

Reinforcement Learning for Engineer-Mathematicians

A Comprehensive Guide to Theory and Applications

Author Name

About This Book

This book provides a comprehensive introduction to reinforcement learning with a focus on mathematical rigor and practical applications. Each chapter includes:

- Intuitive explanations and motivating examples
- Formal mathematical treatment
- Practical algorithms and implementation notes
- Interactive Jupyter notebooks with complete code

figures/rl_diagram.png

Published Year

Institution/Publisher

Copyright © 2024 Author Name
All rights reserved.

*To all engineers and mathematicians who seek to bridge
the gap between theory and practice*

Preface

This book bridges the gap between the mathematical rigor of reinforcement learning theory and its practical applications in engineering systems. Written for engineer-mathematicians, it provides both the theoretical foundations necessary to understand why algorithms work and the practical insights needed to apply them successfully in real-world scenarios.

The field of reinforcement learning has evolved rapidly, with deep learning enabling applications previously thought impossible. However, many engineering applications require understanding the underlying mathematics to ensure safety, reliability, and optimal performance. This book fills that need by providing a comprehensive treatment that balances theory with practice.

Each chapter builds upon previous concepts while maintaining mathematical rigor. Examples are drawn from engineering disciplines including robotics, control systems, power systems, manufacturing, and communications. The goal is to equip readers with both the theoretical understanding and practical skills needed to successfully apply reinforcement learning in their own domains.

Acknowledgments

The author wishes to thank the many researchers, colleagues, and students who contributed to the development of this book through discussions, feedback, and collaboration. Special thanks to the reinforcement learning community for their open sharing of ideas and implementations.

Contents

Preface	v
Acknowledgments	vii
I Mathematical Foundations	1
1 Introduction and Mathematical Prerequisites	5
1.1 Motivation: From Control Theory to Learning Systems	5
1.2 Mathematical Notation and Conventions	6
1.2.1 Sets and Spaces	7
1.2.2 Functions and Operators	7
1.2.3 Probability and Expectation	7
1.3 Probability Theory Refresher	7
1.3.1 Probability Spaces and Random Variables	8
1.3.2 Conditional Expectation and Martingales	8
1.3.3 Concentration Inequalities	8
1.4 Linear Algebra Essentials	9
1.4.1 Vector Spaces and Norms	9
1.4.2 Inner Products and Orthogonality	10
1.4.3 Eigenvalues and Spectral Theory	10
1.5 Optimization Fundamentals	10
1.5.1 Convex Analysis	10
1.5.2 Unconstrained Optimization	11
1.5.3 Gradient Descent and Convergence Analysis	11
1.6 Stochastic Processes and Markov Chains	11
1.6.1 Discrete-Time Stochastic Processes	12
1.6.2 Markov Chain Analysis	12
1.6.3 Mixing Times and Convergence Rates	12
1.7 Chapter Summary	13
2 Markov Decision Processes (MDPs)	15

2.1	Understanding MDPs Through Examples	15
2.2	Formal Definition and Mathematical Properties	16
2.2.1	Assumptions and Regularity Conditions	16
2.2.2	State and Action Spaces	17
2.3	Policies and Value Functions	17
2.3.1	Types of Policies	17
2.3.2	Value Function Theory	18
2.3.3	Bellman Equations	19
2.4	Optimal Policies and Bellman Optimality	19
2.4.1	Partial Ordering on Policies	19
2.4.2	Bellman Optimality Equations	19
2.5	Contraction Mapping Theorem and Fixed Points	20
2.5.1	Bellman Operators	20
2.5.2	Contraction Properties	20
2.5.3	Banach Fixed Point Theorem Application	20
2.6	Policy Improvement and Optimality	21
2.6.1	Policy Improvement Theorem	21
2.6.2	Policy Iteration Algorithm	21
2.7	Computational Complexity Analysis	21
2.7.1	Value Iteration Complexity	21
2.7.2	Policy Iteration Complexity	22
2.7.3	Modified Policy Iteration	22
2.8	Connections to Classical Control Theory	22
2.8.1	Linear Quadratic Regulator (LQR)	22
2.8.2	Hamilton-Jacobi-Bellman Equation	22
2.9	Advanced Topics	23
2.9.1	Partially Observable MDPs (POMDPs)	23
2.9.2	Constrained MDPs	23
2.10	Chapter Summary	23
3	Dynamic Programming Foundations	25
3.1	Principle of Optimality	25
3.1.1	Mathematical Formulation	25
3.1.2	Engineering Interpretation	26
3.2	Value Iteration: Convergence Analysis	26
3.2.1	Algorithm Description	26
3.2.2	Convergence Theory	26
3.2.3	Error Bounds and Stopping Criteria	27
3.2.4	Computational Complexity	27
3.3	Policy Iteration: Mathematical Guarantees	28
3.3.1	Policy Evaluation	28

3.3.2	Iterative Policy Evaluation	28
3.3.3	Policy Improvement Analysis	28
3.3.4	Global Convergence	29
3.4	Modified Policy Iteration	29
3.4.1	Algorithm and Convergence	29
3.4.2	Optimal Choice of Evaluation Steps	29
3.5	Asynchronous Dynamic Programming	30
3.5.1	Gauss-Seidel Value Iteration	30
3.5.2	Prioritized Sweeping	30
3.5.3	Real-Time Dynamic Programming	30
3.6	Linear Programming Formulation	31
3.6.1	Primal LP Formulation	31
3.6.2	Dual LP Formulation	31
3.7	Connections to Classical Control Theory	32
3.7.1	Discrete-Time Optimal Control	32
3.7.2	Stochastic Optimal Control	32
3.7.3	Model Predictive Control (MPC)	32
3.8	Computational Considerations	33
3.8.1	Curse of Dimensionality	33
3.8.2	Approximate Dynamic Programming	33
3.9	Chapter Summary	33
 II Core Algorithms and Theory		35
 4 Monte Carlo Methods		39
4.1	Monte Carlo Estimation Theory	39
4.1.1	Basic Monte Carlo Principle	39
4.1.2	Application to Value Function Estimation	39
4.2	First-Visit vs. Every-Visit Methods	40
4.2.1	First-Visit Monte Carlo	40
4.2.2	Every-Visit Monte Carlo	40
4.3	Variance Reduction Techniques	40
4.3.1	Incremental Implementation	41
4.3.2	Baseline Subtraction	41
4.3.3	Control Variates	41
4.4	Importance Sampling in RL	41
4.4.1	Ordinary Importance Sampling	41
4.4.2	Weighted Importance Sampling	42
4.4.3	Per-Decision Importance Sampling	42
4.5	Off-Policy Monte Carlo Methods	42
4.5.1	Off-Policy Policy Evaluation	42

4.5.2	Off-Policy Monte Carlo Control	42
4.6	Convergence Analysis and Sample Complexity	42
4.6.1	Finite Sample Analysis	42
4.6.2	Asymptotic Convergence Rate	44
4.6.3	Sample Complexity	44
4.7	Practical Considerations	44
4.7.1	Exploration vs. Exploitation	44
4.7.2	Function Approximation	44
4.8	Chapter Summary	45
5	Temporal Difference Learning	47
5.1	TD(0) Algorithm and Mathematical Analysis	47
5.1.1	Basic TD(0) Update	47
5.1.2	Relationship to Bellman Equation	47
5.1.3	Convergence Analysis	48
5.2	Bias-Variance Tradeoff in TD Methods	48
5.2.1	Bias Analysis	48
5.2.2	Variance Analysis	49
5.2.3	Mean Squared Error Decomposition	49
5.3	TD() and Eligibility Traces	49
5.3.1	Forward View: n-step Returns	49
5.3.2	-Return	49
5.3.3	Backward View: Eligibility Traces	50
5.3.4	Equivalence Theorem	50
5.4	Convergence Theory for Linear Function Approximation	51
5.4.1	Projected Bellman Equation	51
5.4.2	Error Bounds	51
5.5	Comparison with Monte Carlo and DP Methods	51
5.5.1	Computational Complexity	51
5.5.2	Sample Efficiency	51
5.5.3	Bootstrapping vs. Sampling	52
5.6	Advanced Topics	52
5.6.1	Multi-step Methods	52
5.6.2	True Online TD()	52
5.6.3	Gradient TD Methods	52
5.7	Chapter Summary	53
6	Chapter 06 Title	55
III	Function Approximation and Deep Learning	57
7	Chapter 07 Title	61

8 Chapter 08 Title	63
9 Chapter 09 Title	65
IV Advanced Topics	67
10 Chapter 10 Title	71
11 Chapter 11 Title	73
12 Chapter 12 Title	75
13 Chapter 13 Title	77
V Implementation and Practice	79
14 Chapter 14 Title	83
15 Chapter 15 Title	85
16 Chapter 16 Title	87
VI Future Directions	89
17 Chapter 17 Title	93
18 Chapter 18 Title	95
Mathematical Reference	97
.1 Matrix Calculus for RL	97
.1.1 Gradients and Jacobians	97
.1.2 Chain Rule	97
.1.3 Common Derivatives	97
.2 Probability Distributions Commonly Used	98
.2.1 Discrete Distributions	98
.2.2 Continuous Distributions	98
.3 Optimization Algorithms Summary	98
.3.1 Gradient Descent Variants	98
.4 Convergence Analysis Techniques	99
.4.1 Lyapunov Functions	99
.4.2 Martingale Convergence Theorem	99
.4.3 Robbins-Monro Conditions	99

Implementation Templates	101
.5 Basic RL Algorithm Implementations	101
.5.1 Value Iteration	101
.5.2 Q-Learning	102
.5.3 Policy Gradient (REINFORCE)	102
.6 Environment Interface Specifications	104
.6.1 OpenAI Gym Compatible Environment	104
.7 Logging and Visualization Code	105
.7.1 Training Logger	105
.8 Performance Benchmarking Utilities	107
.8.1 Algorithm Comparison Framework	107
Case Studies	111
.9 Case Study 1: Autonomous Drone Navigation	111
.9.1 Problem Description	111
.9.2 Algorithm Selection and Implementation	111
.9.3 Training Process and Results	112
.9.4 Lessons Learned	112
.10 Case Study 2: Smart Grid Energy Management	112
.10.1 Problem Description	113
.10.2 MDP Formulation	113
.10.3 Implementation Details	113
.10.4 Results and Performance Analysis	114
.11 Case Study 3: Manufacturing Process Optimization	114
.11.1 Problem Description	114
.11.2 Multi-Objective Optimization	115
.11.3 Implementation and Results	115
.12 Performance Comparisons	115
.12.1 Algorithm Performance Summary	115
.12.2 Common Success Factors	115
.13 Troubleshooting Guide	116
.13.1 Common Training Issues	116
.13.2 Hyperparameter Tuning Guidelines	116
.13.3 Debugging Checklist	117

List of Figures

List of Tables

List of Algorithms

1	Policy Iteration	21
2	Modified Policy Iteration	22
3	Value Iteration	26
4	Modified Policy Iteration	29
5	Gauss-Seidel Value Iteration	30
6	Prioritized Sweeping	31
7	Real-Time Dynamic Programming	31
8	Model Predictive Control	32
9	First-Visit Monte Carlo Policy Evaluation	40
10	Off-Policy Monte Carlo Policy Evaluation	43
11	Off-Policy Monte Carlo Control	43
12	Tabular TD(0) Policy Evaluation	48
13	TD() with Eligibility Traces	50

Part I

Mathematical Foundations

This part establishes the mathematical foundations necessary for understanding reinforcement learning from both theoretical and engineering perspectives. We begin with essential mathematical prerequisites, then develop the formal framework of Markov Decision Processes, and conclude with classical dynamic programming methods that form the basis for modern reinforcement learning algorithms.

The treatment emphasizes mathematical rigor while maintaining practical relevance for engineering applications. Each concept is developed with careful attention to assumptions, proofs, and connections to control theory and optimization.

Chapter 1

Introduction and Mathematical Prerequisites

Chapter Overview

This chapter introduces the fundamental mathematical tools needed for reinforcement learning and provides intuitive motivation for why RL represents a paradigm shift from classical control theory. We'll cover probability theory, linear algebra, optimization, and stochastic processes with practical examples.

1.1 Motivation: From Control Theory to Learning Systems

Why Reinforcement Learning?

Imagine teaching a child to ride a bicycle. You don't give them the equations of motion or tell them exactly how to balance. Instead, they learn through trial and error, gradually improving their balance and control. This is the essence of reinforcement learning.

Reinforcement learning represents a fundamental paradigm shift from classical control theory and optimization. While traditional engineering approaches rely on explicit models and well-defined objectives, reinforcement learning enables systems to learn optimal behavior through interaction with their environment.

Traditional Control vs. Reinforcement Learning

Consider a classic engineering problem: designing a controller for an inverted pendulum.

Traditional Control Approach:

1. Deriving the system dynamics using Lagrangian mechanics
2. Linearizing around the equilibrium point
3. Designing a feedback controller using pole placement or LQR
4. Implementing the controller with known parameters

Reinforcement Learning Approach:

1. Define states (angle, angular velocity) and actions (applied force)
2. Specify a reward function (positive for upright, negative for falling)
3. Allow the agent to explore and learn through trial and error
4. Converge to an optimal policy without explicit knowledge of dynamics

This fundamental difference opens up possibilities for systems where:

- Dynamics are unknown or too complex to model accurately
- Environment conditions change over time
- Multiple conflicting objectives must be balanced
- System parameters vary or degrade over time

Industrial Example: Power Grid Management

Modern power grids face unprecedented challenges with renewable energy integration, electric vehicle charging, and dynamic pricing. Traditional grid control relies on pre-computed lookup tables and heuristic rules. Reinforcement learning enables real-time optimization that adapts to:

- Variable renewable generation
- Changing demand patterns
- Equipment failures and network topology changes
- Market price fluctuations

1.2 Mathematical Notation and Conventions

Notation Guide

Throughout this book, we adopt consistent mathematical notation that aligns with both control theory and machine learning conventions. Don't worry if some symbols are unfamiliar now—we'll introduce them gradually with intuitive explanations.

Throughout this book, we adopt consistent mathematical notation that aligns with both

control theory and machine learning conventions.

