
A Generalized Reinforcement-Learning Model: Convergence and Applications

Michael L. Littman

Department of Computer Science
Brown University
Providence, RI 02912-1910, USA
mlittman@cs.brown.edu

Csaba Szepesvári*

Research Group of Artificial Intelligence
"József Attila" University, Szeged
Szeged 6720, Aradi vrt tere 1. HUNGARY
szepes@math.u-szeged.hu

Abstract

Reinforcement learning is the process by which an autonomous agent uses its experience interacting with an environment to improve its behavior. The Markov decision process (MDP) model is a popular way of formalizing the reinforcement-learning problem, but it is by no means the only way. In this paper, we show how many of the important theoretical results concerning reinforcement learning in MDPs extend to a generalized MDP model that includes MDPs, two-player games and MDPs under a worst-case optimality criterion as special cases. The basis of this extension is a stochastic-approximation theorem that reduces asynchronous convergence to synchronous convergence.

1 INTRODUCTION

Reinforcement learning is the process by which an agent improves its behavior in an environment via experience. A *reinforcement-learning scenario* is defined by the experience presented to the agent at each step, and the criterion for evaluating the agent's behavior.

One particularly well-studied reinforcement-learning scenario is that of a single agent maximizing expected discounted total reward in a finite-state environment; in this scenario experiences are of the form $\langle x, a, y, r \rangle$, with state x , action a , resulting state y , and the agent's scalar immediate reward r . A discount parameter

$0 \leq \gamma < 1$ controls the degree to which future rewards are significant compared to immediate rewards.

The theory of Markov decision processes has been used as a theoretical foundation for important results concerning this reinforcement-learning scenario. A (finite) Markov decision process (MDP) is defined by the tuple $\langle S, A, P, R \rangle$, where S is a finite set of states, A a finite set of actions, P a transition function, and R a reward function. The optimal behavior for an agent in an MDP depends on the optimality criterion; for the infinite-horizon expected discounted total-reward criterion, the optimal behavior can be found by identifying the optimal value function, defined recursively by

$$V^*(x) = \max_a \left(R(x, a) + \gamma \sum_y P(x, a, y) V^*(y) \right),$$

for all states $x \in S$, where $R(x, a)$ is the immediate reward for taking action a from state x , γ the discount factor, and $P(x, a, y)$ the probability that state y is reached from state x when action $a \in A$ is chosen. These simultaneous equations, known as the *Bellman equations*, can be solved using a variety of techniques ranging from successive approximation to linear programming (Puterman, 1994).

In the absence of complete information regarding the transition and reward functions, reinforcement-learning methods can be used to find optimal value functions. Researchers have explored model-free (direct) methods, such as Q-learning (Watkins and Dayan, 1992), and model-based (indirect) methods, such as prioritized sweeping (Moore and Atkeson, 1993), and many converge to optimal value functions under the proper conditions (Tsitsiklis, 1994; Jaakkola et al., 1994; Gullapalli and Barto, 1994).

Not all reinforcement-learning scenarios of interest can

*Also Department of Adaptive Systems, Joint Department of the "József Attila" University, Szeged and the Institute of Isotopes of the Hungarian Academy of Sciences, Budapest 1525, P.O. Box. 77. HUNGARY

be modeled as MDPs. For example, a great deal of reinforcement-learning research has been directed to the problem of solving two-player games (e.g. Tesauro, 1995), and the reinforcement-learning algorithms for solving MDPs and their convergence proofs do not apply directly to games.

In one form of two-player game, experiences are of the form $\langle x, a, y, r \rangle$, where states x and y contain additional information concerning which player (*maximizer* or *minimizer*) gets to choose the action in that state. There are deep similarities between MDPs and this type of game; for example, it is possible to define a set of Bellman equations for the optimal minimax value of a two-player zero-sum game,

$$V^*(x) = \begin{cases} \max_{a \in A} (R(x, a) + \gamma \sum_y P(x, a, y) V^*(y)), & \text{if maximizer moves in } x \\ \min_{a \in A} (R(x, a) + \gamma \sum_y P(x, a, y) V^*(y)), & \text{if minimizer moves in } x, \end{cases}$$

where $R(x, a)$ is the reward to the maximizing player. When $0 \leq \gamma < 1$, these equations have a unique solution and can be solved by successive-approximation methods. In addition, we show that simple extensions of several reinforcement-learning algorithms for MDPs converge to optimal value functions in these games.

In this paper, we introduce a generalized Markov decision process model with applications to reinforcement learning, and list some important results concerning the model. Generalized MDPs provide a foundation for the use of reinforcement learning in MDPs and games, as well as in risk-sensitive reinforcement learning (Heger, 1994), exploration-sensitive reinforcement learning (John, 1995), and reinforcement learning in simultaneous-action games (Littman, 1994). Our main theorem addresses conditions for the convergence of asynchronous stochastic processes and shows how these conditions relate to conditions for convergence of a corresponding synchronous process; it can be used to prove the convergence of model-free and model-based reinforcement-learning algorithms under a variety of reinforcement-learning scenarios.

