

Contents

1	Intr	oduct	ion to Machine Learning 5		
	1.1	What	is Machine Learning?		
		1.1.1	Fundamental Concept		
		1.1.2	Types of Machine Learning 6		
		1.1.3	How Does Machine Learning Work? 6		
		1.1.4	Mathematical Framework		
	1.2	Types	of Machine Learning		
		1.2.1	Supervised Learning		
		1.2.2	Unsupervised Learning		
		1.2.3	Reinforcement Learning		
	1.3	Machi	ne Learning Algorithms		
		1.3.1	Supervised Learning Algorithms		
		1.3.2	Unsupervised Learning Algorithms		
	1.4	Applie	cations of Machine Learning		
		1.4.1	Healthcare		
		1.4.2	Finance and Banking		
		1.4.3	Self-Driving Cars		
		1.4.4	Retail and E-commerce		
		1.4.5	Natural Language Processing (NLP) 20		
	1.5	Challe	enges in Machine Learning		
		1.5.1	Data-Related Challenges		
		1.5.2	Model-Related Challenges		
		1.5.3	Algorithm-Related Challenges		
		1.5.4	Real-World Deployment Challenges		
	1.6				
		1.6.1	Input Representation		
		1.6.2	Hypothesis Class		

4 CONTENTS

		1.6.3	Version Space
2	Intr	oduct	ion to Deep Learning 29
	2.1		luction
		2.1.1	Why Deep Learning?
	2.2	Artific	cial Neural Networks (ANNs)
		2.2.1	Structure of an Artificial Neural Network (ANN) 30
		2.2.2	Summary
		2.2.3	Common Activation Functions
	2.3	Feedfo	orward Neural Networks (FFNs)
		2.3.1	Forward Propagation
		2.3.2	Example Calculation
	2.4	Backp	propagation and Training
		2.4.1	Loss Function
		2.4.2	Gradient Descent
	2.5	Convo	olutional Neural Networks (CNNs)
		2.5.1	Definition and Introduction
		2.5.2	Limitations of Artificial Neural Networks (ANNs) 35
		2.5.3	Significance of CNNs
		2.5.4	Mathematical Formulation of CNN Layers 35
		2.5.5	CNN Architecture Diagram
		2.5.6	Example: Image Classification with CNN 37
	2.6	Recur	rent Neural Networks (RNNs)
		2.6.1	Limitations of Traditional Neural Networks 38
		2.6.2	Why RNNs?
		2.6.3	Mathematical Formulation of RNN
		2.6.4	Backpropagation Through Time (BPTT) 39
		2.6.5	Challenges with RNNs
	2.7	Trans	formers in NLP
		2.7.1	Introduction
		2.7.2	Why Are Transformers Revolutionary? 41
	2.8		encoder and Variational Autoencoder (VAE) 42
		2.8.1	Autoencoder (AE)
		2.8.2	Variational Autoencoder (VAE) 42
		2.8.3	Diagram

Chapter 1

Introduction to Machine Learning

1.1 What is Machine Learning?

Machine Learning (ML) is a branch of Artificial Intelligence (AI) that focuses on creating algorithms that allow computers to learn from data. Instead of being explicitly programmed to perform a task, a machine learning model improves its performance over time by recognizing patterns in data. The key idea is that systems can automatically learn and improve from experience without human intervention, using data to make decisions or predictions.

1.1.1 Fundamental Concept

The fundamental concept of machine learning revolves around learning from data. Traditionally, when we want to perform a task (like recognizing objects in an image or predicting stock prices), we would write a program with explicit rules for each possible scenario. However, in machine learning, we don't write those explicit rules. Instead, we provide the machine with a large dataset of examples, and the machine learns the rules on its own by finding patterns and relationships within the data.

For example, in a spam email filter, instead of writing a program to specify every rule for what makes an email spam, we simply provide a dataset of labeled emails (spam and non-spam) and allow the machine to figure out the patterns that distinguish spam from non-spam.

1.1.2 Types of Machine Learning

Machine Learning can be broadly classified into three main categories based on the type of data and feedback provided to the model:

1. **Supervised Learning**: In supervised learning, the model is trained on labeled data, meaning the input data comes with a known output (or label). The algorithm learns to map inputs to the correct output by minimizing the difference between the predicted and actual outputs.

Example: A dataset of images labeled as "cat" or "dog." The model learns to identify features that differentiate cats from dogs.

2. Unsupervised Learning: Unsupervised learning involves training a model on data without labels. The goal is to discover hidden patterns or structures in the data. The model is expected to infer the underlying structure of the data on its own.

Example: Grouping customers with similar purchasing habits into clusters (known as clustering) without having labels for those groups.

3. **Reinforcement Learning**: In reinforcement learning, the model learns by interacting with an environment. It makes decisions and takes actions that lead to rewards or penalties, and it aims to maximize the cumulative reward over time.

Example: A robot learns to navigate a maze by trial and error, receiving positive or negative feedback (reward or penalty) depending on the actions it takes.

1.1.3 How Does Machine Learning Work?

The machine learning process typically follows these steps:

- 1. **Data Collection**: Gather large amounts of relevant data that represent the problem you're trying to solve. For example, if you're building a spam filter, you need a collection of emails.
- 2. **Data Preparation**: Clean and preprocess the data by removing noise or irrelevant information. This can include normalizing data, handling missing values, or encoding categorical variables.

- 7
- 3. Model Selection: Choose an appropriate machine learning algorithm or model. The choice depends on the type of problem you're trying to solve (e.g., regression, classification, clustering).
- 4. **Training**: Train the model on the data by feeding it input-output pairs and adjusting the model parameters so that the output is as close as possible to the desired result.
- 5. **Evaluation**: Evaluate the model's performance using a separate test dataset that it hasn't seen before. Common evaluation metrics include accuracy, precision, recall, and F1 score for classification tasks.
- 6. **Deployment**: Once the model performs well, deploy it for real-world use, where it can make predictions or decisions based on new, unseen data.

1.1.4 Mathematical Framework

Mathematically, machine learning can be seen as a function approximation problem. Given a set of training data $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where x_i represents input features and y_i represents the corresponding output labels, we want to learn a function $f: X \to Y$ that maps an input x to an output y.