1.2.1 Sets and Spaces

Understanding Spaces

Think of a "space" as the collection of all possible values a variable can take. For example, if we're controlling a robot arm, the state space might include all possible joint angles and velocities.

$$\mathcal{S} = \text{State space (all possible states)} \quad (1.1)$$

$$\mathcal{A} = \text{Action space (all possible actions)} \quad (1.2)$$

$$\mathcal{R} = \text{Reward space (all possible rewards)} \quad (1.3)$$

$$\mathbb{R}^n = n\text{-dimensional real vector space} \quad (1.4)$$

$$\mathbb{R}^{m \times n} = \text{Space of } m \times n \text{ real matrices} \quad (1.5)$$

1.2.2 Functions and Operators

$$\pi : \mathcal{S} \rightarrow \mathcal{A} \quad (\text{Deterministic policy}) \quad (1.6)$$

$$\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A}) \quad (\text{Stochastic policy}) \quad (1.7)$$

$$V^\pi : \mathcal{S} \rightarrow \mathbb{R} \quad (\text{Value function}) \quad (1.8)$$

$$Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \quad (\text{Action-value function}) \quad (1.9)$$

$$T : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}} \quad (\text{Bellman operator}) \quad (1.10)$$

where $\Delta(\mathcal{A})$ denotes the space of probability distributions over \mathcal{A} .

1.2.3 Probability and Expectation

$$\mathbb{P}(s'|s, a) = \text{Transition probability} \quad (1.11)$$

$$\mathbb{E}_\pi[\cdot] = \text{Expectation under policy } \pi \quad (1.12)$$

$$\mathbb{E}_{s \sim \mu}[\cdot] = \text{Expectation over distribution } \mu \quad (1.13)$$

1.3 Probability Theory Refresher

Reinforcement learning is fundamentally about making decisions under uncertainty. A solid understanding of probability theory is essential for analyzing convergence properties, sample complexity, and algorithm performance.

1.3.1 Probability Spaces and Random Variables

Definition 1.1 (Probability Space). A probability space is a triple $(\Omega, \mathcal{F}, \mathbb{P})$ where:

- Ω is the sample space (set of all possible outcomes)
- \mathcal{F} is a σ -algebra on Ω (collection of measurable events)
- $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ is a probability measure satisfying:
 1. $\mathbb{P}(\Omega) = 1$

2. For disjoint events A_1, A_2, \dots : $\mathbb{P}(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \mathbb{P}(A_i)$

Definition 1.2 (Random Variable). A random variable X is a measurable function $X : \Omega \rightarrow \mathbb{R}$ such that for every Borel set $B \subseteq \mathbb{R}$, the preimage $X^{-1}(B) \in \mathcal{F}$.

1.3.2 Conditional Expectation and Martingales

Conditional expectation plays a crucial role in reinforcement learning, particularly in the analysis of temporal difference methods and policy gradient algorithms.

Definition 1.3 (Conditional Expectation). Given random variables X and Y , the conditional expectation $\mathbb{E}[X|Y]$ is the unique (almost surely) random variable that is:

1. Measurable with respect to $\sigma(Y)$
2. Satisfies $\mathbb{E}[\mathbb{E}[X|Y] \cdot \mathbf{1}_A] = \mathbb{E}[X \cdot \mathbf{1}_A]$ for all $A \in \sigma(Y)$

Theorem 1.1 (Law of Total Expectation). For random variables X and Y :

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|Y]] \tag{1.14}$$

Definition 1.4 (Martingale). A sequence of random variables $\{X_t\}_{t=0}^{\infty}$ is a martingale with respect to filtration $\{\mathcal{F}_t\}_{t=0}^{\infty}$ if:

1. X_t is \mathcal{F}_t -measurable for all t
2. $\mathbb{E}[|X_t|] < \infty$ for all t
3. $\mathbb{E}[X_{t+1}|\mathcal{F}_t] = X_t$ almost surely

Martingales are fundamental for proving convergence of stochastic algorithms in reinforcement learning.

1.3.3 Concentration Inequalities

Concentration inequalities provide bounds on the probability that random variables deviate from their expected values. These are essential for finite-sample analysis of RL algorithms.

Theorem 1.2 (Hoeffding's Inequality). Let X_1, \dots, X_n be independent random variables

with $X_i \in [a_i, b_i]$ almost surely. Then for any $t > 0$:

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \right| \geq t \right) \leq 2 \exp \left(- \frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right) \quad (1.15)$$

Theorem 1.3 (Azuma's Inequality). Let $\{X_t\}_{t=0}^\infty$ be a martingale with respect to $\{\mathcal{F}_t\}_{t=0}^\infty$ such that $|X_{t+1} - X_t| \leq c_t$ almost surely. Then:

$$\mathbb{P}(|X_n - X_0| \geq t) \leq 2 \exp \left(- \frac{t^2}{2 \sum_{i=0}^{n-1} c_i^2} \right) \quad (1.16)$$

1.4 Linear Algebra Essentials

Linear algebra provides the foundation for function approximation, policy parameterization, and many algorithmic techniques in reinforcement learning.

1.4.1 Vector Spaces and Norms

Definition 1.5 (Vector Space). A vector space V over field \mathbb{F} (typically \mathbb{R} or \mathbb{C}) is a set equipped with vector addition and scalar multiplication satisfying:

1. Commutativity: $u + v = v + u$
2. Associativity: $(u + v) + w = u + (v + w)$
3. Identity: $\exists 0 \in V$ such that $v + 0 = v$
4. Inverse: $\forall v \in V, \exists -v$ such that $v + (-v) = 0$
5. Scalar associativity: $a(bv) = (ab)v$
6. Scalar identity: $1v = v$
7. Distributivity: $a(u + v) = au + av$ and $(a + b)v = av + bv$

Definition 1.6 (Norm). A norm on vector space V is a function $\|\cdot\| : V \rightarrow \mathbb{R}_{\geq 0}$ satisfying:

1. $\|v\| = 0$ if and only if $v = 0$
2. $\|av\| = |a|\|v\|$ for scalar a
3. $\|u + v\| \leq \|u\| + \|v\|$ (triangle inequality)

Common norms in \mathbb{R}^n :

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad (\ell_1 \text{ norm}) \quad (1.17)$$

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad (\text{Euclidean norm}) \quad (1.18)$$

$$\|x\|_\infty = \max_{i=1,\dots,n} |x_i| \quad (\ell_\infty \text{ norm}) \quad (1.19)$$

1.4.2 Inner Products and Orthogonality

Definition 1.7 (Inner Product). An inner product on real vector space V is a function $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$ satisfying:

1. Symmetry: $\langle u, v \rangle = \langle v, u \rangle$
2. Linearity: $\langle au + bv, w \rangle = a\langle u, w \rangle + b\langle v, w \rangle$
3. Positive definiteness: $\langle v, v \rangle \geq 0$ with equality iff $v = 0$

The induced norm is $\|v\| = \sqrt{\langle v, v \rangle}$.

Theorem 1.4 (Cauchy-Schwarz Inequality). For vectors u, v in an inner product space:

$$|\langle u, v \rangle| \leq \|u\| \|v\| \quad (1.20)$$

with equality if and only if u and v are linearly dependent.

1.4.3 Eigenvalues and Spectral Theory

Definition 1.8 (Eigenvalue and Eigenvector). For matrix $A \in \mathbb{R}^{n \times n}$, scalar λ is an eigenvalue with corresponding eigenvector $v \neq 0$ if:

$$Av = \lambda v \quad (1.21)$$

Theorem 1.5 (Spectral Theorem for Symmetric Matrices). Every real symmetric matrix A has an orthonormal basis of eigenvectors with real eigenvalues. If $A = Q\Lambda Q^T$ where Q is orthogonal and Λ is diagonal, then:

$$A = \sum_{i=1}^n \lambda_i q_i q_i^T \quad (1.22)$$

where λ_i are eigenvalues and q_i are corresponding orthonormal eigenvectors.

1.5 Optimization Fundamentals

Optimization theory underpins virtually all reinforcement learning algorithms, from value iteration to policy gradient methods.

1.5.1 Convex Analysis

Definition 1.9 (Convex Set). A set $C \subseteq \mathbb{R}^n$ is convex if for all $x, y \in C$ and $\lambda \in [0, 1]$:

$$\lambda x + (1 - \lambda)y \in C \quad (1.23)$$

Definition 1.10 (Convex Function). A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if its domain is convex and for all x, y in the domain and $\lambda \in [0, 1]$:

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (1.24)$$

Theorem 1.6 (First-Order Characterization of Convexity). For differentiable function f , the following are equivalent:

1. f is convex
2. $f(y) \geq f(x) + \nabla f(x)^T(y - x)$ for all x, y
3. ∇f is monotone: $(\nabla f(x) - \nabla f(y))^T(x - y) \geq 0$

1.5.2 Unconstrained Optimization

Theorem 1.7 (Necessary Conditions for Optimality). If x^* is a local minimum of differentiable function f , then:

$$\nabla f(x^*) = 0 \quad (1.25)$$

If f is twice differentiable, then additionally:

$$\nabla^2 f(x^*) \succeq 0 \quad (\text{positive semidefinite}) \quad (1.26)$$

Theorem 1.8 (Sufficient Conditions for Optimality). If $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*) \succ 0$ (positive definite), then x^* is a strict local minimum.

1.5.3 Gradient Descent and Convergence Analysis

The gradient descent algorithm iterates:

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) \quad (1.27)$$

Theorem 1.9 (Convergence of Gradient Descent). For convex function f with L -Lipschitz gradient and step size $\alpha \leq 1/L$:

$$f(x_k) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\alpha k} \quad (1.28)$$

where x^* is the optimal solution.

For strongly convex functions, the convergence rate improves to exponential.

1.6 Stochastic Processes and Markov Chains

Understanding stochastic processes is crucial for analyzing the temporal dynamics of reinforcement learning systems.

1.6.1 Discrete-Time Stochastic Processes

Definition 1.11 (Stochastic Process). A discrete-time stochastic process is a sequence of random variables $\{X_t\}_{t=0}^{\infty}$ where each X_t takes values in some state space \mathcal{S} .

Definition 1.12 (Markov Property). A stochastic process $\{X_t\}_{t=0}^{\infty}$ satisfies the Markov property if:

$$\mathbb{P}(X_{t+1} = s' | X_t = s, X_{t-1} = s_{t-1}, \dots, X_0 = s_0) = \mathbb{P}(X_{t+1} = s' | X_t = s) \quad (1.29)$$

for all states $s, s', s_0, \dots, s_{t-1}$ and times $t \geq 0$.

1.6.2 Markov Chain Analysis

For finite state space $\mathcal{S} = \{1, 2, \dots, n\}$, a Markov chain is characterized by its transition matrix $P \in \mathbb{R}^{n \times n}$ where $P_{ij} = \mathbb{P}(X_{t+1} = j | X_t = i)$.

Definition 1.13 (Irreducibility and Aperiodicity). A Markov chain is:

- **Irreducible** if every state is reachable from every other state
- **Aperiodic** if $\gcd\{n \geq 1 : P_{ii}^{(n)} > 0\} = 1$ for some state i

Theorem 1.10 (Fundamental Theorem of Markov Chains). For an irreducible, aperiodic, finite Markov chain:

1. There exists a unique stationary distribution π satisfying $\pi = \pi P$
2. $\lim_{t \rightarrow \infty} P^t = \mathbf{1}\pi^T$ where $\mathbf{1}$ is the vector of ones
3. For any initial distribution μ_0 : $\lim_{t \rightarrow \infty} \|\mu_t - \pi\|_{TV} = 0$

1.6.3 Mixing Times and Convergence Rates

Definition 1.14 (Total Variation Distance). The total variation distance between distributions μ and ν on finite space \mathcal{S} is:

$$\|\mu - \nu\|_{TV} = \frac{1}{2} \sum_{s \in \mathcal{S}} |\mu(s) - \nu(s)| \quad (1.30)$$

Definition 1.15 (Mixing Time). The mixing time of a Markov chain is:

$$t_{mix}(\epsilon) = \min\{t : \max_{i \in \mathcal{S}} \|P^t(i, \cdot) - \pi\|_{TV} \leq \epsilon\} \quad (1.31)$$

Understanding mixing times is essential for analyzing sample complexity in reinforcement learning algorithms that rely on sampling from stationary distributions.

1.7 Chapter Summary

This chapter established the mathematical foundations necessary for rigorous analysis of reinforcement learning algorithms. Key concepts include:

- The paradigm shift from model-based control to learning-based optimization
- Probability theory tools: conditional expectation, martingales, concentration inequalities
- Linear algebra foundations: vector spaces, norms, spectral theory
- Convex optimization and gradient descent convergence analysis
- Markov chain theory and convergence to stationary distributions

These mathematical tools will be applied throughout the book to analyze algorithm convergence, sample complexity, and performance guarantees. The next chapter develops the formal framework of Markov Decision Processes, which provides the mathematical foundation for all subsequent reinforcement learning algorithms.

Chapter 2

Markov Decision Processes (MDPs)

Chapter Overview

This chapter introduces Markov Decision Processes (MDPs), the mathematical foundation of reinforcement learning. We'll build intuition through concrete examples before diving into the formal theory, then explore solution methods like value iteration and policy iteration.

What is an MDP?

Think of an MDP as a mathematical description of a decision-making situation where:

- You observe the current situation (state)
- You choose an action based on what you observe
- The world responds by transitioning to a new state and giving you a reward
- This process repeats over time

The key insight is that the future only depends on the current state, not the entire history—this is the Markov property.

Markov Decision Processes provide the mathematical framework for modeling sequential decision-making under uncertainty. This chapter develops the formal theory of MDPs with particular attention to mathematical rigor and engineering applications.

2.1 Understanding MDPs Through Examples

Before diving into formal definitions, let's build intuition through a concrete example.

Grid World Navigation

Consider a robot navigating a 4×4 grid world:

- **States:** Each cell in the grid (16 total states)
- **Actions:** Move up, down, left, or right
- **Transitions:** Move to adjacent cell (or stay put if hitting a wall)
- **Rewards:** +10 for reaching the goal, -1 for each step, -10 for falling into holes
- **Goal:** Find the shortest path to the target while avoiding obstacles

2.2 Formal Definition and Mathematical Properties

Now that we have intuition, let's formalize these concepts.

