# Reinforcement Learning

An Introduction

**second edition**

Richard S. Sutton and Andrew G. Barto

**Adaptive Computation and Machine Learning**

Francis Bach, series editor

A complete list of books published in the Adaptive Computation and Machine Learning series appears at the back of this book.

# Reinforcement Learning:

# An Introduction

# second edition

Richard S. Sutton and Andrew G. Barto

In memory of A. Harry Klopf

# Contents

# Preface to the Second Edition

The twenty years since the publication of the first edition of this book have seen tremendous progress in artificial intelligence, propelled in large part by advances in machine learning, including advances in reinforcement learning. Although the impressive computational power that became available is responsible for some of these advances, new developments in theory and algorithms have been driving forces as well. In the face of this progress, a second edition of our 1998 book was long overdue, and we finally began the project in 2012. Our goal for the second edition was the same as our goal for the first: to provide a clear and simple account of the key ideas and algorithms of reinforcement learning that is accessible to readers in all the related disciplines. The edition remains an introduction, and we retain a focus on core, online learning algorithms. This edition includes some new topics that rose to importance over the intervening years, and we expanded coverage of topics that we now understand better. But we made no attempt to provide comprehensive coverage of the field, which has exploded in many different directions. We apologize for having to leave out all but a handful of these contributions.

As in the first edition, we chose not to produce a rigorous formal treatment of reinforcement learning, or to formulate it in the most general terms. However, our deeper understanding of some topics since the first edition required a bit more mathematics to explain; we have set off the more mathematical parts in shaded boxes that the non-mathematically-inclined may choose to skip. We also use a slightly different notation than was used in the first edition. In teaching, we have found that the new notation helps to address some common points of confusion. It emphasizes the difference between random variables, denoted with capital letters, and their instantiations, denoted in lower case. For example, the state, action, and reward at time step $t$ are denoted $S_t$, $A_t$, and $R_t$, while their possible values might be denoted $s$, $a$, and $r$. Along with this, it is natural to use lower case for value functions (e.g., $v_\pi$) and restrict capitals to their tabular estimates (e.g., $Q_t(s, a)$). Approximate value functions are deterministic functions of random parameters and are thus also in lower case (e.g., $\hat{v}(s, \mathbf{w}_t) \approx v_\pi(s)$). Vectors, such as the weight vector $\mathbf{w}_t$ (formerly $\boldsymbol{\theta}_t$) and the feature vector $\mathbf{x}_t$ (formerly $\boldsymbol{\phi}_t$), are bold and written in lowercase even if they are random variables. Uppercase bold is reserved for matrices. In the first edition we used special notations, $\mathcal{P}^a_{ss'}$ and $\mathcal{R}^a_{ss'}$, for the transition probabilities and expected rewards. One weakness of that notation is that it still did not fully characterize the dynamics of the rewards, giving only their expectations, which is sufficient for dynamic programming but not for reinforcement learning. Another weakness

is the excess of subscripts and superscripts. In this edition we use the explicit notation of $p(s', r \,|\, s, a)$ for the joint probability for the next state and reward given the current state and action. All the changes in notation are summarized in a table on page xix.

The second edition is significantly expanded, and its top-level organization has been changed. After the introductory first chapter, the second edition is divided into three new parts. The first part (Chapters 2–8) treats as much of reinforcement learning as possible without going beyond the tabular case for which exact solutions can be found. We cover both learning and planning methods for the tabular case, as well as their unification in $n$-step methods and in Dyna. Many algorithms presented in this part are new to the second edition, including UCB, Expected Sarsa, Double learning, tree-backup, $Q(\sigma)$, RTDP, and MCTS. Doing the tabular case first, and thoroughly, enables core ideas to be developed in the simplest possible setting. The second part of the book (Chapters 9–13) is then devoted to extending the ideas to function approximation. It has new sections on artificial neural networks, the fourier basis, LSTD, kernel-based methods, Gradient-TD and Emphatic-TD methods, average-reward methods, true online TD($\lambda$), and policy-gradient methods. The second edition significantly expands the treatment of off-policy learning, first for the tabular case in Chapters 5–7, then with function approximation in Chapters 11 and 12. Another change is that the second edition separates the forward-view idea of $n$-step bootstrapping (now treated more fully in Chapter 7) from the backward-view idea of eligibility traces (now treated independently in Chapter 12). The third part of the book has large new chapters on reinforcement learning's relationships to psychology (Chapter 14) and neuroscience (Chapter 15), as well as an updated case-studies chapter including Atari game playing, Watson's wagering strategy, and the Go playing programs AlphaGo and AlphaGo Zero (Chapter 16). Still, out of necessity we have included only a small subset of all that has been done in the field. Our choices reflect our long-standing interests in inexpensive model-free methods that should scale well to large applications. The final chapter now includes a discussion of the future societal impacts of reinforcement learning. For better or worse, the second edition is about twice as large as the first.

