

Multiagent Reinforcement Learning: Algorithm Converging to Nash Equilibrium in General-Sum Discounted Stochastic Games

Natalia Akchurina
International Graduate School of Dynamic Intelligent Systems
University of Paderborn
100 Warburger Str.
Paderborn, Germany
anatalia@mail.uni-paderborn.de

ABSTRACT

This paper introduces a multiagent reinforcement learning algorithm that converges with a given accuracy to stationary Nash equilibria in general-sum discounted stochastic games. Under some assumptions we formally prove its convergence to Nash equilibrium in self-play. We claim that it is the first algorithm that converges to stationary Nash equilibrium in the general case.

Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning; I.2.11 [Artificial Intelligence]: Distributed Artificial Intelligence—*Multiagent systems*

General Terms

Algorithms, Theory

Keywords

algorithmic game theory, stochastic games, computation of equilibria, multiagent reinforcement learning

1. INTRODUCTION

Reinforcement learning turned out a technique that allowed robots to ride a bicycle, computers to play backgammon on the level of human world masters and solve such complicated tasks of high dimensionality as elevator dispatching. Can it come to rescue in the next generation of challenging problems like playing football or bidding on virtual markets? Reinforcement learning that provides a way of programming agents without specifying how the task is to be achieved could be again of use here but the convergence of reinforcement learning algorithms to optimal policies is only guaranteed under the conditions of stationarity of the environment that is violated in multiagent systems. For reinforcement learning in multiagent environments general-sum discounted stochastic games become a formal framework instead of Markov decision processes. Also the optimal pol-

Cite as: Multiagent Reinforcement Learning: Algorithm Converging to Nash Equilibrium in General-Sum Discounted Stochastic Games, Natalia Akchurina, *Proc. of 8th Int. Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2009)*, Decker, Sichman, Sierra and Castelfranchi (eds.), May, 10–15, 2009, Budapest, Hungary, pp. 725–732
Copyright © 2009, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org), All rights reserved.

icy concept in multiagent systems is different — we can't speak anymore about optimal policy (policy that provides the maximum cumulative reward) without taking into account the policies of other agents that influence our payoffs. In the environment where every agent tries to maximize its cumulative reward it is the most natural to accept Nash equilibrium as the optimal solution concept. In Nash equilibrium each agent's policy is the best response to the other agents' policies. Thus no agent can gain from unilateral deviation.

A number of algorithms [1, 2, 4, 5, 6, 7, 8, 9] were proposed to extend reinforcement learning approach to multiagent systems. The convergence to Nash equilibria was proved for very restricted class of environments: strictly competitive [6], strictly cooperative [2, 7] and 2-agent 2-action iterative game [1]. In [5] convergence to Nash equilibrium has been achieved in self-play for strictly competitive and strictly cooperative games under additional very restrictive condition that all equilibria encountered during learning stage are unique [7].

In this paper we propose a reinforcement learning algorithm that converges to Nash equilibria with some given accuracy in general-sum discounted stochastic games and prove it formally under some assumptions. We claim that it is the first algorithm that finds Nash equilibrium for the general case.

The paper is organized as follows. In section 2 we present formal definitions of stochastic game, Nash equilibrium, as well as prove some theorems that we will need for equilibrium approximation theorem in section 3. Section 3 is devoted to discussion and necessary experimental estimations of the conditions of the equilibrium approximation theorem. In sections 5 and 6 the developed algorithm Nash-DE¹ and analysis of the results of experiments are presented correspondingly.

2. PRELIMINARIES

In this section we recall some definitions from game theory.

Definition 1. A pair of matrices (M^1, M^2) constitute a bimatrix game G , where M^1 and M^2 are of the same size. The rows of M^k correspond to actions of player 1, $a^1 \in A^1$.

¹Maintaining the tradition the name of the algorithm reflects the result — the approximation of Nash-equilibrium as well as the approach — differential equations.

The columns of M^k correspond to actions of player 2, $a^2 \in A^2$. A^1 and A^2 are the sets of discrete actions of players 1 and 2 respectively. The payoff $r^k(a^1, a^2)$ to player k can be found in the corresponding entry of the matrix M^k , $k = 1, 2$.

Definition 2. A pure ε -equilibrium of bimatrix game G is a pair of actions (a_*^1, a_*^2) such that

$$r^1(a_*^1, a_*^2) \geq r^1(a^1, a_*^2) - \varepsilon \text{ for all } a^1 \in A^1$$

$$r^2(a_*^1, a_*^2) \geq r^2(a_*^1, a^2) - \varepsilon \text{ for all } a^2 \in A^2$$

Definition 3. A mixed ε -equilibrium of bimatrix game G is a pair of vectors (ρ_*^1, ρ_*^2) , such that

$$\rho_*^1 M^1 \rho_*^2 \geq \rho^1 M^1 \rho_*^2 - \varepsilon \text{ for all } \rho^1 \in \sigma(A^1)$$

$$\rho_*^1 M^2 \rho_*^2 \geq \rho_*^1 M^2 \rho^2 - \varepsilon \text{ for all } \rho^2 \in \sigma(A^2)$$

where $\sigma(A^k)$ is the set of probability distributions over action space A^k , such that for any $\rho^k \in \sigma(A^k)$, $\sum_{a \in A^k} \rho_a^k = 1$

$$\begin{aligned} \rho^1 M^k \rho^2 &= \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \rho_{a^1}^1 r^k(a^1, a^2) \rho_{a^2}^2 = \\ &= \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} r^k(a^1, a^2) \prod_{i=1}^2 \rho_{a^i}^i \end{aligned}$$

is the expected reward of agent k induced by (ρ^1, ρ^2) .

Definition 4. Nash equilibrium of bimatrix game G is ε -equilibrium with $\varepsilon = 0$.

Obviously definitions 1, 2, 3 and 4 can be generalized for arbitrary number of players.