Definition 2.1 (Markov Decision Process). A Markov Decision Process is a 5-tuple $(\mathcal{S}, \mathcal{A}, P, \mathcal{R}, \gamma)$ where:

- \mathcal{S} is the **state space** (all possible situations)
- \mathcal{A} is the **action space** (all possible actions)
- $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the **transition kernel** (dynamics)
- $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the **reward function** (immediate feedback)
- $\gamma \in [0, 1)$ is the **discount factor** (how much we value future rewards)

Understanding the Components

- **State space \mathcal{S} :** All possible configurations of your system
- **Action space \mathcal{A} :** All decisions you can make in any given state
- **Transition function P :** Describes how actions change states (the "physics" of your world)
- **Reward function \mathcal{R} :** Immediate feedback telling you how good an action was
- **Discount factor γ :** How much you care about future vs. immediate rewards (0 = only care about immediate, close to 1 = care about long-term)

The transition kernel satisfies $\sum_{s' \in \mathcal{S}} P(s, a, s') = 1$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, and we write $P(s'|s, a) = P(s, a, s')$ for the probability of transitioning to state s' from state s under action a .

2.2.1 Assumptions and Regularity Conditions

For mathematical tractability, we typically assume:

[Measurability] The state and action spaces are measurable spaces, and the transition kernel and reward function are measurable with respect to the appropriate σ -algebras.

[Bounded Rewards] The reward function satisfies $\sup_{s,a} |\mathcal{R}(s, a)| \leq R_{max} < \infty$.

[Discount Factor] The discount factor satisfies $\gamma \in [0, 1)$ to ensure convergence of infinite-horizon value functions.

2.2.2 State and Action Spaces

Discrete Spaces

For finite MDPs with $|\mathcal{S}| = n$ and $|\mathcal{A}| = m$, we can represent:

- Transition probabilities as tensors $P^a \in \mathbb{R}^{n \times n}$ for each action a
- Rewards as matrices $R \in \mathbb{R}^{n \times m}$
- Policies as matrices $\Pi \in [0, 1]^{n \times m}$ with $\sum_a \Pi(s, a) = 1$

Continuous Spaces

For continuous state spaces $\mathcal{S} \subseteq \mathbb{R}^d$, the transition kernel becomes a probability measure:

$$P(\cdot | s, a) : \mathcal{B}(\mathcal{S}) \rightarrow [0, 1] \quad (2.1)$$

where $\mathcal{B}(\mathcal{S})$ is the Borel σ -algebra on \mathcal{S} .

Engineering Example: Inverted Pendulum

Consider an inverted pendulum with:

- State: $s = (\theta, \dot{\theta}) \in [-\pi, \pi] \times [-10, 10]$ (angle and angular velocity)
- Action: $a \in [-5, 5]$ (applied torque)
- Dynamics: $\ddot{\theta} = \frac{g}{l} \sin \theta + \frac{a}{ml^2}$ plus noise
- Reward: $r(s, a) = -\theta^2 - 0.1\dot{\theta}^2 - 0.01a^2$ (quadratic cost)

2.3 Policies and Value Functions

What is a Policy?

A policy is simply a decision-making rule. It tells an agent what action to take in each possible state. Think of it as a strategy or game plan.

2.3.1 Types of Policies

Definition 2.2 (Deterministic Policy). A deterministic policy is a function $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maps each state to exactly one action.

Deterministic Policy Example

In our grid world: "Always move towards the goal" could be a deterministic policy where $\pi(\text{state}) = \text{direction_to_goal}$.

Definition 2.3 (Stochastic Policy). A stochastic policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ assigns a probability distribution over actions for each state, where $\Delta(\mathcal{A})$ is the space of probability measures on \mathcal{A} .

Stochastic Policy Example

In grid world: "Move towards goal with probability 0.8, move randomly otherwise" gives $\pi(\text{best_action}|s) = 0.8$ and equal probability to other actions.

Definition 2.4 (History-Dependent Policy). A history-dependent policy depends on the entire sequence of past states and actions: $\pi_t : (\mathcal{S} \times \mathcal{A})^t \times \mathcal{S} \rightarrow \Delta(\mathcal{A})$

Why Focus on Markovian Policies?

While policies could potentially use the entire history, the Markov property means that optimal policies only need to depend on the current state. This greatly simplifies our analysis!

Theorem 2.1 (Sufficiency of Markovian Policies). For any history-dependent policy, there exists a Markovian policy that achieves the same expected discounted reward.

Proof. This follows from the Markov property of the state transitions. The expected future reward depends only on the current state, not the history of how that state was reached. \square

2.3.2 Value Function Theory

Definition 2.5 (State Value Function). The state value function for policy π is:

$$V^\pi(s) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s \right] \quad (2.2)$$

Definition 2.6 (Action Value Function). The action value function (Q-function) for policy π is:

$$Q^\pi(s, a) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a \right] \quad (2.3)$$

Theorem 2.2 (Existence and Uniqueness of Value Functions). Under Assumptions 1-3, the value functions V^π and Q^π exist, are unique, and satisfy $\|V^\pi\|_\infty \leq \frac{R_{max}}{1-\gamma}$.

Proof. The geometric series $\sum_{t=0}^{\infty} \gamma^t R_{max}$ converges to $\frac{R_{max}}{1-\gamma}$ since $\gamma < 1$. Uniqueness follows from the linearity of expectation. \square

2.3.3 Bellman Equations

The fundamental recursive relationships in reinforcement learning are the Bellman equations.

Theorem 2.3 (Bellman Equations for Policy Evaluation). For any policy π :

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left[\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s') \right] \quad (2.4)$$

$$Q^\pi(s, a) = \mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \sum_{a' \in \mathcal{A}} \pi(a'|s') Q^\pi(s', a') \quad (2.5)$$

Proof. By the tower rule of conditional expectation:

$$V^\pi(s) = \mathbb{E}^\pi \left[\mathcal{R}(S_0, A_0) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} \mathcal{R}(S_t, A_t) \mid S_0 = s \right] \quad (2.6)$$

$$= \mathbb{E}^\pi[\mathcal{R}(S_0, A_0) \mid S_0 = s] + \gamma \mathbb{E}^\pi[V^\pi(S_1) \mid S_0 = s] \quad (2.7)$$

Expanding the expectations gives the Bellman equation. \square

2.4 Optimal Policies and Bellman Optimality

2.4.1 Partial Ordering on Policies

Definition 2.7 (Policy Partial Order). Policy π_1 dominates policy π_2 (written $\pi_1 \geq \pi_2$) if:

$$V^{\pi_1}(s) \geq V^{\pi_2}(s) \quad \forall s \in \mathcal{S} \quad (2.8)$$

Theorem 2.4 (Existence of Optimal Policies). There exists an optimal deterministic policy π^* such that:

$$V^{\pi^*}(s) = \max_{\pi} V^\pi(s) \equiv V^*(s) \quad \forall s \in \mathcal{S} \quad (2.9)$$

2.4.2 Bellman Optimality Equations

Theorem 2.5 (Bellman Optimality Equations). The optimal value functions satisfy:

$$V^*(s) = \max_{a \in \mathcal{A}} \left[\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^*(s') \right] \quad (2.10)$$

$$Q^*(s, a) = \mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \max_{a' \in \mathcal{A}} Q^*(s', a') \quad (2.11)$$

Corollary 2.1 (Optimal Policy Extraction). An optimal policy can be extracted from the optimal value functions as:

$$\pi^*(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a) \quad (2.12)$$

2.5 Contraction Mapping Theorem and Fixed Points

The mathematical foundation for proving convergence of dynamic programming algorithms relies on contraction mapping theory.

2.5.1 Bellman Operators

Definition 2.8 (Bellman Operator). For policy π , the Bellman operator $T^\pi : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ is defined by:

$$(T^\pi V)(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left[\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s') \right] \quad (2.13)$$

Definition 2.9 (Bellman Optimality Operator). The Bellman optimality operator $T^* : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ is defined by:

$$(T^* V)(s) = \max_{a \in \mathcal{A}} \left[\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s') \right] \quad (2.14)$$

2.5.2 Contraction Properties

Theorem 2.6 (Contraction Property of Bellman Operators). Under the supremum norm $\|V\|_\infty = \max_{s \in \mathcal{S}} |V(s)|$:

1. T^π is a γ -contraction: $\|T^\pi V_1 - T^\pi V_2\|_\infty \leq \gamma \|V_1 - V_2\|_\infty$
2. T^* is a γ -contraction: $\|T^* V_1 - T^* V_2\|_\infty \leq \gamma \|V_1 - V_2\|_\infty$

Proof. For the policy operator:

$$|(T^\pi V_1)(s) - (T^\pi V_2)(s)| = \left| \sum_a \pi(a|s) \gamma \sum_{s'} P(s'|s, a) [V_1(s') - V_2(s')] \right| \quad (2.15)$$

$$\leq \sum_a \pi(a|s) \gamma \sum_{s'} P(s'|s, a) |V_1(s') - V_2(s')| \quad (2.16)$$

$$\leq \gamma \|V_1 - V_2\|_\infty \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) \quad (2.17)$$

$$= \gamma \|V_1 - V_2\|_\infty \quad (2.18)$$

The proof for T^* follows similarly using the fact that the max operator is non-expansive. \square

2.5.3 Banach Fixed Point Theorem Application

Theorem 2.7 (Banach Fixed Point Theorem). Let (X, d) be a complete metric space and $T : X \rightarrow X$ be a contraction mapping with contraction factor $\gamma < 1$. Then:

1. T has a unique fixed point $x^* \in X$
2. For any $x_0 \in X$, the sequence $x_{n+1} = T(x_n)$ converges to x^*

3. The convergence rate is geometric: $d(x_n, x^*) \leq \gamma^n d(x_0, x^*)$

Corollary 2.2 (Convergence of Value Iteration). The value iteration algorithm $V_{k+1} = T^*V_k$ converges geometrically to the unique optimal value function V^* at rate γ .

2.6 Policy Improvement and Optimality

2.6.1 Policy Improvement Theorem

Theorem 2.8 (Policy Improvement Theorem). Let π be any policy and define the improved policy π' by:

$$\pi'(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^\pi(s, a) \quad (2.19)$$

Then $\pi' \geq \pi$, with strict inequality unless π is optimal.

Proof. For any state s :

$$Q^\pi(s, \pi'(s)) \geq Q^\pi(s, \pi(s)) = V^\pi(s) \quad (2.20)$$

By the policy evaluation equation and induction, this implies $V^{\pi'}(s) \geq V^\pi(s)$. \square

2.6.2 Policy Iteration Algorithm

The policy improvement theorem leads to the policy iteration algorithm:

Algorithm 1 Policy Iteration

Input: Initial policy π_0

Output: Optimal policy π^*

$k \leftarrow 0$

repeat

Policy Evaluation: Solve $V^{\pi_k} = T^{\pi_k} V^{\pi_k}$

Policy Improvement: $\pi_{k+1}(s) \leftarrow \operatorname{argmax}_a Q^{\pi_k}(s, a)$

$k \leftarrow k + 1$

until $\pi_k = \pi_{k-1}$

return $\pi^* = \pi_k$

Theorem 2.9 (Convergence of Policy Iteration). Policy iteration converges to an optimal policy in finitely many iterations for finite MDPs.

2.7 Computational Complexity Analysis

2.7.1 Value Iteration Complexity

For finite MDPs with $|\mathcal{S}| = n$ and $|\mathcal{A}| = m$:

- **Time per iteration:** $O(mn^2)$ operations
- **Iterations to ϵ -accuracy:** $O(\log(\epsilon^{-1}))$ iterations

- **Total complexity:** $O(mn^2 \log(\epsilon^{-1}))$

2.7.2 Policy Iteration Complexity

- **Policy evaluation:** $O(n^3)$ for direct matrix inversion, $O(n^2)$ per iteration for iterative methods
- **Policy improvement:** $O(mn^2)$
- **Number of policy iterations:** At most m^n (typically much smaller)

2.7.3 Modified Policy Iteration

To balance the computational costs, modified policy iteration performs only k steps of policy evaluation:

Algorithm 2 Modified Policy Iteration

Input: Initial policy π_0 , evaluation steps k
 Initialize V_0 arbitrarily
for $i = 0, 1, 2, \dots$ **do**
 for $j = 1, 2, \dots, k$ **do**
 $V_j \leftarrow T^{\pi_i} V_{j-1}$
 end for
 $\pi_{i+1}(s) \leftarrow \operatorname{argmax}_a [r(s, a) + \gamma \sum_{s'} P(s'|s, a) V_k(s')]$
end for

2.8 Connections to Classical Control Theory

2.8.1 Linear Quadratic Regulator (LQR)

For linear dynamics $s_{t+1} = As_t + Ba_t + w_t$ and quadratic costs $r(s, a) = -s^T Qs - a^T Ra$, the optimal value function is quadratic: $V^*(s) = -s^T Ps$ where P satisfies the discrete algebraic Riccati equation:

$$P = Q + A^T P A - A^T P B (R + B^T P B)^{-1} B^T P A \quad (2.21)$$

The optimal policy is linear: $\pi^*(s) = -Ks$ where $K = (R + B^T P B)^{-1} B^T P A$.

2.8.2 Hamilton-Jacobi-Bellman Equation

For continuous-time systems, the Bellman equation becomes the Hamilton-Jacobi-Bellman (HJB) partial differential equation:

$$\frac{\partial V}{\partial t} + \min_a \left[r(s, a) + \frac{\partial V}{\partial s} f(s, a) \right] = 0 \quad (2.22)$$

where $f(s, a)$ is the system dynamics.

2.9 Advanced Topics

2.9.1 Partially Observable MDPs (POMDPs)

Definition 2.10 (POMDP). A POMDP extends an MDP with observations: $(\mathcal{S}, \mathcal{A}, \mathcal{O}, P, \mathcal{R}, \mathcal{Z}, \gamma)$ where:

- \mathcal{O} is the observation space
- $\mathcal{Z} : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow [0, 1]$ is the observation model

The optimal policy depends on the belief state $b(s) = \mathbb{P}(S_t = s | h_t)$ where h_t is the history of observations.