In Section 2, we present generalized MDPs and motivate their form via two detailed examples. In Section 3, we describe a stochastic-approximation theorem, and in Section 4 we show several applications of the theorem that prove the convergence of learning processes in generalized MDPs.

2 THE GENERALIZED MODEL

In this section, we introduce our generalized MDP model. We begin by summarizing some of the more significant results regarding the standard MDP model and some important results for two-player games.

2.1 MARKOV DECISION PROCESSES

To provide a point of departure for our generalization of Markov decision processes, we first describe the use of reinforcement learning in the MDPs; proofs of the unattributed claims can be found in Puterman's (1994) MDP book. The ultimate target of learning is an optimal policy. A *policy* is some function that tells the agent which actions should be chosen under which circumstances. A policy π is *optimal* under the expected discounted total reward criterion if, with respect to the space of all possible policies, π maximizes the expected discounted total reward from all states.

Directly maximizing over the space of all possible policies is impractical. However, MDPs have an important property that makes it unnecessary to consider such a broad space of possibilities. We say a policy π is *stationary* and *deterministic* if it maps directly from states to actions, ignoring everything else, and we write $\pi(x)$ as the action chosen by π when the current state is x . In expected discounted total reward MDP environments, there is always a stationary deterministic policy that is optimal; we will use the word "policy" to mean stationary deterministic policy, unless otherwise stated.

The value function for a policy π , V^π , maps states to their expected discounted total reward under policy π . It can be defined by the simultaneous equations

$$V^\pi(x) = R(x, \pi(x)) + \gamma \sum_y P(x, \pi(x), y) V^\pi(y),$$

for all $x \in S$. The optimal value function V^* is the value function of an optimal policy; it is unique for $0 \leq \gamma < 1$. The *myopic policy* with respect to a value function V is the policy π_V such that

$$\pi_V(x) = \arg \max_a \left(R(x, a) + \gamma \sum_y P(x, a, y) V(y) \right).$$

Any myopic policy with respect to the optimal value function is optimal.

The Bellman equations can be operationalized in the form of the dynamic-programming operator T , which

maps value functions to value functions:

$$[TV](x) = \max_a \left(R(x, a) + \gamma \sum_y P(x, a, y) V(y) \right).$$

For $0 \leq \gamma < 1$, successive applications of T to a value function bring it closer and closer to the optimal value function V^* , which is the unique fixed point of T : $V^* = TV^*$.

In reinforcement learning, R and P are not known in advance. In model-based reinforcement learning, R and P are estimated on-line, and the value function is updated according to the approximate dynamic-programming operator derived from these estimates; this algorithm converges to the optimal value function under a wide variety of choices of the order states are updated (Gullapalli and Barto, 1994).

The method of Q-learning (Watkins and Dayan, 1992) uses experience to estimate the optimal value function without ever explicitly approximating R and P . The algorithm estimates the optimal Q function

$$Q^*(x, a) = R(x, a) + \gamma \sum_y P(x, a, y) V^*(y),$$

from which the optimal value function can be computed via $V^*(x) = \max_a Q^*(x, a)$. Given the experience at step t $\langle x_t, a_t, y_t, r_t \rangle$ and the current estimate $Q_t(x, a)$ of the optimal Q function, Q-learning updates

$$Q_{t+1}(x_t, \mathbf{a}_t) := (1 - \alpha_t(x_t, \mathbf{a}_t)) Q_t(x_t, \mathbf{a}_t) + \alpha_t(x_t, a_t) (r_t + \gamma \max_a Q_t(y_t, a)),$$

where $0 \leq \alpha_t(x, \mathbf{a}) \leq 1$ is a time-dependent learning rate controlling the blending rate of new estimates with old estimates for each state-action pair. The estimated Q function converges to Q^* under the proper conditions (Watkins and Dayan, 1992).

2.2 ALTERNATING MARKOV GAMES

In alternating Markov games, two players take turns issuing actions to maximize their expected discounted total reward. The model is defined by the tuple $\langle S_1, S_2, A, B, P, R \rangle$, where S_1 is the set of states in which player 1 issues actions from the set A , S_2 is the set of states in which player 2 issues actions from the set B , P is the transition function, and R is the reward function for player 1. In the zero-sum games we consider, the rewards to player 2 (the minimizer) are simply the additive inverse of the rewards to player 1 (the maximizer). Markov decision processes are a special case of alternating Markov games in which $S_2 = \emptyset$;

Condon (1992) proves this and the other unattributed results in this section.

