



Value-function reinforcement learning in Markov games

Action editor: Ron Sun

Michael L. Littman

AT&T Labs Research, 180 Park Avenue, Florham Park, NJ 07932-0971, USA

Received 1 October 2000; accepted 10 October 2000

Abstract

Markov games are a model of multiagent environments that are convenient for studying multiagent reinforcement learning. This paper describes a set of reinforcement-learning algorithms based on estimating value functions and presents convergence theorems for these algorithms. The main contribution of this paper is that it presents the convergence theorems in a way that makes it easy to reason about the behavior of simultaneous learners in a shared environment. © 2001 Published by Elsevier Science B.V.

Keywords: Reinforcement learning; Temporal difference learning; Value functions; Game theory; Markov games; Q -learning; Nash equilibria

1. Introduction

Game theory (von Neumann & Morgenstern, 1947) provides a powerful set of conceptual tools for reasoning about behavior in multiagent environments. Markov games (van der Wal, 1981), or stochastic games (Owen, 1982; Shapley, 1953), are a formalization of temporally extended agent interaction.

Reinforcement learning (Kaelbling, Littman & Moore, 1996; Sutton & Barto, 1998) is the problem of an agent learning to behave from experience. One well studied approach to reinforcement learning is

building a value function — a mapping from state to expected reward. Several authors have applied value-function reinforcement learning to Markov games to create agents that learn from experience how to best interact with other agents. This paper presents several value-function reinforcement-learning algorithms and what is known about how they behave when learning simultaneously in different types of games.

Section 2 describes single-agent environments and the basic Q -learning algorithm, which converges to an optimal value function and optimal behavior in this type of environment. Section 3 examines multiagent environments and the Nash Q -learning algorithm. Section 4 looks at one situation in which Nash Q -learning converges — when there are adversarial equilibria; Section 5 examines another — when there

E-mail address: mlittman@research.att.com (M.L. Littman).

are coordination equilibria. Section 6 presents some concluding thoughts.

2. Single-agent environments

Markov decision processes (MDPs) (Bellman, 1957; Howard, 1960) are a descriptive model of single-agent environments. In the MDP framework, it is assumed that, although there may be a great deal of uncertainty about the effects of an agent's actions, there is never any uncertainty about the agent's current state — it has complete and perfect perceptual abilities.

2.1. Markov decision processes

Mathematically, a Markov decision process is a tuple $\langle \mathcal{S}, \mathcal{A}, T, R, \beta \rangle$, where

- \mathcal{S} is a finite set of *states* of the environment;
- \mathcal{A} is a finite set of *actions* available to the agent;
- $T: \mathcal{S} \times \mathcal{A} \rightarrow \Pi(\mathcal{S})$ is the *transition function*, giving for each state and agent action, a probability distribution over states ($T(s, a, s')$ is the probability of ending in state s' , given that the agent starts in state s and takes action a);
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, giving the expected immediate reward gained by the agent for taking each action in each state ($R(s, a)$ is the expected reward for taking action a in state s); and
- $0 \leq \beta < 1$ is a discount factor.

This paper considers only models with finite state and action spaces.

In an MDP, agents should act in such a way as to maximize some measure of their long-run reward received. Under the *discounted objective*, which is the focus of this paper, the discount factor $0 \leq \beta < 1$ controls how much effect future rewards have on the decisions at each moment, with small values of β emphasizing near-term gain and larger values giving significant weight to later situations. Concretely, a reward of r received t steps in the future is worth $\beta^t r$ to the agent now. Mathematically, the discount factor has the desirable property that if all immediate rewards are bounded, then the infinite sum of the

discounted rewards is also bounded. From an applications perspective, the discount factor can be thought of as the probability that the agent will be allowed to continue gathering reward after the current step, or, from an economic perspective, as an inverse interest rate on reward (Puterman, 1994).

A policy is a description of the behavior of an agent. A *stationary policy*, $\pi: \mathcal{S} \rightarrow \Pi(\mathcal{A})$, specifies, for each state, a probability distribution over actions to be taken. Notationally, $\pi(s, a)$ is the probability assigned to action a in state s . A *deterministic* policy is one that assigns probability 1 to some action in each state. Every MDP has a deterministic stationary optimal policy (Bertsekas, 1987).

A policy π for an agent can be evaluated by computing the long-run value the agent can expect to gain. Let $Q^\pi(s, a)$ be the expected discounted future reward to the agent for starting in state s and executing action a for one step, then continuing according to policy π . This can be defined by a set of simultaneous linear equations, one for each state s :

$$Q^\pi(s, a) = R(s, a) + \beta \sum_{s' \in \mathcal{S}} T(s, a, s') \times \sum_{a' \in \mathcal{A}} \pi(s', a') Q^\pi(s', a'). \quad (1)$$

The function Q^π is called the *Q-function* for π .

Given an initial state s , the agent should execute a policy π that maximizes $\sum_a \pi(s, a) Q^\pi(s, a)$. Howard (1960) showed that there exists a stationary deterministic policy π^* that is optimal for every starting state. The *Q-function* for this policy, written Q^* , is defined by the set of equations

$$Q^*(s, a) = R(s, a) + \beta \sum_{s' \in \mathcal{S}} T(s, a, s') \max_{a' \in \mathcal{A}} Q^*(s', a'), \quad (1)$$

and the *greedy policy* that assigns probability one to action $\operatorname{argmax}_a Q^*(s, a)$ in state s is optimal (Puterman, 1994).

The presence of the maximization operator in Eq. (1) means the system of equations is not linear. Methods such as value iteration, policy iteration, linear programming, and modified policy iteration can be used to solve the equations (Puterman, 1994).

