# A Review of *Reinforcement Learning*

## *Sebastian Thrun and Michael L. Littman*

The reinforcement learning problem is the challenge of AI in a microcosm; how can we build an agent that can plan, learn, perceive, and act in a complex world? There's a great new book on the market that lays out the conceptual and algorithmic foundations of this exciting area. Reinforcement learning pioneers Rich Sutton and Andy Barto have published *Reinforcement Learning: An Introduction,* providing a highly accessible starting point for interested students, researchers, and practitioners.

In the reinforcement learning framework, an agent acts in an environment whose state it can sense and occasionally receives some penalty or reward based on its state and action. Its learning task is to find a policy for action selection that maximizes its reward over the long haul; this task requires not only choosing actions that are associated with high reward in the current state but thinking ahead by choosing actions that will lead the agents to more lucrative parts of the state space. Although there are many ways to attack this problem, the paradigm described in the book is to construct a value function that evaluates the "goodness" of different situations. In particular, the *value* of a state is the long-term reward that can be attained if actions are chosen optimally. Recent research has produced a flurry of algorithms for learning value functions, theoretical insights into their power and limitations, and a series of fielded applications. The authors have done a wonderful job of boiling down disparate and complex reinforcement learning algorithms to a set of fundamental components, then showing how these components work together. The differences between dynamic programming,

Monte Carlo methods, and temporal difference learning are teased apart, then tied back together in a unified way. Innovations such as backup diagrams, which decorate the book cover, help convey the power and excitement behind reinforcement learning methods to both novices and veterans like us.

The book consists of three parts, one dedicated to the problem description and two others to a range of reinforcement learning algorithms, their analysis, and related research issues.

We enthusiastically applaud the authors' decision to articulate the problem addressed in the book before talking in length about its various solutions. After all, a thorough discussion of the problem is necessary for

---

Reinforcement Learning: An Introduction, *Richard S. Sutton and Andrew G. Barto, The MIT Press, Cambridge, Massachusetts, 1998, 322 pp., ISBN 0-262-19398-1.*

---

veterans to understand the aims and scope of reinforcement learning research let alone novices in the field. At 85 pages in length, however, one might wonder what it is about the reinforcement learning problem that its description deserves (or requires?) twice as many pages as the typical journal paper. Is the reinforcement learning problem so complicated that it takes that long to describe and discuss it?

In truth, Part 1 does much more

than just pose the problem. Chapter 1 contains a highly informal introduction to the broad problem domain: learning to select actions while interacting with an environment to achieve long-term goals. The example of Tic Tac Toe makes concepts such as reward, value functions, and the exploration-exploitation dilemma feel natural—all concepts that find a more mathematical treatment later in the book. The first chapter also provides an invaluable description of the history of reinforcement learning, placing recent research efforts in context. This history is an early example of a series of detailed literature reviews, found at the end of each chapter, which could alone justify the expense of purchasing the book.

Next, the book dives into a highly restricted instance of the reinforcement learning problem: the *k*-arm bandit problem. This well-researched problem lacks state transitions—there is only a single state—but it otherwise possesses the typical characteristics that set reinforcement learning apart from, say, supervised learning. The placement of the problem is well chosen because it illustrates with rigor the key concepts of the algorithms yet to come: the idea of interaction with an environment, reward, value functions, and the exploration-exploitation dilemma. It is followed by what readers knowledgeable in AI might choose as their starting point into the book: the formal, mathematical definition of the reinforcement learning problem. At this point, the authors state clearly the key assumptions that underlie the methods expounded in the book, including the critical Markov assumption that renders the environment's state fully observable. Their crystal-clear description of Bellman optimality leaves little room for misunderstanding. Trust us; even if you are familiar with reinforcement learning, you will find that Part 1 is insightful and great fun to read! By the end, we could hardly wait to get to Part 2 to learn about reinforcement learning algorithms.

