

# Adaptive Zero-Sum Stochastic Game for Two Finite Markov Chains

A. S. Poznyak\* and K. Najim\*\*

\*CINVESTAV-IPN, Control Automatico, A.P. 14-740, C.P.07300 Mexico D.F., Mexico,  
fax: (52-5)747-70-89 or 747-70-02, e-mail: apoznyak@ctrl.cinvestav.mx

\*\*Process Control Laboratory, E.N.S.I.G.C., Chemin de la loge, 31078 Toulouse cedex, France,  
fax: (+33) 5 62 25 23 18, e-mail: Kaddour.Najim@ensigct.fr

## Abstract

A two finite Markov chains repeated zero-sum stochastic game with unknown transition matrices and payoffs is considered. The control objective is to obtain the equilibrium point based only on current measurements. The behavior of each players is modelled by a finite controlled Markov chain. A novel adaptive policy is developed based on Lagrange multipliers involved into "learning through reinforcement" procedure. A regularized Lagrange function and a new normalization procedure are introduced. The saddle-point of this function is shown to be unique. The convergence properties are proved and the order of almost sure convergence is estimated as  $(n^{-\frac{1}{3}})$ .

*Keywords:* Repeated stochastic game; adaptive control; controlled Markov chains,, learning, reinforcement.

## 1 Introduction

A lot of engineering, business and economic problems can be formulated as stochastic games (see, for example [1]). Finite *stochastic games* can be regarded as competitive *Markov Decision Processes* [2] where there are two or more controllers (players). There is an extensive literature on the existence of optimal strategies [5], [14], [11] and [6]. Several studies have been dedicated to the development of algorithms for the computation of equilibrium points and optimal strategies in stochastic games [10], [12], [4] and [13]. The term "algorithm" was interpreted broadly enough [10]. An *adaptive* or *learning algorithm* (recursive) can be defined as a procedure which forms a new estimate, incorporating new information (realizations), from the old estimate using a fixed amount of computations and memory.

This paper presents a novel algorithm for two finite Markov chains repeated zero-sum stochastic game with incomplete information on the transitions and payoffs. The characteristics (statistics, etc.) of the game are a

*priori* unknown. Only the realizations (states and realized payoffs) of the game are available at each stage. Each player is modelled by a finite controlled irreducible Markov chain [9]. The control objective corresponds to the optimization of the random limiting average payoff which captures the long-run average performance. The algorithm presented in this paper is adaptive in the sense that it provides learning control policies for both players (Markov chains). It is based on the Lagrange multipliers approach and a new normalization procedure in order to preserve the probability measure.

## 2 System Description

Introduce the following notations:

$X_k = (x^k(1), \dots, x^k(K_k))$  is the finite state spaces associated respectively to the first player ( $k = 1$ ) and to the second ( $k = 2$ ) player;

$U_k = (u^k(1), \dots, u^k(N_k))$  is respectively the set of actions;

$\{\eta_n\}$  is a sequence of Borel functions (below, payoffs)

$$\eta_n = \eta_n(\omega, x_n^1, x_n^2, u_n^1, u_n^2) \\ \omega \in \Omega, x_n^1 \in X_1, x_n^2 \in X_2, u_n^1 \in U_1, u_n^2 \in U_2$$

given at a probability space  $(\Omega, \mathcal{F}, P)$  with a fixed flow  $\mathcal{F}_n \subseteq \mathcal{F}_{n+1}$  of  $\sigma$ -algebras. Here  $x_n^k, u_n^k$  ( $k = 1, 2$ ) are the current states and control actions.

The behavior of each player is modelled by a finite homogeneous controlled Markov chain. These Markov chains are assumed to be *irreducible* (completely ergodic, that is, each state can be visited infinitely often) [9].

The game is played in stages ( $n = 1, 2, \dots$ ). The play starts at stage 1 in the initial states  $x_1^1$  and  $x_1^2$ . Each of the players is allowed to randomize over pure action choices ( $u_1^1 \in U_1$  and  $u_1^2 \in U_2$ ). These choices

induce immediate random payoffs  $\eta_1$  and so on. The first player tries to maximize the payoff he can guarantee himself by playing appropriate. The second player is interested in the minimization of the payoff which he at most has to pay.