Barto, Sutton and Watkins (1989) argue that MDPs are an appropriate model for studying re-

137 inforcement learning in single-agent environments.
 138 The next section describes a reinforcement-learning
 139 algorithm with guaranteed performance in MDP
 140 environments.

141 2.2. *Q-learning*

142 *Q-learning* (Watkins, 1989; Watkins & Dayan,
 143 1992) can be viewed as a sampled, asynchronous
 144 method for estimating the optimal *Q*-function for an
 145 unknown MDP. *Q-learning* is a temporal-difference
 146 learning method (Sutton, 1988), and its basic version
 147 keeps a table of values, $Q[s, a]$, with an entry for
 148 each state/action pair. The entry $Q[s, a]$ is an
 149 estimate for $Q^*(s, a)$ as defined in Eq. (1). An agent
 150 uses its experience to improve its estimate, blending
 151 new information into its prior experience according
 152 to a *learning rate* $0 < \alpha < 1$.

153 The *Q*-function is an ideal data structure for
 154 reinforcement learning. There are three fundamental
 155 functions in algorithms for solving MDPs: the value
 156 function V mapping state to value, the *Q*-function Q
 157 mapping state and action to value, and the policy π
 158 mapping state to a probability distribution over
 159 actions. Given a model in the form of transition and
 160 reward functions, any of the mappings can be
 161 computed from any one of the others. Without access
 162 to T and R , however, only the *Q*-function can be
 163 used to reconstruct the other two: $V(s) = \max_a Q(s,$
 164 $a)$ and $\pi(s, a) = 1$ if $a = \operatorname{argmax}_{a'} Q(s, a')$, and 0
 165 otherwise. In addition, the *Q*-function is not difficult
 166 to estimate from experience.

167 The experience available to a reinforcement-learn-
 168 ing agent in an MDP environment can be summa-
 169 rized by a sequence of experience tuples $\langle s, a, r, s' \rangle$.
 170 An experience tuple is a snapshot of a single
 171 transition: the agent starts in state s , takes action a ,
 172 receives reward r and ends up in state s' .

173 Given an experience tuple $\langle s, a, r, s' \rangle$, the *Q*-
 174 learning rule is

$$175 \quad Q[s, a] := (1 - \alpha)Q[s, a] + \alpha(r + \beta \max_{a'} Q[s', a']).$$

176 (2)

177 This creates a new estimate of $Q^*(s, a)$ by adding the
 178 immediate reward to the current discounted estimate
 179 of the discounted reward starting from s' . Because of
 180 the way r and s' are chosen, the average value of this

new estimate is exactly $R(s, a) + \beta \sum_{s'} T(s, a, s') \max_{a'} Q[s', a']$. In a noise-free environment, this
 value would be directly assigned to $Q[s, a]$. How-
 ever, to get an accurate estimate, we need to average
 together many independent samples. The learning
 rate α blends our present estimate with our previous
 estimates to produce a best guess at $Q^*(s, a)$; it
 needs to be decreased slowly for the estimated *Q*-
 function to converge to Q^* (Jaakkola, Jordan &
 Singh, 1994; Szepesvári & Littmann, 1999; Tsitsik-
 lis, 1994; Watkins & Dayan, 1992).

The greedy policy, which assigns probability one
 to action $\operatorname{argmax}_a Q[s, a]$ in state s , receives rewards
 that converge to optimality as Q approaches Q^* .
 However, if this greedy policy is used to choose
 actions throughout the learning process, the agent
 may not explore a sufficient amount to guarantee
 optimal performance. So, the basic *Q-learning* algo-
 rithm identifies optimal behavior, but cannot adopt it.

Singh, Jaakkola, Littman and Szepesvári (2000)
 show that the conflict between learning the optimal
 policy and executing the optimal policy can be
 overcome by selecting actions that are greedy in the
 limit with infinite exploration (GLIE). A concrete
 example of a GLIE policy is decaying ϵ -greedy
 exploration. Let $n(s)$ be the number of times state s
 has been encountered so far during learning. Let
 $\epsilon(s) = c/n(s)$ for $0 < c < 1$. Now, if the agent selects
 the greedy action in state s with probability $1 - \epsilon(s)$
 and a random ‘exploratory’ action in state s with
 probability $\epsilon(s)$, the rewards obtained by the learner
 will converge to optimal.

An agent *converges in behavior* if its action
 distribution becomes stationary (fixed) in the limit. A
 GLIE policy need not converge in behavior, since
 ties in greedy actions are broken arbitrarily. How-
 ever, if there is a unique optimal policy and the
Q-function converges, a *Q-learning* agent will con-
 verge in behavior as well. Convergence in behavior
 is important in multi-agent settings, since an agent’s
 optimal policy may depend on how other agents
 behave. If one agent converges in behavior, we can
 attempt to analyze how other agents in the environ-
 ment respond.

The following theorem is a consequence of the
 results of Singh et al. (2000).

Theorem 1. *In a single-agent environment, an agent*

229 following the Q -learning update rule will converge
 230 to the optimal Q -function with probability one.
 231 Furthermore, if the agent follows a GLIE policy and
 232 the optimal policy is unique, it will converge in
 233 behavior with probability one.

234 3. Multiagent environments

235 In the Markov decision process model, a decision-
 236 making agent interacts with its environment, repre-
 237 sented as a probabilistic transition function. In this
 238 view, secondary agents must be fixed in their be-
 239 havior. The framework of Markov games admits a
 240 wider view that includes multiple adaptive agents
 241 with interacting or competing goals.