2.9.2 Constrained MDPs

Definition 2.11 (Constrained MDP). A constrained MDP adds constraint functions $c_i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and thresholds d_i :

$$\max_{\pi} \mathbb{E}^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \right] \quad (2.23)$$

$$\text{subject to } \mathbb{E}^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t c_i(S_t, A_t) \right] \leq d_i, \quad i = 1, \dots, m \quad (2.24)$$

Solutions typically use Lagrangian methods and primal-dual algorithms.

2.10 Chapter Summary

This chapter established the mathematical foundations of Markov Decision Processes:

- Formal definition of MDPs and regularity assumptions
- Policy and value function theory with existence and uniqueness results
- Bellman equations and optimality conditions
- Contraction mapping theory and convergence guarantees
- Dynamic programming algorithms: value iteration and policy iteration
- Computational complexity analysis
- Connections to classical control theory and advanced extensions

The mathematical framework developed here provides the foundation for all reinforcement learning algorithms. The next chapter examines dynamic programming methods in detail, providing the algorithmic foundation for modern RL techniques.

Chapter 3

Dynamic Programming Foundations

Dynamic programming provides the theoretical and algorithmic foundation for reinforcement learning. This chapter develops the mathematical theory of dynamic programming with emphasis on convergence analysis, computational complexity, and connections to classical optimal control.

3.1 Principle of Optimality

The fundamental insight underlying dynamic programming is Bellman’s principle of optimality, which enables the decomposition of complex sequential decision problems into simpler subproblems.

Theorem 3.1 (Principle of Optimality). An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.

3.1.1 Mathematical Formulation

For an MDP $(\mathcal{S}, \mathcal{A}, P, \mathcal{R}, \gamma)$, consider a finite-horizon problem with horizon T . Define the optimal value function:

$$V_t^*(s) = \max_{\pi} \mathbb{E} \left[\sum_{k=t}^{T-1} \gamma^{k-t} \mathcal{R}(S_k, A_k) \mid S_t = s, \pi \right] \quad (3.1)$$

Theorem 3.2 (Finite-Horizon Optimality). The optimal value function satisfies the recursive relation:

$$V_T^*(s) = 0 \quad \forall s \in \mathcal{S} \quad (3.2)$$

$$V_t^*(s) = \max_{a \in \mathcal{A}} \left[\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) V_{t+1}^*(s') \right] \quad (3.3)$$

for $t = T - 1, T - 2, \dots, 0$.

Proof. The proof follows by backward induction. At time T , no more rewards can be collected, so $V_T^*(s) = 0$. For $t < T$, any optimal policy must choose the action that maximizes immediate reward plus discounted future value, leading to the recursive formulation. \square

3.1.2 Engineering Interpretation

The principle of optimality has direct parallels in engineering optimization:

Optimal Control Example

Consider a spacecraft trajectory optimization problem:

- State: position and velocity $(x, v) \in \mathbb{R}^6$
- Control: thrust vector $u \in \mathbb{R}^3$
- Dynamics: $\dot{x} = v$, $\dot{v} = u/m - \nabla\Phi(x)$ (gravitational field)
- Cost: fuel consumption $\int_0^T \|u(t)\| dt$

The principle of optimality implies that if we have an optimal trajectory from Earth to Mars, then any sub-trajectory (e.g., from lunar orbit to Mars) must also be optimal for the sub-problem.

3.2 Value Iteration: Convergence Analysis

Value iteration is the most fundamental algorithm in dynamic programming, providing a constructive method for computing optimal value functions.

3.2.1 Algorithm Description

Algorithm 3 Value Iteration

Input: MDP $(\mathcal{S}, \mathcal{A}, P, \mathcal{R}, \gamma)$, tolerance $\epsilon > 0$

Output: ϵ -optimal value function V

Initialize $V_0(s)$ arbitrarily for all $s \in \mathcal{S}$

$k \leftarrow 0$

repeat

for each $s \in \mathcal{S}$ **do**

$V_{k+1}(s) \leftarrow \max_{a \in \mathcal{A}} [\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V_k(s')]$

end for

$k \leftarrow k + 1$

until $\|V_k - V_{k-1}\|_\infty < \epsilon(1 - \gamma)/(2\gamma)$

return V_k

3.2.2 Convergence Theory

Theorem 3.3 (Convergence of Value Iteration). For any initial value function V_0 , the value iteration sequence $\{V_k\}_{k=0}^\infty$ defined by $V_{k+1} = T^*V_k$ converges to the unique optimal value function V^* at geometric rate γ .

Specifically:

$$\|V_k - V^*\|_\infty \leq \gamma^k \|V_0 - V^*\|_\infty \quad (3.4)$$

Proof. Since T^* is a γ -contraction in the supremum norm and V^* is the unique fixed point of T^* , the result follows directly from the Banach fixed point theorem. \square

3.2.3 Error Bounds and Stopping Criteria

Theorem 3.4 (Error Bounds for Value Iteration). If $\|V_{k+1} - V_k\|_\infty \leq \delta$, then:

$$\|V_k - V^*\|_\infty \leq \frac{\gamma\delta}{1-\gamma} \quad (3.5)$$

$$\|V_{k+1} - V^*\|_\infty \leq \frac{\delta}{1-\gamma} \quad (3.6)$$

Proof. Using the triangle inequality and contraction property:

$$\|V_k - V^*\|_\infty = \|V_k - T^*V_k + T^*V_k - V^*\|_\infty \quad (3.7)$$

$$\leq \|V_k - T^*V_k\|_\infty + \|T^*V_k - T^*V^*\|_\infty \quad (3.8)$$

$$= \|V_k - V_{k+1}\|_\infty + \gamma\|V_k - V^*\|_\infty \quad (3.9)$$

Solving for $\|V_k - V^*\|_\infty$ gives the first bound. The second follows similarly. \square

Corollary 3.1 (Practical Stopping Criterion). To achieve $\|V_k - V^*\|_\infty \leq \epsilon$, it suffices to stop when:

$$\|V_{k+1} - V_k\|_\infty \leq \epsilon(1-\gamma) \quad (3.10)$$

3.2.4 Computational Complexity

Theorem 3.5 (Sample Complexity of Value Iteration). To achieve ϵ -optimal value function, value iteration requires:

$$O\left(\frac{\log(\epsilon^{-1}) + \log(\|V_0 - V^*\|_\infty)}{1-\gamma}\right) \quad (3.11)$$

iterations.

For each iteration:

- **Time complexity:** $O(|\mathcal{S}|^2|\mathcal{A}|)$ for tabular case
- **Space complexity:** $O(|\mathcal{S}|)$ for storing value function
- **Total operations:** $O(|\mathcal{S}|^2|\mathcal{A}| \log(\epsilon^{-1})/(1-\gamma))$

3.3 Policy Iteration: Mathematical Guarantees

Policy iteration alternates between policy evaluation and policy improvement, providing an alternative approach with different computational characteristics.

3.3.1 Policy Evaluation

Given policy π , policy evaluation solves the linear system:

$$V^\pi = T^\pi V^\pi \quad (3.12)$$

In matrix form for finite MDPs:

$$V^\pi = R^\pi + \gamma P^\pi V^\pi \quad (3.13)$$

where $R^\pi \in \mathbb{R}^{|\mathcal{S}|}$ and $P^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ are policy-specific reward and transition matrices.

Theorem 3.6 (Unique Solution to Policy Evaluation). The linear system $(I - \gamma P^\pi)V^\pi = R^\pi$ has a unique solution:

$$V^\pi = (I - \gamma P^\pi)^{-1} R^\pi \quad (3.14)$$

since $\rho(P^\pi) \leq 1$ and $\gamma < 1$ ensure $(I - \gamma P^\pi)$ is invertible.

3.3.2 Iterative Policy Evaluation

For large state spaces, direct matrix inversion is computationally prohibitive. Iterative policy evaluation uses:

$$V_{k+1}^\pi = T^\pi V_k^\pi \quad (3.15)$$

Theorem 3.7 (Convergence of Iterative Policy Evaluation). The sequence $\{V_k^\pi\}$ converges geometrically to V^π at rate γ :

$$\|V_k^\pi - V^\pi\|_\infty \leq \gamma^k \|V_0^\pi - V^\pi\|_\infty \quad (3.16)$$

3.3.3 Policy Improvement Analysis

Theorem 3.8 (Strict Improvement or Optimality). Given policy π and improved policy π' defined by:

$$\pi'(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^\pi(s, a) \quad (3.17)$$

Then either:

1. $V^{\pi'}(s) > V^\pi(s)$ for some $s \in \mathcal{S}$ (strict improvement), or
2. $V^{\pi'}(s) = V^\pi(s)$ for all $s \in \mathcal{S}$ (optimality)

Proof. By construction, $Q^\pi(s, \pi'(s)) \geq Q^\pi(s, \pi(s)) = V^\pi(s)$ for all s . If inequality is strict

for any state, then by the policy evaluation equations, strict improvement propagates. If equality holds everywhere, then π satisfies the Bellman optimality equation and is optimal. \square

3.3.4 Global Convergence

Theorem 3.9 (Finite Convergence of Policy Iteration). For finite MDPs, policy iteration converges to an optimal policy in finitely many iterations. Specifically, the number of iterations is bounded by $|\mathcal{A}|^{|\mathcal{S}|}$.

Proof. Since each iteration either strictly improves the policy or terminates at optimality, and there are finitely many deterministic policies, convergence must occur in finite time. The bound follows from counting the total number of deterministic policies. \square

3.4 Modified Policy Iteration

Modified policy iteration interpolates between value iteration and policy iteration, providing computational flexibility.

3.4.1 Algorithm and Convergence

Algorithm 4 Modified Policy Iteration

Input: Initial policy π_0 , evaluation steps m

$i \leftarrow 0$

repeat

$V \leftarrow$ arbitrary initialization

for $k = 1, 2, \dots, m$ **do**

$V \leftarrow T^{\pi_i} V$

end for

$\pi_{i+1}(s) \leftarrow \operatorname{argmax}_a [r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s')]$

$i \leftarrow i + 1$

until convergence

Theorem 3.10 (Convergence of Modified Policy Iteration). Modified policy iteration with $m \geq 1$ evaluation steps converges to an optimal policy. The convergence rate depends on m :

- $m = 1$: reduces to value iteration with rate γ
- $m = \infty$: reduces to policy iteration with finite convergence
- $1 < m < \infty$: intermediate convergence rate

3.4.2 Optimal Choice of Evaluation Steps

The computational trade-off between evaluation and improvement can be optimized:

Theorem 3.11 (Optimal Evaluation Steps). For modified policy iteration, the optimal number of evaluation steps m^* minimizes total computational cost:

$$m^* = \underset{m}{\operatorname{argmin}} [\text{cost per iteration} \times \text{number of iterations}] \quad (3.18)$$

Under reasonable assumptions about computational costs, $m^* = O(\log(1/(1 - \gamma)))$.

3.5 Asynchronous Dynamic Programming

Traditional DP algorithms update all states synchronously. Asynchronous variants can offer computational advantages and theoretical insights.

3.5.1 Gauss-Seidel Value Iteration

Algorithm 5 Gauss-Seidel Value Iteration

```

Order states  $s_1, s_2, \dots, s_n$ 
repeat
  for  $i = 1, 2, \dots, n$  do
     $V(s_i) \leftarrow \max_a \left[ r(s_i, a) + \gamma \sum_j P(s_j | s_i, a) V(s_j) \right]$ 
  end for
until convergence
    
```

Theorem 3.12 (Convergence of Gauss-Seidel Value Iteration). Gauss-Seidel value iteration converges to the optimal value function. The convergence rate can be faster than standard value iteration due to more frequent updates.

3.5.2 Prioritized Sweeping

Definition 3.1 (Bellman Error). For state s and value function V , the Bellman error is:

$$\delta(s) = \left| \max_a \left[r(s, a) + \gamma \sum_{s'} P(s' | s, a) V(s') \right] - V(s) \right| \quad (3.19)$$

Prioritized sweeping updates states in order of decreasing Bellman error, focusing computation on states where updates will have the largest impact.

3.5.3 Real-Time Dynamic Programming

Real-time DP focuses updates on states visited by a simulated or actual agent trajectory.

Theorem 3.13 (Convergence of RTDP). Under appropriate exploration conditions, real-time DP converges to optimal values on the states reachable under the optimal policy.

Algorithm 6 Prioritized Sweeping

```

Initialize priority queue  $\mathcal{Q}$  with all states
while  $\mathcal{Q}$  not empty do
     $s \leftarrow$  state with highest priority in  $\mathcal{Q}$ 
    Update  $V(s)$  using Bellman equation
    Remove  $s$  from  $\mathcal{Q}$ 
    for each predecessor  $s'$  of  $s$  do
        if Bellman error of  $s'$  exceeds threshold then
            Add  $s'$  to  $\mathcal{Q}$  with updated priority
        end if
    end for
end while
    
```

Algorithm 7 Real-Time Dynamic Programming

```

Initialize current state  $s$ 
repeat
    Update  $V(s)$  using Bellman equation
    Choose action  $a = \operatorname{argmax}_a Q(s, a)$ 
    Simulate or execute action:  $s \leftarrow s'$  with probability  $P(s'|s, a)$ 
until termination
    
```

3.6 Linear Programming Formulation

Dynamic programming problems can be formulated as linear programs, providing alternative solution methods and theoretical insights.

3.6.1 Primal LP Formulation

The optimal value function can be found by solving:

$$V = \sum_{s \in \mathcal{S}} \alpha(s) V(s) \quad (3.20)$$

$$\text{subject to } V(s) \geq r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s') \quad \forall s, a \quad (3.21)$$

where $\alpha(s) > 0$ represents state weights.

Theorem 3.14 (LP-DP Equivalence). The optimal solution to the linear program equals the optimal value function V^* .