Definition 5. A 2-player discounted stochastic game Γ is a 7-tuple $\langle S, A^1, A^2, \gamma, r^1, r^2, p \rangle$, where S is discrete finite set of states ($|S| = N$), A^k is the discrete action space of player k for $k = 1, 2$ ($|A^k| = m^k$), $\gamma \in [0, 1)$ is the discount factor, $r^k : S \times A^1 \times A^2 \rightarrow \mathbb{R}$ is the reward function for player k bounded in absolute value by R_{\max} , $p : S \times A^1 \times A^2 \rightarrow \Delta$ is the transition probability map, where Δ is the set of probability distributions over state space S .

Discount factor γ reflects the notion that a reward at time $t + 1$ is worth only $\gamma < 1$ of what it is worth at time t .

Every state of a 2-player stochastic game can be regarded as a bimatrix game.

It is assumed that for every $s, s' \in S$ and for every action $a^1 \in A^1$ and $a^2 \in A^2$, transition probabilities $p(s'|s, a^1, a^2)$ are stationary for all $t = 0, 1, 2, \dots$ and $\sum_{s' \in S} p(s'|s, a^1, a^2) = 1$.

Policy of agent $k = 1, 2$ is a vector $x^k = (x_{s_1}^k, x_{s_2}^k, \dots, x_{s_N}^k)$, where $x_s^k = (x_{s a_1^k}^k, x_{s a_2^k}^k, \dots, x_{s a_{m^k}^k}^k)$, $x_{s h}^k \in \mathbb{R}$ being the probability assigned by agent k to its action $h \in A^k$ in state s . A policy x^k is called a *stationary policy* if it is fixed over time. Since all probabilities are nonnegative and sum up to one, the vector $x_s^k \in \mathbb{R}^{m^k}$ belongs to the unit simplex Δ^k in m^k -space defined as

$$\Delta^k = \{x_s^k \in \mathbb{R}_+^{m^k} : \sum_{a \in A^k} x_{s a}^k = 1\}$$

The policy x^k will belong then to policy space of agent k :

$$\Theta^k = \times_{s \in S} \Delta^k$$

Each player k ($k = 1, 2$) strives to learn policy by immediate rewards so as to maximize its expected discounted cumulative reward (players don't know state transition probabilities and payoff functions):

$$v^k(s, x^1, x^2) = \sum_{t=0}^{\infty} \gamma^t E(r_t^k | x^1, x^2, s_0 = s)$$

where x^1 and x^2 are the policies of players 1 and 2 respectively and s is the initial state.

$v^k(s, x^1, x^2)$ is called the discounted value of policies (x^1, x^2) in state s to player k .

Definition 6. A 2-player discounted stochastic game Γ is called zero-sum when $r^1(s, a^1, a^2) + r^2(s, a^1, a^2) = 0$ for all $s \in S$, $a^1 \in A^1$ and $a^2 \in A^2$, otherwise general-sum.

Definition 7. A ε -equilibrium of 2-player discounted stochastic game Γ is a pair of policies (x_*^1, x_*^2) such that for all $s \in S$ and for all policies $x^1 \in \Theta^1$ and $x^2 \in \Theta^2$:

$$v^1(s, x_*^1, x_*^2) \geq v^1(s, x^1, x_*^2) - \varepsilon$$

$$v^2(s, x_*^1, x_*^2) \geq v^2(s, x_*^1, x^2) - \varepsilon$$

Definition 8. Nash equilibrium of 2-player discounted stochastic game Γ is ε -equilibrium with $\varepsilon = 0$.

Definition 9. A n -player discounted stochastic game is a tuple $\langle K, S, A^1, \dots, A^n, \gamma, r^1, \dots, r^n, p \rangle$, where $K = \{1, 2, \dots, n\}$ is the player set, S is the discrete state space ($|S| = N$), A^k is the discrete action space of player k for $k \in K$ ($|A^k| = m^k$), $\gamma \in [0, 1)$ is the discount factor, $r^k : S \times A^1 \times A^2 \times \dots \times A^n \rightarrow \mathbb{R}$ is the reward function for player k bounded in absolute value by R_{\max} , $p : S \times A^1 \times A^2 \times \dots \times A^n \rightarrow \Delta$ is the transition probability map, where Δ is the set of probability distributions over state space S .

Definitions 6, 7 and 8 can be generalized for n -player stochastic game.

Definition 10. A profile is a vector $x = (x^1, x^2, \dots, x^n)$, where each component x^k is a policy for player $k \in K$. The space of all profiles $\Phi = \times_{k \in K} \Theta^k$.

Let's define the probability transition matrix induced by x :

$$p(s'|s, x) = \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^n \in A^n} p(s'|s, a^1, a^2, \dots, a^n) \prod_{i=1}^n x_{s a^i}^i$$

$$P(x) = (p(s'|s, x))_{s, s' \in S}$$

The immediate expected reward of player k in state s induced by x will be:

$$r^k(s, x) = \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^n \in A^n} r^k(s, a^1, a^2, \dots, a^n) \prod_{i=1}^n x_{s a^i}^i$$

Then the immediate expected reward matrix induced by profile x will be:

$$r(x) = (r^k(s, x))_{s \in S, k \in K}$$

The discounted value matrix of x will be [3]:

$$v(x) = [I - \gamma P(x)]^{-1} r(x)$$

where I is $N \times N$ identity matrix.

Note that the following recursive formula will hold for the discounted value matrix [3]:

$$v(x) = r(x) + \gamma P(x)v(x)$$

The k th columns of $r(x)$ and $v(x)$ (the immediate expected reward of player k induced by profile x and the discounted value of x to agent k) let us respectively denote $r^k(x)$ and $v^k(x)$.

2.1 Useful Theorems

In this section we will prove a lemma and a theorem for an arbitrary n -player discounted stochastic game $\Gamma = \langle K, S, A^1, \dots, A^n, \gamma, r^1, \dots, r^n, p \rangle$.

LEMMA 1. If $k \in K$, $x \in \Phi$ and $v, \epsilon \in \mathbb{R}^N$ are such that

$$v \geq r^k(x) + \gamma P(x)v - \epsilon$$

then $v \geq v^k(x) - \sum_{t=0}^{\infty} \gamma^t P^t(x)\epsilon$.