242 3.1. Markov games

243 In its general form, an n -player Markov game is
 244 defined by a tuple $\langle \mathcal{S}, \mathcal{A}_1, \dots, \mathcal{A}_n, T, R_1, \dots, R_n, \beta \rangle$, where
 245

- \mathcal{S} is a finite set of *states* of the environment;
- $\mathcal{A}_1, \dots, \mathcal{A}_n$ is a collection of finite sets of *actions* available to each agent (\mathcal{A}_i is the set of actions for agent i);
- $T: \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_n \rightarrow \Pi(\mathcal{S})$ is the *transition function*, giving for each state and one action from each agent, a probability distribution over states ($T(s, a_1, \dots, a_n, s')$ is the probability of ending in state s' , given that the agents start in state s and agent 1 chooses $a_1 \in \mathcal{A}_1$, agent 2 chooses $a_2 \in \mathcal{A}_2$, etc.);
- $R_i: \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_n \rightarrow \mathbb{R}$ for $1 \leq i \leq n$ is each agent's reward function, giving the expected immediate reward gained by agent i for each set of action choices the group of agents could make in each state ($R_i(s, a_1, \dots, a_n)$ is the expected reward to agent i in state s when agent 1 chooses a_1 , agent 2 chooses a_2 , etc.); and
- $0 \leq \beta < 1$ is a discount factor.

265 Once again, agents attempt to maximize their
 266 expected sum of discounted rewards, $E\{\sum_{j=0}^{\infty} \beta^j r_{i+j}^i\}$, where r_{i+j}^i is the reward received j steps
 267 into the future by agent i .
 268

269 If a set of agents adopts stationary policies

π_1, \dots, π_n , we can define a set of Q -functions for agent $1 \leq i \leq n$ much like in the MDP case as follows:

$$\begin{aligned} Q_i^\pi(s, a_1, \dots, a_n) &= R_i(s, a_1, \dots, a_n) \\ &+ \beta \sum_{s' \in \mathcal{S}} T(s, a_1, \dots, a_n, s') \\ &\cdot \sum_{a'_1 \in \mathcal{A}_1, \dots, a'_n \in \mathcal{A}_n} \pi_1(s', a'_1) \cdot \dots \\ &\cdot \pi_n(s', a'_n) Q_i^\pi(s', a'_1, \dots, a'_n). \end{aligned}$$

Here, Q -functions are defined over joint actions for each of the agents. Each agent receives rewards according to its reward function, with transitions dependent on the actions chosen jointly by the set of agents.

With respect to the same set of stationary policies, we can also define the *best-response* Q -function for each agent i :

$$\begin{aligned} Q_{\Delta i}^\pi(s, a_1, \dots, a_n) &= R_i(s, a_1, \dots, a_n) \\ &+ \beta \sum_{s' \in \mathcal{S}} T(s, a_1, \dots, a_n, s') \cdot \\ &\max_{a'_i \in \mathcal{A}_i} \sum_{a'_1 \in \mathcal{A}_1} \dots \sum_{a'_{i-1} \in \mathcal{A}_{i-1}} \sum_{a'_{i+1} \in \mathcal{A}_{i+1}} \dots \sum_{a'_n \in \mathcal{A}_n} \\ &\times \pi_1(s', a'_1) \cdot \dots \cdot \pi_{i-1}(s', a'_{i-1}) \cdot \pi_{i+1} \\ &\vee (s', a'_{i+1}) \cdot \dots \cdot \pi_n(s', a'_n) Q_{\Delta i}^\pi(s', a'_1, \dots, a'_n). \end{aligned}$$

The idea here is that $Q_{\Delta i}^\pi(s, a_1, \dots, a_n)$ is the Q -function obtained by holding all policies except for π_i fixed, then having agent i choose actions to maximize reward. Note that holding the behavior of all other agents fixed leaves agent i in a single-agent environment — an MDP.

This observation, combined with results of Singh et al. (2000), leads to this theorem.

Theorem 2. *In a multiagent environment, an agent following the Q -learning update rule will converge to the optimal response Q -function with probability one as long as all other agents converge in behavior with probability one. Furthermore, if the agent follows a GLIE policy and its best response policy is unique, it will also converge in behavior with probability one.*

Note that Theorem 2 does not imply that two

simultaneous Q -learners will converge to mutual best responses.

As in Section 2, an optimal policy is one that maximizes the expected sum of discounted reward. There are subtleties in applying this objective to Markov games, however. Firstly, consider the parallel scenario in MDPs.

In an MDP, an optimal policy is one that maximizes the expected sum of discounted reward; it is *undominated*, meaning that there is no state from which any other policy can achieve a better expected sum of discounted reward. Every MDP has at least one optimal policy, and of the optimal policies for a given MDP, at least one is stationary and deterministic.

For many Markov games, there is no policy that is undominated because performance depends critically on the behavior of the other agents in the environment. How, then, can we define an optimal policy? An elegant idea from the game-theory literature is to define an agent's optimal behavior as being its behavior at a *Nash equilibrium*. A set of policies π_1, \dots, π_n is in Nash equilibrium if each is a best response to the others. That is, if for all $1 \leq i \leq n$, the value attained by agent i from any state s ,

$$\sum_{a_1, \dots, a_n} \pi_1(s, a_1) \cdots \pi_n(s, a_n) Q_i^\pi(s, a_1, \dots, a_n),$$

is equal to its best response value

$$\max_{a_i \in \mathcal{A}_i} \sum_{a_1, \dots, a_{i-1}} \sum_{a_{i+1}, \dots, a_n} \pi_1(s, a_1) \cdots$$

$$\pi_{i-1}(s, a_{i-1}) \cdots \pi_{i+1}(s, a_{i+1}) \cdots$$

$$\pi_n(s, a_n) Q_{\Delta i}^\pi(s, a_1, \dots, a_n).$$

At a Nash equilibrium, each agent is maximizing its reward given that all other agents remain fixed.

Filar and Vrieze (1997) show that every Markov game has a Nash equilibrium in stationary policies. However, in contrast to MDPs, these policies are stochastic, in general. A classic example is 'Rock, Paper, Scissors', in which any deterministic policy can be consistently defeated, whereas the optimal stochastic policy always breaks even.