In supervised learning, this can be represented as:

$$y_i = f(x_i) + \epsilon_i$$

Where:

- y_i is the output (or label) for input x_i .
- $f(x_i)$ is the true function that we are trying to learn.
- ϵ_i represents the error or noise in the data.

Our goal is to find the best approximation for the function f based on the data D.

1.2 Types of Machine Learning

Machine learning can be classified into three primary types based on how the model learns from the data: Supervised Learning, Unsupervised Learning, and Reinforcement Learning. Each type has its own unique approach to learning and is suitable for different types of problems. In this section, we will explore each type in detail with relevant examples.

1.2.1 Supervised Learning

Supervised learning is the most common type of machine learning where the model is trained on a labeled dataset. This means that each input in the training set is paired with the correct output. The model learns to map inputs to outputs by minimizing the error between predicted and actual values. The learning process aims to find a function $f: X \to Y$ that best describes the relationship between the input X and the output Y.

Mathematical Representation

In supervised learning, we are given a dataset of input-output pairs $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ where each $x_i \in X$ is an input (or feature vector) and $y_i \in Y$ is the corresponding output (or label). The objective is to find a function f such that:

$$y_i = f(x_i) + \epsilon_i$$

Where ϵ_i is the error or noise in the data, and $f(x_i)$ is the model's prediction.

Example

A typical example of supervised learning is email spam classification. Given a labeled dataset where each email is tagged as either "spam" or "not spam," the model learns to identify the characteristics of spam emails. Once the model is trained, it can classify new, unseen emails into either of the two categories.

1.2.2 Unsupervised Learning

Unsupervised learning is used when the dataset does not contain labeled outputs. The model tries to learn the underlying structure of the data without the guidance of predefined labels. The goal is to explore the data to find patterns, structures, or relationships within the dataset. This can include tasks like clustering, anomaly detection, or dimensionality reduction.

Mathematical Representation

In unsupervised learning, we are only given a set of input data $\{x_1, x_2, \ldots, x_n\}$, and the model needs to find patterns or groupings in the data. Unlike supervised learning, there are no output labels y_i for each input x_i .

For clustering, for example, the goal is to group similar data points together into clusters. A common algorithm used for clustering is the k-means algorithm, which minimizes the within-cluster variance.

Example

A common application of unsupervised learning is customer segmentation in marketing. Given a large set of customer data, unsupervised learning can group customers with similar behaviors into clusters, helping businesses target specific groups with tailored marketing strategies.

1.2.3 Reinforcement Learning

Reinforcement Learning (RL) is a branch of machine learning where an agent learns to make decisions by interacting with its environment. Unlike supervised learning, where the model is trained with labeled data, reinforcement learning is based on the principle of learning from trial and error. In this process, an agent explores different actions, receives feedback (rewards or penalties), and learns to choose actions that maximize the cumulative reward over time.

Key Concepts of Reinforcement Learning

In RL, the agent's objective is to find the optimal policy that defines the best action to take in each state to maximize the total accumulated reward. The

agent learns this policy over time based on interactions with the environment. The basic components of reinforcement learning are:

- 1. **Agent**: The learner or decision-maker that interacts with the environment and performs actions.
- 2. **Environment**: The external system that the agent interacts with, where the agent performs actions and receives feedback.
- 3. **State** (s): A representation of the current situation or configuration of the environment.
- 4. **Action** (a): The decision or move made by the agent to interact with the environment.
- 5. **Reward** (r): The immediate feedback received by the agent after performing an action in a specific state. The reward can be positive (for good actions) or negative (for bad actions).

The goal of reinforcement learning is for the agent to maximize the cumulative sum of rewards, known as the **return**, over time.

Markov Decision Process (MDP)

A Reinforcement Learning problem is often modeled as a Markov Decision Process (MDP), which provides a mathematical framework for modeling decision-making problems. An MDP is defined by the following components:

- 1. State Space (S): The set of all possible states the agent can be in.
- 2. Action Space (A): The set of all possible actions the agent can take.
- 3. Transition Function (P(s'|s, a)): A function that defines the probability of transitioning from state s to state s' after taking action a.
- 4. Reward Function (R(s, a)): A function that provides the immediate reward received after taking action a in state s.
- 5. **Discount Factor** (γ): A factor between 0 and 1 that represents the importance of future rewards relative to immediate rewards.

The goal of the agent in an MDP is to find a policy $\pi: S \to A$ that maximizes the expected cumulative reward, typically over an infinite time horizon. The agent chooses actions based on the policy and updates its behavior based on the rewards it receives.

Q-Learning

Q-learning is a model-free reinforcement learning algorithm that allows an agent to learn the optimal policy without needing to know the environment's dynamics. It is based on learning the **Q-value** function, which represents the expected cumulative reward for an agent starting in state s, taking action a, and following the optimal policy thereafter.

The Q-value function is updated iteratively using the following equation:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

Where:

- Q(s,a) is the Q-value for state s and action a.
- α is the learning rate, which determines how much new information overrides old information.
- r is the immediate reward received after taking action a in state s.
- γ is the discount factor.
- s' is the next state after taking action a in state s.
- $\max_{a'} Q(s', a')$ is the maximum Q-value for the next state s'.

The Q-learning algorithm iterates through the process of updating the Q-values for different state-action pairs until convergence, at which point the Q-values represent the optimal policy.

Examples of Reinforcement Learning

Reinforcement learning has been applied successfully to a variety of real-world problems, including:

- 1. **Game Playing**: RL algorithms like AlphaGo have been used to teach machines to play complex games such as Go, chess, and poker. These algorithms learn strategies by playing against themselves or other agents, receiving rewards based on game outcomes.
- 2. **Robotics**: In robotics, RL is used for tasks such as robotic control, navigation, and manipulation. Robots learn to perform tasks such as picking up objects or navigating through environments by interacting with their surroundings.
- 3. Autonomous Vehicles: RL is used to train self-driving cars to make decisions on the road, such as navigating through traffic, avoiding obstacles, and following traffic rules.
- 4. **Finance and Trading**: RL algorithms are applied to financial trading systems, where agents learn to make buy, sell, or hold decisions based on stock market data to maximize profits.
- 5. **Healthcare**: RL has been used for personalized treatment planning and drug discovery, where agents learn to optimize patient outcomes based on various treatment options.