3.6.2 Dual LP Formulation

The dual problem involves finding an optimal state-action visitation measure:

$$\mu = \sum_{s, a} \mu(s, a) r(s, a) \quad (3.22)$$

$$\text{subject to } \sum_a \mu(s, a) - \gamma \sum_{s', a'} \mu(s', a') P(s|s', a') = \alpha(s) \quad \forall s \quad (3.23)$$

$$\mu(s, a) \geq 0 \quad \forall s, a \quad (3.24)$$

Theorem 3.15 (Strong Duality). Under mild conditions, strong duality holds between the primal and dual formulations, and complementary slackness conditions characterize optimal policies.

3.7 Connections to Classical Control Theory

3.7.1 Discrete-Time Optimal Control

Consider the discrete-time optimal control problem:

$$\sum_{t=0}^{T-1} L(x_t, u_t) + L_T(x_T) \quad (3.25)$$

$$\text{subject to } x_{t+1} = f(x_t, u_t) + w_t \quad (3.26)$$

$$u_t \in \mathcal{U}(x_t) \quad (3.27)$$

The dynamic programming solution gives the Hamilton-Jacobi-Bellman equation:

$$V_t(x) = \min_{u \in \mathcal{U}(x)} [L(x, u) + \mathbb{E}[V_{t+1}(f(x, u) + w)]] \quad (3.28)$$

3.7.2 Stochastic Optimal Control

For stochastic control systems $dx_t = f(x_t, u_t)dt + \sigma(x_t, u_t)dW_t$, the continuous-time HJB equation is:

$$\frac{\partial V}{\partial t} + \min_u \left[L(x, u) + \frac{\partial V}{\partial x} f(x, u) + \frac{1}{2} \text{tr} \left(\sigma(x, u)^T \frac{\partial^2 V}{\partial x^2} \sigma(x, u) \right) \right] = 0 \quad (3.29)$$

3.7.3 Model Predictive Control (MPC)

MPC can be viewed as approximate dynamic programming with receding horizon:

Algorithm 8 Model Predictive Control

repeat

 Measure current state x_t

 Solve optimization problem over horizon $[t, t + H]$:

$u_t^*, \dots, u_{t+H-1}^* = \text{argmin} \sum_{k=0}^{H-1} L(x_{t+k}, u_{t+k}) + L_H(x_{t+H})$

 Apply u_t^* and advance to next time step

until termination

The connection to DP provides stability and performance guarantees for MPC under appropriate conditions.

3.8 Computational Considerations

3.8.1 Curse of Dimensionality

The computational complexity of DP algorithms scales exponentially with state space dimension:

- Memory: $O(|\mathcal{S}|)$ for value function storage
- Computation: $O(|\mathcal{S}|^2|\mathcal{A}|)$ per iteration
- For continuous spaces: requires discretization or function approximation

3.8.2 Approximate Dynamic Programming

To handle large state spaces, approximate DP uses function approximation:

$$V(s) \approx \sum_{i=1}^n w_i \phi_i(s) \quad (3.30)$$

where $\{\phi_i\}$ are basis functions and $\{w_i\}$ are parameters.

Theorem 3.16 (Error Propagation in Approximate DP). If the approximation error is bounded by ϵ in supremum norm:

$$\|V - \hat{V}\|_{\infty} \leq \epsilon \quad (3.31)$$

then the policy derived from \hat{V} satisfies:

$$\|V^{\hat{\pi}} - V^*\|_{\infty} \leq \frac{2\gamma\epsilon}{(1-\gamma)^2} \quad (3.32)$$

3.9 Chapter Summary

This chapter developed the mathematical foundations of dynamic programming:

- Principle of optimality and recursive decomposition
- Value iteration: convergence theory, error bounds, complexity analysis
- Policy iteration: linear algebra formulation, finite convergence
- Modified policy iteration and computational trade-offs
- Asynchronous variants: Gauss-Seidel, prioritized sweeping, real-time DP
- Linear programming formulations and duality theory
- Connections to classical optimal control and MPC

- Computational challenges and approximate methods

These algorithmic foundations provide the basis for understanding modern reinforcement learning methods. The next chapter begins our exploration of learning algorithms that estimate value functions from experience rather than exact knowledge of the MDP.

Part II

Core Algorithms and Theory

This part develops the fundamental learning algorithms that form the core of reinforcement learning. Moving beyond the dynamic programming methods of Part I, which assume complete knowledge of the MDP, we now consider algorithms that learn from experience through interaction with the environment.

We begin with Monte Carlo methods that estimate value functions from complete episodes. We then develop temporal difference learning, which enables learning from individual transitions. Finally, we examine Q-learning and SARSA, which learn action-value functions and form the foundation for more advanced algorithms.

Throughout this part, we emphasize mathematical rigor in convergence analysis while maintaining practical relevance for engineering applications. Each algorithm is developed with careful attention to assumptions, convergence conditions, and sample complexity bounds.

Chapter 4

Monte Carlo Methods

Monte Carlo methods form the foundation of model-free reinforcement learning, enabling value function estimation from sample episodes without requiring knowledge of the environment dynamics. This chapter develops the mathematical theory of Monte Carlo estimation in the RL context, with emphasis on convergence analysis and variance reduction techniques.

4.1 Monte Carlo Estimation Theory

Monte Carlo methods estimate expectations by sampling. In reinforcement learning, we use sample episodes to estimate value functions without requiring the transition probabilities or reward function.

4.1.1 Basic Monte Carlo Principle

Consider estimating the expectation $\mathbb{E}[X]$ of random variable X . The Monte Carlo estimator uses n independent samples X_1, \dots, X_n :

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad (4.1)$$

Theorem 4.1 (Strong Law of Large Numbers). If $\mathbb{E}[|X|] < \infty$, then $\hat{\mu}_n \rightarrow \mathbb{E}[X]$ almost surely as $n \rightarrow \infty$.

Theorem 4.2 (Central Limit Theorem). If $\text{Var}(X) = \sigma^2 < \infty$, then:

$$\sqrt{n}(\hat{\mu}_n - \mathbb{E}[X]) \xrightarrow{d} \mathcal{N}(0, \sigma^2) \quad (4.2)$$

4.1.2 Application to Value Function Estimation

For policy π , the value function is:

$$V^\pi(s) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s \right] \quad (4.3)$$

Monte Carlo estimation uses sample returns $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ from episodes starting in state s to estimate $V^{\pi}(s)$.

4.2 First-Visit vs. Every-Visit Methods

4.2.1 First-Visit Monte Carlo

Algorithm 9 First-Visit Monte Carlo Policy Evaluation

Input: Policy π to evaluate

Initialize $V(s) \in \mathbb{R}$ arbitrarily for all $s \in \mathcal{S}$

Initialize $Returns(s) \leftarrow$ empty list for all $s \in \mathcal{S}$

repeat

 Generate episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

for $t = T - 1, T - 2, \dots, 0$ **do**

$G \leftarrow \gamma G + R_{t+1}$

if S_t not appear in S_0, S_1, \dots, S_{t-1} **then**

 Append G to $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

end if

end for

until convergence

4.2.2 Every-Visit Monte Carlo

Every-visit MC updates the value estimate every time a state is visited in an episode, not just the first time.

Theorem 4.3 (Convergence of First-Visit Monte Carlo). First-visit Monte Carlo converges to $V^{\pi}(s)$ as the number of first visits to state s approaches infinity, assuming:

1. Episodes are generated according to policy π
2. Each state has non-zero probability of being the starting state
3. Returns have finite variance

Proof. Each first visit to state s provides an unbiased sample of the return. By the strong law of large numbers, the sample average converges to the true expectation. \square

Theorem 4.4 (Convergence of Every-Visit Monte Carlo). Every-visit Monte Carlo also converges to $V^{\pi}(s)$ under similar conditions, despite the correlation between visits within the same episode.

4.3 Variance Reduction Techniques

4.3.1 Incremental Implementation

Instead of storing all returns, we can update estimates incrementally:

$$V_{n+1}(s) = V_n(s) + \frac{1}{n+1}[G_n - V_n(s)] \quad (4.4)$$

More generally, with step size α :

$$V(s) \leftarrow V(s) + \alpha[G - V(s)] \quad (4.5)$$

4.3.2 Baseline Subtraction

To reduce variance, we can subtract a baseline $b(s)$ that doesn't depend on the action:

$$G_t - b(S_t) \quad (4.6)$$

The optimal baseline that minimizes variance is:

$$b^*(s) = \frac{\mathbb{E}[G_t^2 \mid S_t = s]}{\mathbb{E}[G_t \mid S_t = s]} = \mathbb{E}[G_t \mid S_t = s] = V^\pi(s) \quad (4.7)$$

4.3.3 Control Variates

For correlated random variable Y with known expectation $\mathbb{E}[Y] = \mu_Y$:

$$\hat{\mu}_{CV} = \hat{\mu}_X - c(\hat{\mu}_Y - \mu_Y) \quad (4.8)$$

The optimal coefficient is:

$$c^* = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)} \quad (4.9)$$

4.4 Importance Sampling in RL

Importance sampling enables off-policy learning by weighting samples according to the ratio of target to behavior policy probabilities.

4.4.1 Ordinary Importance Sampling

To estimate $\mathbb{E}_\pi[X]$ using samples from policy μ :

$$\hat{\mu}_{IS} = \frac{1}{n} \sum_{i=1}^n \rho_i X_i \quad (4.10)$$

where $\rho_i = \frac{\pi(A_i|S_i)}{\mu(A_i|S_i)}$ is the importance sampling ratio.

Theorem 4.5 (Unbiasedness of Importance Sampling). $\mathbb{E}[\hat{\mu}_{IS}] = \mathbb{E}_\pi[X]$ if $\mu(a|s) > 0$ whenever $\pi(a|s) > 0$.

4.4.2 Weighted Importance Sampling

To reduce variance when some importance weights are very large:

$$\hat{\mu}_{WIS} = \frac{\sum_{i=1}^n \rho_i X_i}{\sum_{i=1}^n \rho_i} \quad (4.11)$$

Theorem 4.6 (Bias-Variance Tradeoff). Weighted importance sampling is biased but often has lower variance than ordinary importance sampling:

$$\text{Bias}[\hat{\mu}_{WIS}] \neq 0 \text{ (in general)} \quad (4.12)$$

$$\text{Var}[\hat{\mu}_{WIS}] \leq \text{Var}[\hat{\mu}_{IS}] \text{ (typically)} \quad (4.13)$$

4.4.3 Per-Decision Importance Sampling

For episodic tasks, the importance sampling ratio for a complete episode is:

$$\rho_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)} \quad (4.14)$$

This can have very high variance. Per-decision importance sampling uses only the relevant portion of the trajectory.

4.5 Off-Policy Monte Carlo Methods

4.5.1 Off-Policy Policy Evaluation

4.5.2 Off-Policy Monte Carlo Control

4.6 Convergence Analysis and Sample Complexity

4.6.1 Finite Sample Analysis

Theorem 4.7 (Finite Sample Bound for Monte Carlo). Let $V_n(s)$ be the Monte Carlo estimate after n visits to state s . Under bounded rewards $|R| \leq R_{max}$:

$$\mathbb{P}(|V_n(s) - V^\pi(s)| \geq \epsilon) \leq 2 \exp\left(-\frac{2n\epsilon^2(1-\gamma)^2}{R_{max}^2}\right) \quad (4.15)$$

Algorithm 10 Off-Policy Monte Carlo Policy Evaluation

Input: Target policy π , behavior policy μ
Initialize $V(s) \in \mathbb{R}$ arbitrarily for all $s \in \mathcal{S}$
Initialize $C(s) \leftarrow 0$ for all $s \in \mathcal{S}$
repeat
 Generate episode using μ : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
 for $t = T - 1, T - 2, \dots, 0$ **do**
 $G \leftarrow \gamma G + R_{t+1}$
 $C(S_t) \leftarrow C(S_t) + W$
 $V(S_t) \leftarrow V(S_t) + \frac{W}{C(S_t)}[G - V(S_t)]$
 $W \leftarrow W \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$
 if $W = 0$ **then**
 break
 end if
 end for
until convergence

Algorithm 11 Off-Policy Monte Carlo Control

Initialize $Q(s, a) \in \mathbb{R}$ arbitrarily for all s, a
Initialize $C(s, a) \leftarrow 0$ for all s, a
Initialize $\pi(s) \leftarrow \operatorname{argmax}_a Q(s, a)$ for all s
repeat
 Choose any soft policy μ (e.g., ϵ -greedy)
 Generate episode using μ
 $G \leftarrow 0$
 $W \leftarrow 1$
 for $t = T - 1, T - 2, \dots, 0$ **do**
 $G \leftarrow \gamma G + R_{t+1}$
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)}[G - Q(S_t, A_t)]$
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$
 if $A_t \neq \pi(S_t)$ **then**
 break
 end if
 $W \leftarrow W \frac{1}{\mu(A_t|S_t)}$
 end for
until convergence

4.6.2 Asymptotic Convergence Rate

Theorem 4.8 (Central Limit Theorem for Monte Carlo). If $\text{Var}^\pi[G_t|S_t = s] = \sigma^2(s) < \infty$, then:

$$\sqrt{n}(V_n(s) - V^\pi(s)) \xrightarrow{d} \mathcal{N}(0, \sigma^2(s)) \quad (4.16)$$

This gives the convergence rate $O(n^{-1/2})$, which is slower than the $O(n^{-1})$ rate achievable by temporal difference methods under certain conditions.

4.6.3 Sample Complexity

Theorem 4.9 (Sample Complexity of Monte Carlo). To achieve ϵ -accurate value function estimation with probability $1 - \delta$:

$$n \geq \frac{R_{\max}^2 \log(2/\delta)}{2\epsilon^2(1 - \gamma)^2} \quad (4.17)$$

samples are sufficient.

4.7 Practical Considerations

4.7.1 Exploration vs. Exploitation

Monte Carlo control methods face the exploration-exploitation dilemma. Common approaches:

Exploring Starts: Assume episodes start in randomly selected state-action pairs.

