



Classificazione Decimale Dewey

— 006.30151 (23.) INTELLIGENZA ARTIFICIALE. Matematica

Thema

— Soggetto: UYQ. Intelligenza artificiale

— Qualificatori: 4CT. Per l'istruzione superiore/terziaria/universitaria

GERARDO IOVANE

MATHEMATICAL METHODS FOR ARTIFICIAL INTELLIGENCE





©

ISBN
979-12-218-2687-6

PRIMA EDIZIONE
ROMA 30 APRILE 2026

Table of Contents

PART I – MATHEMATICAL FOUNDATIONS FOR ARTIFICIAL INTELLIGENCE.....	7
1. THE MATHEMATICAL LANGUAGE OF ARTIFICIAL INTELLIGENCE	7
2. LINEAR ALGEBRA AND NUMERICAL METHODS.....	21
3. PROBABILITY, STATISTICS, AND STATISTICAL INFERENCE	35
4. OPTIMIZATION THEORY FOR LEARNING.....	51
PART II – DATA AND THE MATHEMATICAL PIPELINE.....	67
5. DATA THEORY AND DATASET CONSTRUCTION	67
6. REPRESENTATIONS AND FEATURE ENGINEERING.....	81
7. MODEL EVALUATION AND PERFORMANCE ASSESSMENT.....	95
PART III – MACHINE LEARNING.....	109
8. SUPERVISED LEARNING LINEAR AND KERNEL METHODS	109
9. UNSUPERVISED LEARNING CLUSTERING AND LATENT MODELS	121
10. ENSEMBLE METHODS: TREES, BAGGING, AND BOOSTING.....	133
11. SEQUENTIAL AND PROBABILISTIC MODELS.....	143
PART IV – DEEP LEARNING.....	155
12. NEURAL NETWORKS AS UNIVERSAL APPROXIMATORS.....	155
13. BACKPROPAGATION AND COMPUTATIONAL GRAPHS.....	165
14. CONVOLUTIONAL ARCHITECTURES FOR IMAGES AND SIGNALS.....	175
15. SEQUENCES, ATTENTION, AND TRANSFORMERS.....	185
16. GENERATIVE MODELS.....	195
17. LARGE-SCALE TRAINING AND SYSTEMS OPTIMIZATION	205
18. INTERPRETABILITY, RELIABILITY, AND SAFETY	215
PART V – STRONG AI: REASONING, PLANNING, AND AGENCY.....	225
19. FOUNDATIONS OF AGENCY AND DECISION-MAKING.....	225
20. SYMBOLIC REASONING AND LOGICAL INFERENCE	235
21. SELF-SUPERVISED LEARNING AND WORLD MODELS	245
PART VI – TOWARD AGI: GENERALITY, SCALING, ALIGNMENT, AND THEORY	255
22. GENERALIZATION, TRANSFER, AND META-LEARNING.....	255
23. COMPLEXITY THEORY AND FUNDAMENTAL LIMITS	263

24.	ALIGNMENT, SAFETY, AND GOVERNANCE.....	271
25.	AGI ARCHITECTURES AND MULTI-AGENT SYSTEMS	279
PART VII – DEPLOYMENT, PRODUCTION, AND LIFECYCLE MANAGEMENT		287
26.	FROM TRAINING TO PRODUCTION.....	287
27.	CASE STUDIES AND COMPLETE PIPELINES	295
PART VIII – ENCYCLOPEDIA APPENDICES.....		303
28.	REFERENCE TABLES AND QUICK GUIDE.....	303
SUGGESTED READING.....		313

Part I – Mathematical foundations for artificial intelligence

1. The mathematical language of artificial intelligence

The development of artificial intelligence systems rests fundamentally upon a precise mathematical framework. While empirical success has driven much of the recent progress in machine learning and deep learning, the underlying mechanisms can be understood only through rigorous formalization. This chapter establishes the essential mathematical vocabulary and structures that will permeate the entire treatment of AI theory presented in this volume.