At this point, it seems appropriate to make a few comments about the general style of the book. The text is extremely accessible, from the first

---

page to the last. The authors have successfully integrated formal algorithmic descriptions and analysis with intuitions, motivations, numeric examples, and exercises (whose solutions can be obtained from Richard Sutton). Exercises and examples are merged into the floating text, and we recommend spending some time thinking about them! Every exercise has at least one "aha!" insight in it; some have two. Great care was taken in choosing tractable, yet nontrivial examples of reinforcement learning problems; optimal blackjack, soap-bubble shape prediction, and cost-effective rental car management were some of our favorites, in addition to 101 variations on the classic gridworld problem. New ideas are introduced, described, discussed, and evaluated with meticulous care, and often theoretical results or easily replicable experiments accompany the general discourse on the material. The experiments are nicely done because they give the reader a hands-on apprecia-

tion for the underlying ideas. Because of the conversational writing style, algorithms such as TD(0), which encompass a collection of concepts, are not really described at a single location; instead, the book gradually progresses from basic algorithmic methods to today's state-of-the-art reinforcement learning algorithms. This approach makes the book a little difficult to use as a desktop reference or for a course that does not follow the same logical train of thoughts. However, there is a lot of wisdom in the book's development of ideas!

Another aspect of the writing style worth mentioning is that it is a bit dogmatic in places. For example, the book argues that *evolutionary methods* and other methods that directly search the policy space without constructing a value function do not fall into the scope of reinforcement learning. Why not? One of the criticisms the book gives of evolutionary methods is that they are imprecise about credit assignment—the whole policy is
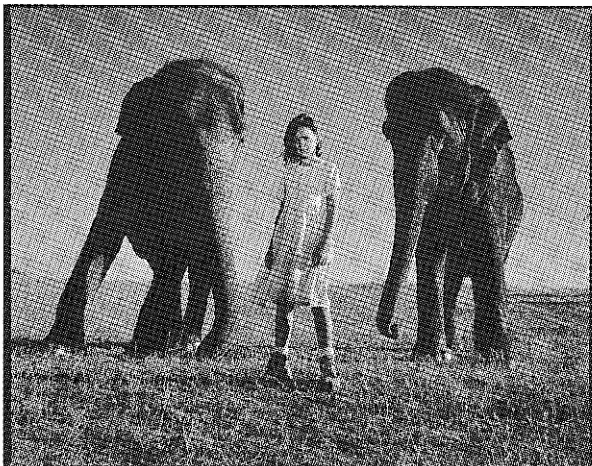
rewarded or punished. However, can't this same criticism be leveled at methods that use function approximation to represent the value function? This and other examples leave us with the feeling that the authors were a bit dismissive of reinforcement learning methods that aren't temporal difference. The reader should keep in mind that the material in the book is extremely valuable but occasionally biased.

Let's move on to the second algorithmic part of the book. Had one of us written this book, we would have been tempted to start its algorithmic part with a description of Q-learning, pointing out that this is by far the most popular reinforcement learning algorithm to date—not so these authors. The second part of the book introduces three families of methods, which form the basis of virtually all reinforcement learning algorithms: (1) dynamic programming, (2) Monte Carlo methods, and (3) temporal difference learning. This way, the reader

not only understands contemporary reinforcement learning algorithms but also appreciates their roots and connections to related (and often older) methods. If you now wonder what the differences are, *dynamic programming methods* compute value functions by backing up values from successor states to predecessor states, and they systematically update one state after another, using a model of the next-state distribution. *Monte Carlo methods* don't require such a model and instead sample entire trajectories to update the value functions based on the episodes' final outcomes. *Temporal difference methods* integrate ideas from both dynamic programming and Monte Carlo methods: Like Monte Carlo methods, they learn from sampled trajectories, but unlike them and like dynamic programming methods, they back up values from state to state (bootstrap). The celebrated Q-learning algorithm is finally introduced in the middle of Chapter 6 as a special case of temporal difference.

