

Artificial Neural Networks: From Theory to Practice

A Comprehensive Textbook for Computer Science Students

September 27, 2025

Contents

1	Introduction to Machine Learning	5
1.1	What is Learning?	5
1.1.1	Core Elements of Learning	5
1.1.2	Key Questions in Learning	5
1.2	What is Reasoning?	5
1.3	Types of Reasoning: Comprehensive Overview	5
1.3.1	Inductive Reasoning	6
1.3.2	Deductive Reasoning	6
1.3.3	Abductive Reasoning (Inference to Best Explanation)	6
1.4	Animal Learning	6
1.4.1	Example: Bait Shyness in Rats	6
1.5	Human Learning: The Cognitive Approach	6
1.5.1	Cognitive Stages of Development	6
1.5.2	The Role of Memory	7
1.5.3	Learning Styles and Strategies	7
1.6	What is Machine Learning?	7
1.6.1	Parallel: Spam Email Filtering	7
1.7	Types of Machine Learning	8
1.7.1	Supervised Learning	8
1.7.2	Unsupervised Learning	8
1.7.3	Reinforcement Learning	8
1.8	The Data and Observation Model	8
1.8.1	The Design Matrix	8
1.9	Probabilistic Model of Learning	8
1.10	Inductive Bias in Machine Learning	9
1.10.1	What is Inductive Bias?	9
1.10.2	Types of Inductive Biases	9
1.10.3	Inductive Bias in Specific Architectures	9
1.11	Mathematical Foundations of Learning	10
1.11.1	Learning as Optimization	10
1.11.2	PAC Learning Framework	10
1.12	Symbolic AI vs Machine Learning	10
1.12.1	What is Symbolic AI?	10
1.12.2	Symbolic AI vs Machine Learning Comparison	10
2	Foundations of Computation	13
2.1	What is Computation?	13
2.2	Computational Models: Theoretical Foundations	13
2.3	Four Fundamental Computational Models	13
2.3.1	Turing Machine (Alan Turing, 1936)	13
2.3.2	Lambda Calculus (Alonzo Church, 1930s)	14
2.3.3	Cellular Automata (Stanislaw Ulam, John von Neumann, later Conway)	14
2.3.4	Biological Computation (Inspired by Nature)	15
2.4	Biological Neural Networks: Nature's Computational Model	15
2.4.1	Characteristics of Biological Neural Networks	15
2.4.2	Neuron Structure	16

2.4.3	Neuron Structure and Components	16
2.4.4	Signal Transmission and Firing	16
2.4.5	Synapses: Chemistry and Types	16
2.4.6	Plasticity and Learning	17
2.5	Artificial Neural Networks	17
2.5.1	Introduction: From Biology to Computation	17
2.5.2	The Abstract Neuron: Building Block of Intelligence	17
2.6	Activation Functions: Mathematical Properties	17
2.6.1	Common Activation Functions	17
2.6.2	Mathematical Requirements for Activation Functions	18
2.6.3	Neural Networks as Function Approximators	18
2.7	Artificial Neural Networks: From Theory to Implementation	19
2.7.1	Fundamental Architecture: Primitive Functions and Composition Rules	19
2.7.2	Neural Networks as Function Approximators	20
2.7.3	Function Approximation: The Classical Problem	20
2.7.4	Learning from Data: The Key Difference	21
2.8	Computational Complexity in Neural Networks	22
2.8.1	Forward Pass Complexity	22
2.8.2	Backward Pass Complexity (Backpropagation)	22
2.8.3	Scalability Considerations	22
3	Linear Classification	23
3.1	Learning Goals	23
3.2	Fundas: Mathematical Foundations	23
3.2.1	Vector Representation	23
3.2.2	Binary Classification Framework	23
3.3	Introduction to Binary Classification	23
3.3.1	Real-World Applications	24
3.4	Binary Linear Classifiers: First Principles	24
3.4.1	Core Concept	24
3.4.2	Mathematical Formulation	24
3.4.3	Eliminating Redundancy: The Bias-Threshold Trick	25
3.5	Geometric Interpretation	25
3.5.1	Input Space Perspective	25
3.5.2	Weight Space Perspective	25
3.6	Hand-Designed Logic Functions	25
3.6.1	AND Function	25
3.6.2	OR Function	26
3.6.3	NOT Function	26
3.7	Limitations of Linear Classifiers	26
3.7.1	Linear Separability	26
3.7.2	Convexity and Separability	26
3.7.3	The XOR Problem	26
3.7.4	Overcoming Limitations: Basis Functions	26
3.8	Practical Considerations	27
3.8.1	Feature Engineering	27
3.8.2	Model Selection	27
3.8.3	Performance Evaluation	27
3.9	Chapter Summary	27
4	The Perceptron	29
4.1	Historical Development of Neural Networks	29
4.1.1	Timeline of Neural Network Evolution	29
4.1.2	Threshold Logic: The Foundation	29
4.2	McCulloch-Pitts Neuron: The First Artificial Neuron	29
4.2.1	Mathematical Model	30
4.2.2	Logic Gate Implementation	30
4.2.3	Limitations of McCulloch-Pitts Neurons	30

4.3	The Perceptron: A Detailed Introduction	30
4.3.1	Definition: The Anatomy of a Perceptron	31
4.4	The Perceptron Learning Rule	31
4.4.1	Mathematical Derivation	31
4.4.2	Convergence Theorem (Rosenblatt, 1962)	31
4.4.3	Learning Rate Analysis	32
4.5	Vectorisation of Perceptron Learning Rule: A Matrix-Based Example	32
4.5.1	Mathematical Foundation of Vectorization	32
4.5.2	Define the Matrices	32
4.5.3	Epoch 1: Mathematical Flow	33
4.5.4	Mathematical Insight: Gradient Descent Connection	34
4.5.5	Computational Complexity Analysis	34
4.6	Geometric Interpretation of the Perceptron	35
4.6.1	Mathematical Foundation of Decision Boundaries	35
4.6.2	Example: The NOT Function	35
4.6.3	Example: The AND Function	35
4.7	Limitations of the Perceptron	36
4.7.1	Linear Separability Constraint	36
4.7.2	Solutions to Linear Separability Limitation	38
4.8	Historical Impact and Legacy	38
4.8.1	The Perceptron Controversy	38
4.8.2	Modern Relevance	39
4.9	Exercises	39

Chapter 1

Introduction to Machine Learning

1.1 What is Learning?

Learning is the process of acquiring new knowledge, skills, or behaviors through experience. This process transforms inputs—such as data, experiences, or information—into useful capabilities like expertise, new skills, or predictive models.

1.1.1 Core Elements of Learning

Every learning process involves these fundamental components:

1. **Input:** Data, experiences, or information that enters the learning system
2. **Processing:** The manipulation or transformation of that data according to learning rules
3. **Output:** The result—new knowledge, skills, or predictive capabilities
4. **Feedback:** Information about the effectiveness of learning outcomes
5. **Memory:** The ability to retain and access learned information

1.1.2 Key Questions in Learning

- What are the essential inputs for the learning process?
- How do we measure the effectiveness and success of learning?
- What are the underlying mechanisms and processes by which learning occurs?
- How can we generalize from specific experiences to handle new situations?

1.2 What is Reasoning?

Reasoning is the ability to draw logical conclusions from known facts or learned knowledge. Unlike learning, reasoning relies on logical inference rather than large amounts of data.

1.3 Types of Reasoning: Comprehensive Overview

Understanding different types of reasoning is crucial for designing effective learning systems. Each type has distinct characteristics and applications in both biological and artificial intelligence systems.

1.3.1 Inductive Reasoning

Definition: Inductive reasoning extracts patterns from observed data to make predictions about future or unseen cases. This approach moves from specific observations to general conclusions, yielding probable rather than certain results.

Key Characteristics:

- Most prevalent form of reasoning in the animal kingdom and primary mode in machine learning
- Forms the basis of most learned behaviors in animals
- Used extensively in deep learning and LLMs
- Enables generalization from limited examples to broader patterns

1.3.2 Deductive Reasoning

Definition: Deductive reasoning moves from general rules and premises to reach specific, guaranteed conclusions. It starts with a general rule and a specific case to reach a logical conclusion.

Key Characteristics:

- Provides certainty when premises are true
- Animals generally lack this capability for abstract reasoning
- LLMs can only mimic this through pattern matching
- Forms the basis of formal logic and mathematical proof

1.3.3 Abductive Reasoning (Inference to Best Explanation)

Definition: Abductive reasoning starts with an observation and seeks to find the simplest and most likely explanation. It's the process of finding a hypothesis that, if true, would best explain the observation.