The mathematical language of AI serves multiple purposes beyond mere notational convenience. It provides a unified framework for reasoning about diverse phenomena including optimization landscapes, statistical generalization, representational capacity, and computational complexity. More importantly, formalization exposes the implicit assumptions embedded in algorithms and models, making explicit the conditions under which theoretical guarantees hold and empirical performance can be expected.

We begin by establishing notational conventions and fundamental mathematical structures. The treatment then progresses through the essential elements of multivariate calculus, culminating in a formal characterization of learning problems. Throughout, we emphasize not merely the definitions themselves but their role in enabling precise statements about learning systems.

1.1. Notation and mathematical structures

The foundation of any rigorous treatment requires careful specification of the mathematical objects under consideration and the notation used to manipulate them. We adopt conventions that balance precision with readability while remaining consistent with the contemporary literature in machine learning theory.

Let \mathbb{N} denote the natural numbers, \mathbb{Z} the integers, \mathbb{R} the real numbers, and \mathbb{C} the complex numbers. We write \mathbb{R}^+ for the non-negative reals. When discussing vector spaces, we typically work over \mathbb{R} unless otherwise specified. For a vector space V over \mathbb{R} , we denote its dimension by $\dim(V)$, with the understanding that V may be infinite-dimensional in functional analysis contexts.

Vectors are represented by lowercase bold letters in prose and standard lowercase in mathematical expressions. Thus a vector $\mathbf{x} \in \mathbb{R}^d$ has components x_1, \dots, x_d . Matrices are denoted by uppercase letters, so $A \in \mathbb{R}^{m \times n}$ has entries A_{ij} for $i = 1, \dots, m$ and $j = 1, \dots, n$. The transpose of a matrix A is written A^T , while A^{-1} denotes the inverse when it exists. The Moore-Penrose pseudoinverse is denoted A^\dagger .

Tensors, which arise naturally in deep learning contexts, are written as $T \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_k}$. Individual entries are accessed through multi-index notation $T_{i_1 i_2 \dots i_k}$. While general tensor operations can be formalized through the language of multilinear algebra, in practice we often work with specific contractions and reshaping operations that have direct computational implementations.

Norms provide essential tools for measuring magnitudes and distances. The Euclidean norm on \mathbb{R}^d is defined as $\|\mathbf{x}\|_2 = (\sum_{i=1}^d x_i^2)^{1/2}$, while the L^1 norm is $\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i|$ and the L^∞ norm is $\|\mathbf{x}\|_\infty = \max_i |x_i|$. For matrices, the induced operator norm is $\|A\| = \sup\{\|A\mathbf{x}\|_2 : \|\mathbf{x}\|_2 \leq 1\}$. These norms generalize to the L^p spaces through the definition $\|f\|_p = (\int |f(x)|^p dx)^{1/p}$ for measurable functions.

The inner product on \mathbb{R}^d is written $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^d x_i y_i = \mathbf{x}^T \mathbf{y}$. This induces the Euclidean norm through $\|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ and satisfies the Cauchy-Schwarz inequality $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$, with equality if and only if \mathbf{x} and \mathbf{y} are linearly dependent. The Cauchy-Schwarz inequality has profound implications throughout learning theory,

appearing in analyses of kernel methods, attention mechanisms, and generalization bounds.

Beyond finite-dimensional Euclidean spaces, we encounter various function spaces. The space $C^k(\Omega)$ consists of k -times continuously differentiable functions on a domain $\Omega \subset \mathbb{R}^d$. When Ω is compact, $C^k(\Omega)$ becomes a Banach space under appropriate norms. The Lebesgue spaces $L^p(\Omega)$ contain measurable functions with finite p -norm, forming the natural setting for statistical learning theory where we integrate with respect to probability measures.

Asymptotic notation provides a language for discussing algorithm complexity and approximation quality. We write $f(n) = O(g(n))$ if there exist constants C and n_0 such that $|f(n)| \leq C|g(n)|$ for all $n \geq n_0$. The notation $f(n) = \Omega(g(n))$ indicates that $g(n) = O(f(n))$, while $f(n) = \Theta(g(n))$ means $f(n) = O(g(n))$ and $f(n) = \Omega(g(n))$ simultaneously. These asymptotic relationships are essential when analyzing sample complexity, computational cost, and approximation error.