This book is designed to be used as the primary text for a one- or two-semester course on reinforcement learning. For a one-semester course, the first ten chapters should be covered in order and form a good core, to which can be added material from the other chapters, from other books such as Bertsekas and Tsitsiklis (1996), Wiering and van Otterlo (2012), and Szepesvári (2010), or from the literature, according to taste. Depending of the students' background, some additional material on online supervised learning may be helpful. The ideas of options and option models are a natural addition (Sutton, Precup and Singh, 1999). A two-semester course can cover all the chapters as well as supplementary material. The book can also be used as part of broader courses on machine learning, artificial intelligence, or neural networks. In this case, it may be desirable to cover only a subset of the material. We recommend covering Chapter 1 for a brief overview, Chapter 2 through Section 2.4, Chapter 3, and then selecting sections from the remaining chapters according to time and interests. Chapter 6 is the most important for the subject and for the rest of the book. A course focusing on machine learning or neural networks should cover Chapters 9 and 10, and a course focusing on artificial intelligence or planning should cover Chapter 8. Throughout the book, sections and chapters that are more difficult and not essential to the rest of the book are marked

with a ∗. These can be omitted on first reading without creating problems later on. Some exercises are also marked with a ∗ to indicate that they are more advanced and not essential to understanding the basic material of the chapter.

Most chapters end with a section entitled "Bibliographical and Historical Remarks," wherein we credit the sources of the ideas presented in that chapter, provide pointers to further reading and ongoing research, and describe relevant historical background. Despite our attempts to make these sections authoritative and complete, we have undoubtedly left out some important prior work. For that we again apologize, and we welcome corrections and extensions for incorporation into the electronic version of the book.

Like the first edition, this edition of the book is dedicated to the memory of A. Harry Klopf. It was Harry who introduced us to each other, and it was his ideas about the brain and artificial intelligence that launched our long excursion into reinforcement learning. Trained in neurophysiology and long interested in machine intelligence, Harry was a senior scientist affiliated with the Avionics Directorate of the Air Force Office of Scientific Research (AFOSR) at Wright-Patterson Air Force Base, Ohio. He was dissatisfied with the great importance attributed to equilibrium-seeking processes, including homeostasis and error-correcting pattern classification methods, in explaining natural intelligence and in providing a basis for machine intelligence. He noted that systems that try to maximize something (whatever that might be) are qualitatively different from equilibrium-seeking systems, and he argued that maximizing systems hold the key to understanding important aspects of natural intelligence and for building artificial intelligences. Harry was instrumental in obtaining funding from AFOSR for a project to assess the scientific merit of these and related ideas. This project was conducted in the late 1970s at the University of Massachusetts Amherst (UMass Amherst), initially under the direction of Michael Arbib, William Kilmer, and Nico Spinelli, professors in the Department of Computer and Information Science at UMass Amherst, and founding members of the Cybernetics Center for Systems Neuroscience at the University, a farsighted group focusing on the intersection of neuroscience and artificial intelligence. Barto, a recent PhD from the University of Michigan, was hired as post doctoral researcher on the project. Meanwhile, Sutton, an undergraduate studying computer science and psychology at Stanford, had been corresponding with Harry regarding their mutual interest in the role of stimulus timing in classical conditioning. Harry suggested to the UMass group that Sutton would be a great addition to the project. Thus, Sutton became a UMass graduate student, whose PhD was directed by Barto, who had become an Associate Professor. The study of reinforcement learning as presented in this book is rightfully an outcome of that project instigated by Harry and inspired by his ideas. Further, Harry was responsible for bringing us, the authors, together in what has been a long and enjoyable interaction. By dedicating this book to Harry we honor his essential contributions, not only to the field of reinforcement learning, but also to our collaboration. We also thank Professors Arbib, Kilmer, and Spinelli for the opportunity they provided to us to begin exploring these ideas. Finally, we thank AFOSR for generous support over the early years of our research, and the NSF for its generous support over many of the following years.