1.4 Animal Learning

1.4.1 Example: Bait Shyness in Rats

Rats demonstrate a fundamental learning principle through their feeding behavior:

- They sample novel food cautiously
- If the food causes illness, they avoid it in the future
- Past experience directly informs future decisions

This natural learning process parallels challenges in machine learning.

1.5 Human Learning: The Cognitive Approach

Human learning is a complex and multifaceted process that has been studied extensively in psychology and neuroscience. It involves a combination of conscious and unconscious processes, leading to the acquisition of knowledge and skills that are both explicit (declarative) and implicit (procedural).

1.5.1 Cognitive Stages of Development

Jean Piaget's theory of cognitive development provides a classic framework for understanding how learning capabilities evolve from infancy to adulthood.

- **Sensorimotor Stage (0-2 years):** Learning occurs through sensory experiences and motor interactions with the environment. Object permanence is a key milestone.
- **Preoperational Stage (2-7 years):** Children begin to think symbolically and use words and pictures to represent objects. Their thinking is egocentric.

- **Concrete Operational Stage (7-11 years):** Children begin to think logically about concrete events. They grasp concepts like conservation.
- **Formal Operational Stage (12+ years):** Abstract reasoning and hypothetical thinking emerge.

This staged progression suggests that the ability to learn and the types of learning that are possible change fundamentally over a lifetime, a concept that has parallels in the development of more sophisticated machine learning models.

1.5.2 The Role of Memory

Memory is central to learning. The Atkinson-Shiffrin model is a classic theory that proposes three stages of memory:

1. **Sensory Memory:** A very brief buffer for sensory information.
2. **Short-Term Memory (Working Memory):** Holds a small amount of information for a short duration. It is where conscious thought and processing occur. This is analogous to the memory (RAM) of a computer.
3. **Long-Term Memory:** The vast, semi-permanent storage of knowledge and skills. This is analogous to a computer's hard drive.

The process of moving information from short-term to long-term memory, known as encoding, is critical for learning. Retrieval is the process of accessing this stored information. In machine learning, the "memory" is stored in the model's parameters (weights).

1.5.3 Learning Styles and Strategies

Humans employ a variety of strategies to learn, which can be broadly categorized:

- **Rote learning:** Memorization through repetition (e.g., flashcards). This is similar to overfitting in machine learning, where a model memorizes the training data.
- **Observational learning:** Learning by watching others.
- **Associative learning:** Connecting stimuli or events that occur together in the environment (classical and operant conditioning).
- **Cognitive learning:** Learning through understanding, reasoning, and problem-solving. This is the goal of more advanced AI systems.

1.6 What is Machine Learning?

Machine learning is a subfield of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. At its core, machine learning involves developing models that can identify patterns in data and make predictions or decisions based on those patterns.

1.6.1 Parallel: Spam Email Filtering

Consider how biological learning principles apply to spam detection:

- **Naive approach:** Memorize all past spam emails
- **Problem:** Cannot classify previously unseen emails
- **Solution:** Extract generalizable patterns (like words, phrases, or sender patterns)
- **Key insight:** Both rats and spam filters must generalize from specific experiences to handle new, similar situations

1.7 Types of Machine Learning

There are three main categories of machine learning algorithms:

1.7.1 Supervised Learning

In supervised learning, the algorithm is trained on a labeled dataset, meaning that each data point is tagged with a correct output. The goal is to learn a mapping function that can predict the output for new, unseen data. Common supervised learning tasks include classification and regression.

1.7.2 Unsupervised Learning

Unsupervised learning deals with unlabeled data. The algorithm tries to learn the underlying structure of the data without any explicit guidance. Common tasks include clustering, dimensionality reduction, and density estimation.

1.7.3 Reinforcement Learning

Reinforcement learning is concerned with how an agent ought to take actions in an environment so as to maximize some notion of cumulative reward. The algorithm learns by trial and error, receiving feedback in the form of rewards or punishments.

1.8 The Data and Observation Model

In machine learning, we typically represent our data as a matrix. Let's denote the dataset as \mathcal{D} . A common convention is to represent the data as a design matrix, X .

1.8.1 The Design Matrix

The design matrix X is an $m \times n$ matrix, where m is the number of training examples (or observations) and n is the number of features (or variables). Each row of the matrix represents a single data point, and each column represents a feature.

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix}$$

Each row i corresponds to a data point \mathbf{x}_i^T , which is a row vector of size n .

1.9 Probabilistic Model of Learning

A powerful way to think about machine learning is from a probabilistic perspective. We can think of the learning process as finding a model that best explains the data. This is often framed as finding the parameters θ of a model that maximize the likelihood of observing the data.

The likelihood function is given by:

$$\mathcal{L}(\theta|X) = P(X|\theta)$$

The goal of learning is to find the parameters $\hat{\theta}$ that maximize this likelihood. This is known as Maximum Likelihood Estimation (MLE).

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(X|\theta)$$

1.10 Inductive Bias in Machine Learning

1.10.1 What is Inductive Bias?

Definition: Inductive bias refers to the set of assumptions that a learning algorithm makes to generalize from limited training data to unseen data.

Why is it Critical? Inductive bias is essential because:

- Machine learning models have limited training data
- Models must generalize from past observations to unseen cases
- Without appropriate bias, models may overfit (memorizing training data without learning generalizable patterns)
- All successful learning algorithms require appropriate assumptions about their domain

1.10.2 Types of Inductive Biases

Preference for Simpler Models (Occam's Razor)

- **Assumption:** Simpler explanations are preferred over complex ones
- **Example:** Decision trees with fewer splits are preferred because they generalize better
- **In Deep Learning:** Regularization techniques (L1, L2) penalize complex models

Smoothness Assumption

- **Assumption:** Data points that are close together should have similar outputs
- **Example:** In image classification, two similar images should belong to the same class
- **In ML:** K-Nearest Neighbors (KNN) assumes nearby data points have the same label

1.10.3 Inductive Bias in Specific Architectures

Convolutional Neural Networks (CNNs)

CNNs are designed for image processing and rely on key inductive biases:

1. Locality Bias (Local Connectivity)

- **Assumption:** Nearby pixels are more relevant than distant pixels
- **Example:** In facial recognition, CNN detects eyes, nose, mouth before recognizing entire face

2. Translation Invariance

- **Assumption:** An object should be recognized regardless of position
- **How it works:** CNNs use shared convolutional filters
- **Example:** Handwritten digit "3" recognized anywhere in the image

Recurrent Neural Networks (RNNs & LSTMs)

RNNs are designed for sequential data and rely on:

1. Temporal Dependency Bias

- **Assumption:** Recent information is more important than distant past
- **Example:** In "The cat sat on the mat", nearby words are more related

Transformers (BERT, GPT)

1. Attention-Based Bias (Self-Attention)

- **Assumption:** Important words can be anywhere in a sentence
- **Example:** In "The dog chased the ball...which was blue", "which" refers to "ball"

1.11 Mathematical Foundations of Learning

1.11.1 Learning as Optimization

Machine learning can be viewed as an optimization problem where we seek to find the best parameters θ that minimize a loss function $L(\theta)$:

$$\theta^* = \arg \min_{\theta} L(\theta)$$

Components of a Learning System

1. **Hypothesis Space \mathcal{H} :** The set of all possible functions the model can represent
2. **Loss Function $L(\theta)$:** Measures how well the model performs on the training data
3. **Optimization Algorithm:** Method to find θ^* (e.g., gradient descent)
4. **Regularization:** Techniques to prevent overfitting and improve generalization

The Bias-Variance Tradeoff

The expected prediction error can be decomposed as:

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

where:

- **Bias:** Error due to simplifying assumptions in the model
- **Variance:** Error due to sensitivity to small fluctuations in training data
- **Irreducible Error:** Inherent noise in the data

1.11.2 PAC Learning Framework

Probably Approximately Correct (PAC) learning provides theoretical foundations for when learning is possible.

A concept class \mathcal{C} is PAC-learnable if there exists an algorithm that, for any distribution \mathcal{D} and any $\epsilon, \delta > 0$, can find a hypothesis h such that:

$$\Pr[\text{error}(h) \leq \epsilon] \geq 1 - \delta$$

using polynomially many samples and computational steps.

1.12 Symbolic AI vs Machine Learning

1.12.1 What is Symbolic AI?

Also known as **Good Old-Fashioned AI (GOFAI)**, it represents knowledge using symbols, rules, and logic. It uses explicitly programmed rules for reasoning.