Probability theory introduces additional notation. A probability space is denoted (Ω, \mathcal{F}, P) where Ω is the sample space, \mathcal{F} a σ -algebra of events, and P a probability measure. Random variables $X: \Omega \rightarrow \mathbb{R}$ have probability distributions characterized by cumulative distribution functions or probability density functions $p(x)$. The expectation of X under P is written $\mathbb{E}[X]$ or $\mathbb{E}_P[X]$ when the measure must be specified. Variance is $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$. Independence of random variables X and Y is denoted $X \perp Y$, while conditional independence given Z is written $X \perp Y \mid Z$.

The notation $:=$ indicates definition rather than equality. Thus $\mu := \mathbb{E}[X]$ defines μ as the expected value of X . Logical quantifiers \forall (for all) and \exists (there exists) allow concise statement of propositions. The symbol \Rightarrow denotes implication while \Leftrightarrow indicates logical equivalence. Set inclusion is written \subset for strict inclusion and \subseteq for inclusion allowing equality.

1.2. Functions, composition, and transformations

Functions form the primary objects of study in machine learning, serving simultaneously as hypotheses to be learned, transformations to be optimized, and mappings describing data generation processes. A function $f: X \rightarrow Y$ maps elements from domain X to codomain Y . When $X \subset \mathbb{R}^d$ and $Y \subset \mathbb{R}^m$, we write $f(x) \in \mathbb{R}^m$ for $x \in \mathbb{R}^d$, with component functions $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ for $i = 1, \dots, m$.

Function composition lies at the heart of deep learning architectures. Given $f: X \rightarrow Y$ and $g: Y \rightarrow Z$, their composition is $(g \circ f)(x) = g(f(x))$. Deep neural networks implement compositions of the form $f = f_L \circ f_{L-1} \circ \dots \circ f_1$ where each layer f_ℓ represents an affine transformation followed by a nonlinear activation. The expressiveness of such compositions derives from the universal approximation properties we shall examine in Chapter 12.

Certain function properties recur throughout learning theory and optimization. A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$, we have $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$. Convexity ensures that local minima are global minima, a property exploited in classical machine learning but largely absent from deep learning. Nonetheless, understanding convex optimization provides essential intuition for more general non-convex problems.

Monotonicity constrains function behavior along orderings. A function $f: \mathbb{R} \rightarrow \mathbb{R}$ is monotone increasing if $x \leq y$ implies $f(x) \leq f(y)$. Many activation functions used in neural networks, such as the rectified linear unit (ReLU) $\sigma(z) = \max(0, z)$, are monotone increasing. Monotonicity interacts with optimization and expressiveness in subtle ways, as we shall explore when discussing initialization and training dynamics. Lipschitz continuity provides a quantitative form of continuity controlling the rate of change. A function $f: X \rightarrow Y$ between metric spaces is L -Lipschitz continuous if for all $x, x' \in X$, we have $d_Y(f(x), f(x')) \leq L d_X(x, x')$. The smallest such L is the Lipschitz constant of f . When $X = Y = \mathbb{R}^d$ with Euclidean metric, this becomes $\|f(x) - f(y)\|_2 \leq L \|x - y\|_2$.

Lipschitz continuity has profound implications for learning theory. It ensures that small perturbations to inputs produce bounded changes in outputs, a form of stability

essential for generalization. Neural networks with Lipschitz-constrained weights exhibit improved robustness to adversarial perturbations. Moreover, Lipschitz constants appear in generalization bounds through Rademacher complexity and related measures of function class capacity.

A function $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ is differentiable at \mathbf{x} if there exists a linear map $Df(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}^m$ such that

$$\lim_{\mathbf{h} \rightarrow 0} \frac{\|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - Df(\mathbf{x})\mathbf{h}\|}{\|\mathbf{h}\|} = 0.$$

The linear map $Df(\mathbf{x})$ is the derivative or differential of f at \mathbf{x} . When $m = 1$, the derivative is identified with the gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^d$ through $Df(\mathbf{x})\mathbf{h} = \langle \nabla f(\mathbf{x}), \mathbf{h} \rangle$. For vector-valued functions, the derivative is represented by the Jacobian matrix $J_f(\mathbf{x}) \in \mathbb{R}^{m \times d}$ with entries $(J_f)_{ij} = \partial f_i / \partial x_j$.