We have very many people to thank for their inspiration and help with this second edition. Everyone we acknowledged for their inspiration and help with the first edition

# Preface to the First Edition

We first came to focus on what is now known as reinforcement learning in late 1979. We were both at the University of Massachusetts, working on one of the earliest projects to revive the idea that networks of neuronlike adaptive elements might prove to be a promising approach to artificial adaptive intelligence. The project explored the "heterostatic theory of adaptive systems" developed by A. Harry Klopf. Harry's work was a rich source of ideas, and we were permitted to explore them critically and compare them with the long history of prior work in adaptive systems. Our task became one of teasing the ideas apart and understanding their relationships and relative importance. This continues today, but in 1979 we came to realize that perhaps the simplest of the ideas, which had long been taken for granted, had received surprisingly little attention from a computational perspective. This was simply the idea of a learning system that *wants* something, that adapts its behavior in order to maximize a special signal from its environment. This was the idea of a "hedonistic" learning system, or, as we would say now, the idea of reinforcement learning.

Like others, we had a sense that reinforcement learning had been thoroughly explored in the early days of cybernetics and artificial intelligence. On closer inspection, though, we found that it had been explored only slightly. While reinforcement learning had clearly motivated some of the earliest computational studies of learning, most of these researchers had gone on to other things, such as pattern classification, supervised learning, and adaptive control, or they had abandoned the study of learning altogether. As a result, the special issues involved in learning how to get something from the environment received relatively little attention. In retrospect, focusing on this idea was the critical step that set this branch of research in motion. Little progress could be made in the computational study of reinforcement learning until it was recognized that such a fundamental idea had not yet been thoroughly explored.

The field has come a long way since then, evolving and maturing in several directions. Reinforcement learning has gradually become one of the most active research areas in machine learning, artificial intelligence, and neural network research. The field has developed strong mathematical foundations and impressive applications. The computational study of reinforcement learning is now a large field, with hundreds of active researchers around the world in diverse disciplines such as psychology, control theory, artificial intelligence, and neuroscience. Particularly important have been the contributions establishing and developing the relationships to the theory of optimal control and dynamic programming.

The overall problem of learning from interaction to achieve goals is still far from being solved, but our understanding of it has improved significantly. We can now place component ideas, such as temporal-difference learning, dynamic programming, and function approximation, within a coherent perspective with respect to the overall problem.

Our goal in writing this book was to provide a clear and simple account of the key ideas and algorithms of reinforcement learning. We wanted our treatment to be accessible to readers in all of the related disciplines, but we could not cover all of these perspectives in detail. For the most part, our treatment takes the point of view of artificial intelligence and engineering. Coverage of connections to other fields we leave to others or to another time. We also chose not to produce a rigorous formal treatment of reinforcement learning. We did not reach for the highest possible level of mathematical abstraction and did not rely on a theorem–proof format. We tried to choose a level of mathematical detail that points the mathematically inclined in the right directions without distracting from the simplicity and potential generality of the underlying ideas.

# Summary of Notation

Capital letters are used for random variables, whereas lower case letters are used for the values of random variables and for scalar functions. Quantities that are required to be real-valued vectors are written in bold and in lower case (even if random variables). Matrices are bold capitals.