ϵ -Greedy Policies: Use soft policies that maintain exploration:

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|} & \text{if } a = \arg\max_a Q(s, a) \\ \frac{\epsilon}{|\mathcal{A}(s)|} & \text{otherwise} \end{cases} \quad (4.18)$$

4.7.2 Function Approximation

For large state spaces, we approximate value functions:

$$V(s) \approx \hat{V}(s, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(s) \quad (4.19)$$

The Monte Carlo update becomes:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha[G_t - \hat{V}(S_t, \mathbf{w})]\nabla_{\mathbf{w}}\hat{V}(S_t, \mathbf{w}) \quad (4.20)$$

Theorem 4.10 (Convergence with Linear Function Approximation). Under linear function approximation with linearly independent features, Monte Carlo methods converge to the best linear approximation in the L^2 norm weighted by the stationary distribution.

4.8 Chapter Summary

This chapter established the foundations of Monte Carlo methods in reinforcement learning:

- Monte Carlo estimation theory and convergence properties
- First-visit vs. every-visit methods with convergence guarantees
- Variance reduction techniques: baselines, control variates, importance sampling
- Off-policy learning through importance sampling with bias-variance analysis
- Sample complexity bounds and convergence rates
- Practical considerations for exploration and function approximation

Monte Carlo methods provide unbiased estimates and are conceptually simple, but they require complete episodes and have slower convergence than temporal difference methods. The next chapter develops temporal difference learning, which enables learning from individual transitions.

Chapter 5

Temporal Difference Learning

Temporal Difference (TD) learning combines ideas from Monte Carlo methods and dynamic programming, enabling learning from incomplete episodes while maintaining the model-free nature of Monte Carlo methods. This chapter develops the mathematical theory of TD learning with emphasis on convergence analysis and the fundamental bias-variance tradeoff.

5.1 TD(0) Algorithm and Mathematical Analysis

5.1.1 Basic TD(0) Update

The core insight of temporal difference learning is to use the current estimate of the successor state's value to update the current state's value:

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \quad (5.1)$$

The TD error is defined as:

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \quad (5.2)$$

5.1.2 Relationship to Bellman Equation

The TD(0) update can be viewed as a stochastic approximation to the Bellman equation. The expected TD update is:

$$\mathbb{E}[\delta_t | S_t = s] = \mathbb{E}[R_{t+1} + \gamma V(S_{t+1}) - V(S_t) | S_t = s] \quad (5.3)$$

$$= \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma V(s') - V(s)] \quad (5.4)$$

$$= (T^\pi V)(s) - V(s) \quad (5.5)$$

Algorithm 12 Tabular TD(0) Policy Evaluation

Input: Policy π to evaluate, step size $\alpha \in (0, 1]$
 Initialize $V(s) \in \mathbb{R}$ arbitrarily for all $s \in \mathcal{S}$, except $V(\text{terminal}) = 0$
repeat
 Initialize S
 repeat
 $A \leftarrow$ action given by π for S
 Take action A , observe R, S'
 $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$
 $S \leftarrow S'$
 until S is terminal
until convergence or sufficient accuracy

where T^π is the Bellman operator for policy π .

5.1.3 Convergence Analysis

Theorem 5.1 (Convergence of TD(0) - Tabular Case). For the tabular case with appropriate step size sequence $\{\alpha_t\}$ satisfying:

$$\sum_{t=0}^{\infty} \alpha_t = \infty \quad (5.6)$$

$$\sum_{t=0}^{\infty} \alpha_t^2 < \infty \quad (5.7)$$

TD(0) converges to V^π with probability 1.

Proof Sketch. The proof uses stochastic approximation theory. Define the ODE:

$$\frac{dV}{dt} = \mathbb{E}[\delta_t | V] = T^\pi V - V \quad (5.8)$$

Since T^π is a contraction, the unique fixed point is V^π . The stochastic approximation theorem ensures convergence of the discrete updates to the ODE solution. \square

5.2 Bias-Variance Tradeoff in TD Methods

5.2.1 Bias Analysis

TD(0) uses the biased estimate $R_{t+1} + \gamma V(S_{t+1})$ as a target for $V(S_t)$, while Monte Carlo uses the unbiased estimate G_t .

Theorem 5.2 (Bias of TD Target). The TD target $R_{t+1} + \gamma V(S_{t+1})$ has bias:

$$\text{Bias}[R_{t+1} + \gamma V(S_{t+1})] = \gamma[\hat{V}(S_{t+1}) - V^\pi(S_{t+1})] \quad (5.9)$$

where \hat{V} is the current estimate.

5.2.2 Variance Analysis

Theorem 5.3 (Variance Comparison). Under the assumption that value function errors are small, the variance of the TD target is approximately:

$$\text{Var}[R_{t+1} + \gamma V(S_{t+1})] \approx \text{Var}[R_{t+1}] \quad (5.10)$$

while the Monte Carlo target has variance:

$$\text{Var}[G_t] = \text{Var} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right] \quad (5.11)$$

which is typically much larger.

5.2.3 Mean Squared Error Decomposition

$$\text{MSE} = \text{Bias}^2 + \text{Variance} + \text{Noise} \quad (5.12)$$

TD methods trade increased bias for reduced variance, often resulting in lower overall MSE and faster convergence.

5.3 TD() and Eligibility Traces

TD() provides a family of algorithms that interpolate between TD(0) and Monte Carlo methods through the use of eligibility traces.

5.3.1 Forward View: n-step Returns

The n-step return combines rewards from the next n steps with the estimated value of the state reached after n steps:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}) \quad (5.13)$$

The n-step TD update is:

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t^{(n)} - V(S_t)] \quad (5.14)$$

5.3.2 λ -Return

The λ -return combines all n-step returns:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \quad (5.15)$$

Theorem 5.4 (-Return Properties). The γ -return satisfies:

$$G_t^\lambda = R_{t+1} + \gamma[(1 - \lambda)V(S_{t+1}) + \lambda G_{t+1}^\lambda] \quad (5.16)$$

$$\lim_{\lambda \rightarrow 0} G_t^\lambda = R_{t+1} + \gamma V(S_{t+1}) \quad (\text{TD}(0)) \quad (5.17)$$

$$\lim_{\lambda \rightarrow 1} G_t^\lambda = G_t \quad (\text{Monte Carlo}) \quad (5.18)$$

5.3.3 Backward View: Eligibility Traces

Eligibility traces provide an online, incremental implementation of TD():

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \quad (5.19)$$

$$e_t(s) = \begin{cases} \gamma \lambda e_{t-1}(s) + 1 & \text{if } s = S_t \\ \gamma \lambda e_{t-1}(s) & \text{if } s \neq S_t \end{cases} \quad (5.20)$$

$$V(s) \leftarrow V(s) + \alpha \delta_t e_t(s) \quad \forall s \quad (5.21)$$

Algorithm 13 TD() with Eligibility Traces

Input: Policy π , step size α , trace decay λ

Initialize $V(s) \in \mathbb{R}$ arbitrarily for all s

repeat

Initialize S , $e(s) = 0$ for all s

repeat

$A \leftarrow$ action given by π for S

Take action A , observe R, S'

$\delta \leftarrow R + \gamma V(S') - V(S)$

$e(S) \leftarrow e(S) + 1$

for all s **do**

$V(s) \leftarrow V(s) + \alpha \delta e(s)$

$e(s) \leftarrow \gamma \lambda e(s)$

end for

$S \leftarrow S'$

until S is terminal

until convergence

5.3.4 Equivalence Theorem

Theorem 5.5 (Forward-Backward Equivalence). Under certain conditions, the forward view (using γ -returns) and backward view (using eligibility traces) produce identical updates when applied offline to a complete episode.

5.4 Convergence Theory for Linear Function Approximation

When the state space is large, we use function approximation:

$$V(s) \approx \hat{V}(s, \mathbf{w}) = \mathbf{w}^T \phi(s) \quad (5.22)$$

The TD(0) update becomes:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [R_{t+1} + \gamma \mathbf{w}_t^T \phi(S_{t+1}) - \mathbf{w}_t^T \phi(S_t)] \phi(S_t) \quad (5.23)$$

5.4.1 Projected Bellman Equation

Under linear function approximation, TD(0) converges to the solution of the projected Bellman equation:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - T^\pi(\Phi \mathbf{w})\|_{\mathbf{D}}^2 \quad (5.24)$$

where Φ is the feature matrix and \mathbf{D} is a diagonal matrix of state visitation probabilities.

Theorem 5.6 (Convergence of Linear TD(0)). Under linear function approximation, TD(0) converges to:

$$\mathbf{w}^* = (\Phi^T \mathbf{D} \Phi)^{-1} \Phi^T \mathbf{D} \mathbf{r}^\pi \quad (5.25)$$

where \mathbf{r}^π is the expected reward vector.

5.4.2 Error Bounds

Theorem 5.7 (Approximation Error Bound). Let V^* be the optimal value function and \hat{V}^* be the best linear approximation. Then:

$$\|V^\pi - \hat{V}^\pi\|_{\mathbf{D}} \leq \frac{1}{1 - \gamma} \min_{\mathbf{w}} \|V^\pi - \Phi \mathbf{w}\|_{\mathbf{D}} \quad (5.26)$$

5.5 Comparison with Monte Carlo and DP Methods

5.5.1 Computational Complexity

Method	Memory	Computation per Step	Episode Completion
DP	$O(\mathcal{S} ^2 \mathcal{A})$	$O(\mathcal{S} ^2 \mathcal{A})$	Not Required
MC	$O(\mathcal{S})$	$O(1)$	Required
TD	$O(\mathcal{S})$	$O(1)$	Not Required

5.5.2 Sample Efficiency

Theorem 5.8 (Sample Complexity Comparison). Under certain regularity conditions:

- TD methods: $O(\frac{1}{\epsilon^2(1-\gamma)^2})$ samples for ϵ -accuracy

- MC methods: $O(\frac{1}{\epsilon^2(1-\gamma)^4})$ samples for ϵ -accuracy

TD methods often have better sample efficiency due to lower variance, despite being biased.

5.5.3 Bootstrapping vs. Sampling

Bootstrapping: Using estimates of successor states (DP, TD) **Sampling:** Using actual experience (MC, TD)

TD methods combine both, leading to:

- Faster learning than MC (bootstrapping)
- Model-free nature (sampling)
- Online learning capability

5.6 Advanced Topics

5.6.1 Multi-step Methods

The n-step TD methods generalize between TD(0) and Monte Carlo:

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t^{(n)} - V(S_t)] \quad (5.27)$$

Theorem 5.9 (Optimal Step Size). For n-step methods, there exists an optimal n that minimizes mean squared error, typically $n \in [3, 10]$ for many problems.

5.6.2 True Online TD()

The classical TD() is not equivalent to the forward view when using function approximation. True online TD() corrects this:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \delta_t \mathbf{z}_t + \alpha (\mathbf{w}_t^T \phi_t - \mathbf{w}_{t-1}^T \phi_t) (\mathbf{z}_t - \phi_t) \quad (5.28)$$

$$\mathbf{z}_{t+1} = \gamma \lambda \mathbf{z}_t + \phi_{t+1} - \alpha \gamma \lambda (\mathbf{z}_t^T \phi_{t+1}) \phi_{t+1} \quad (5.29)$$

5.6.3 Gradient TD Methods

To handle function approximation more rigorously, gradient TD methods minimize the mean squared projected Bellman error:

$$\text{MSPBE}(\mathbf{w}) = \|\Pi(\mathbf{T}^\pi \hat{\mathbf{v}} - \hat{\mathbf{v}})\|_{\mathbf{D}}^2 \quad (5.30)$$

$$\nabla \text{MSPBE}(\mathbf{w}) = 2\Phi^T \mathbf{D}(\Pi(\mathbf{T}^\pi \hat{\mathbf{v}} - \hat{\mathbf{v}})) \quad (5.31)$$

5.7 Chapter Summary

This chapter developed the mathematical foundations of temporal difference learning:

- TD(0) algorithm with convergence analysis using stochastic approximation theory
- Bias-variance tradeoff analysis showing TD's advantage in variance reduction
- TD() and eligibility traces providing a spectrum between TD(0) and Monte Carlo
- Convergence theory for linear function approximation with error bounds
- Comparative analysis with Monte Carlo and dynamic programming methods
- Advanced topics including multi-step methods and gradient TD approaches

Temporal difference learning provides the foundation for many modern RL algorithms, combining the best aspects of Monte Carlo and dynamic programming approaches. The next chapter extends these ideas to action-value methods with Q-learning and SARSA.

Chapter 6

Chapter 06 Title

This chapter will cover...

Part III

Function Approximation and Deep Learning

This part bridges classical reinforcement learning with modern deep learning approaches. When state and action spaces are large or continuous, exact representation of value functions becomes intractable. Function approximation provides the mathematical framework for representing value functions compactly while maintaining theoretical guarantees.

We begin with linear function approximation, which provides strong theoretical foundations and convergence guarantees. We then explore the integration of neural networks into reinforcement learning, leading to the deep reinforcement learning revolution. Finally, we develop policy gradient methods that directly optimize parameterized policies.

The treatment emphasizes both the mathematical foundations that ensure stability and convergence, and the practical considerations necessary for successful implementation in engineering systems.

Chapter 7

Chapter 07 Title

This chapter will cover...

Chapter 8

Chapter 08 Title

This chapter will cover...

Chapter 9

Chapter 09 Title

This chapter will cover...

Part IV

Advanced Topics

This part explores advanced reinforcement learning topics that extend beyond the basic framework to handle more complex scenarios. We examine continuous control problems that arise frequently in engineering applications, multi-agent systems where multiple learners interact, model-based approaches that learn environment dynamics, and the fundamental exploration-exploitation tradeoff.

These advanced topics require sophisticated mathematical tools and careful analysis of convergence properties, stability, and sample complexity. The treatment provides both theoretical insights and practical guidance for tackling complex real-world problems.

Chapter 10

Chapter 10 Title

This chapter will cover...

Chapter 11

Chapter 11 Title

This chapter will cover...

Chapter 12

Chapter 12 Title

This chapter will cover...

Chapter 13

Chapter 13 Title

This chapter will cover...

Part V

Implementation and Practice

This part focuses on the practical aspects of implementing and deploying reinforcement learning systems. Moving from theory to practice requires careful consideration of computational efficiency, software engineering best practices, validation methodologies, and deployment considerations.