**Definition 1**  $\{\eta_n\}$  ( $n = 1, 2, \dots$ ) is said to be the **pay-off function sequence** for the first player (loss function for the second one) if  $|\eta_n| \stackrel{a.s.}{\leq} \sigma^+ < \infty$  and

$$E \left\{ \eta_n \mid x_n^1 = x^1(i), x_n^2 = x^2(j); \right. \\ \left. u_n^1 = u^1(l), u_n^2 = u^2(r) \right\} = \frac{v_{ij}^{lr}}{i = \overline{1, K_1}, j = \overline{1, K_2}, l = \overline{1, N_1}, r = \overline{1, N_2}}$$

The play moves to new states ( $x_2^1$  and  $x_2^2$ ) according to the transition probabilities

$$\Pi^k = \left[ \pi_{ij}^{k,l} \right] \quad (k = 1, 2)$$

The element

$$\pi_{ij}^{k,l} \quad (i = 1, \dots, K_k; j = 1, \dots, K_k \text{ and } l = 1, \dots, N_k)$$

represents at time  $n$  ( $n = 1, 2, \dots$ ) the probability of transition from state  $x^k(i)$  to state  $x^k(j)$  under the action  $u^k(l) \in U^k$ , that is,

$$\pi_{ij}^{k,l} = P \{ x_{n+1}^k = x^k(j) \mid x_n^k = x^k(i), u_n^k = u^k(l) \}$$

A sequence of random stochastic matrices  $\{d_n^k\}$  is said to be an **admissible randomized (mixed) control strategy**  $D^k$  for the  $k^{th}$  player if

1) it is causal (independent on the future), that is,

$$d_n^k = \left[ d_n^{k,il} \right]_{i=1, \dots, K_k; l=1, \dots, N_k}$$

is  $F_{n-1}$ -measurable, where

$$F_{n-1} := \sigma \left( x_1^k, u_1^k; \dots; x_{n-1}^k, u_{n-1}^k; k = 1, 2 \right)$$

is the  $\sigma$ -algebra generated by the random variables  $(x_1^k, u_1^k; \dots; x_{n-1}^k, u_{n-1}^k)$  up to time  $n - 1$ ;

2) the random variables  $(u_1^k, \dots, u_{n-1}^k)$  represent the realizations of the control actions, taking values on the finite sets  $U^k = \{u^k(1), \dots, u^k(N_k)\}$ , and satisfying the following property:

$$d_n^{k,il} = \Pr \{ u_n^k = u^k(l) \mid x_n^k = x^k(i) \wedge F_{n-1} \}$$

Denote by  $\Sigma^k$  the class of all **randomized strategies for the  $k^{th}$  player**, that is,  $\Sigma^k = \{ \{d_n^k\} \}$ . The conditional transition probability matrix  $\Pi^k(d_n^k)$  can be defined as

$$\Pi^k(d_n^k) = \left[ \pi_n^{k,ij}(d_n^k) \right]_{i=1, \dots, K_k; j=1, \dots, K_k} \\ \pi_n^{k,ij}(d_n^k) := P \{ x_{n+1}^k = x^k(j) \mid x_n^k = x^k(i) \wedge F_{n-1} \} \\ \stackrel{a.s.}{=} \sum_{l=1}^{N_k} \pi_{ij}^{k,l} d_n^{k,il}$$

represents the law of motion among states. The *average payoff*  $v_n^1$  of the  $1^{th}$  player at the  $n^{th}$  stage is given by

$$v_n^1 = V(d_n^1, d_n^2) \\ := \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} \sum_{l=1}^{N_1} \sum_{r=1}^{N_2} v_{ij}^{lr} p_n^1(i) p_n^2(j) d_n^{1,il} d_n^{2,jr}$$

where  $p_n^k(i_k)$  is the probability for the  $k^{th}$  player to be in the state  $x^k(i_k)$  at time  $n$  verifying

$$p_{n+1}^k(j) = \sum_{i=1}^{K_k} \pi_n^{k,ij}(d_n^k) p_n^k(i) \\ = \sum_{i=1}^{K_k} \left( \sum_{l=1}^{N_k} \pi_{ij}^{k,l} d_n^{k,il} \right) p_n^k(i)$$

Introduce the following functions

$$\Phi_n := \frac{1}{n} \sum_{t=1}^n \xi_t$$

and

$$V_n := \frac{1}{n} \sum_{t=1}^n v_t = \frac{1}{n} \sum_{t=1}^n V(d_t^1, d_t^2)$$

**Lemma 2** Under the accepted assumptions

$$\Phi_n = V_n + o_\omega(n^{1-\varepsilon}), \quad \varepsilon > 0$$

with probability one.

**Proof:** It follows immediately from Lemma 3 for  $v_t := \eta_n$ , and  $a_n := n$ , and Lemma 4 for  $\eta_n := n^{1-\varepsilon}$  (Appendix A, in [9]). ■

The *strategy* of the  $k^{th}$  player consists of any admissible matrix sequences  $D^k = \{d_n^k\}$  ( $k = 1, 2$ ) with  $F_{n-1}$ -measurable elements.