PROOF.

$$\begin{aligned} v &\geq r^k(x) + \gamma P(x)[r^k(x) + \gamma P(x)v - \epsilon] - \epsilon \\ &= r^k(x) + \gamma P(x)r^k(x) + \gamma^2 P^2(x)v \\ &\quad - \epsilon - \gamma P(x)\epsilon \end{aligned}$$

Upon substituting the above inequality into itself i times we obtain:

$$\begin{aligned} v &\geq r^k(x) + \gamma P(x)r^k(x) + \gamma^2 P^2(x)r^k(x) + \\ &\quad + \dots + \\ &\quad + \gamma^{i-1} P^{i-1}(x)r^k(x) + \gamma^i P^i(x)v \\ &\quad - \epsilon - \gamma P(x)\epsilon - \gamma^2 P^2(x)\epsilon - \dots - \gamma^{i-1} P^{i-1}(x)\epsilon \end{aligned}$$

which upon taking the limit as $i \rightarrow \infty$, yields $v \geq v^k(x) - \sum_{t=0}^{\infty} \gamma^t P^t(x)\epsilon$. \square

THEOREM 1. From 1 \Rightarrow 2

1. For each $s \in S$, the vector $(x_s^1, x_s^2, \dots, x_s^n)$ constitutes an ϵ -equilibrium in the n -matrix game $(B_s^1, B_s^2, \dots, B_s^n)$ with equilibrium payoffs $(v_s^1, v_s^2, \dots, v_s^n)$, where for $k \in K$ and $(a^1, a^2, \dots, a^n) \in A^1 \times A^2 \times \dots \times A^n$ entry (a^1, a^2, \dots, a^n) of B_s^k equals

$$\begin{aligned} b^k(s, a^1, a^2, \dots, a^n) &= r^k(s, a^1, a^2, \dots, a^n) \\ &\quad + \gamma \sum_{s' \in S} (p(s'|s, a^1, a^2, \dots, a^n) \\ &\quad + \zeta(s'|s, a^1, a^2, \dots, a^n)) \\ &\quad \cdot (v_{s'}^k + \sigma_{s'}^k) \end{aligned}$$

where $-\sigma < \sigma_s^k < \sigma$, $-\zeta < \zeta(s'|s, a^1, a^2, \dots, a^n) < \zeta$.

2. x is an ϵ -equilibrium in the discounted stochastic game Γ where $\epsilon = [2\gamma(\sigma + \zeta N \max_{k,s} |v_s^k| + N\zeta\sigma) + \epsilon] \sum_{t=0}^{\infty} \gamma^t$.

PROOF.

$$\begin{aligned} b^k(s, a^1, a^2, \dots, a^n) &= r^k(s, a^1, a^2, \dots, a^n) + \\ &\quad + \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) v_{s'}^k + \\ &\quad + \gamma \sum_{s' \in S} \zeta(s'|s, a^1, a^2, \dots, a^n) v_{s'}^k + \\ &\quad + \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) \sigma_{s'}^k + \\ &\quad + \gamma \sum_{s' \in S} \zeta(s'|s, a^1, a^2, \dots, a^n) \sigma_{s'}^k \\ &= r^k(s, a^1, a^2, \dots, a^n) + \xi^k(s, a^1, a^2, \dots, a^n) + \\ &\quad + \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) v_{s'}^k \end{aligned}$$

where

$$\begin{aligned} \xi^k(s, a^1, a^2, \dots, a^n) &= \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) \sigma_{s'}^k \\ &\quad + \gamma \sum_{s' \in S} \zeta(s'|s, a^1, a^2, \dots, a^n) v_{s'}^k + \\ &\quad + \gamma \sum_{s' \in S} \zeta(s'|s, a^1, a^2, \dots, a^n) \sigma_{s'}^k \end{aligned}$$

Let's estimate the worst case:

$$\begin{aligned} &- \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) \sigma \\ &- \gamma \sum_{s' \in S} \zeta \max_{s'} |v_{s'}^k| - \gamma \sum_{s' \in S} \zeta \sigma \\ &< \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) \sigma_{s'}^k \\ &\quad + \gamma \sum_{s' \in S} \zeta(s'|s, a^1, a^2, \dots, a^n) v_{s'}^k \\ &\quad + \gamma \sum_{s' \in S} \zeta(s'|s, a^1, a^2, \dots, a^n) \sigma_{s'}^k \\ &< \gamma \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) \sigma \\ &\quad + \gamma \sum_{s' \in S} \zeta \max_{s'} |v_{s'}^k| + \gamma \sum_{s' \in S} \zeta \sigma \end{aligned}$$

Let's denote $\omega = \gamma\sigma + \gamma\zeta N \max_{k,s} |v_s^k| + \gamma N\zeta\sigma$

$$-\omega < \xi^k(s, a^1, a^2, \dots, a^n) < \omega$$

Let's take some arbitrary $f \in \Theta^1$. If (1) is true, then for each state $s \in S$ by definition of ϵ -equilibrium:

$$\begin{aligned} &r^1(s, f, x^2, \dots, x^n) + \zeta^1(s, f, x^2, \dots, x^n) \\ &\quad + \gamma \sum_{s' \in S} p(s'|s, f, x^2, \dots, x^n) v_{s'}^1 \leq v_s^1 + \epsilon \end{aligned}$$

where

$$\zeta^k(s, x) = \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^n \in A^n} \xi^k(s, a^1, a^2, \dots, a^n) \prod_{i=1}^n x_{sa^i}$$

In the worst case:

$$\begin{aligned}
& - \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^n \in A^n} \omega \prod_{i=1}^n x_{sa^i}^i < \\
& < \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^n \in A^n} \xi^k(s, a^1, a^2, \dots, a^n) \prod_{i=1}^n x_{sa^i}^i < \\
& < \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^n \in A^n} \omega \prod_{i=1}^n x_{sa^i}^i \\
& \quad - \omega < \zeta^k(s, x) < \omega
\end{aligned}$$

$$r^1(s, f, x^2, \dots, x^n) + \gamma \sum_{s' \in S} p(s'|s, f, x^2, \dots, x^n) v_{s'}^1$$

$$\leq v_s^1 + \omega + \epsilon$$

Applying the lemma 1 we get

$$v^1(f, x^2, \dots, x^n) - (\omega + \epsilon) \sum_{t=0}^{\infty} \gamma^t \leq v^1$$

for all $f \in \Theta^1$, and by symmetrical arguments it follows that

$$v^k(x^1, \dots, x^{k-1}, f, x^{k+1}, \dots, x^n) - (\omega + \epsilon) \sum_{t=0}^{\infty} \gamma^t \leq v^k$$

for all $k \in K$ and for all $f \in \Theta^k$.