1.12.2 Symbolic AI vs Machine Learning Comparison

Feature	Symbolic AI	Machine Learning
Knowledge Source	Rules & logic	Data & patterns
Interpretability	Highly explainable	Often a black box
Adaptability	Rigid (manual updates)	Can generalize from data
Data Requirements	Minimal	Requires large datasets
Best Use Cases	Theorem proving	NLP, computer vision

Table 1.1: Comparison of Symbolic AI and Machine Learning.

Chapter 2

Foundations of Computation

2.1 What is Computation?

Computation is the process of performing calculations, manipulating data, or executing a sequence of operations to solve problems or transform inputs into desired outputs. It encompasses both the theoretical and practical aspects of processing information.

2.2 Computational Models: Theoretical Foundations

A **computational model** is a mathematical or conceptual framework that defines how computation is carried out. For an arbitrary computing model, the following metaphoric expression has been proposed:

$$\text{computation} = \text{storage} + \text{transmission} + \text{processing}$$

2.3 Four Fundamental Computational Models

2.3.1 Turing Machine (Alan Turing, 1936)

The Foundation of Algorithmic Computation

Core Characteristics

- **Style:** Imperative / mechanical model of computation
- **Core Idea:** A machine reads/writes symbols on an infinite tape with a finite set of rules
- **Representation:**
 - Infinite tape divided into cells
 - Head that can read/write and move left or right
 - Finite state machine controlling transitions

Strengths and Limitations

Strengths:

- Canonical model for algorithmic computability
- Basis of the Church—Turing Thesis
- Directly models sequential execution

Limitations:

- Low-level, not efficient
- Sequential by nature, doesn't capture parallelism well

Example: A Turing Machine can simulate any algorithm you'd run on a modern computer (given enough tape).

2.3.2 Lambda Calculus (Alonzo Church, 1930s)

The Foundation of Functional Computation

Core Characteristics

- **Style:** Functional model of computation
- **Core Idea:** Everything is a function. Computation = function application + substitution
- **Representation:**
 - Variables (x)
 - Function definitions ($\lambda x. \text{ expression}$)
 - Function application ($((f\ x))$)

Strengths and Limitations

Strengths:

- Basis of functional programming (Haskell, Lisp)
- Good for reasoning about higher-order functions, abstraction, recursion

Limitations:

- Abstract and symbolic; not naturally tied to hardware
- Efficiency is not modeled, just computability

Example: Addition can be defined entirely in terms of functions (Church numerals).

2.3.3 Cellular Automata (Stanislaw Ulam, John von Neumann, later Conway)

The Foundation of Distributed/Parallel Computation

Core Characteristics

- **Style:** Spatial / distributed model of computation
- **Core Idea:** Computation arises from simple local rules applied to a grid of cells over time
- **Representation:**
 - Infinite (or finite) grid of cells
 - Each cell has a finite state (e.g., alive/dead)
 - Transition rules depend only on the local neighborhood

Strengths and Limitations

Strengths:

- Good for modeling parallel, distributed, physical systems
- Supports universal computation (Conway's Life is Turing-complete)

Limitations:

- Not natural for symbolic or algebraic computation
- More suited for simulating dynamics

Example: Conway's *Game of Life* shows how simple rules produce complex, even universal, behaviors.

2.3.4 Biological Computation (Inspired by Nature)

The Foundation of Adaptive/Learning Computation

Core Characteristics

Biological models are inspired by living systems and emphasize parallelism, adaptability, and learning.

- **Learns from data** (training) rather than using fixed rules
- Massive parallelism and fault tolerance
- Self-organization and adaptation
- Pattern recognition and generalization

Examples of Biological Computation

Neural Networks:

- Inspired by the brain's neurons and synapses
- Computation happens through weighted sums and nonlinear activations
- Foundation of modern AI (deep learning for vision, NLP, etc.)

DNA Computing:

- Uses DNA strands and biochemical reactions to encode and solve problems
- Enables massive parallelism (billions of molecules interacting at once)
- Example: Adleman (1994) solved a small Hamiltonian Path problem with DNA

Swarm Intelligence:

- Inspired by ants, bees, and bird flocks
- Simple agents interacting lead to complex global solutions
- Example: Ant Colony Optimization for shortest path problems

2.4 Biological Neural Networks: Nature's Computational Model

The following are key characteristics that make biological neural networks powerful computational systems:

2.4.1 Characteristics of Biological Neural Networks

- **Highly interconnected:** Neurons form a complex web of connections
- **Robustness and Fault Tolerance:** The decay of nerve cells does not affect the overall function of the network significantly
- **Flexibility:** The ability to reorganize and adapt to new situations
- **Handling incomplete information:** Ability to infer appropriate outputs even when some inputs are missing or noisy
- **Parallel processing:** Multiple neurons can process information simultaneously

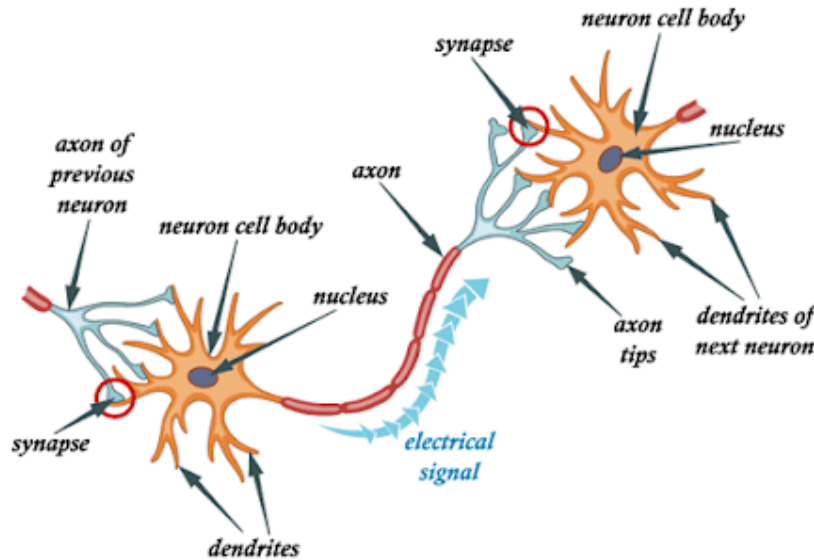


Figure 2.1: Structure of a biological neuron.

2.4.2 Neuron Structure

2.4.3 Neuron Structure and Components

- **Fundamental unit:** neuron (cell body / soma, dendrites, axon, synapses)
- **Dendrites** receive inputs; **axon** transmits output and branches to many synapses (often thousands)
- **Synapse:** junction between axon terminal and target cell
- **Synaptic junctions** form between presynaptic axon terminals and postsynaptic dendrites or the cell body

Typical Sizes

- soma $\sim 10\text{--}80\ \mu\text{m}$
- synaptic gap $\sim 200\ \text{nm}$
- neuron length from $0.01\ \text{mm}$ to $1\ \text{m}$

2.4.4 Signal Transmission and Firing

- **Resting potential** $\sim -70\ \text{mV}$; depolarization above threshold (roughly $\sim 10\ \text{mV}$) triggers firing
- **Action potentials** are all-or-none pulses sent down the axon; information is encoded in firing rate ($\sim 1\text{--}100\ \text{Hz}$)
- **Propagation speed** in brain tissue $\sim 0.5\text{--}2\ \text{m/s}$; synaptic transmission delay $\sim 0.5\ \text{ms}$
- After firing the membrane recovers (**refractory period**); synaptic effects decay with time constant $\sim 5\text{--}10\ \text{ms}$

2.4.5 Synapses: Chemistry and Types

- **Transmission** across synapse is chemical: neurotransmitters released from presynaptic terminal
- **Postsynaptic effect** can be excitatory (depolarizing) or inhibitory (hyperpolarizing)
- All endings of a given axon are typically either excitatory or inhibitory
- **Synaptic strength** depends on activity and can change over time (basis for learning)

2.4.6 Plasticity and Learning

Active synapses that repeatedly contribute to postsynaptic firing tend to strengthen; inactive ones weaken. **Hebb's rule** ("cells that fire together, wire together") describes this activity-dependent plasticity. Continuous modification of synaptic strengths underlies learning and memory formation.

2.5 Artificial Neural Networks

2.5.1 Introduction: From Biology to Computation

Artificial Neural Networks (ANNs) represent one of the most successful attempts to harness the computational principles observed in biological neural systems for solving complex problems.