**Definition 3** The pair of strategies  $\bar{D}^1$  and  $\bar{D}^2$  is said to be the *non-cooperative equilibrium strategy* (in the Nash sense) if for any admissible  $D^1$  and  $D^2$

$$\tilde{V}(D^1, \bar{D}^2) \stackrel{a.s.}{\leq} \tilde{V}(\bar{D}^1, \bar{D}^2) \\ := \min_{x_1^1, x_2^2} \liminf_{n \rightarrow \infty} V_n \stackrel{a.s.}{\leq} \tilde{V}(\bar{D}^1, D^2)$$

**Remark 4** It is evident (Flesch et al., 1997) that in the case of irreducible controlled finite Markov chains (Poznyak et al., 1999), the set of all equilibrium strategies  $\bar{D}^1$  and  $\bar{D}^2$  contains the set of stationary strategies  $\{\bar{d}^k\}$  ( $k = 1, 2$ ) where  $(\bar{d}^1, \bar{d}^2)$  represents a saddle-point (may be not unique) of the function  $V(d^1, d^2)$  which satisfies

$$V(d^1, \bar{d}^2) \leq V(\bar{d}^1, \bar{d}^2) \leq V(\bar{d}^1, d^2)$$

for any  $d^k \in \mathbf{D}_\varepsilon^k$  where the simplex  $\mathbf{D}_\varepsilon^k$  is defined as follows

$$\mathbf{D}_\varepsilon^k := \left\{ d^k = \|d^{k,il}\| \left( i = \overline{1, K_k}; l = \overline{1, N_k} \right) : \right. \\ \left. d^{k,il} \geq \varepsilon \geq 0, \sum_{l=1}^{N_k} d^{k,il} = 1 \right\}$$

This fact follows from: i) the continuity property of the function  $V(d^1, d^2)$ , ii) the compactness of the sets  $\mathbf{D}_\varepsilon^k$ , and iii) Von Neumann theorem.

### 3 Problem Setting

Introduce the following variable change [9]

$$c^{k,il} := p^k(i) d^{k,il} \quad (k = 1, 2)$$

We are concerned with a constrained optimization problem:  $c^{k,il}$  must verify (for each  $k = 1, 2$  and  $j = \overline{1, K_k}$ ) the constraint:

$$\sum_{l=1}^{N_k} c^{k,jl} = \sum_{i=1}^{K_k} \sum_{l=1}^{N_k} \pi_{ij}^{k,l} c^{k,il}$$

Consider the following regularized Lagrange function ( $\delta > 0$ )

$$L_\delta := L_\delta(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) \\ = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} \sum_{l=1}^{N_1} \sum_{r=1}^{N_2} v_{ij}^{lr} c^{1,il} c^{2,il} \\ - \sum_{k=1}^2 (-1)^k \sum_{j=1}^{K_k} \lambda_j^k \left[ \sum_{l=1}^{N_k} c^{k,jl} - \sum_{i=1}^{K_k} \sum_{l=1}^{N_k} \pi_{ij}^{k,l} c^{k,il} \right] \\ + (\delta/2) \sum_{k=1}^2 (-1)^k \left( \sum_{i=1}^{K_k} \sum_{l=1}^{N_k} (c^{k,il})^2 - \sum_{j=1}^{K_k} (\lambda_j^k)^2 \right) \\ \mathbf{c}^k := (c^{k,11}, \dots, c^{k,1N_k}; \dots; c^{k,K_k 1}, \dots, c^{k,K_k N_k})^T \\ \boldsymbol{\lambda}^k := (\lambda_1^k, \dots, \lambda_{K_k}^k)^T \in R^{K_k}$$

given in the set  $(S_0^{K_1 N_1} \times R^{K_1}) \times (S_0^{K_2 N_2} \times R^{K_2})$  with

$$S_{\varepsilon=0}^{K_k N_k} := \left\{ \begin{array}{l} \mathbf{c}^k \in R^{N_k K_k} \mid c^{k,il} \geq \varepsilon \geq 0 \\ \sum_{i=1}^{K_k} \sum_{l=1}^{N_k} c^{k,il} = 1 \quad (i = \overline{1, K_k}; l = \overline{1, N_k}) \end{array} \right\}$$