The third and final part of the book contains five quite different chapters in which the authors explore issues that go beyond the basic reinforcement learning paradigms. It begins with a discussion of eligibility traces and TD($\lambda$). A proof of the equivalence of the forward view and the backward view of TD($\lambda$) provides a refreshing mathematical insight in a book that otherwise contains only a few formal results. A thoughtful discussion on function approximation in reinforcement learning follows. This topic is currently an active research area, and the authors pay tribute to it by carefully describing methods that have been found to work well along with their pitfalls and limitations. This section lacks a description of backpropagation, which plays a key role in several successful applications of neural networks to reinforcement learning, but otherwise provides a range of useful examples. Subsequently, a chapter on reinforcement learning and planning forges ties to more traditional work in AI and psychology. After a final, brief summary of the entire book, the authors conclude their monologue with a description of six successful applications of reinforcement learn-

ing. Included here is an in-depth description of TD-GAMMON, one of the most captivating successes of the entire field because it learns to play the game of Backgammon on par with the best human players.

Now, what is there to be said about the book as a whole? It is a gentle, insightful, and balanced introduction into a topic that over the last decade has been subject to intense research. More than that, it ties together the basic ideas in an astonishingly clear way, with many new insights into the relation of different reinforcement learning algorithms. Should you buy it? If you are a novice to AI and want to learn something about a highly active field in AI, the answer is definitely yes. If you are a teacher who wants to design a course on reinforcement learning, then yes, buy it, and also purchase *Neuro-Dynamic Programming* by Bertsekas and Tsitsiklis (1996) to avoid running out of material after the first few weeks. If you are an active reinforcement learning researcher who'd like to know the latest and best in reinforcement learning, buy it anyhow—be it just for the bibliographic

remarks and clever examples—but be aware that the book isn't intended as a guide to current research. It does, however, highlight intriguing open problems and point out promising lines of future research—indicators that the field of reinforcement learning is healthy and changing and on the verge of new and important discoveries!

### References

Bertsekas, D. P., and Tsitsiklis, J. N. 1996. *Neuro-Dynamic Programming,* Belmon, Mass.: Athena Scientific.

**Sebastian Thrun** is an assistant professor of computer science and robotics at Carnegie Mellon University, with research interests in AI, robotics, and machine learning.

**Michael L. Littman** just joined AT&T Labs-Research and is affiliated with Duke University. He applies probability theory and machine learning to planning and language.

# Android Epistemology

## Edited by Kenneth M. Ford, Clark Glymour, & Patrick J. Hayes

Epistemology has traditionally been the study of human knowledge and rational change of human belief. *Android epistemology* is the exploration of the space of possible machines and their capacities for knowledge, beliefs, attitudes, desires, and action in accord with their mental states. From the perspective of android epistemology, artificial intelligence and computational cognitive psychology form a unified endeavor: artificial intelligence explores any possible way of engineering machines with intelligent features, while cognitive psychology focuses on reverse engineering the most intelligent system we know, us. The editors argue that contemporary android epistemology is the fruition of a long tradition in philosophical theories of knowledge and mind.

The sixteen essays by computer scientists and philosophers collected in this volume include substantial contributions to android epistemology as well as examinations, defenses, elaborations, and challenges to the very idea.

Contributors include Kalyan Basu, Margaret Boden, Selmer Bringsjord, Ronald Chrisley, Paul Churchland, Cary deBessonet, Ken Ford, James Gips, Clark Glymour, Antoni Gomila, Pat Hayes, Umar Khan, Henry Kyburg, Marvin Minsky, Anatol Rapoport, Herbert Simon, Christian Stary, and Lynn Stein.

# Network&Netplay

## Virtual Groups on the Internet

Edited by Fay Sudweeks,
Margaret McLaughlin and Sheizaf Rafaeli
Foreword by Ronald Rice

*Network and Netplay* addresses the mutual influences between information technology and group information and development, to assess the impact of computer-mediated communications on both work and play. Areas discussed include the growth and features of the Internet, network norms and experiences, and the essential nature of network communications.

6 x 9, 320 pp. ISBN 0-262-69206-5

*Prices higher outside the U.S. and subject to change without notice.*

**To order, call 800-356-0343 (US and Canada) or (617) 625-8569. Distributed by The MIT Press, Cambridge, MA 02142**