| | |
|---|---|
| $\doteq$ | equality relationship that is true by definition |
| $\approx$ | approximately equal |
| $\propto$ | proportional to |
| $\Pr\{X\!=\!x\}$ | probability that a random variable $X$ takes on the value $x$ |
| $X \sim p$ | random variable $X$ selected from distribution $p(x) \doteq \Pr\{X\!=\!x\}$ |
| $\mathbb{E}[X]$ | expectation of a random variable $X$, i.e., $\mathbb{E}[X] \doteq \sum_x p(x)x$ |
| $\operatorname{argmax}_a f(a)$ | a value of $a$ at which $f(a)$ takes its maximal value |
| $\ln x$ | natural logarithm of $x$ |
| $e^x$, $\exp(x)$ | the base of the natural logarithm, $e \approx 2.71828$, carried to power $x$; $e^{\ln x} = x$ |
| $\mathbb{R}$ | set of real numbers |
| $f : \mathcal{X} \to \mathcal{Y}$ | function $f$ from elements of set $\mathcal{X}$ to elements of set $\mathcal{Y}$ |
| $\leftarrow$ | assignment |
| $(a, b]$ | the real interval between $a$ and $b$ including $b$ but not including $a$ |
| | |
| $\varepsilon$ | probability of taking a random action in an $\varepsilon$-greedy policy |
| $\alpha, \beta$ | step-size parameters |
| $\gamma$ | discount-rate parameter |
| $\lambda$ | decay-rate parameter for eligibility traces |
| $\mathbb{1}_{predicate}$ | indicator function ($\mathbb{1}_{predicate} \doteq 1$ if the *predicate* is true, else 0) |

In a multi-arm bandit problem:

| | |
|---|---|
| $k$ | number of actions (arms) |
| $t$ | discrete time step or play number |
| $q_*(a)$ | true value (expected reward) of action $a$ |
| $Q_t(a)$ | estimate at time $t$ of $q_*(a)$ |
| $N_t(a)$ | number of times action $a$ has been selected up prior to time $t$ |
| $H_t(a)$ | learned preference for selecting action $a$ at time $t$ |
| $\pi_t(a)$ | probability of selecting action $a$ at time $t$ |
| $\bar{R}_t$ | estimate at time $t$ of the expected reward given $\pi_t$ |

In a Markov Decision Process:

| | |
|---|---|
| $s, s'$ | states |
| $a$ | an action |
| $r$ | a reward |
| $\mathcal{S}$ | set of all nonterminal states |
| $\mathcal{S}^+$ | set of all states, including the terminal state |
| $\mathcal{A}(s)$ | set of all actions available in state $s$ |
| $\mathcal{R}$ | set of all possible rewards, a finite subset of $\mathbb{R}$ |
| $\subset$ | subset of (e.g., $\mathcal{R} \subset \mathbb{R}$) |
| $\in$ | is an element of; e.g. $(s \in \mathcal{S},\ r \in \mathcal{R})$ |
| $|\mathcal{S}|$ | number of elements in set $\mathcal{S}$ |

| | |
|---|---|
| $t$ | discrete time step |
| $T, T(t)$ | final time step of an episode, or of the episode including time step $t$ |
| $A_t$ | action at time $t$ |
| $S_t$ | state at time $t$, typically due, stochastically, to $S_{t-1}$ and $A_{t-1}$ |
| $R_t$ | reward at time $t$, typically due, stochastically, to $S_{t-1}$ and $A_{t-1}$ |
| $\pi$ | policy (decision-making rule) |
| $\pi(s)$ | action taken in state $s$ under *deterministic* policy $\pi$ |
| $\pi(a|s)$ | probability of taking action $a$ in state $s$ under *stochastic* policy $\pi$ |

| | |
|---|---|
| $G_t$ | return following time $t$ |
| $h$ | horizon, the time step one looks up to in a forward view |
| $G_{t:t+n}, G_{t:h}$ | $n$-step return from $t+1$ to $t+n$, or to $h$ (discounted and corrected) |
| $\bar{G}_{t:h}$ | flat return (undiscounted and uncorrected) from $t+1$ to $h$ (Section 5.8) |
| $G_t^\lambda$ | $\lambda$-return (Section 12.1) |
| $G_{t:h}^\lambda$ | truncated, corrected $\lambda$-return (Section 12.3) |
| $G_t^{\lambda s}, G_t^{\lambda a}$ | $\lambda$-return, corrected by estimated state, or action, values (Section 12.8) |

| | |
|---|---|
| $p(s', r|s, a)$ | probability of transition to state $s'$ with reward $r$, from state $s$ and action $a$ |
| $p(s'|s, a)$ | probability of transition to state $s'$, from state $s$ taking action $a$ |
| $r(s, a)$ | expected immediate reward from state $s$ after action $a$ |
| $r(s, a, s')$ | expected immediate reward on transition from $s$ to $s'$ under action $a$ |