But

$$v^k = r^k(x) + \gamma P(x)v^k + \zeta^k(x)$$

$$\begin{aligned}
v^k &= r^k(x) + \gamma P(x)[r^k(x) + \gamma P(x)v^k + \zeta^k(x)] + \zeta^k(x) \\
&= r^k(x) + \gamma P(x)r^k(x) + \gamma^2 P^2(x)v^k \\
&+ \zeta^k(x) + \gamma P(x)\zeta^k(x)
\end{aligned}$$

Upon substituting the above inequality into itself and taking the limit, we get:

$$v^k = v^k(x) + \sum_{t=0}^{\infty} \gamma^t P^t(x)\zeta^k(x)$$

$$v^k(x) - \omega \sum_{t=0}^{\infty} \gamma^t < v^k < v^k(x) + \omega \sum_{t=0}^{\infty} \gamma^t$$

$$v^k(x^1, \dots, x^{k-1}, f, x^{k+1}, \dots, x^n) - (\omega + \epsilon) \sum_{t=0}^{\infty} \gamma^t$$

$$\leq v^k(x) + \omega \sum_{t=0}^{\infty} \gamma^t$$

for all $k \in K$ and for all $f \in \Theta^k$.

$$v^k(x^1, \dots, x^{k-1}, f, x^{k+1}, \dots, x^n) \leq v^k(x) + (2\omega + \epsilon) \sum_{t=0}^{\infty} \gamma^t$$

for all $k \in K$ and for all $f \in \Theta^k$.

So the condition of theorem 1 is satisfied with $\epsilon = [2(\gamma\sigma + \gamma\zeta N \max_{k,s} |v_s^k| + \gamma N \zeta\sigma) + \epsilon] \sum_{t=0}^{\infty} \gamma^t$ and we get (2). \square

3. ϵ -EQUILIBRIUM THEOREM

Let

$$\Gamma = \langle K, S, A^1, \dots, A^n, \gamma, r^1, \dots, r^n, p \rangle$$

and

$$\tilde{\Gamma} = \langle K, S, A^1, \dots, A^n, \gamma, r^1, \dots, r^n, \tilde{p} \rangle$$

be n -player discounted stochastic games such that:

for all $s' \in S$ and $(s, a^1, a^2, \dots, a^n) \in S \times A^1 \times A^2 \times \dots \times A^n$:

$$\begin{aligned}
\tilde{p}(s'|s, a^1, a^2, \dots, a^n) &= p(s'|s, a^1, a^2, \dots, a^n) \\
&+ \zeta(s'|s, a^1, a^2, \dots, a^n)
\end{aligned}$$

where $-\zeta < \zeta(s'|s, a^1, a^2, \dots, a^n) < \zeta$.

Since the agents in reinforcement learning don't know transition probabilities they have only an approximation $\tilde{\Gamma}$ of the actual game Γ at each learning stage. We suppose though that the agents have already learned the reward functions.

Further on $\tilde{\cdot}$ will indicate that we are using an approximation of transitions \tilde{p} instead of the actual transitions p for calculation of the corresponding values.

Let's consider the following system of differential equations:

$$\frac{dx_{sh}^k}{dt} = [\vartheta_{sh}^k(x) - \tilde{v}_s^k(x)]x_{sh}^k \quad k \in K, s \in S, h \in A^k \quad (1)$$

where

$$\vartheta_{sh}^k(x) = r^k(s, \chi) + \gamma \sum_{s' \in S} \tilde{p}(s'|s, \chi) \tilde{v}_{s'}^k(x)$$

and χ denotes the profile equal to x but where player k plays pure strategy h in state s .

Let $x_{sh}^k(t)$ for $k \in K, s \in S$ and $h \in A^k$ be the solution of the system of differential equations 1 satisfying some initial conditions:

$$x_{sh}^k(0) = x_{sh}^{k(0)} \in (0, 1), \quad \sum_{h \in A^k} x_{sh}^k(0) = 1$$

Let y denote the profile where $y_{sh}^k = \frac{1}{T} \int_0^T x_{sh}^k(t) dt$ for some T .

And let $v = (v_s^k)_{s \in S, k \in K}$ denote the matrix where $v_s^k = \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt$.

Let

$$\nu_{sh}^k(y) = r^k(s, \varphi) + \gamma \sum_{s' \in S} \tilde{p}(s'|s, \varphi) v_{s'}^k$$

where φ stands for the profile equal to y but where player k plays pure strategy h in state s .

THEOREM 2. *If $T, \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_4 \in \mathbb{R}$ are such that:*

1. *for all $k \in K$ and $s \in S \exists C1_s^k$ and $C2_s^k : C1_s^k \cup C2_s^k = A^k$ and $C1_s^k \cap C2_s^k = \emptyset$ and such that*

$$(a) \text{ for all } h \in C1_s^k : \left| \frac{1}{T} \ln x_{sh}^k(T) - \frac{1}{T} \ln x_{sh}^k(0) \right| < \epsilon_1$$

$$(b) \text{ for all } h \in C2_s^k : \frac{1}{T} \ln x_{sh}^k(T) - \frac{1}{T} \ln x_{sh}^k(0) \leq -\epsilon_1$$

$$(c) \sum_{h \in C2_s^k} y_{sh}^k (\max_{i \in A^k} \nu_{si}^k - \nu_{sh}^k) < \epsilon_2$$