A popular optimality criterion for alternating Markov games is discounted minimax optimality. Under this criterion, the maximizer chooses actions to maximize its reward against the minimizer's best possible counter-policy. A pair of policies is in *equilibrium* if neither player has any incentive to change policies if the other player's policy remains fixed. The value function for a pair of equilibrium policies is the optimal value function for the game; it is unique when $0 \leq \gamma < 1$, and can be found by successive approximation. For both players, there is always a deterministic stationary optimal policy. Any myopic policy with respect to the optimal value function is optimal.

Dynamic-programming operators, Bellman equations, and reinforcement-learning algorithms can be defined for alternating Markov games by starting with the definitions used in MDPs and changing the maximum operators to either maximums or minimums conditioned on the state. We show below that the resulting algorithms share their convergence properties with the analogous algorithms for MDPs.

2.3 GENERALIZED MDPs

In alternating Markov games and MDPs, optimal behavior can be specified by the Bellman equations; any myopic policy with respect to the optimal value function is optimal. In this section, we generalize the Bellman equations to define optimal behavior for a broad class of reinforcement-learning models. The objective criterion used in these models is additive in that the value of a policy is some measure of the *total* reward received.

The generalized Bellman equations can be written

$$V^*(x) = \bigotimes_a^x \left(R(x, a) + \gamma \bigoplus_y^{x, a} V^*(y) \right). \quad (1)$$

Here " \bigotimes_a^x " is an operator that summarizes values over actions as a function of the state, and " $\bigoplus_y^{x, a}$ " is an operator that summarizes values over next states as a function of the state and action. For Markov decision processes, $\bigotimes_a^x f(x, \mathbf{a}) = \max_a f(x, \mathbf{a})$ and $\bigoplus_y^{x, a} g(y) = \sum_y P(x, a, y) g(y)$. For alternating Markov games, $\bigoplus_y^{x, a}$ is the same and $\bigotimes_a^x f(x, a) = \max_a f(x, \mathbf{a})$ or $\min_a f(x, a)$ depending whether x is in S_1 or S_2 . Many models can be represented in this framework; see Section 4.

From a reinforcement-learning perspective, the value

functions defined by the generalized MDP model can be interpreted as the total value of the rewards received by an agent selecting actions in a stochastic environment. The agent begins in state x , takes action \mathbf{a} , and ends up in state y . The $\bigoplus_y^{x,\mathbf{a}}$ operator defines how the value of the next state should be used in assigning value to the current state. The \bigotimes_a^x operator defines how an optimal agent should choose actions.

When $0 \leq \gamma < 1$ and \bigotimes^x and $\bigoplus^{x,\mathbf{a}}$ are non-expansions, the generalized Bellman equations have a unique optimal solution, and therefore, the optimal value function is well defined. The \bigotimes^x operator is a non-expansion if

$$\left| \bigotimes_a^x f_1(x, \mathbf{a}) - \bigotimes_a^x f_2(x, \mathbf{a}) \right| \leq \max_a |f_1(x, \mathbf{a}) - f_2(x, \mathbf{a})|$$

for all f_1, f_2 , and x . An analogous condition defines when $\bigoplus^{x,\mathbf{a}}$ is a non-expansion.

Many natural operators are non-expansions, such as max, min, midpoint, median, mean, and fixed weighted averages of these operations. Several previously described reinforcement-learning scenarios are special cases of this generalized MDP model including computing the expected return of a fixed policy (Sutton, 1988), finding the optimal risk-averse policy (Heger, 1994), and finding the optimal exploration-sensitive policy (John, 1995).

As with MDPs, we can define a dynamic-programming operator

$$[TV](x) = \bigotimes_a^x \left(R(x, a) + \gamma \bigoplus_y^{x,\mathbf{a}} V(y) \right). \quad (2)$$

The operator T is a contraction mapping for $0 \leq \gamma < 1$. This means

$$\sup_x |[TV_1](x) - [TV_2](x)| \leq \gamma \sup_x |V_1(x) - V_2(x)|$$

where V_1 and V_2 are arbitrary functions and $0 \leq \gamma < 1$ is the index of contraction.

We can define a notion of stationary myopic policies with respect to a value function V ; it is any (stochastic) policy π_V for which $T^\pi V = TV$ where

$$[T^\pi V](x) = \sum_a \pi_a(x) \left(R(x, a) + \gamma \bigoplus_y^{x,\mathbf{a}} V(y) \right).$$

Here $\pi_a(x)$ represents the probability that an agent following π would choose action a in state x . To be certain that every value function possesses a myopic

policy, we require that the operator \bigotimes^x satisfy the following property: for all functions f and states x , $\min_a f(x, \mathbf{a}) \leq \bigotimes_a^x f(x, \mathbf{a}) \leq \max_a f(x, \mathbf{a})$.