Differentiability extends naturally to normed vector spaces. A map $f: V \rightarrow W$ between normed spaces is Fréchet differentiable at $x \in V$ if there exists a bounded linear operator $Df(x): V \rightarrow W$ such that

$$\lim_{\|\mathbf{h}\|_V \rightarrow 0} \frac{\|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - Df(\mathbf{x})\mathbf{h}\|_W}{\|\mathbf{h}\|_V} = 0.$$

This framework encompasses both finite-dimensional calculus and infinite-dimensional functional analysis, providing a unified language for discussing neural networks as functions on finite-dimensional parameter spaces and as elements of infinite-dimensional function classes.

Higher-order derivatives capture curvature information essential for optimization. The Hessian of a twice-differentiable function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is the matrix of second partial derivatives $H_f(\mathbf{x}) \in \mathbb{R}^{d \times d}$ with entries $H_{ij} = \partial^2 f / (\partial x_i \partial x_j)$. When f is twice continuously differentiable, Schwarz's theorem ensures that the Hessian is symmetric. The Hessian's eigenvalues determine local curvature properties, with positive definiteness characterizing strict local minima and negative definiteness characterizing strict local maxima.

1.3. Multivariate calculus and computational graphs

The calculus of vector-valued functions provides the foundation for understanding gradient-based optimization, the workhorse of modern machine learning. We develop the essential results carefully, as their computational implementation through automatic differentiation powers all contemporary deep learning systems.

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable scalar function. The gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^d$ is the vector of partial derivatives:

$$\nabla f(\mathbf{x}) = (\partial f / \partial x_1, \dots, \partial f / \partial x_d)^T.$$

The gradient points in the direction of steepest ascent at \mathbf{x} , and its negative indicates the direction of steepest descent. This geometric interpretation underlies gradient descent optimization: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$, where $\eta_t > 0$ is the learning rate or step size.

For a vector-valued function $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ with components f_1, \dots, f_m , the Jacobian matrix assembles the gradients of the components:

$$J_f(\mathbf{x}) = [\nabla f_1(\mathbf{x})^T, \dots, \nabla f_m(\mathbf{x})^T]^T = [\partial f_i / \partial x_j]_{i,j}$$

The Jacobian captures how perturbations to the input propagate through the function. Indeed, for small $\mathbf{h} \in \mathbb{R}^d$, we have the first-order approximation $f(\mathbf{x} + \mathbf{h}) \approx f(\mathbf{x}) + J_f(\mathbf{x})\mathbf{h}$. This linear approximation forms the basis for understanding how errors propagate through computational graphs.

The chain rule extends differentiation to composed functions, enabling backpropagation through deep networks. Consider functions $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^n$. Their composition $h = g \circ f$ satisfies

$$J_h(\mathbf{x}) = J_g(f(\mathbf{x}))J_f(\mathbf{x}).$$

In the scalar case where $g: \mathbb{R}^m \rightarrow \mathbb{R}$, this becomes

$$\nabla(g \circ f)(\mathbf{x}) = J_f(\mathbf{x})^T \nabla g(f(\mathbf{x})).$$

This formulation reveals the fundamental structure of backpropagation: gradients propagate backward through transposed Jacobians. Given a loss $L = g(f(\mathbf{x}))$ and the gradient $\nabla_f L$ with respect to the intermediate representation $f(\mathbf{x})$, we compute the

gradient with respect to the input as $\nabla_x L = J_f(\mathbf{x})^T \nabla_f L$. This vector-Jacobian product constitutes the basic operation of reverse-mode automatic differentiation.