| | |
|---|---|
| $v_\pi(s)$ | value of state $s$ under policy $\pi$ (expected return) |
| $v_*(s)$ | value of state $s$ under the optimal policy |
| $q_\pi(s, a)$ | value of taking action $a$ in state $s$ under policy $\pi$ |
| $q_*(s, a)$ | value of taking action $a$ in state $s$ under the optimal policy |

| | |
|---|---|
| $V, V_t$ | array estimates of state-value function $v_\pi$ or $v_*$ |
| $Q, Q_t$ | array estimates of action-value function $q_\pi$ or $q_*$ |
| $\bar{V}_t(s)$ | expected approximate action value; for example, $\bar{V}_t(s) \doteq \sum_a \pi(a|s) Q_t(s, a)$ |
| $U_t$ | target for estimate at time $t$ |

| | |
|---|---|
| $\delta_t$ | temporal-difference (TD) error at $t$ (a random variable) (Section 6.1) |
| $\delta_t^s, \delta_t^a$ | state- and action-specific forms of the TD error (Section 12.9) |
| $n$ | in $n$-step methods, $n$ is the number of steps of bootstrapping |
| | |
| $d$ | dimensionality—the number of components of $\mathbf{w}$ |
| $d'$ | alternate dimensionality—the number of components of $\boldsymbol{\theta}$ |
| $\mathbf{w}, \mathbf{w}_t$ | $d$-vector of weights underlying an approximate value function |
| $w_i, w_{t,i}$ | $i$th component of learnable weight vector |
| $\hat{v}(s,\mathbf{w})$ | approximate value of state $s$ given weight vector $\mathbf{w}$ |
| $v_{\mathbf{w}}(s)$ | alternate notation for $\hat{v}(s,\mathbf{w})$ |
| $\hat{q}(s,a,\mathbf{w})$ | approximate value of state–action pair $s,a$ given weight vector $\mathbf{w}$ |
| $\nabla\hat{v}(s,\mathbf{w})$ | column vector of partial derivatives of $\hat{v}(s,\mathbf{w})$ with respect to $\mathbf{w}$ |
| $\nabla\hat{q}(s,a,\mathbf{w})$ | column vector of partial derivatives of $\hat{q}(s,a,\mathbf{w})$ with respect to $\mathbf{w}$ |
| | |
| $\mathbf{x}(s)$ | vector of features visible when in state $s$ |
| $\mathbf{x}(s,a)$ | vector of features visible when in state $s$ taking action $a$ |
| $x_i(s), x_i(s,a)$ | $i$th component of vector $\mathbf{x}(s)$ or $\mathbf{x}(s,a)$ |
| $\mathbf{x}_t$ | shorthand for $\mathbf{x}(S_t)$ or $\mathbf{x}(S_t, A_t)$ |
| $\mathbf{w}^\top\mathbf{x}$ | inner product of vectors, $\mathbf{w}^\top\mathbf{x} \doteq \sum_i w_i x_i$; for example, $\hat{v}(s,\mathbf{w}) \doteq \mathbf{w}^\top\mathbf{x}(s)$ |
| $\mathbf{v}, \mathbf{v}_t$ | secondary $d$-vector of weights, used to learn $\mathbf{w}$ (Chapter 11) |
| $\mathbf{z}_t$ | $d$-vector of eligibility traces at time $t$ (Chapter 12) |
| | |
| $\boldsymbol{\theta}, \boldsymbol{\theta}_t$ | parameter vector of target policy (Chapter 13) |
| $\pi(a|s,\boldsymbol{\theta})$ | probability of taking action $a$ in state $s$ given parameter vector $\boldsymbol{\theta}$ |
| $\pi_{\boldsymbol{\theta}}$ | policy corresponding to parameter $\boldsymbol{\theta}$ |
| $\nabla\pi(a|s,\boldsymbol{\theta})$ | column vector of partial derivatives of $\pi(a|s,\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$ |
| $J(\boldsymbol{\theta})$ | performance measure for the policy $\pi_{\boldsymbol{\theta}}$ |
| $\nabla J(\boldsymbol{\theta})$ | column vector of partial derivatives of $J(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$ |
| $h(s,a,\boldsymbol{\theta})$ | preference for selecting action $a$ in state $s$ based on $\boldsymbol{\theta}$ |
| | |
| $b(a|s)$ | behavior policy used to select actions while learning about target policy $\pi$ |
| $b(s)$ | a baseline function $b : \mathcal{S} \mapsto \mathbb{R}$ for policy-gradient methods |
| $b$ | branching factor for an MDP or search tree |
| $\rho_{t:h}$ | importance sampling ratio for time $t$ through time $h$ (Section 5.5) |
| $\rho_t$ | importance sampling ratio for time $t$ alone, $\rho_t \doteq \rho_{t:t}$ |
| $r(\pi)$ | average reward (reward rate) for policy $\pi$ (Section 10.3) |
| $\bar{R}_t$ | estimate of $r(\pi)$ at time $t$ |
| | |
| $\mu(s)$ | on-policy distribution over states (Section 9.2) |
| $\boldsymbol{\mu}$ | $|\mathcal{S}|$-vector of the $\mu(s)$ for all $s \in \mathcal{S}$ |
| $\|v\|_\mu^2$ | $\mu$-weighted squared norm of value function $v$, i.e., $\|v\|_\mu^2 \doteq \sum_{s\in\mathcal{S}} \mu(s)v(s)^2$ |
| $\eta(s)$ | expected number of visits to state $s$ per episode (page 199) |
| $\Pi$ | projection operator for value functions (page 268) |
| $B_\pi$ | Bellman operator for value functions (Section 11.4) |