2.5.2 The Abstract Neuron: Building Block of Intelligence

The output of a neuron can be expressed as:

$$Y = f\left(\sum_{i=1}^n W_i X_i + b\right) = f(\mathbf{W}^T \mathbf{X} + b)$$

Where:

- Y : Output of the neuron
- f : Activation function (primitive function that introduces non-linearity)
- W_i : Weight associated with input i (learnable parameter)
- X_i : Value of input i
- b : Bias term (learnable parameter that shifts the activation function)
- $\mathbf{W} = [W_1, W_2, \dots, W_n]^T$: Weight vector
- $\mathbf{X} = [X_1, X_2, \dots, X_n]^T$: Input vector

Mathematical Foundation: Linear Combination and Affine Transformation

The computation $\mathbf{W}^T \mathbf{X} + b$ represents an **affine transformation** of the input space. This can be broken down as:

1. **Linear transformation:** $\mathbf{W}^T \mathbf{X}$ scales and rotates the input vector
2. **Translation:** Adding bias b shifts the result by a constant

The activation function f then introduces non-linearity, enabling the neuron to model complex, non-linear relationships between inputs and outputs.

2.6 Activation Functions: Mathematical Properties

Activation functions are crucial for introducing non-linearity into neural networks. Different activation functions have distinct mathematical properties that affect learning dynamics.

2.6.1 Common Activation Functions

Step Function (Heaviside)

$$\phi(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

Properties: Non-differentiable, binary output, historically important for perceptron.

Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Properties:

- Smooth, differentiable: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Range: $(0, 1)$
- Problem: Vanishing gradients for large $|z|$

Hyperbolic Tangent

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = 2\sigma(2z) - 1$$

Properties:

- Range: $(-1, 1)$
- Zero-centered output
- Derivative: $\tanh'(z) = 1 - \tanh^2(z)$

Rectified Linear Unit (ReLU)

$$\text{ReLU}(z) = \max(0, z)$$

Properties:

- Computationally efficient
- Alleviates vanishing gradient problem
- Non-differentiable at $z = 0$
- Can suffer from "dying ReLU" problem

2.6.2 Mathematical Requirements for Activation Functions

For universal approximation, activation functions should be:

1. **Non-linear:** Otherwise, multiple layers collapse to a single linear transformation
2. **Differentiable:** Enables gradient-based optimization (almost everywhere is sufficient)
3. **Monotonic:** Helps with optimization landscape (not strictly required)
4. **Bounded or unbounded:** Different properties affect convergence behavior

2.6.3 Neural Networks as Function Approximators

With sufficient neurons and appropriate activation functions, neural networks can approximate any continuous function to arbitrary precision.

Universal Approximation Theorem

Theorem (Cybenko, 1989; Hornik, 1991): Let ϕ be a continuous sigmoid-type function. Then finite sums of the form:

$$F(x) = \sum_{j=1}^N \alpha_j \phi(y_j^T x + \theta_j)$$

are dense in $C(I_n)$, the space of continuous functions on the unit hypercube $I_n = [0, 1]^n$.

Implications

- **Existence:** There exists a neural network that can approximate any continuous function
- **No constructive proof:** Doesn't tell us how to find the network
- **Width vs. Depth:** Original theorem about width; depth can be more efficient
- **Approximation vs. Learning:** Says nothing about learnability from data

Modern Extensions

- **ReLU networks:** Also have universal approximation properties
- **Deep vs. Wide:** Deep networks can be exponentially more efficient than wide ones
- **Smooth functions:** Require fewer neurons than general continuous functions

2.7 Artificial Neural Networks: From Theory to Implementation

2.7.1 Fundamental Architecture: Primitive Functions and Composition Rules

To understand artificial neural networks, we must first examine their core computational elements. Every computational model requires:

1. **Primitive Functions:** Basic operations that cannot be decomposed further
2. **Composition Rules:** Ways to combine primitive functions to create complex behaviors

Primitive Functions in Neural Networks

In artificial neural networks, **primitive functions are located in the nodes (neurons) of the network**. Each node implements a specific mathematical transformation that processes incoming information and produces an output.

Composition Rules in Neural Networks

The **composition rules are contained implicitly in:**

- **Interconnection pattern of the nodes:** How neurons are connected determines information flow
- **Synchrony or asynchrony of information transmission:** Whether neurons update simultaneously or in sequence
- **Presence or absence of cycles:** Whether information can flow in loops (recurrent networks) or only forward (feedforward networks)

This differs fundamentally from traditional computing models:

Computing Model	Primitive Functions	Composition Rules
von Neumann Processor	Machine instructions (ADD, MOVE, JUMP)	Program sequence + control flow
Artificial Neural Networks	Neuron activation functions	Network topology + connection weights + timing

Table 2.1: Comparison of primitive functions and composition rules across computing models.

2.7.2 Neural Networks as Function Approximators

Networks of Primitive Functions

Artificial neural networks are nothing but networks of primitive functions. Each node transforms its input into a precisely defined output, and the combination of these transformations creates complex computational behaviors.

The Network Function

Consider a neural network that takes inputs (x, y, z) and produces an output through nodes implementing primitive functions f_1, f_2, f_3, f_4 . The network can be thought of as implementing a **network function** ϕ :

$$\phi(x, y, z) = f_4(a_4 \cdot f_3(a_3 \cdot f_2(a_2 \cdot f_1(a_1 \cdot x)))) + \dots$$

Where a_1, a_2, \dots, a_5 are the weights of the network. **Different selections of weights produce different network functions.**

Three Critical Elements

Different models of artificial neural networks differ mainly in three fundamental aspects:

1. Structure of the Nodes

- Choice of activation function (sigmoid, ReLU, tanh, etc.)
- Input integration method (weighted sum, product, etc.)
- Presence of bias terms

2. Topology of the Network

- Feedforward vs. recurrent connections
- Number of layers and neurons per layer
- Connection patterns (fully connected, sparse, convolutional)

3. Learning Algorithm

- Method for finding optimal weights
- Supervised vs. unsupervised vs. reinforcement learning
- Optimization techniques (gradient descent, evolutionary algorithms)

2.7.3 Function Approximation: The Classical Problem

Historical Context

Function approximation is a classical problem in mathematics: **How can we reproduce a given function $F : \mathbb{R} \rightarrow \mathbb{R}$ either exactly or approximately using a given set of primitive functions?**

Traditional approaches include:

- **Polynomial approximation:** Using powers of x (Taylor series)
- **Fourier approximation:** Using trigonometric functions (sine and cosine)
- **Spline approximation:** Using piecewise polynomials

Neural Networks as Universal Approximators

Neural networks provide a revolutionary approach to function approximation:

Key Insight: With sufficient neurons and appropriate activation functions, neural networks can approximate any continuous function to arbitrary precision (Universal Approximation Theorem).

Advantages of Neural Network Approximation

1. **Adaptive:** Networks learn the approximation from data rather than requiring explicit mathematical formulation
2. **Flexible:** Can handle high-dimensional inputs and complex, non-linear relationships
3. **Robust:** Can generalize to unseen data and handle noise
4. **Parallel:** Multiple neurons can process different aspects of the input simultaneously

2.7.4 Learning from Data: The Key Difference

The main difference between Taylor or Fourier series and artificial neural networks is, however, that **the function F to be approximated is given not explicitly but implicitly through a set of input-output examples**. We know F only at some points but we want to generalize as well as possible. This means that we try to adjust the parameters of the network in an optimal manner to reflect the information known and to extrapolate to new input patterns which will be shown to the network afterwards. This is the task of the learning algorithm used to adjust the network's parameters.