3.2. Nash Q -learning

Define $\text{Nash}_i(s, Q_1, \dots, Q_n)$ to be a one-stage Nash equilibrium policy for agent i in state s , where the total payoff to agent j is defined by Q -function Q_j in state s . Define $\text{Val}_i(s, Q_1, \dots, Q_n)$ to be the value obtained by agent i at this Nash equilibrium:

$$\begin{aligned} \text{Val}_i(s, Q_1, \dots, Q_n) = & \sum_{a_1, \dots, a_n} \text{Nash}_i(s, Q_1, \dots, Q_n) \\ & [a_1] \cdots \\ & = \cdot \text{Nash}_n(s, Q_1, \dots, Q_n) \\ & [a_n] Q_i[s, a_1, \dots, a_n]. \end{aligned}$$

In words, the value received by agent i is the expected value of the agent's future reward (in the Q function). The expected value is taken over all possible joint actions of the n agents, where we expect each agent to select actions according to the Nash equilibrium policy it chooses. Note that Val_i need not be unique, in general, as games can have multiple Nash equilibria with different values. Further, even if a game admits only one value of Val_i , it need not be the case that the policy Nash_i achieves this value, because other agents can choose other Nash equilibrium policies that don't fit with Nash_i . Finally, even if there is a unique Nash equilibrium, Nash_i need not achieve Val_i because other agents might not choose equilibrium strategies at all. Nonetheless, a Nash equilibrium policy is still one sensible choice of policy given the complexity of multiagent environments.

Since the Nash and Val functions extend the notion of greedy action choice with argmax and max in MDPs to Markov games, a natural extension to the update rule from Q -learning (Eq. (2)) is to use the Nash equilibrium value in place of the max to estimate each agent i 's Q -function:

$$\begin{aligned} Q_i[s, a_1, \dots, a_n] := & (1 - \alpha) Q_i[s, a_1, \dots, a_n] + \alpha(r_i + \\ & \beta \text{Val}_i(s, Q_1, \dots, Q_n)), \end{aligned} \quad (3)$$

given an experience tuple $\langle s, a_1, \dots, a_n, r_1, \dots, r_n, s' \rangle$. This update rule is due to Hu and Wellman (1998).

This learning algorithm is not known to converge in general, even if there is a unique value of the game. Some conditions to guarantee convergence are

390 given in the next two sections. Note, however, that
 391 Hu (1999) and Hu and Wellman (2000) found that
 392 the rule sometimes converged in simulations even
 393 when the strict assumptions needed to guarantee
 394 convergence were not satisfied.

395 4. Adversarial equilibria

396 Markov games can have a wide variety of reward
 397 structures and learning algorithms can display differ-
 398 ent dynamics depending on this structure. This
 399 section examines payoff structures that, to some
 400 degree, are in conflict with each other.

401 Define an *adversarial equilibrium* in an n -player
 402 game as one that is a saddle point. This means that if
 403 one agent deviates from the equilibrium, it not only
 404 hurts the agent, but it helps all other agents. More
 405 formally, let π_1, π_2 , etc. be policies in equilibrium.
 406 Then, it must be the case that for all states s , for all
 407 alternative policies π'_1, π'_2 , etc.

$$\begin{aligned}
 & \sum_{a_1, \dots, a_n} \pi_1(s, a_1) \cdot \dots \cdot \pi_n(s, a_n) Q_i[s, a_1, \dots, a_n] \\
 & \geq \sum_{a_1, \dots, a_n} \pi_1(s, a_1) \cdot \dots \cdot \pi_{i-1}(s, a_{i-1}) \pi'_i(s, a_i) \pi_{i+1}(s, a_{i+1}) \cdot \dots \\
 & \quad \cdot \pi_n(s, a_n) Q_i[s, a_1, \dots, a_n] \quad (4)
 \end{aligned}$$

412 for all i . Thus, agent i prefers π_i to π'_i ; each prefers
 413 to stay than switch.

414 In addition, if π_1, π_2 , etc. are in adversarial
 415 equilibrium, it must be the case that for all states s ,
 416 for all alternative policies π'_1, π'_2 , etc.

$$\begin{aligned}
 & \sum_{a_1, \dots, a_n} \pi_1(s, a_1) \cdot \dots \cdot \pi_n(s, a_n) Q_i[s, a_1, \dots, a_n] \\
 & \leq \sum_{a_1, \dots, a_n} \pi'_1(s, a_1) \cdot \dots \cdot \pi'_{i-1}(s, a_{i-1}) \pi_i(s, a_i) \pi'_{i+1}(s, a_{i+1}) \cdot \dots \\
 & \quad \cdot \pi'_n(s, a_n) Q_i[s, a_1, \dots, a_n]. \quad (5)
 \end{aligned}$$

421 This means that each agent would prefer that the
 422 other agents switch. In a sense, what is good for one
 423 agent is bad for the others.

424 The following theorem is a consequence of the
 425 results of Hu and Wellman (1998) and Bowling
 426 (2000).

Theorem 3. *In a multiagent environment, an agent following the Nash Q -learning update rule will converge to the optimal Q -function with probability one as long as all Q -functions encountered have adversarial equilibria and these are used in the update rule. Furthermore, if the agent follows a GLIE policy and the limit equilibrium is unique, it will converge in behavior with probability one.*

For convergence to the game's optimal Q -function, all combinations of actions must be executed infinitely often. In a sense, this puts a restriction on the behavior of the other agents in the environment. In general, the learning algorithm converges to the optimal Q -function for the game defined by the set of actions that are executed infinitely often in combination with all other actions.