Challenges in Reinforcement Learning

Reinforcement learning faces several challenges that can hinder its practical application:

- 1. **Exploration vs Exploitation**: The agent must balance between exploring new actions to discover better strategies and exploiting known actions that have already yielded good results. Striking the right balance is critical for optimal learning.
- 2. Sample Efficiency: RL algorithms often require a large number of interactions with the environment to learn effectively, which can be computationally expensive and time-consuming.
- 3. **Delayed Rewards**: In many RL tasks, rewards are delayed, meaning the agent does not receive immediate feedback after taking an action. This makes it challenging for the agent to associate actions with outcomes.

- 4. **Sparse Rewards**: In some environments, rewards are sparse and only provided after a long sequence of actions. This can make it difficult for the agent to learn which actions lead to the desired outcome.
- 5. **Generalization**: RL agents are often trained in specific environments, and they may struggle to generalize their learned policies to new or unseen environments.

Reinforcement learning is a powerful and versatile approach to decision-making and problem-solving, with applications ranging from gaming to robotics and healthcare. While it holds great promise, RL also presents unique challenges related to exploration, sample efficiency, delayed rewards, and generalization. By understanding the core concepts such as MDPs and Q-learning, as well as the challenges and real-world applications, students can gain a deeper appreciation of the potential and limitations of reinforcement learning in AI systems.

1.3 Machine Learning Algorithms

Machine learning algorithms are methods or techniques that enable a machine to learn from data. The objective of these algorithms is to identify patterns in data, learn from them, and make predictions or decisions without being explicitly programmed. There are several algorithms used in different types of machine learning, such as supervised learning, unsupervised learning, and reinforcement learning. This section explores some popular machine learning algorithms and their mathematical underpinnings.

1.3.1 Supervised Learning Algorithms

In supervised learning, the model is trained on labeled data, where the input data is associated with the correct output. The goal is to learn a mapping from the input to the output, which can then be used to predict the output for new, unseen inputs. Some common supervised learning algorithms include:

Linear Regression

Linear regression is one of the simplest and most widely used supervised learning algorithms. It is used to predict a continuous target variable based on one or more input features. The model assumes a linear relationship between the input variables and the target variable.

Mathematical Formulation: Given a dataset of n data points $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ the linear regression model predicts the target value y as a linear combination of the input features x:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m + \epsilon$$

Where:

- y is the target variable.
- x_1, x_2, \ldots, x_m are the input features.
- $\beta_0, \beta_1, \dots, \beta_m$ are the coefficients to be learned.
- ϵ is the error term.

The goal is to find the coefficients $\beta_0, \beta_1, \dots, \beta_m$ that minimize the residual sum of squares, which represents the difference between the predicted and actual values.

Logistic Regression

Logistic regression is used for binary classification problems, where the output variable is categorical with two classes (e.g., "spam" vs. "not spam"). It models the probability that a given input belongs to a particular class.

Mathematical Formulation: The logistic regression model predicts the probability P(y = 1|x) of class 1 given the input x using the logistic (sigmoid) function:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m)}}$$

Where:

- $\beta_0, \beta_1, \ldots, \beta_m$ are the coefficients to be learned.
- x_1, x_2, \ldots, x_m are the input features.
- \bullet e is Euler's number.

The output is a probability value between 0 and 1, which can then be used to classify the input into one of the two classes.

Support Vector Machines (SVM)

Support vector machines are a powerful supervised learning algorithm used for both classification and regression tasks. The SVM algorithm works by finding the hyperplane that best separates the data into different classes while maximizing the margin between the classes.

Mathematical Formulation: Given a dataset of points $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$, where each $x_i \in \mathbb{R}^m$ is a feature vector and $y_i \in \{+1, -1\}$ is the class label, the SVM classifier aims to find a hyperplane $w^Tx + b = 0$ that separates the data into two classes such that the margin between the closest points of each class is maximized. The optimization problem can be written as:

$$\min_{w,b} \frac{1}{2} ||w||^2$$

Subject to the constraint:

$$y_i(w^T x_i + b) \ge 1 \quad \forall i = 1, 2, \dots, n$$

1.3.2 Unsupervised Learning Algorithms

Unsupervised learning algorithms are used when the data is not labeled, and the goal is to uncover patterns or structures in the data. Some common unsupervised learning algorithms include:

K-Means Clustering

K-means is a clustering algorithm used to partition the data into K clusters. The goal is to minimize the variance within each cluster by finding cluster centers that best represent the data points.

Mathematical Formulation: Given a dataset $\{x_1, x_2, \ldots, x_n\}$, the K-means algorithm assigns each point to the nearest of K cluster centers, minimizing the within-cluster sum of squared distances. The objective function to minimize is:

$$J = \sum_{i=1}^{K} \sum_{x_j \in C_i} ||x_j - \mu_i||^2$$

Where:

- C_i is the set of points in cluster i.
- μ_i is the centroid of cluster i.
- x_i is a data point.

Principal Component Analysis (PCA)

Principal component analysis is a dimensionality reduction algorithm used to project high-dimensional data into a lower-dimensional space while retaining as much variance as possible. The goal is to find a set of orthogonal basis vectors (principal components) that capture the most variance in the data.

Mathematical Formulation: PCA finds the eigenvectors and eigenvalues of the data covariance matrix. The data is projected onto the eigenvectors corresponding to the largest eigenvalues, which represent the directions of maximum variance.

Let $X \in \mathbb{R}^{n \times m}$ be the data matrix, where n is the number of samples and m is the number of features. The covariance matrix C is given by:

$$C = \frac{1}{n-1} X^T X$$

The principal components are the eigenvectors v_1, v_2, \ldots, v_m of the covariance matrix C, and the data is projected onto the top k eigenvectors to reduce the dimensionality.

1.4 Applications of Machine Learning

Machine learning has transformed various industries by providing solutions to complex problems that were previously impossible or time-consuming to solve with traditional approaches. The applications of machine learning span many fields, from healthcare to finance, and even to self-driving cars. In this section, we explore some of the key real-world applications of machine learning and their impact on different industries.

1.4.1 Healthcare

In healthcare, machine learning is being used to revolutionize diagnostics, treatment recommendations, and personalized medicine. By analyzing vast amounts of medical data, machine learning algorithms can identify patterns and predict outcomes with a high degree of accuracy. Some notable applications of machine learning in healthcare include:

Disease Diagnosis

Machine learning models, such as deep learning networks, are trained to diagnose diseases from medical images (e.g., X-rays, MRIs, and CT scans) with comparable or even superior accuracy to that of human doctors. For instance, convolutional neural networks (CNNs) are used to detect tumors in medical images.