Classical Series vs. Neural Networks: A Fundamental Distinction

Classical Mathematical Series (Taylor/Fourier):

- **Explicit Function Definition:** The function $F(x)$ is mathematically defined and known
- **Analytical Coefficients:** Series coefficients can be computed directly using calculus
 - Taylor: $a_n = F^{(n)}(x_0)/n!$ (nth derivative at expansion point)
 - Fourier: a_n, b_n computed via integration over the function's period
- **Perfect Representation:** Given enough terms, the series can represent the function exactly
- **No Learning Required:** Coefficients are determined mathematically, not learned

Artificial Neural Networks:

- **Implicit Function Definition:** The function F is unknown but represented by data points
- **Learned Parameters:** Network weights and biases are learned from examples
- **Approximation from Samples:** Must generalize from finite training data to unknown inputs
- **Adaptive Learning:** Parameters adjust through iterative optimization algorithms

Mathematical Formulation of the Learning Problem

Given a training dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, we seek to find parameters $\boldsymbol{\theta}$ that minimize the empirical risk:

$$\mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$$

where:

- $L(\cdot, \cdot)$: Loss function measuring prediction error
- $f(\mathbf{x}; \boldsymbol{\theta})$: Neural network function with parameters $\boldsymbol{\theta}$
- Goal: Minimize true risk $\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim P}[L(f(\mathbf{x}; \boldsymbol{\theta}), y)]$

Generalization Gap

The fundamental challenge is the generalization gap:

$$\text{Generalization Gap} = \mathcal{R}(\boldsymbol{\theta}) - \mathcal{R}_{\text{emp}}(\boldsymbol{\theta})$$

This gap can be controlled through:

1. **Regularization:** Adding penalty terms to control model complexity
2. **Cross-validation:** Using held-out data to estimate generalization performance
3. **Early stopping:** Halting training before overfitting occurs
4. **Data augmentation:** Artificially increasing training set size

2.8 Computational Complexity in Neural Networks

2.8.1 Forward Pass Complexity

For a neural network with L layers, where layer l has n_l neurons:

- **Matrix multiplication:** $O(n_{l-1} \times n_l)$ for each layer
- **Total forward pass:** $O(\sum_{l=1}^L n_{l-1} \times n_l)$
- **Activation functions:** $O(n_l)$ per layer (typically much smaller than matrix operations)

2.8.2 Backward Pass Complexity (Backpropagation)

- **Gradient computation:** Same order as forward pass $O(\sum_{l=1}^L n_{l-1} \times n_l)$
- **Parameter updates:** $O(\text{total parameters})$
- **Memory complexity:** $O(\text{total activations})$ to store intermediate values

2.8.3 Scalability Considerations

- **Batch processing:** Process multiple examples simultaneously for efficiency
- **Parallelization:** Matrix operations are highly parallelizable on GPUs
- **Memory-computation tradeoff:** Can reduce memory by recomputing activations

Chapter 3

Linear Classification

3.1 Learning Goals

By the end of this chapter, you will:

- **Know what is meant by binary linear classification** and understand its fundamental concepts
- **Understand why an explicit threshold for a classifier is redundant** and how bias terms can be eliminated using dummy features
- **Be able to specify weights and biases by hand** to represent simple logical functions (AND, OR, NOT)
- **Be familiar with input space and weight space**, including:
 - Plotting training cases and classification weights in both spaces
 - Understanding the geometric interpretation of linear classifiers
- **Be aware of the limitations of linear classifiers**, including:
 - Understanding convexity and its role in linear separability
 - Knowing how basis function representations can overcome some limitations

3.2 Fundas: Mathematical Foundations

3.2.1 Vector Representation

- **Input vectors:** Each data point is represented as a D -dimensional vector $\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_D^{(i)}]$
- **Weight vectors:** Classification parameters represented as $\mathbf{w} = [w_1, w_2, \dots, w_D]$
- **Linear combination:** $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, where b is the bias term
- **Decision boundary:** The hyperplane where $\mathbf{w}^T \mathbf{x} + b = 0$

3.2.2 Binary Classification Framework

- **Target values:** $t^{(i)} \in \{0, 1\}$ where 0 = negative class, 1 = positive class
- **Classification rule:** $\hat{y} = 1$ if $\mathbf{w}^T \mathbf{x} + b > \tau$, else $\hat{y} = 0$ (where τ is threshold)
- **Training set:** $\{(\mathbf{x}^{(i)}, t^{(i)})\}_{i=1}^N$ where N is the number of examples

3.3 Introduction to Binary Classification

Binary classification represents one of the most fundamental problems in machine learning, where the goal is to predict a binary-valued target from input features. This forms the foundation for understanding more complex classification scenarios.

3.3.1 Real-World Applications

Medical Diagnosis Systems

- **Problem:** Predict whether a patient has a specific disease
- **Features:** Symptoms, test results, patient history
- **Target:** Disease present (1) or absent (0)
- **Example:** Diagnosing diabetes from glucose levels, BMI, and family history

Email Spam Detection

- **Problem:** Classify emails as spam or legitimate
- **Features:** Word frequencies, sender information, email metadata
- **Target:** Spam (1) or not spam (0)
- **Example:** Using keywords like “free money” as indicators

Fraud Detection

- **Problem:** Identify fraudulent transactions
- **Features:** Transaction amount, time, location, merchant type
- **Target:** Fraudulent (1) or legitimate (0)
- **Example:** Detecting unusual spending patterns

3.4 Binary Linear Classifiers: First Principles

3.4.1 Core Concept

A binary linear classifier makes decisions by computing a linear function of the input features and comparing the result to a threshold. This approach assumes that the two classes can be separated by a linear decision boundary in the feature space.

3.4.2 Mathematical Formulation

Step 1: Linear Combination

$$z = w_1x_1 + w_2x_2 + \cdots + w_Dx_D + b \quad (3.1)$$

Step 2: Threshold Decision

$$\hat{y} = \begin{cases} 1 & \text{if } z > \tau \\ 0 & \text{otherwise} \end{cases} \quad (3.2)$$

Where:

- w_i : Weight for feature i (determines importance and direction)
- x_i : Value of feature i
- b : Bias term (shifts the decision boundary)
- τ : Threshold value

3.4.3 Eliminating Redundancy: The Bias-Threshold Trick

Problem: Having both bias (b) and threshold (τ) is redundant.

Solution 1: Absorb threshold into bias term:

- Set new bias: $b' = b - \tau$
- Set threshold to zero: $\tau = 0$
- Decision rule becomes: $\hat{y} = 1$ if $\mathbf{w}^T \mathbf{x} + b' > 0$

Solution 2: Add dummy feature:

- Extend input: $\mathbf{x} \rightarrow [\mathbf{x}, 1]$
- Extend weights: $\mathbf{w} \rightarrow [\mathbf{w}, b]$
- Decision rule: $\hat{y} = 1$ if $\mathbf{w}_{\text{extended}}^T \mathbf{x}_{\text{extended}} > 0$

3.5 Geometric Interpretation

3.5.1 Input Space Perspective

Decision Boundary: The hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ divides the input space into two regions:

- **Positive region:** $\mathbf{w}^T \mathbf{x} + b > 0$ (predicted class 1)
- **Negative region:** $\mathbf{w}^T \mathbf{x} + b < 0$ (predicted class 0)

Properties:

- The weight vector \mathbf{w} is perpendicular to the decision boundary
- The bias b determines the distance of the boundary from the origin
- Moving along \mathbf{w} increases the classifier's output

3.5.2 Weight Space Perspective

Constraint Regions: Each training example imposes a constraint on the weight space:

- **Positive examples:** Require $\mathbf{w}^T \mathbf{x}^{(i)} + b > 0$
- **Negative examples:** Require $\mathbf{w}^T \mathbf{x}^{(i)} + b < 0$

Feasible Region: The intersection of all constraints defines the region of weight vectors that correctly classify all training examples.

3.6 Hand-Designed Logic Functions

3.6.1 AND Function

Truth Table: $(0, 0) \rightarrow 0, (0, 1) \rightarrow 0, (1, 0) \rightarrow 0, (1, 1) \rightarrow 1$

Implementation:

$$w_1 = 1, \quad w_2 = 1 \quad (3.3)$$

$$b = -1.5 \quad (3.4)$$

$$\text{Decision : } \hat{y} = 1 \text{ if } x_1 + x_2 - 1.5 > 0 \quad (3.5)$$

Verification:

$$(0, 0) : \quad 0 + 0 - 1.5 = -1.5 < 0 \rightarrow \hat{y} = 0 \quad \checkmark \quad (3.6)$$

$$(1, 1) : \quad 1 + 1 - 1.5 = 0.5 > 0 \rightarrow \hat{y} = 1 \quad \checkmark \quad (3.7)$$

3.6.2 OR Function

Implementation:

$$w_1 = 1, \quad w_2 = 1 \quad (3.8)$$

$$b = -0.5 \quad (3.9)$$

$$\text{Decision : } \hat{y} = 1 \text{ if } x_1 + x_2 - 0.5 > 0 \quad (3.10)$$

3.6.3 NOT Function

Implementation:

$$w_1 = -1 \quad (3.11)$$

$$b = 0.5 \quad (3.12)$$

$$\text{Decision : } \hat{y} = 1 \text{ if } -x_1 + 0.5 > 0 \quad (3.13)$$

3.7 Limitations of Linear Classifiers

3.7.1 Linear Separability

Definition 1 (Linear Separability). *A dataset is linearly separable if there exists a linear decision boundary that perfectly separates the two classes.*

3.7.2 Convexity and Separability

Definition 2 (Convex Set). *A set is convex if the line segment between any two points in the set lies entirely within the set.*

Theorem 1 (Linear Separability Condition). *If the positive examples form a convex set and the negative examples form a convex set, and these sets do not overlap, then the data is linearly separable.*

3.7.3 The XOR Problem

XOR Truth Table: $(0, 0) \rightarrow 0$, $(0, 1) \rightarrow 1$, $(1, 0) \rightarrow 1$, $(1, 1) \rightarrow 0$

Why XOR is not linearly separable:

- Positive examples: $(0, 1)$ and $(1, 0)$
- Negative examples: $(0, 0)$ and $(1, 1)$
- No single line can separate these points correctly

3.7.4 Overcoming Limitations: Basis Functions

Strategy: Transform inputs using basis functions $\phi(x)$ to create a new feature space where linear separation becomes possible.

