch inference processes data periodically, rather than in real-time.
5. What are characteristics of time-series data? (Choose two.)
 - A. Time-series data captures sequential observations over time intervals.
 - B. Time-series data typically does not exhibit autocorrelation.
 - C. Time-series models are often deployed using synchronous inference.
 - D. Time-series data is commonly used in predictive maintenance applications.

6. Which of the following statements are correct about time-series and image data? (Choose two)
- A. Time-series data is often used for tasks like anomaly detection and predictive maintenance.
 - B. Image data models typically rely on autocorrelation to influence future predictions.
 - C. Convolutional neural networks (CNNs) are well-suited for analyzing time-series data.
 - D. CNNs are commonly used for tasks like defect detection in manufacturing pipelines.
7. What is a key distinction between word embeddings and transformer models in natural language processing? (Choose two.)
- A. Word embeddings capture semantic relationships between words, while transformers handle context within sentences.
 - B. Transformers primarily utilize statistical methods to represent text data.
 - C. Word embeddings use self-attention to weigh the importance of words in a sentence.
 - D. Transformers can capture long-range dependencies, unlike traditional word embeddings.
8. Which of the following techniques allows a model to weigh the importance of each word in a sentence relative to all others?
- A. Term frequency-inverse document frequency (TF-IDF)
 - B. Bag-of-words (BoW)
 - C. Self-attention
 - D. Word2Vec
9. What are two key differences between neural networks and traditional algorithms in AI? (Choose two.)
- A. Neural networks have fixed weights that do not change during training.
 - B. Traditional algorithms rely on explicit programming rather than data-driven learning.
 - C. Neural networks use layers of interconnected nodes, resembling the human brain.
 - D. Traditional algorithms always require large datasets to function effectively.
10. Which of the following statements accurately describe the function of deep learning and neural networks? (Choose two.)
- A. Deep learning involves training models that are based on single-layer neural networks.
 - B. Neural networks are organized into multiple layers and can handle complex tasks.
 - C. Deep learning is solely focused on supervised learning tasks.
 - D. Neural networks can form the foundation for many modern AI systems.