2. for all $k \in K$, $s, s' \in S$ and $(a^1, a^2, \dots, a^n) \in A^1 \times A^2 \times \dots \times A^n$ the following holds:

$$(a) \left| \frac{1}{T} \int_0^T \prod_{i=1, i \neq k}^n x_{sa^i}^i(t) dt - \prod_{i=1, i \neq k}^n y_{sa^i}^i \right| < \epsilon_3$$

$$(b) \left| \frac{1}{T} \int_0^T \tilde{v}_{s'}^k \prod_{i=1, i \neq k}^n x_{sa^i}^i dt - v_{s'}^k \prod_{i=1, i \neq k}^n y_{sa^i}^i \right| < \epsilon_4$$

then

$y_{sh}^k = \frac{1}{T} \int_0^T x_{sh}^k(t) dt$ constitute ϵ -equilibrium of discounted stochastic game Γ where $\epsilon = [2\gamma(\sigma + \zeta N \max_{k,s} |v_s^k| + N\zeta\sigma) + \epsilon] \sum_{t=0}^{\infty} \gamma^t$ and $\epsilon = 2\epsilon_1 + 2R_{\max}\epsilon_3 \prod_{k=1}^n m^k + 2\gamma\epsilon_4 \prod_{k=1}^n m^k + \epsilon_2$ and $\sigma = 3\epsilon_1 + 3R_{\max}\epsilon_3 \prod_{k=1}^n m^k + 3\gamma\epsilon_4 \prod_{k=1}^n m^k + \epsilon_2$.

PROOF. Let

$$b_{sh}^k = \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \dots \sum_{a^n \in A^n} r^k(s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \prod_{i=1, i \neq k}^n y_{sa^i}^i + \gamma \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \dots \sum_{a^n \in A^n} \sum_{s' \in S} \tilde{p}(s'|s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) v_{s'}^k \prod_{i=1, i \neq k}^n y_{sa^i}^i$$

and

$$b_s^k = \sum_{h \in A^k} b_{sh}^k y_{sh}^k$$

If we could show that for some ϵ and σ and for all $k \in K$, $s \in S$ and $h \in A^k$:

- $b_{sh}^k \leq b_s^k + \epsilon$ (in other words that $y_{sh}^k = \frac{1}{T} \int_0^T x_{sh}^k(t) dt$ constitute ϵ -equilibrium of corresponding games with equilibrium payoffs $(b_s^1, b_s^2, \dots, b_s^n)$)
- $|b_s^k - v_s^k| < \sigma$

then by applying the theorem 1 we could get the implication in question.

Let's consider arbitrary agent $k \in K$ and state $s \in S$.

Without losing generality let $C1_s^k = \{a_1^k, \dots, a_l^k\}$ and $C2_s^k = \{a_{l+1}^k, \dots, a_m^k\}$.

Let's consider the first case:

Since $T < \infty$ $x_{sh}^k(t) > 0$ on $[0, T]$ and

$$(\ln x_{sh}^k)' = \frac{x_{sh}^k'}{x_{sh}^k} = \vartheta_{sh}^k(x) - \tilde{v}_s^k(x)$$

which gives, if one integrates from 0 to T and then divides by T :

$$\frac{1}{T} \ln x_{sh}^k(T) - \frac{1}{T} \ln x_{sh}^k(0) = \frac{1}{T} \int_0^T \vartheta_{sh}^k(x(t)) dt - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt$$

Thus

$$\left| \frac{1}{T} \int_0^T \vartheta_{sa_1^k}^k(x(t)) dt - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt \right| < \epsilon_1$$

⋮

$$\left| \frac{1}{T} \int_0^T \vartheta_{sa_l^k}^k(x(t)) dt - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt \right| < \epsilon_1$$

Upon using the properties of integral we get:

$$\begin{aligned} & \frac{1}{T} \int_0^T \vartheta_{sh}^k(x(t)) dt \\ &= \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \dots \sum_{a^n \in A^n} r^k(s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \\ & \cdot \frac{1}{T} \int_0^T \prod_{i=1, i \neq k}^n x_{sa^i}^i(t) dt \\ &+ \gamma \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \dots \sum_{a^n \in A^n} \sum_{s' \in S} \tilde{p}(s'|s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \\ & \cdot \frac{1}{T} \int_0^T \tilde{v}_{s'}^k(x(t)) \prod_{i=1, i \neq k}^n x_{sa^i}^i(t) dt \end{aligned}$$

Let

$$\frac{1}{T} \int_0^T \prod_{i=1, i \neq k}^n x_{sa^i}^i(t) dt - \prod_{i=1, i \neq k}^n \frac{1}{T} \int_0^T x_{sa^i}^i(t) dt = \epsilon_{3sq}$$

$$q = 1, \dots, m^1 m^2 \dots m^n$$

and

$$\begin{aligned} & \frac{1}{T} \int_0^T \tilde{v}_{s'}^k(x(t)) \prod_{i=1, i \neq k}^n x_{sa^i}^i(t) dt - \frac{1}{T} \int_0^T \tilde{v}_{s'}^k(x(t)) dt \\ & \cdot \prod_{i=1, i \neq k}^n \frac{1}{T} \int_0^T x_{sa^i}^i(t) dt = \epsilon_{4s's'j} \end{aligned}$$

$$j = 1, \dots, m^1 m^2 \dots m^n$$

And according to prerequisite 2:

$$-\epsilon_3 < \epsilon_{3sq} < \epsilon_3 \text{ and } -\epsilon_4 < \epsilon_{4s's'j} < \epsilon_4.$$