Example: A deep learning model can be trained on a dataset of labeled medical images to classify whether a patient has a certain type of cancer. After training, the model can analyze new images and assist radiologists in making accurate diagnoses.

Predictive Healthcare Analytics

Machine learning is also used to predict the onset of diseases or complications based on patient data, such as electronic health records (EHRs). Predictive models can identify patients at high risk of developing conditions like diabetes, heart disease, and stroke, enabling preventive measures and early interventions.

Example: Machine learning models can analyze a patient's medical history, lifestyle factors, and genetic data to predict the likelihood of them developing diabetes in the next five years.

1.4.2 Finance and Banking

Machine learning is widely used in the finance and banking sectors for risk assessment, fraud detection, algorithmic trading, and customer service. By analyzing large volumes of financial data, machine learning models can uncover hidden patterns and make real-time decisions. Some applications include:

Fraud Detection

Machine learning algorithms are employed to detect fraudulent activities by analyzing transaction data. For example, a model can be trained on historical transaction data to recognize patterns associated with fraud, such as unusual spending behavior or rapid changes in account activity.

Example: A credit card company uses machine learning algorithms to identify transactions that deviate from a customer's typical spending habits. The model flags suspicious transactions, which are then reviewed by a fraud detection team.

Algorithmic Trading

In algorithmic trading, machine learning models are used to analyze stock market data and identify trends or signals that can inform trading decisions. These models can analyze vast amounts of data in real time, executing trades at high speeds to take advantage of small price movements.

Example: A trading algorithm uses reinforcement learning to optimize buying and selling decisions based on the past performance of stock prices, aiming to maximize returns while minimizing risks.

Credit Scoring

Machine learning models are used to predict the creditworthiness of individuals or businesses based on historical financial data, such as past loans, income, and spending patterns. These models provide more accurate and personalized credit scores compared to traditional credit scoring systems.

Example: A machine learning model predicts whether a borrower will repay a loan based on features such as employment history, income, debt-to-income ratio, and credit card usage.

1.4.3 Self-Driving Cars

Self-driving cars are one of the most prominent applications of machine learning. Autonomous vehicles rely on machine learning algorithms to process data from various sensors (such as cameras, radar, and LIDAR) and make real-time decisions that allow them to navigate the road safely. Key machine learning techniques used in self-driving cars include:

Object Detection and Tracking

Machine learning algorithms, particularly deep learning models, are used to identify and track objects such as pedestrians, other vehicles, traffic signs,

and road markers in the environment. Convolutional neural networks (CNNs) are often used for object detection in images or video streams from cameras.

Example: A self-driving car uses a CNN to analyze real-time images from its cameras and detect pedestrians on the road. The model then tracks the pedestrians' movements and makes decisions to avoid collisions.

Path Planning and Decision Making

Reinforcement learning and other machine learning algorithms are used to make decisions related to vehicle control, such as steering, braking, and acceleration. The car learns the best actions to take based on the environment and the goal of safely reaching its destination.

Example: A reinforcement learning model is trained to drive a car in a simulated environment, learning to navigate through traffic, obey traffic rules, and avoid accidents.

1.4.4 Retail and E-commerce

Machine learning is widely used in the retail and e-commerce industries to enhance customer experiences, optimize supply chains, and increase sales. Some applications include:

Recommendation Systems

Recommendation systems use machine learning to provide personalized product recommendations to users based on their past behavior, preferences, and similar users' behavior. These systems are widely used in platforms like Amazon, Netflix, and Spotify.

Example: A recommendation system on an e-commerce website suggests products to customers based on their browsing history, purchases, and ratings given to similar products by other users.

Customer Service Automation

Machine learning algorithms are used to automate customer service tasks, such as answering customer queries, processing returns, and providing recommendations. Natural language processing (NLP) models, such as chatbots, can handle customer interactions without human intervention.

Example: An e-commerce website uses a chatbot powered by NLP to assist customers with questions about shipping, returns, and product details.

1.4.5 Natural Language Processing (NLP)

Machine learning techniques are widely used in natural language processing (NLP) tasks, which involve the interaction between computers and human language. Some applications include:

Sentiment Analysis

Sentiment analysis is used to determine the sentiment or opinion expressed in a piece of text, such as product reviews, social media posts, or customer feedback. Machine learning models can classify text as positive, negative, or neutral.

Example: A company uses sentiment analysis to analyze customer reviews of its products. The model classifies reviews into categories like "positive," "neutral," or "negative," providing valuable insights for improving products.

Machine Translation

Machine learning models, particularly neural machine translation (NMT) models, are used for automatic translation between languages. These models are trained on large parallel corpora and can translate text with high accuracy.

Example: Google Translate uses machine learning models to translate text between various languages, improving translation quality over time by learning from user feedback.

Machine learning has widespread applications across various industries, from healthcare and finance to self-driving cars and e-commerce. These applications are transforming how we interact with technology and solving complex problems more efficiently. As machine learning continues to evolve, it is expected to play an even larger role in shaping the future of many industries. In the following chapters, we will delve deeper into the mathematical foundations of machine learning and explore its underlying techniques in greater detail.

1.5 Challenges in Machine Learning

While machine learning has made significant advancements in recent years, it still faces a variety of challenges that impact its effectiveness, reliability, and scalability. In this section, we will discuss some of the key challenges in machine learning and how they affect the development and deployment of machine learning models.

1.5.1 Data-Related Challenges

Data Quality and Availability

The quality and quantity of data are crucial factors in training effective machine learning models. In many real-world scenarios, obtaining large and high-quality datasets is challenging. Incomplete, inconsistent, or biased data can lead to poor model performance and inaccurate predictions. Moreover, obtaining labeled data for supervised learning tasks can be costly and time-consuming.

Example: In healthcare, acquiring labeled medical data such as annotated images for disease diagnosis is often difficult due to privacy concerns, and manual annotation by experts is time-consuming and expensive.

Data Privacy and Security

Machine learning models often require large amounts of data, some of which may contain sensitive information such as personal, financial, or medical data. Protecting user privacy and ensuring the security of this data is a major concern. The introduction of regulations such as the GDPR (General Data Protection Regulation) has made it necessary to adopt privacy-preserving machine learning techniques.