Example for XOR:

- Original features: x_1, x_2
- Basis functions: $\phi_1 = x_1$, $\phi_2 = x_2$, $\phi_3 = x_1x_2$
- In the new space $[x_1, x_2, x_1x_2]$, XOR becomes linearly separable

The transformed XOR problem becomes:

$$(0, 0, 0) \rightarrow 0 \quad (3.14)$$

$$(0, 1, 0) \rightarrow 1 \quad (3.15)$$

$$(1, 0, 0) \rightarrow 1 \quad (3.16)$$

$$(1, 1, 1) \rightarrow 0 \quad (3.17)$$

With weights $w_1 = 1$, $w_2 = 1$, $w_3 = -2$, and bias $b = -0.5$:

$$\hat{y} = 1 \text{ if } x_1 + x_2 - 2x_1x_2 - 0.5 > 0 \quad (3.18)$$

3.8 Practical Considerations

3.8.1 Feature Engineering

- **Normalization:** Scale features to similar ranges
- **Interaction terms:** Add products of features (x_1x_2)
- **Polynomial features:** Add powers of features (x_1^2, x_2^2)

3.8.2 Model Selection

- **Bias-variance tradeoff:** Simple models (few features) vs. complex models (many features)
- **Interpretability:** Linear classifiers provide clear feature importance through weights
- **Computational efficiency:** Linear models are fast to train and evaluate

3.8.3 Performance Evaluation

- **Training accuracy:** Percentage of training examples correctly classified
- **Generalization:** Performance on unseen test data
- **Decision boundary visualization:** Plot boundaries in 2D for intuition

3.9 Chapter Summary

Linear classification provides a fundamental framework for understanding binary classification problems. Key takeaways:

1. **Geometric intuition:** Linear classifiers create hyperplane decision boundaries in input space
2. **Mathematical elegance:** Simple linear combinations with threshold decisions
3. **Practical utility:** Many real-world problems are approximately linearly separable
4. **Clear limitations:** Not all problems (like XOR) can be solved with linear classifiers
5. **Extension possibilities:** Basis functions can transform non-linearly separable problems

The concepts developed here—particularly the geometric interpretation in both input and weight spaces—form the foundation for understanding more sophisticated classification algorithms like the perceptron, support vector machines, and neural networks.

The next chapter will introduce the **perceptron algorithm**, which provides a systematic method for learning the weights of a linear classifier from training data, moving beyond hand-designed solutions to automated learning procedures.

Chapter 4

The Perceptron

4.1 Historical Development of Neural Networks

4.1.1 Timeline of Neural Network Evolution

The development of neural networks has proceeded through several distinct phases, each marked by significant theoretical breakthroughs and practical applications.

Period	Year	Key Development	Contributors	Description
Early Foundations	1943	McCulloch-Pitts Neuron	Warren McCulloch, Walter Pitts	First mathematical model of artificial neuron using threshold logic
	1949	Hebbian Learning Rule	Donald Hebb	"Cells that fire together, wire together" - synaptic plasticity principle
First Generation	1957	Perceptron	Frank Rosenblatt	First trainable neural network with learning algorithm
	1960	ADALINE/MADALINE	Bernard Widrow, Marcian Hoff	Adaptive linear neurons with delta rule learning
Winter Period	1969	Perceptron Limitations	Marvin Minsky, Seymour Papert	Proved perceptrons cannot solve XOR problem
Revival Era	1982	Hopfield Networks	John Hopfield	Recurrent networks for associative memory
	1986	Backpropagation	Rumelhart, Hinton, Williams	Efficient algorithm for training multi-layer networks
Modern Era	2012	AlexNet	Alex Krizhevsky, Geoffrey Hinton	Deep CNN wins ImageNet competition
	2017	Transformer Architecture	Vaswani et al. (Google)	Attention-based model for sequences

Table 4.1: Key milestones in neural network development.

4.1.2 Threshold Logic: The Foundation

The simplest kind of computing units used to build artificial neural networks are based on threshold logic. These computing elements are a generalization of the common logic gates used in conventional computing and, since they operate by comparing their total input with a threshold, this field of research is known as **threshold logic**.

4.2 McCulloch-Pitts Neuron: The First Artificial Neuron

The McCulloch-Pitts neuron, introduced in 1943, was the first mathematical model of an artificial neuron. It established the theoretical foundation for neural computation using threshold logic.

4.2.1 Mathematical Model

The McCulloch-Pitts neuron computes its output according to:

$$y = \begin{cases} 1 & \text{if } \sum_{i=1}^n w_i x_i \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

Where:

- x_i are the input values (binary: 0 or 1)
- w_i are the corresponding weights
- θ is the threshold value
- y is the binary output (0 or 1)

4.2.2 Logic Gate Implementation

The McCulloch-Pitts model can implement basic logic functions:

AND Gate

For an AND gate with two inputs:

- Weights: $w_1 = w_2 = 1$
- Threshold: $\theta = 2$
- Result: Output is 1 only when both inputs are 1

OR Gate

For an OR gate with two inputs:

- Weights: $w_1 = w_2 = 1$
- Threshold: $\theta = 1$
- Result: Output is 1 when at least one input is 1

4.2.3 Limitations of McCulloch-Pitts Neurons

- **Fixed weights:** No learning mechanism
- **Binary inputs only:** Cannot handle continuous values
- **Synchronous operation:** All neurons fire simultaneously
- **No adaptation:** Cannot modify behavior based on experience

These limitations led to the development of the Perceptron, which introduced learning capabilities.

4.3 The Perceptron: A Detailed Introduction

The Perceptron is a simple binary classifier that serves as the foundational building block for more complex neural networks.

4.3.1 Definition: The Anatomy of a Perceptron

For an input vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$, the Perceptron computes a single output y . This is done in two steps:

1. **Compute a Weighted Sum:** The model calculates a weighted sum of the inputs, adding a bias. This is the net input z .

$$z = (w_1x_1 + w_2x_2 + \dots + w_nx_n) + b = \mathbf{w} \cdot \mathbf{x} + b$$

2. **Apply an Activation Function:** The output z is passed through a Heaviside step function.

$$y = \phi(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

4.4 The Perceptron Learning Rule

The Perceptron learns by adjusting its weights \mathbf{w} and bias b based on the errors it makes. This learning rule has strong theoretical foundations.

4.4.1 Mathematical Derivation

For a given training example (\mathbf{x}, t) , where t is the true target label, the error ϵ is calculated as $\epsilon = t - y$. The weights and bias are then updated:

$$w_i(\text{new}) = w_i(\text{old}) + \eta \cdot \epsilon \cdot x_i$$

$$b(\text{new}) = b(\text{old}) + \eta \cdot \epsilon$$

Where η is the learning rate.

4.4.2 Convergence Theorem (Rosenblatt, 1962)

Theorem: If the training data is linearly separable, the perceptron learning algorithm will converge to a solution in a finite number of steps.

Proof Sketch

Let \mathbf{w}^* be a weight vector that correctly classifies all training examples with margin $\gamma > 0$:

$$y_i(\mathbf{w}^{*T} \mathbf{x}_i) \geq \gamma \quad \forall i$$

The proof shows that:

1. The dot product $\mathbf{w}_t \cdot \mathbf{w}^*$ grows linearly with updates
2. The norm $\|\mathbf{w}_t\|^2$ grows at most linearly with updates
3. This leads to a contradiction if the algorithm doesn't converge

Convergence Bound

The perceptron will make at most $\left(\frac{R}{\gamma}\right)^2$ mistakes, where:

- R is the maximum norm of any training example: $R = \max_i \|\mathbf{x}_i\|$
- γ is the margin of the optimal separating hyperplane

4.4.3 Learning Rate Analysis

Effect of Learning Rate η

- **Large η :** Faster convergence but may overshoot optimal solution
- **Small η :** More stable learning but slower convergence
- **Theoretical Result:** For linearly separable data, any $\eta > 0$ guarantees convergence

Adaptive Learning Rates

Common strategies include:

- **Time decay:** $\eta_t = \frac{\eta_0}{1+\alpha t}$
- **Step decay:** Reduce η by factor every few epochs
- **Performance-based:** Reduce η when performance plateaus

4.5 Vectorisation of Perceptron Learning Rule: A Matrix-Based Example

To see how vectorization works in practice, let's walk through the process using a full matrix-based approach for the AND function. This method calculates the updates for all samples in the dataset (a "batch") and then applies a single, consolidated update at the end of the epoch.