Theorem 1.1 (Multivariate Chain Rule). Let $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^n$ be differentiable functions. Then $h := g \circ f: \mathbb{R}^d \rightarrow \mathbb{R}^n$ is differentiable with Jacobian

$$J_h(\mathbf{x}) = J_g(f(\mathbf{x}))J_f(\mathbf{x})$$

for all $\mathbf{x} \in \mathbb{R}^d$. Equivalently, in component form,

$$\partial h_i / \partial x_j = \sum_{k=1}^m (\partial g_i / \partial y_k)|_{\mathbf{y}=f(\mathbf{x})} \cdot \partial f_k / \partial x_j$$

for $i = 1, \dots, n$ and $j = 1, \dots, d$.

Proof. Fix $\mathbf{x} \in \mathbb{R}^d$ and let $\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^m$. By differentiability of f at \mathbf{x} , for any $\varepsilon > 0$ there exists $\delta_1 > 0$ such that for $\|\mathbf{h}\| < \delta_1$,

$$\|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - J_f(\mathbf{x})\mathbf{h}\| \leq \varepsilon\|\mathbf{h}\|.$$

Similarly, by differentiability of g at \mathbf{y} , there exists $\delta_2 > 0$ such that for $\|\mathbf{k}\| < \delta_2$,

$$\|g(\mathbf{y} + \mathbf{k}) - g(\mathbf{y}) - J_g(\mathbf{y})\mathbf{k}\| \leq \varepsilon\|\mathbf{k}\|.$$

Set $\mathbf{k} = f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x})$. For sufficiently small \mathbf{h} , we have

$$\|\mathbf{k}\| \leq \|J_f(\mathbf{x})\mathbf{h}\| + \varepsilon\|\mathbf{h}\| \leq (\|J_f(\mathbf{x})\| + \varepsilon)\|\mathbf{h}\|.$$

Thus for $\|\mathbf{h}\|$ sufficiently small, $\|\mathbf{k}\| < \delta_2$. We then compute

$$h(\mathbf{x} + \mathbf{h}) - h(\mathbf{x}) = g(f(\mathbf{x} + \mathbf{h})) - g(f(\mathbf{x})) = g(\mathbf{y} + \mathbf{k}) - g(\mathbf{y}) = J_g(\mathbf{y})\mathbf{k} + r(\mathbf{k})$$

where $\|r(\mathbf{k})\| \leq \varepsilon\|\mathbf{k}\|$. Substituting $\mathbf{k} = J_f(\mathbf{x})\mathbf{h} + s(\mathbf{h})$ with $\|s(\mathbf{h})\| \leq \varepsilon\|\mathbf{h}\|$,

$$h(\mathbf{x} + \mathbf{h}) - h(\mathbf{x}) = J_g(\mathbf{y})J_f(\mathbf{x})\mathbf{h} + J_g(\mathbf{y})s(\mathbf{h}) + r(\mathbf{k}).$$

The error terms satisfy

$$\begin{aligned} \|J_g(\mathbf{y})s(\mathbf{h}) + r(\mathbf{k})\| &\leq \|J_g(\mathbf{y})\|\|s(\mathbf{h})\| + \|r(\mathbf{k})\| \leq (\|J_g(\mathbf{y})\| + 1)\varepsilon\|\mathbf{h}\| + \varepsilon\|\mathbf{k}\| \\ &\leq (\|J_g(\mathbf{y})\| + 1 + \varepsilon(\|J_f(\mathbf{x})\| + \varepsilon))\varepsilon\|\mathbf{h}\|. \end{aligned}$$

Since ε was arbitrary, we conclude that h is differentiable at \mathbf{x} with derivative $J_h(\mathbf{x}) = J_g(f(\mathbf{x}))J_f(\mathbf{x})$.

The chain rule extends to arbitrary compositions. For $h = f_L \circ f_{L-1} \circ \dots \circ f_1$, we have

$$J_h(\mathbf{x}) = J_{f_L}(z_{L-1})J_{f_{L-1}}(z_{L-2}) \dots J_{f_1}(\mathbf{x})$$

where $z_\ell = (f_\ell \circ \dots \circ f_1)(\mathbf{x})$ represents the intermediate activations. This telescoping product reveals both the power and the challenge of deep learning. Each Jacobian modulates the gradient signal, potentially causing exponential decay (vanishing gradients) or exponential growth (exploding gradients) as gradients propagate through many layers.