Example: When using machine learning for analyzing financial data, it is essential to ensure that sensitive customer information is protected and that the model is not violating privacy laws.

1.5.2 Model-Related Challenges

Overfitting and Underfitting

Overfitting occurs when a machine learning model is too complex and learns not only the underlying patterns in the data but also the noise or random fluctuations. This leads to poor generalization to new data. On the other hand, underfitting happens when the model is too simple to capture the underlying patterns in the data, resulting in poor performance.

Example: A deep learning model trained on a small dataset of images might overfit, learning irrelevant details and not generalizing well to new, unseen images. In contrast, a linear regression model on a complex dataset might underfit if it cannot capture the non-linear relationships in the data.

Model Interpretability

Many machine learning models, especially deep learning models, are often considered "black boxes" because it is difficult to understand how they make predictions. Lack of interpretability can be a problem in applications where explainability is important, such as healthcare, finance, and law. Understanding the reasoning behind a model's predictions is crucial for building trust in the system.

Example: In medical diagnosis, doctors need to understand why a machine learning model has classified a patient as high-risk for a certain disease. If the model's decision-making process is not interpretable, it may be difficult for doctors to trust its recommendations.

Scalability

Many machine learning algorithms, particularly deep learning models, require large amounts of computational resources and storage. As the size of the data grows, training models becomes more computationally expensive and time-consuming. This presents a challenge when trying to scale machine learning applications to handle big data in real-time.

Example: Training a large-scale deep learning model on millions of images may require powerful GPUs or cloud computing resources, which may not be available for all organizations.

1.5.3 Algorithm-Related Challenges

Bias and Fairness

Machine learning models can inherit biases present in the training data. If the training data reflects societal biases, the model may learn these biases and make unfair or discriminatory decisions. This is particularly problematic in sensitive areas such as hiring, lending, and criminal justice, where biased algorithms can perpetuate discrimination and inequality.

Example: A hiring algorithm trained on historical employee data may inadvertently favor candidates from certain demographic groups, leading to biased hiring decisions that disadvantage other groups.

Adversarial Attacks

Machine learning models, particularly those used in image recognition, have been found to be vulnerable to adversarial attacks. In these attacks, small, carefully crafted perturbations to the input data can cause the model to make incorrect predictions with high confidence. Adversarial attacks pose a significant security risk, especially in applications like self-driving cars and facial recognition.

Example: An attacker could introduce slight noise to an image of a stop sign, causing a self-driving car's model to misinterpret it as a yield sign, leading to potentially dangerous consequences.

1.5.4 Real-World Deployment Challenges

Model Deployment and Maintenance

Once a machine learning model is developed, deploying it into a real-world environment can be challenging. Models may need to be updated or retrained periodically as new data becomes available, and deployment in dynamic, changing environments can lead to performance degradation over time. Ensuring that the model remains effective and adapts to changing conditions is a key challenge.

Example: A recommendation system used by an e-commerce platform may need to be updated regularly to account for changing customer preferences, product availability, and seasonal trends.

Ethical Considerations

The ethical implications of machine learning are an ongoing area of concern. Decisions made by machine learning models can impact individuals and communities in profound ways, and it is essential to ensure that models are developed and used in an ethical manner. This includes issues such as transparency, accountability, and the potential for unintended consequences.

Example: In predictive policing, machine learning models may be used to determine where to allocate law enforcement resources. If the model is trained on biased data, it may lead to unfair targeting of certain communities, raising ethical concerns.

Machine learning presents numerous challenges that need to be addressed in order to ensure its successful application and deployment in real-world scenarios. These challenges range from issues related to data quality and privacy to difficulties with model interpretability and algorithm fairness. Overcoming these challenges requires ongoing research, development, and the adoption of best practices to create machine learning systems that are not only effective but also ethical, secure, and transparent. In the next section, we will explore some of the techniques and strategies used to address these challenges in the field of machine learning.

1.6 Supervised Learning

Supervised learning is a type of machine learning where the model is trained on a labeled dataset. The objective is to learn a mapping function from input variables (features) to output variables (labels or target values) based on the given training data. The algorithm attempts to generalize the relationships within the data in such a way that it can make accurate predictions on new, unseen data.

1.6.1 Input Representation

In supervised learning, the **input representation** refers to how the input data is structured and formatted to be fed into the machine learning algorithm. This representation plays a crucial role in the performance of the learning algorithm. Proper input representation should capture all the relevant information needed for making accurate predictions.

25

The key aspects of input representation include:

- **Features**: Features are individual measurable properties or characteristics of the data. These could be numerical or categorical. For example, in a housing price prediction model, the features could include the size of the house, the number of rooms, or the proximity to essential services.
- Feature Vectors: A feature vector is the mathematical representation of an input. In most algorithms, an input is represented as a vector of features, denoted by $\mathbf{x} = (x_1, x_2, \dots, x_n)$, where each x_i represents a specific feature of the input data. These vectors allow the algorithm to process and analyze the input efficiently.
- Data Normalization: Many machine learning algorithms assume that all features have the same scale. If features have different units or ranges, it may lead to an imbalanced influence of features on the model. Data normalization ensures that all features are scaled to a similar range, often between 0 and 1.
- Categorical Features: Sometimes features are categorical (e.g., color, type). These need to be encoded into numerical form using techniques like one-hot encoding or label encoding so that they can be processed by the algorithm.
- Missing Data Handling: It's common for datasets to have missing values in certain features. Methods like imputation or dropping missing values are used to handle this issue and ensure the model is trained on complete data.

Mathematically, input representation can be written as:

$$X = \{x_1, x_2, \ldots, x_m\}$$

where $\mathbf{x_i}$ represents a feature vector corresponding to the *i*-th sample, and m is the number of training samples.

1.6.2 Hypothesis Class

The **hypothesis class**, denoted H, is the set of all possible models or functions that the learning algorithm can use to map inputs to outputs. The hypothesis class defines the space of potential mappings that the learning algorithm will explore to find the best possible model.