We examine computational techniques for scaling RL algorithms, software frameworks and tools, and methodologies for validating and deploying RL systems in production environments. The treatment emphasizes engineering best practices while maintaining theoretical rigor.

Chapter 14

Chapter 14 Title

This chapter will cover...

Chapter 15

Chapter 15 Title

This chapter will cover...

Chapter 16

Chapter 16 Title

This chapter will cover...

Part VI

Future Directions

This final part explores emerging paradigms and future directions in reinforcement learning. We examine meta-learning approaches that enable rapid adaptation to new tasks, integration with other fields such as classical optimization and control theory, and discuss open challenges and future research directions.

The treatment provides perspective on the current state of the field and identifies promising directions for future research and development, particularly relevant for engineer-mathematicians working at the intersection of theory and practice.

Chapter 17

Chapter 17 Title

This chapter will cover...

Chapter 18

Chapter 18 Title

This chapter will cover...

Mathematical Reference

This appendix provides a comprehensive mathematical reference for concepts used throughout the book. It serves as a quick reference for key mathematical tools and theorems.

.1 Matrix Calculus for RL

.1.1 Gradients and Jacobians

For scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the gradient is:

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \quad (1)$$

For vector function $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$, the Jacobian is:

$$J_F(x) = \begin{bmatrix} \frac{\partial F_1}{\partial x_1} & \cdots & \frac{\partial F_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial F_m}{\partial x_1} & \cdots & \frac{\partial F_m}{\partial x_n} \end{bmatrix} \quad (2)$$

.1.2 Chain Rule

For composite functions $h(x) = f(g(x))$:

$$\nabla h(x) = J_g(x)^T \nabla f(g(x)) \quad (3)$$

.1.3 Common Derivatives

$$\frac{\partial}{\partial x} x^T A x = (A + A^T)x \quad (4)$$

$$\frac{\partial}{\partial x} a^T x = a \quad (5)$$

$$\frac{\partial}{\partial X} \text{tr}(AXB) = A^T B^T \quad (6)$$

$$\frac{\partial}{\partial X} \log \det(X) = (X^{-1})^T \quad (7)$$

.2 Probability Distributions Commonly Used

.2.1 Discrete Distributions

Bernoulli Distribution: $X \sim \text{Bernoulli}(p)$

$$P(X = 1) = p, \quad P(X = 0) = 1 - p \quad (8)$$

$$\mathbb{E}[X] = p, \quad \text{Var}(X) = p(1 - p) \quad (9)$$

Categorical Distribution: $X \sim \text{Categorical}(\mathbf{p})$

$$P(X = k) = p_k, \quad \sum_{k=1}^K p_k = 1 \quad (10)$$

$$\mathbb{E}[X] = \sum_{k=1}^K k p_k \quad (11)$$

.2.2 Continuous Distributions

Normal Distribution: $X \sim \mathcal{N}(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (12)$$

$$\mathbb{E}[X] = \mu, \quad \text{Var}(X) = \sigma^2 \quad (13)$$

Multivariate Normal: $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (14)$$

$$\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}, \quad \text{Cov}(\mathbf{X}) = \Sigma \quad (15)$$

.3 Optimization Algorithms Summary

.3.1 Gradient Descent Variants

Vanilla Gradient Descent:

$$\theta_{t+1} = \theta_t - \alpha \nabla f(\theta_t) \quad (16)$$

Momentum:

$$v_{t+1} = \beta v_t + \nabla f(\theta_t) \quad (17)$$

$$\theta_{t+1} = \theta_t - \alpha v_{t+1} \quad (18)$$

Adam:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla f(\theta_t) \quad (19)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) [\nabla f(\theta_t)]^2 \quad (20)$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (21)$$

$$\theta_{t+1} = \theta_t - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad (22)$$

.4 Convergence Analysis Techniques

.4.1 Lyapunov Functions

A function $V : \mathcal{X} \rightarrow \mathbb{R}_+$ is a Lyapunov function for dynamical system $x_{t+1} = f(x_t)$ if:

1. $V(x) > 0$ for $x \neq x^*$ and $V(x^*) = 0$
2. $V(f(x)) - V(x) \leq 0$ for all x

.4.2 Martingale Convergence Theorem

Theorem .1 (Martingale Convergence). Let $\{X_t\}$ be a supermartingale that is bounded below. Then X_t converges almost surely to a finite random variable.

.4.3 Robbins-Monro Conditions

For stochastic approximation algorithm $\theta_{t+1} = \theta_t + \alpha_t H(\theta_t, \xi_t)$, convergence occurs under:

$$\sum_{t=0}^{\infty} \alpha_t = \infty \quad (23)$$

$$\sum_{t=0}^{\infty} \alpha_t^2 < \infty \quad (24)$$

$$\mathbb{E}[H(\theta, \xi)] = h(\theta) \text{ has unique zero at } \theta^* \quad (25)$$

Implementation Templates

This appendix provides implementation templates and code examples for key reinforcement learning algorithms. The code is provided in Python and follows best practices for numerical stability and computational efficiency.

.5 Basic RL Algorithm Implementations

.5.1 Value Iteration

```
1 import numpy as np
2
3 def value_iteration(P, R, gamma, tol=1e-6, max_iter=1000):
4     """
5     Value iteration algorithm for finite MDPs.
6
7     Parameters:
8     P: transition probability tensor [S x A x S]
9     R: reward matrix [S x A]
10    gamma: discount factor
11    tol: convergence tolerance
12    max_iter: maximum iterations
13
14    Returns:
15    V: optimal value function
16    policy: optimal policy
17    """
18    S, A = R.shape
19    V = np.zeros(S)
20
21    for i in range(max_iter):
22        V_old = V.copy()
23
24        # Bellman optimality operator
25        Q = R + gamma * np.sum(P * V[None, None, :], axis=2)
26        V = np.max(Q, axis=1)
27
28        # Check convergence
29        if np.max(np.abs(V - V_old)) < tol:
30            break
```

```
31
32     # Extract optimal policy
33     Q = R + gamma * np.sum(P * V[None, None, :], axis=2)
34     policy = np.argmax(Q, axis=1)
35
36     return V, policy
```

Listing 1: Value Iteration Implementation

.5.2 Q-Learning

```
1  import numpy as np
2  from collections import defaultdict
3
4  class QLearning:
5      def __init__(self, n_states, n_actions, alpha=0.1, gamma=0.99, epsilon
        =0.1):
6          self.n_states = n_states
7          self.n_actions = n_actions
8          self.alpha = alpha
9          self.gamma = gamma
10         self.epsilon = epsilon
11         self.Q = np.zeros((n_states, n_actions))
12
13     def select_action(self, state):
14         """Epsilon-greedy action selection"""
15         if np.random.random() < self.epsilon:
16             return np.random.randint(self.n_actions)
17         else:
18             return np.argmax(self.Q[state])
19
20     def update(self, state, action, reward, next_state, done):
21         """Q-learning update rule"""
22         if done:
23             target = reward
24         else:
25             target = reward + self.gamma * np.max(self.Q[next_state])
26
27         td_error = target - self.Q[state, action]
28         self.Q[state, action] += self.alpha * td_error
29
30     def return td_error
31
32     def get_policy(self):
33         """Extract greedy policy"""
34         return np.argmax(self.Q, axis=1)
```

Listing 2: Q-Learning Implementation

.5.3 Policy Gradient (REINFORCE)

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.nn.functional as F
5 from torch.distributions import Categorical
6
7 class PolicyNetwork(nn.Module):
8     def __init__(self, state_dim, action_dim, hidden_dim=128):
9         super(PolicyNetwork, self).__init__()
10        self.fc1 = nn.Linear(state_dim, hidden_dim)
11        self.fc2 = nn.Linear(hidden_dim, hidden_dim)
12        self.fc3 = nn.Linear(hidden_dim, action_dim)
13
14        def forward(self, x):
15            x = F.relu(self.fc1(x))
16            x = F.relu(self.fc2(x))
17            x = F.softmax(self.fc3(x), dim=-1)
18            return x
19
20 class REINFORCE:
21     def __init__(self, state_dim, action_dim, lr=1e-3, gamma=0.99):
22         self.policy = PolicyNetwork(state_dim, action_dim)
23         self.optimizer = optim.Adam(self.policy.parameters(), lr=lr)
24         self.gamma = gamma
25
26     def select_action(self, state):
27         state = torch.FloatTensor(state).unsqueeze(0)
28         probs = self.policy(state)
29         m = Categorical(probs)
30         action = m.sample()
31         return action.item(), m.log_prob(action)
32
33     def update(self, log_probs, rewards):
34         """Update policy using REINFORCE algorithm"""
35         # Compute discounted returns
36         returns = []
37         G = 0
38         for r in reversed(rewards):
39             G = r + self.gamma * G
40             returns.insert(0, G)
41
42         returns = torch.FloatTensor(returns)
43         returns = (returns - returns.mean()) / (returns.std() + 1e-8)
44
45         # Compute policy loss
46         policy_loss = []
47         for log_prob, G in zip(log_probs, returns):
48             policy_loss.append(-log_prob * G)
49
50         policy_loss = torch.stack(policy_loss).sum()

```

```
51
52     # Update policy
53     self.optimizer.zero_grad()
54     policy_loss.backward()
55     self.optimizer.step()
56
57     return policy_loss.item()
```

Listing 3: REINFORCE Implementation

.6 Environment Interface Specifications

.6.1 OpenAI Gym Compatible Environment

```
1  import gym
2  from gym import spaces
3  import numpy as np
4
5  class CustomEnvironment(gym.Env):
6      """Template for custom RL environment"""
7
8      def __init__(self):
9          super(CustomEnvironment, self).__init__()
10
11         # Define action and observation spaces
12         self.action_space = spaces.Discrete(4) # Example: 4 discrete
13         # actions
14         self.observation_space = spaces.Box(
15             low=-np.inf, high=np.inf, shape=(4,), dtype=np.float32
16         )
17
18         # Initialize state
19         self.state = None
20         self.episode_length = 0
21         self.max_episode_length = 1000
22
23     def reset(self):
24         """Reset environment to initial state"""
25         self.state = self._get_initial_state()
26         self.episode_length = 0
27         return self.state
28
29     def step(self, action):
30         """Execute one step in the environment"""
31         # Validate action
32         assert self.action_space.contains(action), f"Invalid action: {
33             action}"
34
35         # Update state based on action
36         self.state = self._update_state(self.state, action)
```

```

35
36     # Compute reward
37     reward = self._compute_reward(self.state, action)
38
39     # Check if episode is done
40     done = self._is_done()
41
42     # Additional info
43     info = {}
44
45     self.episode_length += 1
46
47     return self.state, reward, done, info
48
49     def render(self, mode='human'):
50         """Render the environment"""
51         print(f"State: {self.state}")
52
53     def _get_initial_state(self):
54         """Get initial state (implement based on problem)"""
55         return np.random.randn(4)
56
57     def _update_state(self, state, action):
58         """Update state based on action (implement based on problem)"""
59         # Example: simple dynamics
60         next_state = state + 0.1 * action
61         return next_state
62
63     def _compute_reward(self, state, action):
64         """Compute reward (implement based on problem)"""
65         # Example: negative squared distance from origin
66         return -np.sum(state**2)
67
68     def _is_done(self):
69         """Check if episode should terminate"""
70         return self.episode_length >= self.max_episode_length

```

Listing 4: Custom Environment Template

7 Logging and Visualization Code

7.1 Training Logger

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from collections import deque
4 import json
5
6 class TrainingLogger:
7     def __init__(self, window_size=100):

```

```

8         self.metrics = {}
9         self.window_size = window_size
10        self.episode_rewards = deque(maxlen=window_size)
11        self.episode_lengths = deque(maxlen=window_size)
12
13    def log_episode(self, episode, reward, length, **kwargs):
14        """Log episode statistics"""
15        self.episode_rewards.append(reward)
16        self.episode_lengths.append(length)
17
18        # Log additional metrics
19        for key, value in kwargs.items():
20            if key not in self.metrics:
21                self.metrics[key] = deque(maxlen=self.window_size)
22                self.metrics[key].append(value)
23
24        # Print progress
25        if episode % 100 == 0:
26            avg_reward = np.mean(self.episode_rewards)
27            avg_length = np.mean(self.episode_lengths)
28            print(f"Episode_{episode}: Avg_Reward={avg_reward:.2f}, "
29                  f"Avg_Length={avg_length:.2f}")
30
31    def plot_training_curves(self, save_path=None):
32        """Plot training curves"""
33        fig, axes = plt.subplots(2, 2, figsize=(12, 8))
34
35        # Episode rewards
36        axes[0, 0].plot(self.episode_rewards)
37        axes[0, 0].set_title('Episode_Rewards')
38        axes[0, 0].set_xlabel('Episode')
39        axes[0, 0].set_ylabel('Reward')
40
41        # Episode lengths
42        axes[0, 1].plot(self.episode_lengths)
43        axes[0, 1].set_title('Episode_Lengths')
44        axes[0, 1].set_xlabel('Episode')
45        axes[0, 1].set_ylabel('Length')
46
47        # Moving averages
48        if len(self.episode_rewards) > 10:
49            window = min(50, len(self.episode_rewards) // 4)
50            moving_avg = np.convolve(self.episode_rewards,
51                                     np.ones(window)/window, mode='valid')
52            axes[1, 0].plot(moving_avg)
53            axes[1, 0].set_title(f'Moving_Average_Reward_(window={window})')
54            axes[1, 0].set_xlabel('Episode')
55            axes[1, 0].set_ylabel('Reward')
56
57        # Additional metrics

```

```

58         if self.metrics:
59             metric_name = list(self.metrics.keys())[0]
60             axes[1, 1].plot(self.metrics[metric_name])
61             axes[1, 1].set_title(metric_name)
62             axes[1, 1].set_xlabel('Episode')
63             axes[1, 1].set_ylabel('Value')
64
65         plt.tight_layout()
66
67         if save_path:
68             plt.savefig(save_path)
69         plt.show()
70
71     def save_metrics(self, filepath):
72         """Save metrics to JSON file"""
73         data = {
74             'episode_rewards': list(self.episode_rewards),
75             'episode_lengths': list(self.episode_lengths),
76             'metrics': {k: list(v) for k, v in self.metrics.items()}
77         }
78         with open(filepath, 'w') as f:
79             json.dump(data, f, indent=2)