The value function with respect to a policy π , V^π , can be defined by the simultaneous equations $V^\pi = T^\pi V^\pi$; it is unique when T is a contraction mapping. A policy π is optimal if it is myopic with respect to its own value function. If π^* is an optimal policy, then $V^{\pi^*} = V^*$ because it solves the Bellman equation: $V^{\pi^*} = T^{\pi^*} V^{\pi^*} = TV^{\pi^*}$.

The next section describes a general theorem that can be used to prove the convergence of several reinforcement-learning algorithms for these and other models.

3 CONVERGENCE THEOREM

The process of finding an optimal value function can be viewed in the following general way. At any moment in time, there is a set of values representing the current approximation of the optimal value function. On each iteration, we apply some dynamic-programming operator, perhaps modified by experience, to the current approximation to generate a new approximation. Over time, we would like the approximation to tend toward the optimal value function.

In this process, there are two types of approximation going on simultaneously. The first is an approximation of the dynamic-programming operator for the underlying model, and the second is the use of the approximate dynamic-programming operator to find the optimal value function. This section presents a theorem that gives a set of conditions under which this type of simultaneous stochastic approximation converges to an optimal value function.

First, we need to define the general stochastic process. Let the set X be the states of the model, and the set $\mathcal{B}(X)$ of bounded, real-valued functions over X be the set of value functions. Let $T : \mathcal{B}(X) \rightarrow \mathcal{B}(X)$ be an arbitrary contraction mapping with fixed point V^* .

If we had direct access to the contraction mapping T , we could use it to successively approximate V^* . In most reinforcement-learning scenarios, T is not available and we must use experience to construct approximations of T . Consider a sequence of random operators $T_t : \mathcal{B}(X) \rightarrow (\mathcal{B}(X) \rightarrow \mathcal{B}(X))$ and define $U_{t+1} = [T_t U_t]V$ where V and $U_0 \in \mathcal{B}(X)$ are arbitrary value functions. We say T_t approximates T at V , if U_t converges to TV with probability 1 uniformly

over X^1 . The idea is that T_t is a randomized version of T that uses U_t as “memory” to converge to TV .

The following theorem shows that, under the proper conditions, we can use the sequence T_t to estimate the fixed point V^* of T .

Theorem 1 *Let T be an arbitrary mapping with fixed point V^* , and let T_t approximate T at V^* . Let V_t be an arbitrary value function, and define $V_{t+1} = [T_t V_t] V_t$. If there exist functions $0 \leq F_t(x) \leq 1$ and $0 \leq G_t(x) \leq 1$ satisfying the conditions below with probability one, then V_t converges to V^* with probability 1 uniformly over X :*

1. for all U_1 and $U_2 \in \mathcal{B}(X)$, and all $x \in X$,

$$\begin{aligned} & |([T_t U_1] V^*)(x) - ([T_t U_2] V^*)(x)| \\ & \leq G_t(x) |U_1(x) - U_2(x)|; \end{aligned}$$

2. for all U and $V \in \mathcal{B}(X)$, and all $x \in X$,

$$\begin{aligned} & |([T_t U] V^*)(x) - ([T_t U] V)(x)| \\ & \leq F_t(x) \sup_{x'} |V^*(x') - V(x')|; \end{aligned}$$

3. for all $k > 0$, $\prod_{t=k}^n G_t(x)$ converges to zero uniformly in x as n increases; and,

4. there exists $0 \leq \gamma < 1$ such that for all $x \in X$ and large enough t ,

$$F_t(x) \leq \gamma(1 - G_t(x)).$$

Note that from the conditions of the theorem, it follows that T is a contraction operator at V^* with index of contraction γ . The theorem is proven in a more detailed version of this paper (Szepesvári and Littman, 1996). We next describe some of the intuition behind the statement of the theorem and its conditions.

The iterative approximation of V^* is performed by computing $V_{t+1} = [T_t V_t] V_t$, where T_t approximates T with the help of the “memory” present in V_t . Because of Conditions 1 and 2, $G_t(x)$ is the extent to which the estimated value function depends on its present value and $F_t(x) \approx 1 - G_t(x)$ is the extent to which the estimated value function is based on “new” information (this reasoning becomes clearer in the context of the applications in Section 4).

¹A sequence of functions f_n converges to f^* with probability 1 uniformly over X if, for the events w for which $f_n(w, x) \rightarrow f^*$, the convergence is uniform in x .

In some applications, such as Q-learning, the contribution of new information needs to decay over time to insure that the process converges. In this case, $G_t(x)$ needs to converge to one; Condition 3 allows this as long as the convergence is slow enough to incorporate sufficient information for the process to converge.