Thus we get:

$$\begin{aligned} & \frac{1}{T} \int_0^T \vartheta_{sh}^k(x(t)) dt \\ &= \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \dots \sum_{a^n \in A^n} r^k(s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \\ & \cdot \left(\prod_{i=1, i \neq k}^n \frac{1}{T} \int_0^T x_{sa^i}^i(t) dt + \epsilon_{3sq} \right) \\ &+ \gamma \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \dots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \dots \sum_{a^n \in A^n} \sum_{s' \in S} \tilde{p}(s'|s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \\ & \cdot \left(\frac{1}{T} \int_0^T \tilde{v}_{s'}^k(x(t)) dt \cdot \prod_{i=1, i \neq k}^n \frac{1}{T} \int_0^T x_{sa^i}^i(t) dt + \epsilon_{4s's'j} \right) \end{aligned}$$

$$\begin{aligned}
& \frac{1}{T} \int_0^T \vartheta_{sh}^k(x(t)) dt \\
= & \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \cdots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \cdots \sum_{a^n \in A^n} \\
& r^k(s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \\
& \cdot \prod_{i=1, i \neq k}^n \frac{1}{T} \int_0^T x_{sa^i}^i(t) dt \\
+ & \gamma \sum_{a^1 \in A^1} \sum_{a^2 \in A^2} \cdots \sum_{a^{k-1} \in A^{k-1}} \sum_{a^{k+1} \in A^{k+1}} \cdots \sum_{a^n \in A^n} \sum_{s' \in S} \\
& \tilde{p}(s' | s, a^1, a^2, \dots, a^{k-1}, h, a^{k+1}, \dots, a^n) \\
& \cdot \frac{1}{T} \int_0^T \tilde{v}_{s'}^k(x(t)) dt \cdot \prod_{i=1, i \neq k}^n \frac{1}{T} \int_0^T x_{sa^i}^i(t) dt + \epsilon_{5^k_{sh}}
\end{aligned}$$

where $-\epsilon_5 < \epsilon_{5^k_{sh}} < \epsilon_5$ and $\epsilon_5 = m^1 m^2 \dots m^n (R_{\max} \epsilon_3 + \gamma \epsilon_4) = (R_{\max} \epsilon_3 + \gamma \epsilon_4) \prod_{k=1}^n m^k$

Apparently $\frac{1}{T} \int_0^T \vartheta_{sh}^k(x(t)) dt = b_{sh}^k + \epsilon_{5^k_{sh}}$.

Hence we will have the following inequalities:

$$-\epsilon_6 < b_{sa_1^k}^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt < \epsilon_6$$

⋮

$$-\epsilon_6 < b_{sa_l^k}^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt < \epsilon_6$$

where $\epsilon_6 = \epsilon_1 + \epsilon_5$

So the difference between b_{sh}^k for $h \in C1_s^k$ and $\frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt$ won't exceed ϵ_6 . Hence the difference between any two $b_{sh_1}^k$ and $b_{sh_2}^k$, $h_1, h_2 \in C1_s^k$ won't be more than $2\epsilon_6$.

For all $h \in C2_s^k$ the following condition holds:

$$\frac{1}{T} \ln x_{sh}^k(T) - \frac{1}{T} \ln x_{sh}^k(0) \leq -\epsilon_1$$

$$\begin{aligned}
& \frac{1}{T} \int_0^T \vartheta_{sh}^k(x(t)) dt - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt \\
= & \frac{1}{T} \ln x_{sh}^k(T) - \frac{1}{T} \ln x_{sh}^k(0) \leq -\epsilon_1
\end{aligned}$$

$$b_{sh}^k + \epsilon_{5^k_{sh}} - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt \leq -\epsilon_1$$

$$b_{sh}^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt < \epsilon_5 - \epsilon_1 < \epsilon_5 + \epsilon_1 = \epsilon_6$$

Let $b_{sh_*}^k = \max_{h \in A^k} b_{sh}^k$.

If $h_* \in C1_s^k$ then for any $h \in C1_s^k$ the difference between corresponding b_{sh}^k and $b_{sh_*}^k$ won't exceed $2\epsilon_6$ (as we have already demonstrated the difference between any two $b_{sh_1}^k$ and $b_{sh_2}^k$, $h_1, h_2 \in C1_s^k$ won't be more than $2\epsilon_6$). If $h_* \in C2_s^k$ then for any $h \in C1_s^k$ the difference between corresponding b_{sh}^k and $b_{sh_*}^k$ also won't exceed $2\epsilon_6$ (b_{sh}^k from $C2_s^k$ that deviates from $\frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt$ on more than ϵ_6 can't for sure be

the maximal for the whole A^k because it will be less than any b_{sh}^k for $h \in C1_s^k$).

The condition

$$\sum_{h \in C2_s^k} y_{sh}^k (\max_{i \in A^k} \nu_{si}^k - \nu_{sh}^k) < \epsilon_2$$

we can rewrite as

$$\sum_{h \in C2_s^k} b_{sh_*}^k y_{sh}^k - \sum_{h \in C2_s^k} b_{sh}^k y_{sh}^k < \epsilon_2$$

$$\sum_{h \in C2_s^k} b_{sh}^k y_{sh}^k > \sum_{h \in C2_s^k} b_{sh_*}^k y_{sh}^k - \epsilon_2$$

$$\begin{aligned}
b_s^k &= \sum_{h \in A^k} b_{sh}^k y_{sh}^k = \sum_{h \in C1_s^k} b_{sh}^k y_{sh}^k + \sum_{h \in C2_s^k} b_{sh}^k y_{sh}^k > \\
& \sum_{h \in C1_s^k} (b_{sh_*}^k - 2\epsilon_6) y_{sh}^k + \sum_{h \in C2_s^k} b_{sh_*}^k y_{sh}^k - \epsilon_2 = \\
& b_{sh_*}^k \sum_{h \in A^k} y_{sh}^k - 2\epsilon_6 \sum_{h \in C1_s^k} y_{sh}^k - \epsilon_2 > b_{sh_*}^k - 2\epsilon_6 - \epsilon_2
\end{aligned}$$

Thus the first inequality $b_{sh}^k \leq b_s^k + \epsilon$ that we must prove will hold with $\epsilon = 2\epsilon_6 + \epsilon_2$ for all $h \in A^k$.