Theorem 3 is interesting, but hard to apply in general. The condition that all Q -functions encountered during learning have adversarial equilibria is difficult to verify for an arbitrary Markov game. For example, even if the immediate reward functions R have adversarial equilibria and the Q -functions are initialized to have adversarial equilibria, none of the intermediate Q -functions need have adversarial equilibria.

Also, note that this theorem puts limitations on the Nash Q -learning update itself. In particular, it is not sufficient that the intermediate Q -functions simply possess adversarial equilibria, these equilibria must be used in the update rule. This point is important because of the non-uniqueness of equilibria (a game can have both an adversarial *and* a coordination equilibrium, defined later) and seems not to have been noted by Hu and Wellman (1998) or Bowling (2000).

The next section describes a common class of games that are guaranteed to have adversarial equilibria throughout the learning process.

4.1. Zero-sum Markov games

Zero-sum Markov games are a well-studied specialization of Markov games in which two agents have diametrically opposed goals. In particular, for all $a_1 \in \mathcal{A}_1, a_2 \in \mathcal{A}_2$, and $s \in \mathcal{S}$, $R_1(s, a_1, a_2) = -R_2(s, a_1, a_2)$. In a sense, therefore, there is only a single reward function R_1 , which agent 1 tries to

maximize and agent 2 tries to minimize. Zero-sum games can also be called adversarial or fully competitive for this reason.

The notion of a Nash equilibrium takes on special meaning in zero-sum games: each policy is evaluated with respect to the opposing policy that makes it look the worst. This performance measure prefers conservative strategies that can force any opponent to a stalemate over more daring ones that accrue a great deal of reward against some opponents and lose a great deal to others. This is the essence of minimax: behave so as to maximize your reward in the worst case.

Zero-sum Markov games were first studied by Shapley (1953), who showed that each such game has a unique value function and gave an algorithm that converges to this value function. Thus, even though zero-sum Markov games are a strict generalization of MDPs, they predate MDPs by several years.

Because every reward received by agent 1 in a zero-sum game is received with a sign flip by agent 2, we have that $Q_2 = -Q_1$; therefore, only one Q -function needs to be learned. Therefore, the Nash equilibrium conditions (Eq. (4)) can be written simply as:

$$\sum_{a_1, a_2} \pi_1(s, a_1) \pi_2(s, a_2) Q_1[s, a_1, a_2] \geq \sum_{a_1, a_2} \pi'_1(s, a_1) \pi_2(s, a_2) Q_1[s, a_1, a_2]$$

and

$$\sum_{a_1, a_2} \pi_1(s, a_1) \pi_2(s, a_2) Q_1[s, a_1, a_2] \leq \sum_{a_1, a_2} \pi_1(s, a_1) \pi'_2(s, a_2) Q_1[s, a_1, a_2].$$

Note that these instantly satisfy the conditions for being an adversarial equilibrium (Eq. (5)).

The preceding facts imply that if we apply Nash Q -learning to a zero-sum game, all Q -functions encountered during learning will have adversarial equilibria. Thus, Theorem 3 applies and convergence is guaranteed. The next section explains how to simplify the Nash Q -learning algorithm in the context of zero-sum games.

4.2. Minimax Q -learning

Minimax Q -learning is a value-function reinforcement-learning algorithm specifically designed for zero-sum games. It is described in an earlier paper (Littman, 1994), which includes empirical results on a simple zero-sum Markov game version of soccer. Other researchers have carried out similar studies (Uther & Veloso, 1997).

Note that the definition of the value of a game can be simplified in the zero-sum case as:

$$\text{Val}_1(s, Q_1) = \max_{\pi_1(s, \cdot) \in \Pi(\mathcal{A}_1)} \min_{a_2 \in \mathcal{A}_2} \sum_{a_1 \in \mathcal{A}_1} \pi_1(s, a_1) Q_1[s, a_1, a_2].$$

This equation identifies the probability distribution that maximizes the expected value in the face of the worst-possible action choice of the opponent. This calculation can be carried out via a small linear program. Note that the min can also be defined over stochastic policies, but, since it is ‘inside’ the max, the minimum is achieved for a deterministic action choice.

Using this revised definition of game value, the update rule for minimax Q -learning can be written:

$$Q_1[s, a_1, a_2] := (1 - \alpha) Q_1[s, a_1, a_2] + \alpha(r_1 + \beta \text{Val}_1(s, Q_1)). \quad (6)$$

The convergence of this approach follows from the convergence of the generalized Q -learning algorithm (Littman & Szepesvári, 1996; Szepesvári & Littman, 1999) as well as the results in Section 3.2.

Theorem 4. *In a two-player zero-sum multiagent environment, an agent following the minimax Q -learning update rule will converge to the optimal Q -function with probability one. Furthermore, if the agent follows a GLIE policy and the limit equilibrium is unique, it will converge in behavior with probability one.*

The special structure of zero-sum games makes it possible to provide several additional guarantees about minimax Q -learning. Even if the limit equilibrium is not unique, a minimax Q -learning agent with a GLIE policy in a zero-sum Markov game con-

verges to a policy that always achieves at least its optimal value regardless of its opponent. So, the policy learned by a minimax Q -learning agent is *safe* in that it can be executed in total ignorance of its opponent and still exhibit its intended effects.

Note that adversarial equilibrium policies in general achieve their learned value (or more) regardless of the opponent's action. The policy used in the minimax Q -learning update rule has a stronger guarantee — that it achieves the *largest* value possible in the absence of knowledge of the opponent's policy. This suggests that, even in a non-zero-sum game with adversarial equilibria, minimax Q -learning is preferable to Nash Q -learning — if the agent ignores the opponent's payoffs and assumes it is in a zero-sum game, it will do no worse, and possibly better, in terms of expected reward.