Condition 4 links the values of $G_t(x)$ and $F_t(x)$ through some quantity $\gamma < 1$. If it were somehow possible to update the values synchronously over the entire state space, the process would converge to V^* even when $\gamma = 1$. In the more interesting asynchronous case, when $\gamma = 1$ the long-term behavior of V_t is not immediately clear; it may even be that V_t converges to something other than V^* . The requirement that $\gamma < 1$ insures that the use of outdated information in the asynchronous updates does not cause a problem in convergence.

One of the most noteworthy aspects of this theorem is that it shows how to reduce the problem of approximating V^* to the problem of approximating T at a particular point V (in particular, it is enough if T can be approximated at V^*); in many cases, the latter is much easier to achieve and also to prove. For example, the theorem makes the convergence of Q-learning a consequence of the classical Robbins-Monro theorem (Robbins and Monro, 1951).

4 APPLICATIONS

This section makes use of Theorem 1 to prove the convergence of various reinforcement-learning algorithms.

4.1 GENERALIZED Q-LEARNING FOR EXPECTED VALUE MODELS

Consider the family of finite state and action generalized MDPs defined by the Bellman equations

$$V^*(x) = \bigotimes_a^x \left(R(x, a) + \gamma \sum_y P(x, a, y) V^*(y) \right)$$

where the definition of \bigotimes^x does not depend on R or P . A Q-learning algorithm for this class of models can be defined as follows. Given experience $\langle x_t, a_t, y_t, r_t \rangle$ at time t and an estimate $Q_t(x, a)$ of the optimal Q function, let

$$\begin{aligned} Q_{t+1}(x_t, a_t) & := (1 - \alpha_t(x_t, a_t)) Q_t(x_t, a_t) \\ & + \alpha_t(x_t, a_t) \left(r_t + \gamma \bigotimes_a^x Q_t(y_t, a) \right). \end{aligned}$$

We can derive the assumptions necessary for this learning algorithm to satisfy the conditions of Theorem 1 and therefore converge to the optimal Q values. The dynamic-programming operator defining the optimal Q function is

$$[TQ](x, a) = R(x, a) + \gamma \sum_y P(x, a, y) \bigotimes_{a'}^x Q(y, a').$$

The randomized approximate dynamic-programming operator that gives rise to the Q-learning rule is

$$([T_t Q']Q)(x, a) = \begin{cases} (1 - \alpha_t(x, a))Q'(x, a) + \\ \alpha_t(x, a)(r_t + \gamma \bigotimes_{a'}^x Q(y_t, a')), \\ \text{if } x = x_t \text{ and } a = a_t \\ Q'(x, a), \text{ otherwise.} \end{cases}$$

If

- y_t is randomly selected according to the probability distribution defined by $P(x_t, a_t, \cdot)$,
- \bigotimes^x is a non-expansion, and both the expected value and the variance of $\bigotimes_{a'}^x Q(y_t, a)$ exist given the way y_t is sampled,
- r_t has finite variance and expected value given x_t and a_t equal to $R(x_t, a_t)$,
- the learning rates are decayed so that

$$\sum_t \chi(x_t = x, a_t = a) \alpha_t(x, a) = \infty$$

and $\sum_t \chi(x_t = x, a_t = a) \alpha_t(x, a)^2 < \infty$ with probability 1 uniformly over $X \times A^2$,

then a standard result from the theory of stochastic approximation (Robbins and Monro, 1951) states that T_t approximates T everywhere. That is, this method of using a decayed, exponentially weighted average correctly computes the average one-step reward.

$$\text{Let } G_t(x, a) = \begin{cases} 1 - \alpha_t(x, a), & \text{if } x = x_t \text{ and } a = a_t; \\ 1, & \text{otherwise,} \end{cases}$$

$$\text{and } F_t(x, a) = \begin{cases} \gamma \alpha_t(x, a), & \text{if } x = x_t \text{ and } a = a_t; \\ 0, & \text{otherwise.} \end{cases}$$

These functions satisfy the conditions of Theorem 1 (Condition 3 is implied by the restrictions placed on the sequence of learning rates α_t).

²This condition implies, among other things, that every state-action pair is updated infinitely often. Here, χ denotes the characteristic function.

Theorem 1 therefore implies that this generalized Q-learning algorithm converges to the optimal Q function with probability 1 uniformly over $X \times A$. The convergence of Q-learning for discounted MDPs and alternating Markov games follows trivially from this. Extensions of this result for undiscounted “all-policies-proper” MDPs (Bertsekas and Tsitsiklis, 1989), a soft state aggregation learning rule (Singh et al., 1995), and a “spreading” learning rule are given in a more detailed version of this paper (Szepesvári and Littman, 1996).