As we have just demonstrated for all $h \in A^k$:

$$b_{sh}^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt < \epsilon_6$$

If we multiply each inequality by y_{sh}^k accordingly and sum up we will get:

$$\sum_{h \in A^k} b_{sh}^k y_{sh}^k - \sum_{h \in A^k} \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt \cdot y_{sh}^k < \sum_{h \in A^k} \epsilon_6 y_{sh}^k$$

$$b_s^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt < \epsilon_6$$

From the first inequality that we have already proved and the estimations of b_{sh}^k for $h \in C1_s^k$ we can derive:

$$-\epsilon_6 < b_{sh}^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt \leq b_s^k + \epsilon - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt$$

$$-\epsilon_6 - \epsilon < b_s^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt$$

The second inequality:

$|b_s^k - \frac{1}{T} \int_0^T \tilde{v}_s^k(x(t)) dt| < \sigma$ will hold with $\sigma = 3\epsilon_6 + \epsilon_2$. Let's calculate ϵ and σ .

$$\epsilon = 2\epsilon_6 + \epsilon_2 = 2(\epsilon_1 + \epsilon_5) + \epsilon_2 =$$

$$= 2(\epsilon_1 + (R_{\max} \epsilon_3 + \gamma \epsilon_4) \prod_{k=1}^n m^k) + \epsilon_2 =$$

$$= 2\epsilon_1 + 2R_{\max} \epsilon_3 \prod_{k=1}^n m^k + 2\gamma \epsilon_4 \prod_{k=1}^n m^k + \epsilon_2$$

$$\sigma = 3\epsilon_6 + \epsilon_2 = 3(\epsilon_1 + \epsilon_5) + \epsilon_2 =$$

$$= 3(\epsilon_1 + (R_{\max} \epsilon_3 + \gamma \epsilon_4) \prod_{k=1}^n m^k) + \epsilon_2 =$$

$$= 3\epsilon_1 + 3R_{\max} \epsilon_3 \prod_{k=1}^n m^k + 3\gamma \epsilon_4 \prod_{k=1}^n m^k + \epsilon_2$$

Applying the theorem 1 we get the the implication in question. \square

4. DISCUSSION AND EXPERIMENTAL ESTIMATIONS

Let's consider the conditions of the theorem 2 in detail. For each orbit $x_{sh}^k(t)$ there are only two possibilities:

1. for any $t \in [0, \infty)$ the orbit $x_{sh}^k(t)$ remains bounded from 0 on some value $\delta > 0$
2. $x_{sh}^k(t)$ comes arbitrarily close to 0

In the first case we can reduce ϵ_1 arbitrarily by increasing T (k belongs to $C1_s^k$ in this case).

In the second case if the condition on ϵ_1 for class $C1_s^k$ holds — k belongs to $C1_s^k$ otherwise to $C2_s^k$ ($\frac{1}{T} \ln x_{sh}^k(T) - \frac{1}{T} \ln x_{sh}^k(0) > 0$ will never be true for big enough T).

We can arbitrarily decrease ϵ_2 by increasing T in the second case since ν_{sh}^k is a bounded function.

ϵ_3 and ϵ_4 are much more difficult to deal with. . .

In general the systems of differential equation can be solved:

1. analytically (solution in explicit form)
2. qualitatively (with the use of vector fields)
3. numerically (numerical methods, e.g., Runge-Kutta)

It is hopeless to try to solve the system of such complexity as 1 by the first two approaches and therefore a proof that its solutions satisfy the prerequisites of the theorem 2 seems to us non-trivial. Till now we have managed to find ϵ_3 and ϵ_4 estimations only experimentally.

In table 1 estimations of average relative $\overline{\epsilon_{s_{sh}^k}}$ and average relative $\overline{\epsilon_5}$ are presented for different game classes (with different number of states, agents and actions). The averages are calculated for 100 games of each class and $T = 1000$. The initial conditions for the system of differential equations 1 were chosen quite randomly. The games are generated with uniformly distributed payoffs. Transition probabilities were also derived from uniform distribution. As we can see the preconditions of the theorem 2 hold with a quite acceptable accuracy for all the classes.

5. NASH-DE ALGORITHM

To propose an algorithm we have to make one more assumption that we have managed to confirm only experimentally till now, namely:

The more accurate approximation of Nash-equilibrium we choose as an initial condition for our system 1 the more precisely the prerequisites of the theorem 2 hold.

So now we can propose an iterative algorithm for calculating ϵ -equilibria of discounted stochastic games with some given accuracy ϵ (see algorithm 1).

An example of its work on a 2-state 2-agent 2-action discounted stochastic game is presented on the figure 1 (because of the space restrictions we are illustrating convergence only for state s_1 but no state can be examined in isolation for analysis). On each figure the probabilities assigned to the first actions of the first and the second agents are presented as xy -plot (it is quite descriptive since the probabilities of the second actions are equal to one minus probabilities of

Algorithm 1 Nash-DE algorithm for the player k

Input: accuracy ϵ, T
for all $s \in S, k \in K$ and $h \in A^k$ **do**
 $x_{sh}^k(0) \leftarrow 1/|A^k|$
end for
while $x(0)$ doesn't constitute ϵ -equilibrium **do**
 Find solution of the system 1 through the point $x(0)$ on the interval $[0, T]$ (updating model in parallel)
 Let the initial point be $x_{sh}^k(0) = \frac{1}{T} \int_0^T x_{sh}^k(t) dt$
end while

the first ones). The solutions are lighter at the end of $[0, T]$ interval. The precise Nash-equilibrium is designated by a star and the average $\frac{1}{T} \int_0^T x_{sh}^k(t) dt$ for each iteration — by a cross. Since the agents in reinforcement learning don't know either transition probabilities or reward functions and they learn them online the first policies are quite random. The algorithm converges in self-play to Nash-equilibrium with the given relative accuracy $\epsilon = 1\%$ in two iterations.

6. EXPERIMENTAL RESULTS

Since the agents in reinforcement learning don't know either transition probabilities or reward functions they have to approximate the model somehow. We tested our algorithm as an off-policy version (the agents pursue the best learned policy so far in the most of cases (we chose — 90% of cases) and explore the environment in 10% of cases). The results of the experiments are presented in table 1. The number of independent transitions to be learned can be calculated by the formula $Tr = N(N - 1) \prod_{k=1}^n m^k$ and is presented in the corresponding column for each game class. In column "Iterations" the average number of iterations necessary to find Nash-equilibrium with relative accuracy $\epsilon = 1\%$ is presented. And in the last column — the percentage of games for which we managed to find Nash-equilibrium with the given relative accuracy $\epsilon = 1\%$ after 500 iterations.