The learning process involves selecting the hypothesis $h^* \in H$ that minimizes the error between the predicted and actual outputs. The choice of hypothesis class depends on the type of algorithm used. For example:

- Linear Regression: The hypothesis class is the set of all linear functions of the form $h(\mathbf{x}) = \mathbf{w}^{\top}\mathbf{x} + b$, where \mathbf{w} is the weight vector and b is the bias term.
- **Decision Trees**: The hypothesis class consists of all possible decision tree structures that split the data based on feature values.
- **Neural Networks**: The hypothesis class includes all possible network structures with different activation functions, layer configurations, and weights.

The learning algorithm searches through the hypothesis class to find the best model that generalizes well to new data. The model maps an input vector \mathbf{x} to the output y:

$$h(\mathbf{x}) = y$$

where h is a hypothesis in the hypothesis class H, \mathbf{x} is the input feature vector, and y is the predicted output.

1.6.3 Version Space

The **version space** is the subset of the hypothesis class that is consistent with the training data. It contains all the hypotheses that make correct predictions on the given labeled data. In other words, the version space is the set of hypotheses that fit the training data exactly.

As the algorithm learns, it narrows down the version space to find the best hypothesis. The version space may start out large, but it will shrink as the model learns from more examples. A key goal of supervised learning

is to narrow down the version space to a single hypothesis or a small set of hypotheses that accurately predict the output.

Mathematically, if D represents the training dataset (a set of input-output pairs), the version space is defined as:

$$VS(D) = \{ h \in H \mid h(\mathbf{x_i}) = y_i \text{ for all } (\mathbf{x_i}, y_i) \in D \}$$

where h is a hypothesis, and $(\mathbf{x_i}, y_i)$ represents an input-output pair from the training dataset. The version space contains all hypotheses that correctly classify every training example.

The concept of the version space is particularly important in the context of inductive learning, as it helps to describe the process of hypothesis selection. If the version space is too large, the model might struggle to make accurate predictions, and if it is too small, it might lead to overfitting.

Example

Consider a simple example of a supervised learning task: predicting whether a customer will buy a product based on their age and income. The features are "age" and "income," and the target variable is "buy" (whether the customer buys the product or not).

- The input representation is a feature vector $\mathbf{x} = (\text{age}, \text{income})$. - The hypothesis class could be a set of linear decision functions that separate customers who buy the product from those who do not. - The version space will consist of all decision boundaries (lines in the 2D space of age and income) that correctly classify the training examples.

The learning algorithm will search through the hypothesis class and select the best model that minimizes the classification error on the training data.

In supervised learning, the algorithm learns a mapping from input to output based on labeled training data. The process involves representing the input in a suitable way, selecting a hypothesis class, and narrowing down the version space to find the best hypothesis that generalizes well to new data.

The effective representation of data and the appropriate choice of hypothesis class and version space are key factors that influence the success of supervised learning algorithms. By understanding these concepts, one can design algorithms that learn from data and make accurate predictions.

Chapter 2

Introduction to Deep Learning

2.1 Introduction

Deep Learning is a subset of Machine Learning (ML), which utilizes Artificial Neural Networks (ANNs) with multiple layers to learn complex patterns from data. It is particularly useful for tasks that involve large datasets and require feature extraction, such as image recognition, speech processing, and natural language processing (NLP).

2.1.1 Why Deep Learning?

- Handles large-scale data efficiently.
- Learns complex representations automatically without manual feature engineering.
- Provides state-of-the-art performance in various domains such as health-care, finance, and autonomous driving.
- Enables transfer learning, where pre-trained models can be fine-tuned for different tasks.

2.2 Artificial Neural Networks (ANNs)

An Artificial Neural Network (ANN) is inspired by the human brain and consists of neurons connected by weighted edges. It is capable of learning

from input data and making predictions based on learned patterns.

2.2.1 Structure of an Artificial Neural Network (ANN)

An Artificial Neural Network (ANN) is a computational model inspired by the structure and function of the human brain. It consists of multiple layers of interconnected neurons, where each layer has a specific role in processing information.

Layers of an ANN

An ANN consists of three main types of layers:

• Input Layer:

- The first layer of the ANN that receives the raw data.
- Each neuron in this layer corresponds to a feature in the dataset.
- If the dataset has n features, the input layer contains n neurons.

Mathematically, the input layer is represented as:

$$X = [x_1, x_2, ..., x_n]$$

where x_i represents an individual feature.

• Hidden Layers:

- Hidden layers perform intermediate computations.
- Each neuron receives inputs, applies a weighted sum, and passes it through an activation function.
- The number of hidden layers and neurons per layer depends on the complexity of the problem.

The computation at a hidden layer neuron is given by:

$$z = WX + b$$

where:

-W is the weight matrix.

- -X is the input vector.
- b is the bias term.
- -z is the weighted sum before applying an activation function.

The activation function introduces non-linearity:

$$a = f(z)$$

• Output Layer:

- The final layer that produces the ANN's output.
- The number of neurons depends on the type of task:
 - * Binary Classification (Spam/Not Spam) \rightarrow 1 neuron (Sigmoid activation).
 - * Multi-class Classification (Digit recognition: 0-9) $\rightarrow n$ neurons (Softmax activation).
 - * Regression (House Price Prediction) \rightarrow 1 neuron (Linear activation).

For classification tasks, the output neuron applies an activation function like Sigmoid:

$$y = \frac{1}{1 + e^{-z}}$$

where y is the probability of the class label.

2.2.2 Summary

The layers of an ANN and their purposes are summarized as follows:

Layer	Purpose
Input Layer	Receives raw features as input.
Hidden Layers	Processes data using weights, biases, and activation
	functions.
Output Layer	Produces final output (classification label or numerical
	value).

Table 2.1: Layers of an Artificial Neural Network

An ANN transforms input data step by step to generate meaningful predictions. The hidden layers extract important patterns, making deep learning models powerful for various tasks like image recognition, speech processing, and medical diagnosis.

2.2.3 Common Activation Functions

- Sigmoid: $f(z) = \frac{1}{1+e^{-z}}$ Used in binary classification.
- ReLU (Rectified Linear Unit): $f(z) = \max(0, z)$ Prevents vanishing gradient problems in deep networks.
- Tanh: $f(z) = \frac{e^z e^{-z}}{e^z + e^{-z}}$ Centers the output around zero, making optimization easier.
- **Softmax**: Converts logits into probability distributions for multi-class classification problems.