4.2 Q-LEARNING FOR MARKOV GAMES

Markov games are a generalization of MDPs and alternating Markov games in which both players simultaneously choose actions at each step. The basic model is defined by the tuple $\langle S, A, B, P, R \rangle$ and discount factor γ . As in alternating Markov games, the optimality criterion is one of discounted minimax optimality, but because the players move simultaneously, the Bellman equations take on a more complex form:

$$V^*(x) = \max_{\rho \in \Pi(A)} \min_{b \in B} \sum_{a \in A} \rho(a) \cdot \left(R(x, a, b) + \gamma \sum_{y \in S} P(x, a, b, y) V^*(y) \right).$$

In these equations, $R(x, a, b)$ is the immediate reward for the maximizer for taking action a in state x at the same time the minimizer takes action b , $P(x, a, b, y)$ is the probability that state y is reached from state x when the maximizer takes action a and the minimizer takes action b , and $\Pi(A)$ represents the set of discrete probability distributions over the set A . The sets S , A , and B are finite.

Once again, optimal policies are policies that are in equilibrium, and there is always a pair of optimal policies that are stationary. Unlike MDPs and alternating Markov games, the optimal policies are sometimes stochastic; there are Markov games in which no deterministic policy is optimal. The stochastic nature of optimal policies explains the need for the optimization over probability distributions in the Bellman equations, and stems from the fact that players must avoid being “second guessed” during action selection. An equivalent set of equations can be written with a stochastic choice for the minimizer, and also with the roles of the maximizer and minimizer reversed.

The Q-learning update rule for Markov games (Littman, 1994) given step t experience $\langle x_t, a_t,$

b_t, y_t, r_t has the form

$$Q_{t+1}(x_t, a_t, \mathbf{b}_t) := (1 - \alpha_t(x_t, a_t, \mathbf{b}_t))Q_t(x_t, a_t, \mathbf{b}_t) + \alpha_t(x_t, a_t, \mathbf{b}_t) \left(r_t + \gamma \bigotimes_{a,b}^x Q_t(y_t, a, b) \right),$$

where

$$\bigotimes_{a,b}^x f(x, a, b) = \max_{\rho \in \Pi(A)} \min_{b \in \mathcal{B}} \sum_{a \in \mathcal{A}} \rho(a) f(x, a, b).$$

The results of the previous section prove that this rule converges to the optimal Q function under the proper conditions.

4.3 RISK-SENSITIVE MODELS

Heger (1994) described an optimality criterion for MDPs in which only the *worst* possible value of the next state makes a contribution to the value of a state. An optimal policy under this criterion is one that avoids states for which a bad outcome is possible, even if it is not probable; for this reason, the criterion has a risk-averse quality to it. The generalized Bellman equations for this criterion are

$$V^*(x) = \bigotimes_a^x \left(R(x, a) + \gamma \min_{y: P(x, a, y) > \bullet} V^*(y) \right).$$

The argument in Section 4.5 shows that model-based reinforcement learning can be used to find optimal policies in risk-sensitive models, as long as \bigotimes^x does not depend on R or P , and P is estimated in a way that preserves its zero vs. non-zero nature in the limit.

For the model in which $\bigotimes_a^x f(x, a) = \max_a f(x, a)$, Heger defined a Q-learning-like algorithm that converges to optimal policies without estimating R and P online. In essence, the learning algorithm uses an update rule analogous to the rule in Q-learning with the additional requirement that the initial Q function be set optimistically; that is, $Q_0(x, a)$ must be larger than $Q^*(x, a)$ for all x and a .

Using Theorem 1 it is possible to prove the convergence of a generalization of Heger's algorithm to models where $\bigotimes_a^x f(x, a) = f(x, a^*(f, x))$ for some function $a^*(\cdot)$; that is, as long as the summary value of $f(x, a)$ is equal to $f(x, a^*)$ for some a^* . The proof is based on estimating the Q-learning algorithm from above by an appropriate process where the Q function is updated only if the received experience tuple is an extremity according to the optimality equation; details are given elsewhere (Szepesvári and Littman, 1996).

4.4 EXPLORATION-SENSITIVE MODELS

John (1995) considered the implications of insisting that reinforcement-learning agents keep exploring forever; he found that better learning performance can be achieved if the Q-learning rule is changed to incorporate the condition of persistent exploration. In John's formulation, the agent is forced to adopt a policy from a restricted set; in one example, the agent must choose a stochastic stationary policy that selects actions at random 5% of the time.

This approach requires that the definition of optimality be changed to reflect the restriction on policies. The optimal value function is given by $V^*(x) = \sup_{\pi \in \mathcal{P}_\bullet} V^\pi(x)$, where \mathcal{P}_\bullet is the set of permitted (stationary) policies, and the associated Bellman equations are

$$V^*(x) = \sup_{\pi \in \mathcal{P}_\bullet} \sum_a \pi_a(x) (R(x, a) + \gamma \sum_y P(x, a, y) V^*(y)),$$

which corresponds to a generalized MDP model with $\bigoplus_y^{x,a} g(y) = \sum_y P(x, a, y) g(y)$ and $\bigotimes_a^x f(x, a) = \sup_{\pi \in \mathcal{P}_\bullet} \sum_a \pi_a(x) f(x, a)$. Because $\pi_a(x)$ is a probability distribution over a for any given state x , \bigotimes^x is a non-expansion and, thus, the convergence of the associated Q-learning algorithm follows from the arguments in Section 4.1. As a result, John's learning rule gives the optimal policy under the revised optimality criterion.