In general one can see the following trend: the larger is the model the more iterations the agents require to find a 1%-equilibrium, and the oftener they fail to come to this equilibrium during 500 iterations.

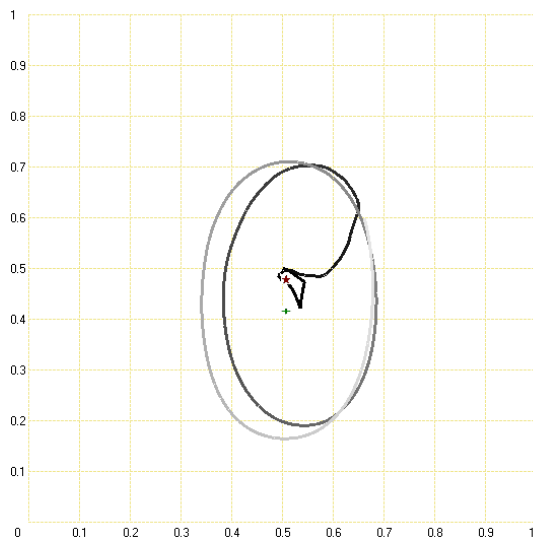
The main reason for it is their incapability to approximate large models to the necessary accuracy (their approximations of transition probabilities are too imprecise — they explore the environment only in 10% of cases each and the transition probabilities of some combinations of actions remain very poorly estimated) and as a result they can't find an equilibrium or converge to it more slowly (let us not forget that the accuracy of transition probabilities acts as a relative factor and comes to ϵ estimation of theorem 2 multiplied by the maximal discounted value). In order to decrease the average number of iterations and to increase the percentage of solved games it appears promising to test a version of the algorithm with a more intensive exploration stage (first learn the model to some given precision and only then act according to the policy found by the algorithm and keep on learning in parallel). For instance, it can be achieved by setting ϵ to larger values at the beginning.

7. CONCLUSION

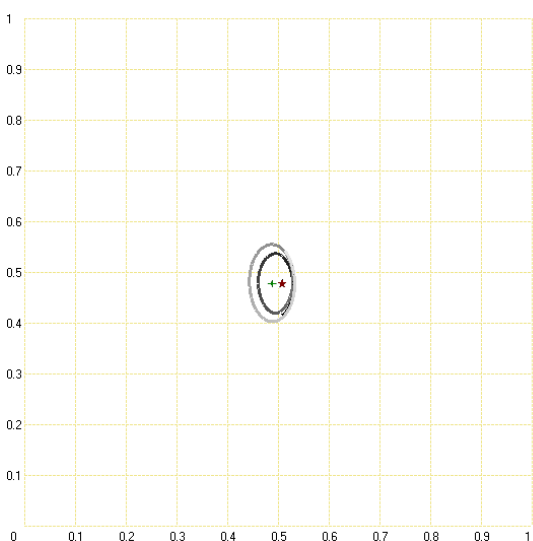
This paper is devoted to an actual topic of extending reinforcement learning approach for multiagent systems. An

Table 1: Estimations and Results of Experiments

States	Agents	Actions	Tr	$\overline{\epsilon_{sh}^k}$	$\overline{\epsilon_5}$	Iterations	%
2	2	2	8	0.08%	0.24%	11.23	98%
2	2	3	18	0.20%	0.36%	9.43	95%
2	2	5	50	0.16%	0.25%	18.60	90%
2	2	10	200	0.48%	0.73%	38.39	94%
2	3	2	16	0.18%	0.85%	16.03	87%
2	3	3	54	0.68%	1.74%	30.64	91%
2	5	2	64	1.80%	4.36%	27.79	87%
5	2	2	80	0.00%	0.04%	31.60	83%
5	2	3	180	0.14%	0.22%	52.26	93%
5	2	5	500	0.10%	0.14%	62.74	91%
5	3	3	540	0.35%	1.58%	85.83	75%
10	2	2	360	0.02%	0.06%	69.68	82%



(a) Iteration 1 State 1



(b) Iteration 2 State 1

Figure 1: Convergence of Algorithm 1

algorithm based on system of differential equations of special type is developed. A formal proof of its convergence with a given accuracy to Nash equilibrium for environments represented as general-sum discounted stochastic games is given under some assumptions. We claim that it is the first algorithm that converges to Nash equilibrium in general case. Thorough testing showed that the assumptions necessary for the formal convergence hold with quite a good accuracy that allowed the proposed algorithm to find Nash equilibrium with relative accuracy of 1% in approximately 90% of randomly generated games.

8. REFERENCES

- [1] M. H. Bowling and M. M. Veloso. Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136(2):215–250, 2002.
- [2] C. Claus and C. Boutilier. The dynamics of reinforcement learning in cooperative multiagent systems. In *AAAI '98/IAAI '98: Proceedings of the fifteenth national/tenth conference on Artificial intelligence/Innovative applications of artificial intelligence*, pages 746–752, Menlo Park, CA, USA, 1998. American Association for Artificial Intelligence.
- [3] J. Filar and K. Vrieze. *Competitive Markov decision processes*. Springer-Verlag New York, Inc., New York, NY, USA, 1996.
- [4] A. Greenwald and K. Hall. Correlated-q learning. In *Proceedings of the Twentieth International Conference on Machine Learning*, pages 242–249, 2003.
- [5] J. Hu and M. P. Wellman. Multiagent reinforcement learning: theoretical framework and an algorithm. In *Proc. 15th International Conf. on Machine Learning*, pages 242–250. Morgan Kaufmann, San Francisco, CA, 1998.
- [6] M. L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *ICML*, pages 157–163, 1994.
- [7] M. L. Littman. Friend-or-foe q-learning in general-sum games. In C. E. Brodley and A. P. Danyluk, editors, *ICML*, pages 322–328. Morgan Kaufmann, 2001.
- [8] G. Tesauro. Extending q-learning to general adaptive multi-agent systems. In S. Thrun, L. K. Saul, and B. Schölkopf, editors, *NIPS*. MIT Press, 2003.
- [9] M. Zinkevich, A. R. Greenwald, and M. L. Littman. Cyclic equilibria in markov games. In *NIPS*, 2005.