In spite of its strong asymptotic guarantees, minimax Q -learning can be slow to learn. Uther and Veloso (1997) argue that single-agent Q -learning is a sensible alternative in zero-sum games because it appears to learn more quickly. However, because Q -learning's update and action selection are deterministic, it can be 'tricked' into playing suboptimally. Consider the behavior of Q -learning in the game 'Rock, Paper, Scissors'. An opponent algorithm could be written that simulates the updates performed by Q -learning and always selects the optimal response to what Q -learning is about to choose. Thus, Q -learning will score considerably worse than minimax Q -learning (–1 at every stage as opposed to an average score of 0 for minimax Q -learning).

The next section examines the behavior of Nash Q -learning and variants in coordination settings.

5. Coordination equilibria

This section examines payoff structures that make it so that the agents, to some degree, are working toward a common goal.

Define a *coordination equilibrium* in an n -player game as one for which all agents achieve their maximum possible payoff. If π_1, π_2 , etc. are in coordination equilibrium, we have that

$$\sum_{a_1, \dots, a_n} \pi_1(s, a_1) \cdots \pi_n(s, a_n) Q_i[s, a_1, \dots, a_n] = \max_{a_1, \dots, a_n} Q_i[s, a_1, \dots, a_n]$$

for all $1 \leq i \leq n$ and states s . One observation is that if a game has a coordination equilibrium, it has a deterministic coordination equilibrium. This follows from the fact that the value for each agent is a convex combination of the values in the Q -functions.

One consequence of a game possessing a coordination equilibrium is that no agent has any incentive to switch because no other set of policies could result in a higher score. In a sense, what is best for one agent is also best for the others.

The following theorem is a consequence of the results of Hu and Wellman (1998) and Bowling (2000).

Theorem 5. *In a multiagent environment, an agent following the Nash Q -learning update rule will converge to the optimal Q -function with probability one as long as all Q -functions encountered have coordination equilibria and these are used in the update rule. Furthermore, if the agent follows a GLIE policy and the limit equilibrium is unique, it will converge in behavior with probability one.*

Like Theorem 3, this theorem is hard to apply because the conditions are difficult to verify in advance. The next section addresses a special case in which the conditions are easily verifiable — team games.

5.1. Team Markov games

In team Markov games, agents have precisely the same goals. In particular, for all $a_1 \in \mathcal{A}_1, a_2 \in \mathcal{A}_2$, etc., and $s \in \mathcal{S}$, $R_1(s, a_1, \dots, a_n) = R_2(s, a_1, \dots, a_n) = \dots$. In a sense, therefore, there is only a single reward function R_1 , which all agents try to maximize together. Team games can also be called coordination games or fully cooperative games for this reason. Boutilier (1996) refers to team Markov games as multiagent decision processes, since, acting together, the group of agents is faced with an MDP.

Because every reward received by agent 1 in a team game is received by all agents, we have that

639 $Q_1 = \dots = Q_n$; therefore, only one Q -function
 640 needs to be learned. Let $a_1^*, \dots, a_n^* =$
 641 $\operatorname{argmax}_{a_1, \dots, a_n} Q_1[s, a_1, \dots, a_n]$. The set of
 642 policies in which each agent i assigns probability 1
 643 to a_i^* is a coordination equilibrium.

644 The above facts imply that if we apply Nash
 645 Q -learning to a team game, all Q -functions en-
 646 countered during learning will have coordination
 647 equilibria. Thus, Theorem 5 applies and convergence
 648 is guaranteed. The next section explains how to
 649 simplify the Nash Q -learning algorithm in the con-
 650 text of team games.

651 5.2. Team Q -learning

652 Team Q -learning is a value-function rein-
 653 forcement-learning algorithm specifically designed
 654 for team games. The definition of the value of a
 655 game can be simplified in the team case as:

$$656 \operatorname{Val}_1(s, Q_1) = \max_{a_1, \dots, a_n} Q_1[s, a_1, \dots, a_n].$$

657 This equation simply returns the largest entry in the
 658 Q table for the given state.

659 Using this revised definition of game value, the
 660 update rule for team Q -learning can be written:

$$661 Q_1[s, a_1, \dots, a_n] := (1 - \alpha)Q_1[s, a_1, \dots, a_n] + \\ 662 \alpha(r_1 + \beta \operatorname{Val}_1(s, Q_1)). \quad (7)$$

663 This update is slightly different from the joint-
 664 action learner (JAL) algorithm of Claus and Boutilier
 665 (1998), since it doesn't use an opponent model.

666 The convergence of this approach follows from
 667 the convergence of the generalized Q -learning algo-
 668 rithm (Littman & Szepesvári, 1996; Szepesvári &
 669 Littman, 1999) as well as the results in Section 3.2.

670 **Theorem 6.** *In a multiagent environment where*
 671 *agents have identical reward functions, an agent*
 672 *following the team Q -learning update rule will*
 673 *converge to the optimal Q -function with probability*
 674 *one. Furthermore, if the agent follows a GLIE policy*
 675 *and the limit equilibrium is unique, it will converge*
 676 *in behavior with probability one.*

677 Note that this algorithm learns precisely the same

678 values as Nash Q -learning for games with coordina-
 679 tion equilibria, even in a non-team game. The
 680 calculation for the team Q -learning update is simpler,
 681 though, and more easily applied to n -player games.

682 An important open problem for team Markov
 683 games is finding a robust and general way of
 684 selecting an equilibrium when there are multiple
 685 coordination equilibria. Boutilier (1996) argues for
 686 establishing a tie-breaking scheme (e.g. lexicog-
 687 raphic ordering) to pick the equilibrium to play.
 688 However, in the presence of noise, there is no
 689 guarantee there will ever be ties and the learners
 690 could perceive different unique equilibria and there-
 691 fore not converge on optimal play.