4.5 MODEL-BASED METHODS

The defining assumption in reinforcement learning is that the reward and transition functions, R and P , are not known in advance. Although Q-learning shows that optimal value functions can be estimated without ever explicitly learning R and P , learning R and P makes more efficient use of experience at the expense of additional storage and computation (Moore and Atkeson, 1993). The parameters of R and P can be gleaned from experience by keeping statistics for each state-action pair on the expected reward and the proportion of transitions to each next state. In model-based reinforcement learning, R and P are estimated on-line, and the value function is updated according to the approximate dynamic-programming operator derived from these estimates. Theorem 1 implies the convergence of a wide variety of model-based reinforcement-learning methods.

The dynamic-programming operator defining the optimal value for generalized MDPs is given in Equation 2. Here we assume that $\bigoplus^{x,a}$ may depend on P and/or

R , but \otimes^x may not. It is possible to extend the following argument to allow \otimes^x to depend on P and R as well. In model-based reinforcement learning, R and P are estimated by the quantities R_t and P_t , and $\oplus^{x,a,t}$ is an estimate of the $\oplus^{x,y}$ operator defined using R_t and P_t . As long as every state-action pair is visited infinitely often, there are a number of simple methods for computing R_t and P_t that converge to R and P . A bit more care is needed to insure that $\oplus^{x,a,t}$ converges to $\oplus^{x,a}$, however. For example, in expected-reward models, $\oplus_y^{x,a} g(y) = \sum_y P(x,a,y)g(y)$ and the convergence of P_t to P guarantees the convergence of $\oplus^{x,a,t}$ to $\oplus^{x,a}$. On the other hand, in a risk-sensitive model, $\oplus_y^{x,a} g(y) = \min_{y:P(x,a,y)>0} g(y)$ and it is necessary to approximate P in a way that insures that the set of y such that $P_t(x,a,y) > \epsilon$ converges to the set of y such that $P(x,a,y) > 0$. This can be accomplished easily, for example, by setting $P_t(x,a,y) = \epsilon$ if no transition from x to y under a has been observed.

Assuming P and R are estimated in a way that results in the convergence of $\oplus^{x,a,t}$ to $\oplus^{x,a}$, the sequence of dynamic-programming operators T_t defined by

$$([T_t U]V)(x) = \begin{cases} \otimes_{\bullet}^x (R_t(x,a) + \gamma \oplus_y^{x,a,t} V(y)), & \text{if } x \in \tau_t \\ U(x), & \text{otherwise,} \end{cases}$$

approximates T for all value functions. The set $\tau_t \subseteq S$ represents the set of states whose values are updated on step t ; one popular choice is to set $\tau_t = \{x_t\}$.

The functions

$$G_t(x) = \begin{cases} 0, & \text{if } x \in \tau_t; \\ 1, & \text{otherwise,} \end{cases}$$

and

$$F_t(x) = \begin{cases} \gamma, & \text{if } x \in \tau_t; \\ \bullet, & \text{otherwise,} \end{cases}$$

satisfy the conditions of Theorem 1 as long as each x is in infinitely many τ_t sets (Condition 3) and the discount factor γ is less than 1 (Condition 4).

As a consequence of this argument and Theorem 1, model-based methods can be used to find optimal policies in MDPs, alternating Markov games, Markov games, risk-sensitive MDPs, and exploration-sensitive MDPs. Also, letting $R_t = R$ and $P_t = P$ for all t , this result implies that real-time dynamic programming (Barto et al., 1995) converges to the optimal value function.

5 CONCLUSIONS

In this paper, we presented a generalized model of Markov decision processes, and proved the convergence of several reinforcement-learning algorithms in the generalized model.

Other Results We have derived a collection of results (Szepesvári and Littman, 1996) for the generalized MDP model that demonstrate its general applicability: the Bellman equations can be solved by value iteration; a myopic policy with respect to an approximately optimal value function gives an approximately optimal policy; when \otimes^x has a particular ‘‘maximization’’ property, policy iteration converges to the optimal value function, and, for models with the maximization property and finite state and action spaces, both value iteration and policy iteration identify optimal policies in pseudopolynomial time.

Related Work The work presented here is closely related to several previous research efforts. Szepesvári (1995) described a related generalized reinforcement-learning model and presented conditions under which there is an optimal (stationary) policy that is myopic with respect to the optimal value function.