4.5.1 Mathematical Foundation of Vectorization

Vectorization allows us to process multiple training examples simultaneously using matrix operations instead of iterating through samples one by one. This approach is:

- **Computationally efficient:** Modern hardware (GPUs) excels at matrix operations
- **Mathematically elegant:** Compact representation of batch operations
- **Numerically stable:** Reduces accumulated floating-point errors

4.5.2 Define the Matrices

For the AND function, we use an augmented Input Matrix X , a Weight Vector W , and a Target Vector T .

Input Matrix X (Augmented Design Matrix)

The input matrix includes a bias column (first column of 1s) followed by the feature columns:

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}_{4 \times 3}$$

where:

- Rows represent training examples (4 samples)
- First column represents bias input (always 1)
- Remaining columns represent feature inputs x_1, x_2

Target Vector T

$$T = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}_{4 \times 1}$$

Weight Vector W

Let's initialize the weight vector W to zeros and use a learning rate $\eta = 0.1$.

$$W_0 = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}_{3 \times 1}$$

where w_0 is the bias weight, w_1 and w_2 are feature weights.

4.5.3 Epoch 1: Mathematical Flow

Step 1: Compute Net Input Z

The net input is computed as $Z = X \cdot W$, where we multiply the 4×3 input matrix by the 3×1 weight vector to get a 4×1 output vector.

$$Z = X \cdot W_0 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}_{4 \times 3} \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}_{3 \times 1} = \begin{pmatrix} 1 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 \\ 1 \cdot 0 + 0 \cdot 0 + 1 \cdot 0 \\ 1 \cdot 0 + 1 \cdot 0 + 0 \cdot 0 \\ 1 \cdot 0 + 1 \cdot 0 + 1 \cdot 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}_{4 \times 1}$$

Matrix Dimension Check: $(4 \times 3) \times (3 \times 1) = (4 \times 1)$

Step 2: Apply Activation Function to get Output Y

Apply the Heaviside step function element-wise: $\phi(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$

$$Y = \phi(Z) = \phi \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}_{4 \times 1} = \begin{pmatrix} \phi(0) \\ \phi(0) \\ \phi(0) \\ \phi(0) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}_{4 \times 1}$$

Since all net inputs are 0, and $\phi(0) = 1$ by our step function definition, all outputs are 1.

Step 3: Calculate the Error Vector E

$$E = T - Y = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \\ -1 \\ 0 \end{pmatrix}$$

Step 4: Calculate the Total Weight Update ΔW

The weight update is computed as $\Delta W = \eta \cdot (X^T \cdot E)$, where we multiply the transpose of the input matrix by the error vector.

First, let's compute X^T :

$$X^T = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}^T = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix}_{3 \times 4}$$

Now compute the weight update:

$$\begin{aligned}\Delta W &= \eta \cdot (X^T \cdot E) = 0.1 \cdot \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix}_{3 \times 4} \begin{pmatrix} -1 \\ -1 \\ -1 \\ 0 \end{pmatrix}_{4 \times 1} \\ &= 0.1 \cdot \begin{pmatrix} 1(-1) + 1(-1) + 1(-1) + 1(0) \\ 0(-1) + 0(-1) + 1(-1) + 1(0) \\ 0(-1) + 1(-1) + 0(-1) + 1(0) \end{pmatrix} = 0.1 \cdot \begin{pmatrix} -3 \\ -1 \\ -1 \end{pmatrix} = \begin{pmatrix} -0.3 \\ -0.1 \\ -0.1 \end{pmatrix}_{3 \times 1}\end{aligned}$$

Matrix Dimension Check: $(3 \times 4) \times (4 \times 1) = (3 \times 1)$

Step 5: Update the Weight Vector \mathbf{W}

$$W_1 = W_0 + \Delta W = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} -0.3 \\ -0.1 \\ -0.1 \end{pmatrix} = \begin{pmatrix} -0.3 \\ -0.1 \\ -0.1 \end{pmatrix}$$

After the first epoch, our new weight vector is $W_1 = (-0.3, -0.1, -0.1)^T$.

4.5.4 Mathematical Insight: Gradient Descent Connection

The vectorized perceptron learning rule is actually a special case of gradient descent. The weight update $\Delta W = \eta \cdot X^T \cdot E$ can be derived from minimizing the perceptron loss function:

Loss Function

For a single misclassified example, the perceptron loss is:

$$L = \max(0, -y_{\text{true}} \cdot z)$$

where $z = \mathbf{w}^T \mathbf{x}$ is the net input.

Gradient of the Loss

The gradient with respect to weights is:

$$\frac{\partial L}{\partial \mathbf{w}} = -y_{\text{true}} \cdot \mathbf{x}$$

for misclassified examples, and 0 for correctly classified ones.

Batch Update Rule

For a batch of examples, the total gradient is:

$$\frac{\partial L_{\text{total}}}{\partial \mathbf{w}} = \sum_i \frac{\partial L_i}{\partial \mathbf{w}} = X^T \cdot (Y_{\text{pred}} - Y_{\text{true}})$$

This shows that our vectorized update rule $\Delta W = \eta \cdot X^T \cdot E$ is exactly gradient descent with the perceptron loss function.

4.5.5 Computational Complexity Analysis

Vectorized vs. Sequential Processing

- **Sequential:** $O(m \cdot n)$ time for m examples and n features, but cannot leverage parallel hardware
- **Vectorized:** Same $O(m \cdot n)$ time complexity, but:
 - Can utilize SIMD (Single Instruction, Multiple Data) operations
 - Reduces Python interpreter overhead
 - Enables GPU acceleration for large matrices
 - Better cache locality and memory access patterns

4.6 Geometric Interpretation of the Perceptron

Understanding a Perceptron involves grasping the geometry of how it makes decisions. We can visualize this geometry in two primary ways: the **Input Space** and the **Weight Space**.

4.6.1 Mathematical Foundation of Decision Boundaries

The perceptron's decision boundary is defined by the hyperplane equation:

$$\mathbf{w}^T \mathbf{x} + b = 0$$

Hyperplane Properties

- **Normal Vector:** The weight vector \mathbf{w} is perpendicular to the decision boundary
- **Distance from Origin:** $\frac{|b|}{\|\mathbf{w}\|_2}$ gives the perpendicular distance from the hyperplane to the origin
- **Classification Rule:**
 - Points where $\mathbf{w}^T \mathbf{x} + b > 0$ are classified as positive (class 1)
 - Points where $\mathbf{w}^T \mathbf{x} + b < 0$ are classified as negative (class 0)
 - Points on the boundary satisfy $\mathbf{w}^T \mathbf{x} + b = 0$

Margin and Support

The **functional margin** of a point \mathbf{x}_i with true label $y_i \in \{-1, +1\}$ is:

$$\gamma_i = y_i(\mathbf{w}^T \mathbf{x}_i + b)$$

The **geometric margin** is the normalized version:

$$\hat{\gamma}_i = \frac{y_i(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|_2}$$

This represents the perpendicular distance from the point to the decision boundary.

4.6.2 Example: The NOT Function

The Input Space

The Input Space is a geometric representation of the data points. The goal of the Perceptron is to find a **decision boundary**—a line in 2D, or a hyperplane in higher dimensions—that perfectly separates the positive and negative data points.

The Weight Space

While the input space plots the data, the **Weight Space** plots the possible solutions. The axes of this space are the weights themselves. Every point in this space represents a different Perceptron model.

4.6.3 Example: The AND Function

The Input Space

The Input Space shows our data points and the decision boundary that separates them. For the AND function, we have three "negative" points (target=0) and one "positive" point (target=1).

The Weight Space

The Weight Space represents the set of all possible solutions. Each of our four data points imposes a constraint on the possible values of the weights.