Taylor's theorem provides higher-order approximations essential for understanding optimization landscapes. For $f: \mathbb{R}^d \rightarrow \mathbb{R}$ twice continuously differentiable, we have

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{h} \rangle + \frac{1}{2} \langle \mathbf{h}, H_f(\mathbf{x}) \mathbf{h} \rangle + o(\|\mathbf{h}\|^2)$$

as $\mathbf{h} \rightarrow 0$. The quadratic term involving the Hessian $H_f(\mathbf{x})$ captures local curvature. When $H_f(\mathbf{x})$ is positive definite, \mathbf{x} is a strict local minimum. The condition number $\kappa(H_f(\mathbf{x})) = \lambda_{\max}/\lambda_{\min}$, the ratio of largest to smallest eigenvalue, governs the difficulty of optimization near \mathbf{x} .

Theorem 1.2 (Second-Order Taylor Expansion). Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be twice continuously differentiable on an open convex set $U \subset \mathbb{R}^d$. Then for any $\mathbf{x}, \mathbf{x} + \mathbf{h} \in U$,

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{h} \rangle + \frac{1}{2} \langle \mathbf{h}, H_f(\mathbf{x} + \theta \mathbf{h}) \mathbf{h} \rangle$$

for some $\theta \in (0,1)$, where H_f denotes the Hessian matrix of f .

This result follows from repeated application of the mean value theorem along the line segment $[\mathbf{x}, \mathbf{x} + \mathbf{h}]$. The first-order approximation $f(\mathbf{x} + \mathbf{h}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{h} \rangle$ underlies gradient-based optimization, while the second-order term reveals how curvature affects convergence rates. Newton's method exploits second-order information through the update $\mathbf{x}_{t+1} = \mathbf{x}_t - H_f(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$, achieving quadratic convergence near minima but at the cost of computing and inverting the Hessian.

Modern deep learning typically operates with first-order methods, computing only gradients rather than full Hessians. Stochastic gradient descent and its variants (momentum, Adam, RMSprop) use gradient estimates from minibatches to perform scalable optimization. The success of these methods despite ignoring second-order information reflects both the structure of neural network loss landscapes and implicit regularization through architectural choices and training procedures.

Computational graphs formalize the structure of composed functions as directed acyclic graphs where nodes represent operations and edges represent data flow. Each node corresponds to a function $f_i: \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{m_i}$, and edges indicate which outputs feed into which inputs. The overall computation $h: \mathbb{R}^d \rightarrow \mathbb{R}$ can be expressed as a composition of these primitives.

Automatic differentiation exploits the chain rule to compute gradients efficiently through computational graphs. Forward-mode differentiation computes Jacobian-vector products by propagating perturbations forward through the graph, while reverse-mode (backpropagation) computes vector-Jacobian products by propagating gradient information backward. For functions $\mathbb{R}^d \rightarrow \mathbb{R}$, reverse-mode requires only $O(1)$ backward passes regardless of d , making it optimal for the typical machine learning setting where d is large (millions or billions of parameters) but the output is a scalar loss.

The graph structure imposes an evaluation order respecting dependencies. A topological sort provides a valid sequence for forward evaluation. The reverse of this sequence gives a valid order for backpropagation. Modern automatic differentiation frameworks construct these computational graphs either statically (define-then-run) or dynamically (define-by-run), with the latter offering greater flexibility for control flow at some computational overhead.

1.4. The formal structure of learning problems

Having established the mathematical language, we can now formalize what constitutes a learning problem and situate optimization within this framework. This formulation provides the conceptual scaffolding for the entire treatment of machine learning to follow.

A learning problem begins with a data-generating distribution D defined over a space $X \times Y$, where X is the input or feature space and Y the output or label space. In supervised learning, each sample $(x, y) \sim D$ represents an input-output pair. The

distribution D is typically unknown, reflecting our epistemic position: we observe only finite samples from D , not D itself.