692 6. Conclusions

693 This section explores some of the implications of
 694 the results described in the paper.

695 6.1. Value-function reinforcement learning

696 The theorems described above can be used to
 697 reason about the dynamics of simultaneously learn-
 698 ing agents.

699 It follows from Theorem 4 that two independent
 700 minimax Q -learning agents will learn Q -functions
 701 that converge. In a zero-sum Markov game, the
 702 resulting behavior will be optimal for both agents
 703 (mutual best responses). This follows from the fact
 704 that equilibria in zero-sum games are 'interchange-
 705 able'. In other environments, even with more than
 706 two agents, all minimax Q -learners behave in a way
 707 that maximizes their reward given no assumption
 708 whatsoever about the behavior of the other agents.

709 In a zero-sum Markov game between a minimax
 710 Q -learner and a Q -learner, the Q -learner will con-
 711 verge to optimal behavior (Theorem 2) if the equilib-
 712 rium is unique.

713 In a sense, team Markov games are harder. A set
 714 of team Q -learners will converge to optimal play if
 715 the limit Q -functions have a unique maximum for
 716 each state. If not, even though all agents will
 717 converge on an equilibrium, if they converge on
 718 different equilibria, they can score arbitrarily poorly.

720 This is also true of a Q -learner among team Q -
721 learners.

722 Claus and Boutilier (1998) examined simultaneous
723 Q -learners in team games. They argue that two such
724 learners converge to a Nash pair, although it need
725 not be the optimal one. As far as I know, no formal
726 proof has been presented. Nonetheless, Sen et al.
727 (1994), Mundhe and Sen (2000) and others have
728 successfully applied this technique to particular
729 games, validating it as a reasonable approach.

730 In general games, if a Nash Q -learner's behavior
731 converges, a Q -learner will learn a best response.
732 However, as in the case with two team Q -learners,
733 two Nash Q -learners need not converge to compat-
734 ible equilibria and can score arbitrarily badly. A
735 Nash Q -learner with a minimax Q -learner need not
736 score well, but the minimax Q -learner will at least
737 learn a 'safe' strategy.

738 6.2. Nash Q -learning

739 Nash Q -learning is a promising algorithm and has
740 been executed with positive results on some small,
741 but interesting, environments. The results described
742 in this paper, however, undercut its theoretical
743 justifications. As detailed next, any situation in
744 which the use of Nash Q -learning is justified,
745 minimax Q -learning or team Q -learning would ap-
746 pear to be a more sensible choice.

747 Bowling (2000) pointed out that Nash Q -learning
748 is only guaranteed to converge if the Q -functions at
749 every state either always has an adversarial equilib-
750 rium or always has a coordination equilibrium. As
751 mentioned in Section 4, because of the possibility of
752 games having both types of equilibria, the learner
753 must know in advance which type of equilibrium to
754 pick in each state. If a state is of the 'adversarial'
755 type, using minimax Q -learning (Section 4.2) gives a
756 stronger guarantee than Nash Q -learning on the
757 amount of reward the agent will receive. If a state is
758 of the 'coordination' type, using team Q -learning
759 (Section 5.2) results in learning the same values as
760 Nash Q -learning, but with a simpler computation.

761 Thus, even though Nash Q -learning can be applied
762 in very general settings, in any case in which its use
763 is formally justified — at present — minimax Q -
764 learning or team Q -learning is preferable. This

indicates that more theoretical work on Nash Q -
learning could be quite beneficial.

6.3. Policy modeling

As is evident from the theorems cited in this
paper, multiagent environments with multiple
equilibria pose difficult problems for the conver-
gence of learning algorithms. As Bowling and Vel-
oso (2000) point out, if a game has multiple equilib-
ria, the optimal policy must depend on the policies of
the other agents. In this type of environment, op-
ponent-independent algorithms like Nash Q -learning,
minimax Q -learning, and team Q -learning are not
sufficient to identify optimal behavior. Some depen-
dence on the policies of other agents is necessary,
and to do this, some assumptions must be made
about how other agents will behave.

One reasonable assumption is that past behavior is
a strong indicator of future behavior. Q -learning, as
described earlier, implicitly uses this idea by choos-
ing actions that maximize payoff as estimated by
past interactions. It is also possible to model other
agents more directly, for example as probabilistic
stationary policies. Claus and Boutilier (1998) de-
fined a joint-action Q -learning algorithm (JAL) that
uses this idea to play team games, and Uther and
Veloso (1997) used this approach in zero-sum games
to good effect. The fictitious play approach in game
theory, surveyed by Vrieze and Tijs (1982), also
stems from this insight. Agent modeling would seem
to be even more important in environments with
more complex payoff structures. However, there is
also some evidence that opponent modeling can be
problematic in some circumstances (Sun & Qi,
2000).

A drawback of straightforward agent modeling is
that it is purely reactive. In a sense, each agent is
letting the other agents pick the equilibrium and then
each learns a best response. It would be beneficial to
identify a method that is partly opponent-independ-
ent, like Nash Q -learning, and partly opponent-
sensitive, like Q -learning. The latter attribute could
encourage equilibrium behavior, so the actions
chosen fit well with the actions chosen by other
agents. The former attribute could help encourage
agents to select an equilibrium that results in high
payoff. The hope is that this could reduce the

probability of selecting ‘optimal’ actions that score poorly because agents don’t agree on an equilibrium or jointly selecting actions that form a highly suboptimal equilibrium.

In sum, reinforcement learning is a powerful framework for studying learning in multiagent scenarios. Adversarial environments are well behaved in that optimal play can be guaranteed against an arbitrary opponent. Coordination environments are less well behaved, as strong assumptions need to be made about other agents to guarantee convergence to optimal behavior. In other types of environments, no value-function reinforcement-learning algorithms with guaranteed convergence properties are known. Nevertheless, the last few years have expanded our understanding of reinforcement-learning approaches and have helped clarify areas in which further study is needed.