| | |
|---|---|
| $\mathbf{A}$ | $d \times d$ matrix $\mathbf{A} \doteq \mathbb{E}\left[\mathbf{x}_t\left(\mathbf{x}_t - \gamma\mathbf{x}_{t+1}\right)^\top\right]$ |
| $\mathbf{b}$ | $d$-dimensional vector $\mathbf{b} \doteq \mathbb{E}[R_{t+1}\mathbf{x}_t]$ |
| $\mathbf{w}_{\text{TD}}$ | TD fixed point $\mathbf{w}_{\text{TD}} \doteq \mathbf{A}^{-1}\mathbf{b}$ (a $d$-vector, Section 9.4) |
| $\mathbf{I}$ | identity matrix |
| $\mathbf{P}$ | $|\mathcal{S}| \times |\mathcal{S}|$ matrix of state-transition probabilities under $\pi$ |
| $\mathbf{D}$ | $|\mathcal{S}| \times |\mathcal{S}|$ diagonal matrix with $\boldsymbol{\mu}$ on its diagonal |
| $\mathbf{X}$ | $|\mathcal{S}| \times d$ matrix with the $\mathbf{x}(s)$ as its rows |
| | |
| $\bar{\delta}_{\mathbf{w}}(s)$ | Bellman error (expected TD error) for $v_{\mathbf{w}}$ at state $s$ (Section 11.4) |
| $\bar{\delta}_{\mathbf{w}}$, BE | Bellman error vector, with components $\bar{\delta}_{\mathbf{w}}(s)$ |
| $\overline{\text{VE}}(\mathbf{w})$ | mean square value error $\overline{\text{VE}}(\mathbf{w}) \doteq \|v_{\mathbf{w}} - v_\pi\|_\mu^2$ (Section 9.2) |
| $\overline{\text{BE}}(\mathbf{w})$ | mean square Bellman error $\overline{\text{BE}}(\mathbf{w}) \doteq \left\|\bar{\delta}_{\mathbf{w}}\right\|_\mu^2$ |
| $\overline{\text{PBE}}(\mathbf{w})$ | mean square projected Bellman error $\overline{\text{PBE}}(\mathbf{w}) \doteq \left\|\Pi\bar{\delta}_{\mathbf{w}}\right\|_\mu^2$ |
| $\overline{\text{TDE}}(\mathbf{w})$ | mean square temporal-difference error $\overline{\text{TDE}}(\mathbf{w}) \doteq \mathbb{E}_b\left[\rho_t\delta_t^2\right]$ (Section 11.5) |
| $\overline{\text{RE}}(\mathbf{w})$ | mean square return error (Section 11.6) |