2.3 Feedforward Neural Networks (FFNs)

A Feedforward Neural Network (FFN) is the simplest type of ANN where information moves in one direction (forward) from input to output. It does not have cycles or loops.

2.3.1 Forward Propagation

Forward propagation involves calculating the output of each neuron in every layer by applying weights, biases, and activation functions. The mathematical formulation for layer l is:

$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]} (2.1)$$

$$a^{[l]} = f(z^{[l]}) (2.2)$$

where l is the layer index, $a^{[l-1]}$ represents activations from the previous layer, and $W^{[l]}$ and $b^{[l]}$ are the weights and biases for layer l.

33

2.3.2 Example Calculation

Consider a simple neural network with:

- Input: $X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- Weights: $W = \begin{bmatrix} w_1 & w_2 \end{bmatrix}$
- \bullet Bias: b

The neuron computes:

$$z = w_1 x_1 + w_2 x_2 + b (2.3)$$

Applying ReLU activation:

$$a = \max(0, z) \tag{2.4}$$

2.4 Backpropagation and Training

Once the output is generated, the network needs to learn optimal weights through **backpropagation**. Backpropagation uses **gradient descent** to minimize a loss function.

2.4.1 Loss Function

A loss function quantifies the error of the network's predictions. Common loss functions include:

• Mean Squared Error (MSE) for regression:

$$L = \frac{1}{n} \sum (y - \hat{y})^2 \tag{2.5}$$

• Cross-Entropy Loss for classification:

$$L = -\sum y \log(\hat{y}) \tag{2.6}$$

2.4.2 Gradient Descent

Gradient Descent updates weights iteratively to minimize the loss function:

$$W = W - \alpha \frac{\partial L}{\partial W} \tag{2.7}$$

where α is the learning rate, which determines the step size in the direction of the gradient.

The gradient of the loss function with respect to weight W is computed as:

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial W} \tag{2.8}$$

Using Stochastic Gradient Descent (SGD), the update rule becomes:

$$W^{(t+1)} = W^{(t)} - \alpha \nabla L(W^{(t)})$$
(2.9)

where t is the iteration step.

Example: If the true value is y = 1 and the predicted value is $\hat{y} = 0.8$, then for MSE:

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) = 2(0.8 - 1) = -0.4 \tag{2.10}$$

which is then used to update the weights.

2.5 Convolutional Neural Networks (CNNs)

2.5.1 Definition and Introduction

A Convolutional Neural Network (CNN) is a class of deep learning models specifically designed for processing structured grid-like data, such as images. Unlike traditional Artificial Neural Networks (ANNs), CNNs exploit spatial hierarchies in data using convolutional layers. They are widely used in image recognition, object detection, and other computer vision tasks.

CNNs consist of multiple layers, including convolutional layers, activation functions, pooling layers, and fully connected layers, which work together to extract and learn important features from the input data.

2.5.2 Limitations of Artificial Neural Networks (ANNs)

Traditional Artificial Neural Networks (ANNs) have several limitations when applied to image processing and spatial data:

- ANNs treat each input feature independently, ignoring spatial relationships.
- The number of parameters grows significantly with image size, leading to high computational cost.
- Fully connected layers require large datasets and extensive training.

These limitations necessitate the use of Convolutional Neural Networks (CNNs), which preserve spatial structures and reduce parameter complexity.

2.5.3 Significance of CNNs

CNNs are designed to process data with grid-like topology, such as images, by capturing spatial hierarchies. They achieve this through three fundamental concepts:

- Local Connectivity: Instead of fully connecting all neurons, CNNs connect only to local regions.
- Shared Weights: Filters (kernels) slide over the input, reducing the number of parameters.
- Pooling: Downsampling layers reduce spatial dimensions while retaining important features.

2.5.4 Mathematical Formulation of CNN Layers

Convolutional Layer

The convolution operation applies a filter (kernel) to the input image:

$$S(i,j) = \sum_{m} \sum_{n} I(i+m, j+n) K(m, n)$$
 (2.11)

where:

- I(i,j) is the input image matrix,
- K(m, n) is the kernel,
- S(i, j) is the output feature map.

Activation Function (ReLU)

CNNs use non-linear activation functions such as ReLU (Rectified Linear Unit):

$$f(x) = \max(0, x) \tag{2.12}$$

This ensures non-linearity, allowing CNNs to learn complex patterns.

Pooling Layer

Pooling reduces spatial dimensions while retaining critical information. The most common operation is Max Pooling:

$$P(i,j) = \max\{S(i+m,j+n)\}, \quad \forall m,n \in \text{pool size}$$
 (2.13)

Pooling reduces overfitting and computation.

Fully Connected Layer (FC)

After convolution and pooling, feature maps are flattened and passed through a fully connected layer:

$$y = Wx + b \tag{2.14}$$

where:

- W is the weight matrix,
- x is the input vector,
- b is the bias term,
- y is the final output.

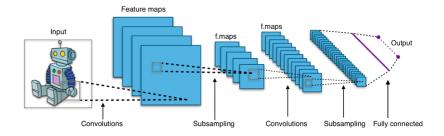


Figure 2.1: Basic CNN Architecture: Convolution, ReLU, Pooling, Fully Connected Layers

2.5.5 CNN Architecture Diagram

2.5.6 Example: Image Classification with CNN

For an image classification task, CNN processes an input image as follows:

- 1. The convolutional layers extract features such as edges and textures.
- 2. Pooling layers reduce dimensionality and retain essential features.
- 3. Fully connected layers classify the image into predefined categories.
- 4. The final output is a probability distribution over possible labels.

This hierarchical feature extraction makes CNNs effective for visual recognition tasks.

CNNs overcome the limitations of traditional ANNs by leveraging local connectivity, weight sharing, and hierarchical feature extraction. They are widely used in computer vision applications, including image classification, object detection, and segmentation.

2.6 Recurrent Neural Networks (RNNs)

Recurrent Neural Networks (RNNs) are a class of artificial neural networks specifically designed for sequential data. Unlike traditional feedforward networks, RNNs introduce connections that allow information to persist across time steps, making them well-suited for tasks involving time-series data, natural language processing, and speech recognition.