Acknowledgements

Craig Boutilier, Michael Bowling, Junling Hu, Satinder Singh, Peter Stone, Ron Sun and Michael Wellman provided pointers to papers and other valuable help.

References

Barto, A. G., Sutton, R. S., & Watkins, C. J. C. H. (1991). Learning and sequential decision making. In: Gabriel, M., & Moore, J. (Eds.), *Learning and computational neuroscience: foundations of adaptive networks*, MIT Press, Cambridge, MA, Tech. Rep. 89-95, Department of Computer and Information Science, University of Massachusetts, Amherst, MA, 1989.

Bellman, R. (1957). *Dynamic Programming*, Princeton University Press, Princeton, NJ.

Bertsekas, D. P. (1987). *Dynamic programming: deterministic and stochastic models*, Prentice-Hall, Englewood Cliffs, NJ.

Boutilier, C. (1996). Planning, learning and coordination in multiagent decision processes. In: Proceedings of the Sixth Conference on Theoretical Aspects of Rationality and Knowledge (TARK-96).

Bowling, M. (2000). Convergence problems of general-sum multiagent reinforcement learning. In: Proceedings of the Seventeenth International Conference on Machine Learning.

Bowling, M., & Veloso, M. M. (2000). An analysis of stochastic game theory for multiagent reinforcement learning. In: Technical Report CMU-CS-00-165, Computer Science Department, Carnegie Mellon University.

Claus, C., & Boutilier, C. (1998). The dynamics of reinforcement learning in cooperative multiagent systems. In: Proceedings of the Fifteenth National Conference on Artificial Intelligence.

Filar, J., & Vrieze, K. (1997). *Competitive Markov decision processes*, Springer-Verlag.

Howard, R. A. (1960). *Dynamic programming and Markov processes*, MIT Press, Cambridge, MA.

Hu, J. (1999). *Learning in Dynamic Noncooperative Multiagent Systems*, Department of Computer Science, University of Michigan, Ph.D. thesis.

Hu, J., & Wellman, M. P. (1998). Multiagent reinforcement learning: theoretical framework and an algorithm. In: Shavlik, J. (Ed.), *Proceedings of the Fifteenth International Conference on Machine Learning*, Morgan Kaufmann.

Hu, J., & Wellman, M. P. (2000). Experimental results on Q -learning for general-sum stochastic games. In: Langley, P. (Ed.), *Proceedings of the Seventeenth International Conference on Machine Learning*, Morgan Kaufmann.

Jaakkola, T., Jordan, M. I., & Singh, S. P. (1994). On the convergence of stochastic iterative dynamic programming algorithms. *Neural Comput.* 6(6), 1185–1201.

Kaelbling, L. P., Littman, M. L., & Moore, A. W. (1996). Reinforcement learning: a survey. *J. Artificial Intell. Res.* 4, 237–285.

Littman, M. L. (1994). Markov games as a framework for multiagent reinforcement learning. In: Proceedings of the Eleventh International Conference on Machine Learning, Morgan Kaufmann, San Francisco, CA.

Littman, M. L., & Szepesvári, C. (1996). A generalized reinforcement-learning model: convergence and applications. In: Saitta, L. (Ed.), *Proceedings of the Thirteenth International Conference on Machine Learning*.

Mundhe, M., & Sen, S. (2000). Evaluating concurrent reinforcement learners. In: Proceedings of the Fourth International Conference on Multiagent Systems, IEEE Press.

Owen, G. (1982). In: 2nd Edition, *Game Theory*, Academic Press, Orlando, FL.

Puterman, M. L. (1994). *Markov decision processes — discrete stochastic dynamic programming*, John Wiley, New York.

Sen, S., Sekaran, M., & Hale, J. (1994). Learning to coordinate without sharing information. In: Proceedings of the Twelfth National Conference on Artificial Intelligence.

Shapley, L. (1953). Stochastic games. *Proc. Natl. Acad. Sci. USA* 39, 1095–1100.

Singh, S., Jaakkola, T., Littman, M. L., & Szepesvári, C. (2000). Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine Learn.* 39, 287–308.

Sun, R., & Qi, D. (2000). Rationality assumptions and optimality of co-learning. In: Proceedings of PRIMA’2000, Lecture Notes in Artificial Intelligence, Springer-Verlag.

Sutton, R. S. (1988). Learning to predict by the method of temporal differences. *Machine Learn.* 3(1), 9–44.

Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: an introduction*, MIT Press.

Szepesvári, C., & Littman, M. L. (1999). A unified analysis of value-function-based reinforcement-learning algorithms. *Neural Comput.* 11(8), 2017–2059.

- 916 Tsitsiklis, J. N. (1994). Asynchronous stochastic approximation and Q -learning. *Machine Learn.* 16(3), 185–202. 926
- 917 Uther, W., & Veloso, M. (1997). *Adversarial reinforcement learning*, Unpublished manuscript. 927
- 918 van der Wal, J. (1981). In: Stochastic dynamic programming, Mathematical centre tracts, Vol. 139, Mathematisch Centrum, Amsterdam. 928
- 919 von Neumann, J., & Morgenstern, O. (1947). *Theory of games and economic behavior*, Princeton University Press, Princeton, NJ. 929
- 920 Vrieze, O. J., & Tijs, S. H. (1982). Fictitious play applied to sequences of games and discounted stochastic games. *Int. J. Game Theory* 11(2), 71–85. 930
- 921 Watkins, C. J. C. H. (1989). *Learning from delayed rewards*, King's College, Cambridge, UK, Ph.D. thesis. 931
- 922 Watkins, C. J. C. H., & Dayan, P. (1992). Q -learning. *Machine Learn.* 8(3), 279–292. 932

UNCORRECTED PROOF