Given a hypothesis space H consisting of functions $h: X \rightarrow Y$, we seek to find $h^* \in H$ that minimizes the risk or expected loss

$$R(h) := \mathbb{E}_{(x,y) \sim D} [L(h(x), y)]$$

where $L: Y \times Y \rightarrow \mathbb{R}^+$ is a loss function quantifying the cost of predicting $h(x)$ when the true label is y . Common losses include squared error $L(\hat{y}, y) = (\hat{y} - y)^2$ for regression and cross-entropy $L(\hat{y}, y) = -\sum_k y_k \log(\hat{y}_k)$ for classification with probabilistic predictions.

The minimizer $h^* = \operatorname{argmin}_{h \in H} R(h)$, when it exists, is called the optimal hypothesis or Bayes optimal predictor within H . This represents the best possible performance achievable by the hypothesis class under the given loss. However, $R(h)$ cannot be computed directly without knowing D . This fundamental limitation drives the central strategy of empirical risk minimization.

Observing a training sample $S = \{(x_i, y_i)\}_{i=1}^n$ drawn i.i.d. from D , we define the empirical risk

$$R_{\text{emp}}(h; S) := \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i).$$

This empirical average provides an estimate of the true risk $R(h)$. The law of large numbers ensures that $R_{\text{emp}}(h; S) \rightarrow R(h)$ as $n \rightarrow \infty$ for any fixed h , provided appropriate regularity conditions hold. Empirical risk minimization seeks

$$\hat{h}_S = \operatorname{argmin}_{h \in H} R_{\text{emp}}(h; S),$$

selecting the hypothesis that performs best on the training data.

The fundamental question of learning theory concerns generalization: does \hat{h}_S , which minimizes empirical risk, achieve near-optimal true risk? Formally, we want to bound the generalization gap $R(\hat{h}_S) - R(h^*)$. This gap decomposes as $R(\hat{h}_S) - R(h^*) = [R(\hat{h}_S) - R_{\text{emp}}(\hat{h}_S; S)] + [R_{\text{emp}}(\hat{h}_S; S) - R_{\text{emp}}(h^*; S)] + [R_{\text{emp}}(h^*; S) - R(h^*)]$

. The second term is non-positive by definition of \hat{h}_S . The third term concentrates

around zero by the law of large numbers. The first term, measuring how empirical and true risk diverge for the selected hypothesis, requires more sophisticated analysis involving the complexity of H .

Capacity measures like VC dimension, Rademacher complexity, and covering numbers quantify how large H is in a sense relevant for generalization. Roughly, more complex hypothesis classes can fit training data better (lower empirical risk) but suffer larger generalization gaps for fixed sample size n . This creates the fundamental trade-off between expressiveness and generalization captured by the bias-variance decomposition and related frameworks.

Optimization algorithms implement empirical risk minimization approximately. Gradient descent and its stochastic variants iterate toward local minima of R_{emp} . When H is parameterized as $H = \{h_\theta : \theta \in \Theta \subset \mathbb{R}^p\}$, the empirical risk becomes a function of parameters:

$$J(\theta) := R_{\text{emp}}(h_\theta; S) = \frac{1}{n} \sum_{i=1}^n L(h_\theta(\mathbf{x}_i), y_i).$$

Gradient-based optimization seeks to minimize J through updates $\theta_{t+1} = \theta_t - \eta_t \nabla_\theta J(\theta_t)$. In stochastic gradient descent, gradients are estimated from random minibatches rather than the full dataset, providing computational efficiency and implicit regularization benefits.

Deep learning introduces additional structure: the hypothesis h_θ is a composition $h_\theta = f_L \circ \dots \circ f_1$ of parameterized layers. Each layer f_ℓ has parameters θ_ℓ , and the overall parameter vector θ concatenates all layer parameters. Backpropagation computes $\nabla_\theta J$ efficiently by applying the chain rule layer by layer in reverse topological order. The gradient with respect to layer ℓ parameters depends on the gradient from layer $\ell + 1$, creating a recursive structure that mirrors the forward computation.

Constraints and regularization modify the basic optimization problem. Ridge regression adds a penalty $\lambda \|\theta\|^2$ to the empirical risk, encouraging simpler models through parameter shrinkage. Constrained optimization enforces explicit requirements such as $\|\theta\|_2 \leq B$, maintaining parameters within a bounded set. These techniques

implement Occam's razor: among hypotheses fitting the data equally well, prefer the simplest.