```

Listing 5: Training Logger

.8 Performance Benchmarking Utilities

.8.1 Algorithm Comparison Framework

```

1  import time
2  import numpy as np
3  from typing import Dict, List, Callable
4
5  class AlgorithmComparison:
6      def __init__(self, environment_factory: Callable):
7          self.environment_factory = environment_factory
8          self.results = {}
9
10     def run_algorithm(self, algorithm_name: str, algorithm_class,
11                      algorithm_params: Dict, n_runs: int = 5,
12                      n_episodes: int = 1000):
13         """Run algorithm multiple times and collect statistics"""
14         print(f"Running {algorithm_name}...")
15
16         run_results = []
17
18         for run in range(n_runs):
19             print(f"Run {run+1}/{n_runs}")
20
21             # Create fresh environment and algorithm

```

```
22     env = self.environment_factory()
23     algorithm = algorithm_class(**algorithm_params)
24
25     # Track performance
26     episode_rewards = []
27     episode_lengths = []
28     start_time = time.time()
29
30     for episode in range(n_episodes):
31         state = env.reset()
32         episode_reward = 0
33         episode_length = 0
34         done = False
35
36         while not done:
37             action = algorithm.select_action(state)
38             next_state, reward, done, _ = env.step(action)
39
40             # Update algorithm
41             if hasattr(algorithm, 'update'):
42                 algorithm.update(state, action, reward, next_state,
43                                 done)
44
45             state = next_state
46             episode_reward += reward
47             episode_length += 1
48
49             episode_rewards.append(episode_reward)
50             episode_lengths.append(episode_length)
51
52         training_time = time.time() - start_time
53
54         run_results.append({
55             'episode_rewards': episode_rewards,
56             'episode_lengths': episode_lengths,
57             'training_time': training_time,
58             'final_performance': np.mean(episode_rewards[-100:])
59         })
60
61     self.results[algorithm_name] = run_results
62
63     def compare_algorithms(self):
64         """Generate comparison statistics"""
65         comparison = {}
66
67         for alg_name, results in self.results.items():
68             final_perfs = [r['final_performance'] for r in results]
69             training_times = [r['training_time'] for r in results]
70
71             comparison[alg_name] = {
72                 'mean_performance': np.mean(final_perfs),
```



```

72         'std_performance': np.std(final_perfs),
73         'mean_training_time': np.mean(training_times),
74         'std_training_time': np.std(training_times)
75     }
76
77     return comparison
78
79 def plot_comparison(self):
80     """Plot algorithm comparison"""
81     plt.figure(figsize=(15, 5))
82
83     # Performance comparison
84     plt.subplot(1, 3, 1)
85     alg_names = list(self.results.keys())
86     performances = []
87     errors = []
88
89     for alg_name in alg_names:
90         final_perfs = [r['final_performance'] for r in self.results[
91             alg_name]]
92         performances.append(np.mean(final_perfs))
93         errors.append(np.std(final_perfs))
94
95     plt.bar(alg_names, performances, yerr=errors, capsize=5)
96     plt.title('Final Performance Comparison')
97     plt.ylabel('Average Return')
98     plt.xticks(rotation=45)
99
100    # Learning curves
101    plt.subplot(1, 3, 2)
102    for alg_name in alg_names:
103        all_rewards = []
104        for result in self.results[alg_name]:
105            all_rewards.append(result['episode_rewards'])
106
107        mean_rewards = np.mean(all_rewards, axis=0)
108        std_rewards = np.std(all_rewards, axis=0)
109        episodes = np.arange(len(mean_rewards))
110
111        plt.plot(episodes, mean_rewards, label=alg_name)
112        plt.fill_between(episodes, mean_rewards - std_rewards,
113                        mean_rewards + std_rewards, alpha=0.3)
114
115    plt.title('Learning Curves')
116    plt.xlabel('Episode')
117    plt.ylabel('Episode Reward')
118    plt.legend()
119
120    # Training time comparison
121    plt.subplot(1, 3, 3)
122    times = []

```

```
122     time_errors = []
123
124     for alg_name in alg_names:
125         training_times = [r['training_time'] for r in self.results[
126             alg_name]]
127         times.append(np.mean(training_times))
128         time_errors.append(np.std(training_times))
129
130     plt.bar(alg_names, times, yerr=time_errors, capsize=5)
131     plt.title('Training Time Comparison')
132     plt.ylabel('Time (seconds)')
133     plt.xticks(rotation=45)
134
135     plt.tight_layout()
136     plt.show()
```

Listing 6: Algorithm Comparison

Case Studies

This appendix presents detailed case studies that demonstrate the application of reinforcement learning to real-world engineering problems. Each case study includes problem formulation, algorithm selection and tuning, implementation details, and lessons learned.

.9 Case Study 1: Autonomous Drone Navigation

.9.1 Problem Description

A quadrotor drone must navigate through a complex environment with obstacles while minimizing energy consumption and flight time. The drone has continuous state and action spaces, making this a challenging continuous control problem.

State Space: $s = (x, y, z, \dot{x}, \dot{y}, \dot{z}, \phi, \theta, \psi, \dot{\phi}, \dot{\theta}, \dot{\psi}) \in \mathbb{R}^{12}$

where (x, y, z) is position, $(\dot{x}, \dot{y}, \dot{z})$ is velocity, (ϕ, θ, ψ) are Euler angles, and $(\dot{\phi}, \dot{\theta}, \dot{\psi})$ are angular velocities.

Action Space: $a = (T, \tau_\phi, \tau_\theta, \tau_\psi) \in \mathbb{R}^4$

where T is thrust and $(\tau_\phi, \tau_\theta, \tau_\psi)$ are torques about each axis.

Dynamics: The quadrotor dynamics are governed by:

$$m\ddot{\mathbf{r}} = T\mathbf{R}\mathbf{e}_3 - mg\mathbf{e}_3 \quad (26)$$

$$\mathbf{I}\dot{\boldsymbol{\omega}} = \boldsymbol{\tau} - \boldsymbol{\omega} \times \mathbf{I}\boldsymbol{\omega} \quad (27)$$

where \mathbf{R} is the rotation matrix and \mathbf{I} is the inertia tensor.

.9.2 Algorithm Selection and Implementation

Algorithm Choice: Deep Deterministic Policy Gradient (DDPG) was selected for its ability to handle continuous action spaces and its sample efficiency.

Network Architecture:

- Actor: $[12] \rightarrow [256] \rightarrow [256] \rightarrow [4]$ with tanh output activation
- Critic: $[16] \rightarrow [256] \rightarrow [256] \rightarrow [1]$ (state-action input)

- Target networks with soft updates ($\tau = 0.001$)

Reward Function:

$$r(s, a) = -\|s_{pos} - s_{target}\|^2 - 0.1\|a\|^2 - 10 \cdot \mathbf{1}_{collision} \quad (28)$$

.9.3 Training Process and Results

Training Configuration:

- Episodes: 2000
- Steps per episode: 1000
- Replay buffer size: 10^6
- Batch size: 256
- Learning rates: Actor 10^{-4} , Critic 10^{-3}
- Exploration noise: Ornstein-Uhlenbeck process

Performance Metrics:

- Success rate: 94% (reaching target within 1m)
- Average flight time: 12.3 seconds
- Energy efficiency: 15% improvement over PID controller
- Collision rate: 2%

.9.4 Lessons Learned

1. **Reward Shaping Critical:** Initial sparse rewards led to poor exploration. Dense reward with distance-based terms significantly improved learning.
2. **Curriculum Learning Effective:** Starting with simple environments and gradually increasing complexity improved sample efficiency.
3. **Simulation-to-Reality Gap:** Robust training with domain randomization was essential for real-world transfer.
4. **Safety Considerations:** Emergency safety controller was necessary during real-world testing.

.10 Case Study 2: Smart Grid Energy Management

.10.1 Problem Description

A microgrid with renewable energy sources, battery storage, and flexible loads must optimize energy dispatch to minimize costs while maintaining reliability constraints.

State Space:

- Battery state of charge: $SOC \in [0, 1]$
- Renewable generation forecast: $P_{ren} \in [0, P_{max}]$
- Load demand forecast: $P_{load} \in [0, P_{max}]$
- Electricity price: $\lambda \in [\lambda_{min}, \lambda_{max}]$
- Time of day: $t \in [0, 23]$

Action Space:

- Battery charge/discharge power: $P_{bat} \in [-P_{bat,max}, P_{bat,max}]$
- Grid import/export: $P_{grid} \in [-P_{grid,max}, P_{grid,max}]$
- Load curtailment: $P_{curt} \in [0, P_{load}]$

.10.2 MDP Formulation

Transition Dynamics:

$$SOC_{t+1} = SOC_t + \frac{\eta P_{bat,t} \Delta t}{E_{bat,max}} \quad (29)$$

$$P_{balance} = P_{ren} + P_{grid} + P_{bat} - P_{load} + P_{curt} \quad (30)$$

Constraints:

$$SOC_{min} \leq SOC_t \leq SOC_{max} \quad (31)$$

$$|P_{balance}| \leq \epsilon \quad (\text{power balance}) \quad (32)$$

Reward Function:

$$r_t = -(\lambda_t P_{grid,t} \Delta t + C_{curt} P_{curt,t} + C_{deg} |P_{bat,t}|) \quad (33)$$

.10.3 Implementation Details

Algorithm: Soft Actor-Critic (SAC) for its sample efficiency and robustness.

Feature Engineering:

- Moving averages of renewable generation and demand

- Seasonal and diurnal patterns encoded as sinusoidal features
- Price forecasts using historical patterns

Constraint Handling: Projection method to ensure feasible actions:

$$a_{feasible} = \Pi_{\mathcal{A}}(a_{proposed}) \quad (34)$$

.10.4 Results and Performance Analysis

Performance Comparison:

Method	Daily Cost (\$)	Renewable Utilization (%)	Constraint Violations
Rule-based	142.50	78.3	0
MPC	138.20	82.1	0
SAC	134.80	85.7	0.02%

Key Insights:

1. 5.4% cost reduction compared to rule-based controller
2. 2.6% cost reduction compared to model predictive control
3. Learned to anticipate price patterns and pre-charge batteries
4. Robust performance under forecast uncertainties

.11 Case Study 3: Manufacturing Process Optimization

.11.1 Problem Description

A chemical batch process must optimize temperature and pressure profiles to maximize product yield while minimizing energy consumption and processing time.

Process Dynamics:

$$\frac{dC_A}{dt} = -k_1(T)C_A \quad (35)$$

$$\frac{dC_B}{dt} = k_1(T)C_A - k_2(T)C_B \quad (36)$$

$$\frac{dT}{dt} = \frac{Q - Q_{loss}(T)}{mC_p} \quad (37)$$

where $k_i(T) = A_i \exp(-E_i/RT)$ are temperature-dependent rate constants.

.11.2 Multi-Objective Optimization

Objectives:

$$J_1 = \text{maximize } C_B(t_f) \quad (\text{yield}) \quad (38)$$

$$J_2 = \text{minimize } \int_0^{t_f} Q(t)dt \quad (\text{energy}) \quad (39)$$

$$J_3 = \text{minimize } t_f \quad (\text{time}) \quad (40)$$

Scalarization: Weighted sum approach:

$$r(s, a) = w_1 \frac{C_B(t_f)}{C_{B,max}} - w_2 \frac{Q(t)}{Q_{max}} - w_3 \frac{1}{t_{max}} \quad (41)$$

.11.3 Implementation and Results

Algorithm: Twin Delayed DDPG (TD3) for stability in continuous control.

Results:

- 12% improvement in product yield
- 18% reduction in energy consumption
- 8% reduction in batch time
- Consistent performance across different initial conditions

.12 Performance Comparisons

.12.1 Algorithm Performance Summary

Case Study	Problem Type	Algorithm	Sample Efficiency	Final Performance
Drone Navigation	Continuous Control	DDPG	Medium	94% success
Smart Grid	Constrained Control	SAC	High	5.4% improvement
Manufacturing	Multi-objective	TD3	Medium	12% yield gain

.12.2 Common Success Factors

1. **Domain Knowledge Integration:** Incorporating engineering insights into reward design and feature engineering
2. **Simulation Fidelity:** High-fidelity simulators were crucial for initial learning
3. **Constraint Handling:** Explicit constraint enforcement prevented unsafe exploration

4. **Robust Training:** Domain randomization and robust optimization improved real-world performance

.13 Troubleshooting Guide

.13.1 Common Training Issues

Poor Convergence:

- Check reward function scaling and normalization
- Verify network initialization and learning rates
- Ensure sufficient exploration during early training
- Monitor gradient norms for vanishing/exploding gradients

Unstable Training:

- Reduce learning rates, especially for critic networks
- Use target networks with appropriate update rates
- Implement gradient clipping
- Check for numerical instabilities in environment dynamics

Poor Real-World Transfer:

- Increase simulation realism and randomization
- Implement domain adaptation techniques
- Use conservative policy updates
- Include safety constraints and emergency controllers

.13.2 Hyperparameter Tuning Guidelines

Learning Rates:

- Actor: typically 10^{-4} to 10^{-3}
- Critic: typically 10^{-3} to 10^{-2}
- Use learning rate schedules for long training runs

Network Architecture:

- Start with 2-3 hidden layers of 256-512 units
- Use batch normalization for deep networks
- Consider layer normalization for recurrent policies

Exploration:

- Gaussian noise: $\sigma = 0.1$ to 0.2 of action range
- Ornstein-Uhlenbeck: $\theta = 0.15$, $\sigma = 0.2$
- Decay exploration over training

.13.3 Debugging Checklist**1. Environment Sanity Checks:**

- Verify state and action space definitions
- Test random policy performance
- Check reward function computation
- Validate episode termination conditions

2. Algorithm Implementation:

- Verify gradient computation and backpropagation
- Check replay buffer implementation
- Validate target network updates
- Monitor loss functions and metrics

3. Training Diagnostics:

- Plot learning curves and moving averages
- Monitor exploration statistics
- Track gradient norms and weight distributions
- Analyze action distributions over time