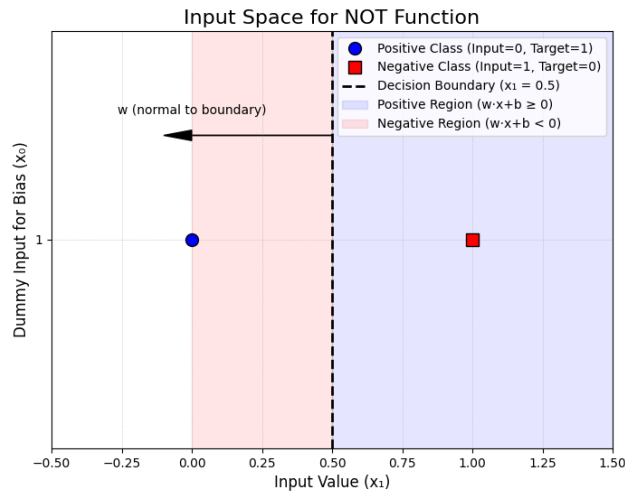


Figure 4.1: Input Space for the NOT Function.

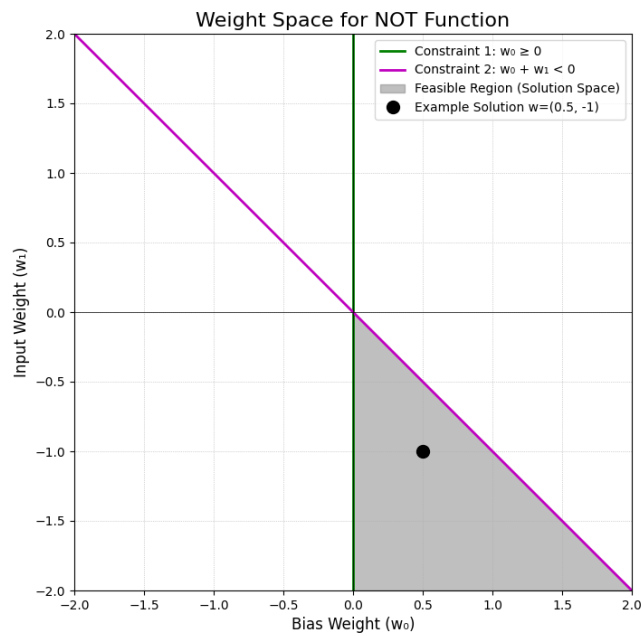


Figure 4.2: Weight Space for the NOT Function.

4.7 Limitations of the Perceptron

4.7.1 Linear Separability Constraint

The fundamental limitation of the perceptron is that it can only learn linearly separable functions.

Definition: Linear Separability

A dataset is **linearly separable** if there exists a hyperplane that perfectly separates the positive and negative examples:

$$\exists \mathbf{w}, b : y_i(\mathbf{w}^T \mathbf{x}_i + b) > 0 \quad \forall i$$

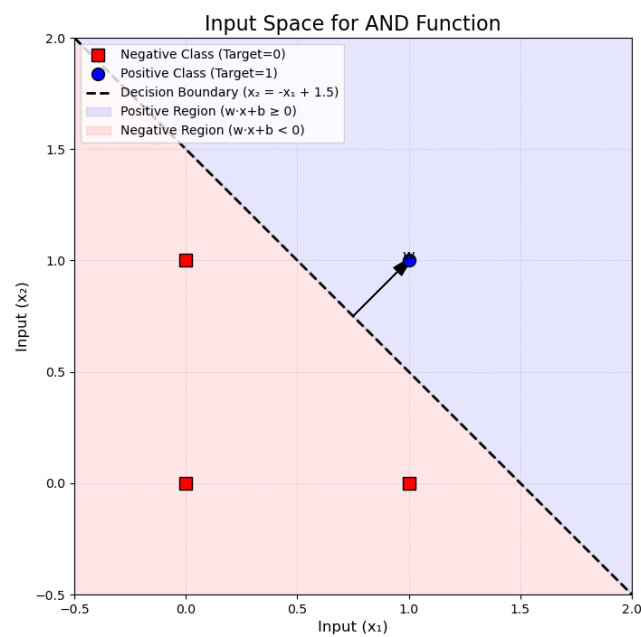
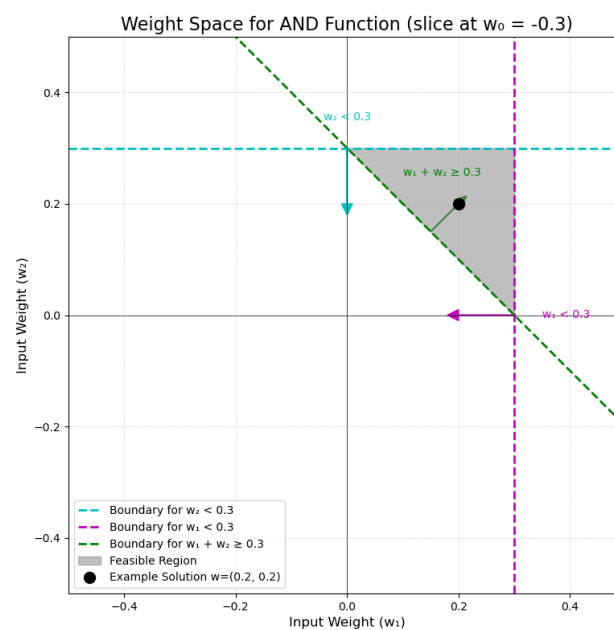


Figure 4.3: Input Space for the AND Function.

Figure 4.4: A 2D slice of the Weight Space for the AND Function, with $w_0 = -0.3$.

The XOR Problem (Minsky & Papert, 1969)

Consider the XOR function:

x_1	x_2	XOR
0	0	0
0	1	1
1	0	1
1	1	0

Mathematical Proof of Non-Linear Separability: Assume there exists a linear classifier $\mathbf{w}^T \mathbf{x} + b = 0$ that separates XOR data. For the four points, we need:

$$w_1 \cdot 0 + w_2 \cdot 0 + b < 0 \quad (\text{point } (0,0)) \quad (4.1)$$

$$w_1 \cdot 0 + w_2 \cdot 1 + b > 0 \quad (\text{point } (0,1)) \quad (4.2)$$

$$w_1 \cdot 1 + w_2 \cdot 0 + b > 0 \quad (\text{point } (1,0)) \quad (4.3)$$

$$w_1 \cdot 1 + w_2 \cdot 1 + b < 0 \quad (\text{point } (1,1)) \quad (4.4)$$

From equations (1) and (2): $b < 0$ and $w_2 + b > 0 \Rightarrow w_2 > -b > 0$ From equations (1) and (3): $b < 0$ and $w_1 + b > 0 \Rightarrow w_1 > -b > 0$ From equation (4): $w_1 + w_2 + b < 0$

But this contradicts $w_1 > -b$ and $w_2 > -b$, since:

$$w_1 + w_2 + b > -b + (-b) + b = -b > 0$$

Therefore, no linear separator exists for XOR.

4.7.2 Solutions to Linear Separability Limitation

Multi-Layer Perceptrons (MLPs)

- Add hidden layers with non-linear activation functions
- Can approximate any continuous function (Universal Approximation Theorem)
- Require more sophisticated training algorithms (backpropagation)

Feature Engineering

- Transform input space to make data linearly separable
- For XOR: Add feature $x_3 = x_1 \oplus x_2$ (though this requires knowing the solution)
- Kernel methods: Implicitly map to higher-dimensional spaces

Ensemble Methods

- Combine multiple linear classifiers
- Voting or weighted combination schemes
- Can learn non-linear decision boundaries

4.8 Historical Impact and Legacy

4.8.1 The Perceptron Controversy

The 1969 book "Perceptrons" by Minsky and Papert highlighted the limitations of single-layer perceptrons, leading to:

- **AI Winter:** Reduced funding and interest in neural networks
- **Focus shift:** Emphasis moved to symbolic AI and expert systems
- **Delayed progress:** Multi-layer networks existed but lacked efficient training methods

4.8.2 Modern Relevance

Despite limitations, perceptrons remain important because:

- **Building blocks:** Neurons in modern deep networks are perceptron variants
- **Theoretical foundation:** Understanding linear classifiers is crucial
- **Computational efficiency:** Still useful for linearly separable problems
- **Online learning:** Perceptron learning rule works in streaming settings

4.9 Exercises

1. Determine the weights of a network with 4 input and 2 output units using Perceptron learning law

$$\text{Input: } \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \text{ Output: } (1 \quad 1 \quad 1 \quad 0 \quad 0)$$