The learning problem admits various generalizations. In unsupervised learning, Y is absent and we seek structure in the distribution of X alone, such as clustering or density estimation. Reinforcement learning replaces the fixed distribution D with an interactive environment, and the goal becomes maximizing cumulative reward through sequential decision-making. Multi-task learning shares representations across related problems, while meta-learning optimizes for rapid adaptation to new tasks.

Despite this diversity, the core structure remains: define a hypothesis space, specify an objective functional over that space, and employ optimization to select a good hypothesis. The mathematical language established in this chapter provides the foundation for making each of these components precise and for analyzing their interactions rigorously.

1.5. Concluding remarks

This chapter has established the mathematical vocabulary and fundamental structures that pervade artificial intelligence research. The notational system introduced here will remain consistent throughout this volume, enabling precise statements about models, algorithms, and theoretical guarantees.

We have emphasized several recurring themes. First, formalization exposes implicit assumptions and enables rigorous reasoning about learning systems. Second, the chain rule and its manifestations in backpropagation form the computational foundation of modern deep learning. Third, the statistical learning framework decomposes the problem of intelligence into manageable mathematical components: distributions, hypothesis classes, loss functions, and optimization algorithms.

The subsequent chapters build upon this foundation, developing the algebraic, probabilistic, and optimization-theoretic machinery required for a comprehensive treatment of AI. Chapter 2 examines linear algebra and numerical methods, providing tools for understanding dimensionality, spectral properties, and iterative algorithms.

Chapter 3 introduces probability and statistics, formalizing uncertainty quantification and inference. Chapter 4 develops optimization theory, analyzing the convergence and complexity of algorithms that power modern machine learning systems.

Throughout, we maintain the standard of rigor established here: precise definitions, complete proofs of fundamental results, and explicit connections between theory and practice. The goal is not merely to present algorithms but to understand why they work, when they fail, and how they can be improved. This mathematical perspective, while demanding, provides the deepest insight into the nature of learning and intelligence.

EXERCISES

1. Prove that the function $f(\mathbf{x}) = \|\mathbf{x}\|_2^2$ is convex on \mathbb{R}^d by verifying the definition directly.
2. Show that if $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and differentiable, then $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.
3. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be twice continuously differentiable. Prove that if $H_f(\mathbf{x})$ is positive definite for all \mathbf{x} , then f is strictly convex.
4. Consider the composition $h = g \circ f$ where $f: \mathbb{R} \rightarrow \mathbb{R}^2$ is defined by $f(t) = (\cos(t), \sin(t))$ and $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ by $g(x_1, x_2) = x_1^2 + x_2^2$. Compute $J_h(t)$ using the chain rule and verify by direct computation.
5. Let $\sigma(z) = 1/(1 + e^{-z})$ be the logistic sigmoid. Show that σ is Lipschitz continuous and compute its Lipschitz constant.
6. Prove that for any norm $\|\cdot\|$ on \mathbb{R}^d , the function $f(\mathbf{x}) = \|\mathbf{x}\|$ is convex.
7. Consider the neural network layer $f(\mathbf{x}) = \sigma(W\mathbf{x} + \mathbf{b})$ where σ is applied element-wise. Derive the formula for the Jacobian $J_f(\mathbf{x})$ in terms of W and the diagonal matrix $\text{Diag}(\sigma'(W\mathbf{x} + \mathbf{b}))$.
8. Show that if $f, g: \mathbb{R}^d \rightarrow \mathbb{R}$ are convex, then their pointwise maximum $h(\mathbf{x}) = \max(f(\mathbf{x}), g(\mathbf{x}))$ is also convex.

9. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and L -Lipschitz continuous. Prove that $\|\nabla f(\mathbf{x})\|_2 \leq L$ for all \mathbf{x} .
10. Verify that the empirical risk $R_{\text{emp}}(h; S)$ is an unbiased estimator of the true risk $R(h)$ when samples are drawn i.i.d. from D .