Tsitsiklis (1994) developed the connection between stochastic-approximation theory and reinforcement learning in MDPs. Our work is similar in spirit to that of Jaakkola, Jordan, and Singh (1994). We believe the form of Theorem 1 makes it particularly convenient for proving the convergence of reinforcement-learning algorithms; our theorem reduces the proof of the convergence of an asynchronous process to a simpler proof of convergence of a corresponding synchronized one. This idea enables us to prove the convergence of asynchronous stochastic processes whose underlying synchronous process is not of the Robbins-Monro type (e.g., risk-sensitive MDPs, model-based algorithms, etc.).

Future Work There are many areas of interest in the theory of reinforcement learning that we would like to address in future work. The results in this paper primarily concern reinforcement-learning in contractive models ($\gamma < 1$ or all-policies-proper), and there are important non-contractive reinforcement-learning scenarios, for example, reinforcement learning under an average-reward criterion (Mahadevan, 1996). It would be interesting to develop a TD(λ) algorithm (Sutton, 1988) for generalized MDPs. Theorem 1 is not re-

stricted to finite state spaces, and it might be valuable to prove the convergence of a reinforcement-learning algorithm for a infinite state-space model.

Conclusion By identifying common elements among several reinforcement-learning scenarios, we created a new class of models that generalizes existing models in an interesting way. In the generalized framework, we replicated the established convergence proofs for reinforcement learning in Markov decision processes, and proved new results concerning the convergence of reinforcement-learning algorithms in game environments, under a risk-sensitive assumption, and under an exploration-sensitive assumption. At the heart of our results is a new stochastic-approximation theorem that is easy to apply to new situations.

Acknowledgements

Research supported by PHARE H9305-02/1022 and OTKA Grant no. F020132 and by Bellcore's Support for Doctoral Education Program.

References

- Barto, A. G., Bradtke, S. J., and Singh, S. P. (1995). Learning to act using real-time dynamic programming. *Artificial Intelligence*, 72(1):81–138.
- Bertsekas, D. P. and Tsitsiklis, J. N. (1989). *Parallel and Distributed Computation: Numerical Methods*. Prentice-Hall, Englewood Cliffs, NJ.
- Condon, A. (1992). The complexity of stochastic games. *Information and Computation*, 96(2):203–224.
- Gullapalli, V. and Barto, A. G. (1994). Convergence of indirect adaptive asynchronous value iteration algorithms. In Cowan, J. D., Tesauro, G., and Alspector, J., editors, *Advances in Neural Information Processing Systems 6*, pages 695–702, San Mateo, CA. Morgan Kaufmann.
- Heger, M. (1994). Consideration of risk in reinforcement learning. In *Proceedings of the Eleventh International Conference on Machine Learning*, pages 105–111, San Francisco, CA. Morgan Kaufmann.
- Jaakkola, T., Jordan, M. I., and Singh, S. P. (1994). On the convergence of stochastic iterative dynamic programming algorithms. *Neural Computation*, 6(6).
- John, G. H. (1995). When the best move isn't optimal: Q-learning with exploration. Unpublished manuscript.
- Littman, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the Eleventh International Conference on Machine Learning*, pages 157–163, San Francisco, CA. Morgan Kaufmann.
- Mahadevan, S. (1996). Average reward reinforcement learning: Foundations, algorithms, and empirical results. *Machine Learning*, 22(1/2/3):159–196.
- Moore, A. W. and Atkeson, C. G. (1993). Prioritized sweeping: Reinforcement learning with less data and less real time. *Machine Learning*, 13.
- Puterman, M. L. (1994). *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY.
- Robbins, H. and Monro, S. (1951). A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407.
- Singh, S., Jaakkola, T., and Jordan, M. (1995). Reinforcement learning with soft state aggregation. In Tesauro, G., Touretzky, D. S., and Leen, T. K., editors, *Advances in Neural Information Processing Systems 7*, Cambridge, MA. The MIT Press.
- Sutton, R. S. (1988). Learning to predict by the method of temporal differences. *Machine Learning*, 3(1):9–44.
- Szepesvári, C. (1995). General framework for reinforcement learning. In *Proceedings of ICANN'95 Paris*.
- Szepesvári, C. and Littman, M. L. (1996). Generalized Markov decision processes: Dynamic programming and reinforcement-learning algorithms. Technical Report CS-96-11, Brown University, Providence, RI.
- Tesauro, G. (1995). Temporal difference learning and TD-Gammon. *Communications of the ACM*, pages 58–67.
- Tsitsiklis, J. N. (1994). Asynchronous stochastic approximation and Q-learning. *Machine Learning*, 16(3).
- Watkins, C. J. C. H. and Dayan, P. (1992). Q-learning. *Machine Learning*, 8(3):279–292.