2.6.1 Limitations of Traditional Neural Networks

Traditional methods such as Artificial Neural Networks (ANNs) and Convolutional Neural Networks (CNNs) have limitations when dealing with sequential data:

- Fixed Input Size: ANNs and CNNs assume fixed-length inputs and cannot handle variable-length sequences.
- Lack of Memory: These networks treat each input independently and do not retain past information, making them ineffective for sequential dependencies.
- Scalability Issues: Increasing the size of the network to capture sequential relationships leads to an explosion in parameters, making training computationally expensive.

To address these issues, RNNs introduce the concept of hidden states and loops, enabling the network to retain information over multiple time steps.

2.6.2 Why RNNs?

RNNs overcome the limitations of ANNs and CNNs by incorporating **temporal dependencies**. They achieve this by:

- Maintaining a hidden state that evolves over time.
- Sharing parameters across different time steps.
- Being adaptable to variable-length sequences, making them ideal for tasks like language modeling and speech processing.

2.6.3 Mathematical Formulation of RNN

An RNN processes sequential data by maintaining a hidden state h_t that captures past information. The update rule for an RNN is given by:

$$h_t = f(W_h h_{t-1} + W_x x_t + b_h) (2.15)$$

where:

• h_t is the hidden state at time step t,

- x_t is the input at time step t,
- W_h and W_x are weight matrices,
- b_h is the bias term,
- f is a non-linear activation function (commonly tanh or ReLU).

The output at each time step is computed as:

$$y_t = f(W_y h_t + b_y) (2.16)$$

where W_y is the weight matrix for the output layer.

2.6.4 Backpropagation Through Time (BPTT)

Training an RNN involves a modified version of backpropagation known as **Backpropagation Through Time (BPTT)**. The loss function is computed over all time steps, and gradients are propagated backward through time.

The total loss is:

$$L = \sum_{t=1}^{T} L_t \tag{2.17}$$

where L_t is the loss at each time step. Gradients are computed recursively using the chain rule.

2.6.5 Challenges with RNNs

Despite their advantages, RNNs suffer from:

- Vanishing Gradient Problem: When gradients become too small, earlier layers learn very slowly, making it difficult to capture long-term dependencies.
- Exploding Gradients: When gradients become too large, training becomes unstable.
- Short-Term Memory: Standard RNNs struggle to remember information over long sequences.

These challenges led to the development of advanced architectures such as Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU).

Recurrent Neural Networks provide a powerful framework for handling sequential data by maintaining hidden states and capturing temporal dependencies. However, they face challenges like vanishing gradients, which are addressed by LSTMs and GRUs. RNNs are widely used in applications such as language modeling, speech recognition, and time-series forecasting.

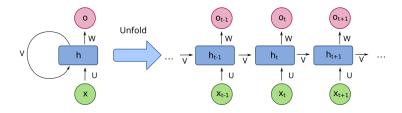


Figure 2.2: Recurrent Neural Network Structure

2.7 Transformers in NLP

2.7.1 Introduction

Natural Language Processing (NLP) has undergone a significant transformation with the advent of **Transformers**, a deep learning architecture that has revolutionized how machines understand and generate human language. Prior to transformers, NLP relied on **Recurrent Neural Networks** (RNNs), including their variants such as **Long Short-Term Memory** (LSTM) and **Gated Recurrent Units** (GRUs). However, these models had several limitations:

- Sequential Processing Bottleneck RNNs and LSTMs process text sequentially, making training and inference slower.
- Difficulty Handling Long Dependencies They struggle to retain information over long sequences.
- Vanishing and Exploding Gradients Learning long-range dependencies becomes challenging as gradients diminish.

• Limited Parallelization – Due to their sequential nature, they cannot fully utilize modern hardware capabilities.

To overcome these challenges, **Vaswani et al.** introduced the Transformer model in their 2017 paper "Attention is All You Need". Unlike previous approaches, transformers leverage the **self-attention mechanism** to process entire sequences in parallel, significantly improving efficiency and capturing long-range dependencies more effectively.

Since their introduction, transformers have become the foundation for modern NLP applications, including:

- Machine Translation (e.g., Google Translate using Transformer models)
- Text Summarization (e.g., BART, T5)
- Sentiment Analysis and Classification (e.g., BERT, RoBERTa)
- Conversational AI and Chatbots (e.g., ChatGPT, Claude, Gemini)
- Search Engine Optimization and Information Retrieval (e.g., BERT in Google Search)
- Text Generation (e.g., GPT-3, GPT-4, LLaMA, Mixtral)

2.7.2 Why Are Transformers Revolutionary?

Transformers introduce several key innovations that differentiate them from previous models:

- 1. **Self-Attention Mechanism** Instead of processing words sequentially, the model attends to all words in a sequence simultaneously, capturing complex dependencies.
- 2. **Positional Encoding** Since transformers lack inherent recurrence, positional encodings are added to maintain word order.
- 3. **Parallel Processing** Unlike RNNs, transformers process the entire input sequence at once, significantly improving training speed.

4. **Scalability** – Transformers are highly scalable and can handle large datasets, making them ideal for modern AI applications.

These innovations have led to the rise of state-of-the-art transformer models, including **BERT**, **GPT**, **T5**, and **LLaMA**, which are now at the forefront of AI-driven NLP systems.

2.8 Autoencoder and Variational Autoencoder (VAE)

2.8.1 Autoencoder (AE)

An autoencoder is a neural network used for unsupervised learning, primarily for dimensionality reduction and feature extraction. It consists of two main components:

- Encoder: Maps the input x to a lower-dimensional latent representation z.
- **Decoder:** Reconstructs the input x from the latent space representation z.

The objective is to minimize the reconstruction error:

$$L_{AE} = ||x - \hat{x}||^2 \tag{2.18}$$

2.8.2 Variational Autoencoder (VAE)

A VAE introduces a probabilistic framework by modeling the latent space as a distribution. Instead of encoding a single latent vector, the encoder outputs the mean (μ) and variance (σ^2) of a Gaussian distribution:

$$z = \mu + \sigma \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$
 (2.19)

The loss function consists of two terms:

• Reconstruction Loss: Measures the difference between input and reconstructed output.

2.8. AUTOENCODER AND VARIATIONAL AUTOENCODER (VAE) 43

• **KL Divergence Loss:** Ensures the latent distribution is close to a standard normal distribution.

Total loss function:

$$L_{VAE} = \mathbb{E}[||x - \hat{x}||^2] + D_{KL}(q(z|x)||p(z))$$
 (2.20)