The saddle-point of this regularized (augmented) Lagrange function will be denoted by

$$\left( \mathbf{c}_\delta^{1,*}, \mathbf{c}_\delta^{2,*}, \boldsymbol{\lambda}_\delta^{1,*}, \boldsymbol{\lambda}_\delta^{2,*} \right) := \\ \arg \max_{\mathbf{c}^1 \in S_0^{K_1 N_1}} \min_{\mathbf{c}^2 \in S_0^{K_2 N_2}} \min_{\boldsymbol{\lambda}^1 \in R^{K_1}} \max_{\boldsymbol{\lambda}^2 \in R^{K_2}} L_\delta \quad (1)$$

The function  $L_{\delta>0}$  is strictly concave with respect to  $\mathbf{c}_\delta^{1,*}$  and  $\boldsymbol{\lambda}_\delta^{2,*}$ , and is strictly convex with respect to  $\mathbf{c}_\delta^{2,*}$  and  $\boldsymbol{\lambda}_\delta^{1,*}$ . It follows that its saddle-point is unique

and possess the Lipschitz property with respect to the parameter  $\delta$ :

$$\sum_{k=1}^2 \left( \|\mathbf{c}_{\delta_1}^{k,*} - \mathbf{c}_{\delta_2}^{k,*}\| + \|\boldsymbol{\lambda}_{\delta_1}^{k,*} - \boldsymbol{\lambda}_{\delta_2}^{k,*}\| \right) \\ \leq Const |\delta_1 - \delta_2|, \quad k = 1, 2$$

It has been shown (see Theorem 2 of section 4.3 in [8]) that if  $\delta \rightarrow 0$ , the saddle-point  $(\mathbf{c}_\delta^{1,*}, \mathbf{c}_\delta^{2,*}, \boldsymbol{\lambda}_\delta^{1,*}, \boldsymbol{\lambda}_\delta^{2,*})$  converges to the solution  $(\mathbf{c}_0^{1,*}, \mathbf{c}_0^{2,*}, \boldsymbol{\lambda}_0^{1,*}, \boldsymbol{\lambda}_0^{2,*})$  of the optimization problem (1) with  $\delta = 0$ , which has the minimal norm (in the case of an irreducible chain this point is also unique):

$$\mathbf{c}_\delta^{k,*} \xrightarrow{\delta \rightarrow 0} \mathbf{c}_0^{k,**} := \arg \min_{\substack{\mathbf{c}_0^{1,*} \in S_0^{K_1 N_1}, \boldsymbol{\lambda}^{k,*} \in R^{K_k} \\ \left( \|\mathbf{c}_0^{k,*}\|^2 + \|\boldsymbol{\lambda}_0^{k,*}\|^2 \right)}} \quad (2)$$

(the minimization is done over all saddle-points of the nonregularized Lagrange functions).

Now the main problem considered in this paper can be formulated as follows: for the given zero-sum stochastic game with incomplete information, find an adaptive non-cooperative equilibrium strategies  $\bar{\mathbf{D}}^1$  and  $\bar{\mathbf{D}}^2$  using only the on-line (current) information. To achieve this objective, the recursive procedures (adaptive policies)

$$\mathbf{c}_{n+1}^k = \mathbf{c}_{n+1}^k(x_n^k, u_n^k, \eta_n, x_{n+1}^k, \mathbf{c}_n^k; k = 1, 2)$$

should be constructed generating the sequences  $\{\mathbf{c}_n^k\}$  which converges in some probability sense to the solutions  $\mathbf{c}_0^{k,**}$  defined by (2). In this process the player's actions are selected according to the randomized (mixed) rule

$$d_n^{k,il} = c_n^{k,il} / \sum_{l=1}^{N_k} c_n^{k,il}$$

(the denominator must be strictly positive). The need of randomization is intuitively obvious when playing against an intelligent opponent.

### 4 Adaptive Strategy

Consider the following iterative procedure:

*Step 1.* Use the available data

$$x_n^k = x^k(\alpha_k), u_n^k = u^k(\beta_k), \eta_n \\ x_{n+1}^k = x^k(\gamma_k), \mathbf{c}_n^k (c_n^{k,il} > 0), \boldsymbol{\lambda}_n$$

to build the following functions

$$\xi_n^k := \eta_n - (-1)^k (\lambda_{\alpha_k}^k - \lambda_{\gamma_k}^k) + (-1)^k \delta_n c_n^{k, \alpha_k \gamma_k}$$

and use the following normalization procedure

$$\begin{aligned} \zeta_n^1 &:= 1 - (a_n^1 \xi_n^1 + b_n^1) / c_n^{1, \alpha_k \gamma_k} \\ \zeta_n^2 &:= (a_n^2 \xi_n^2 + b_n^2) / c_n^{2, \alpha_k \gamma_k} \end{aligned} \quad (3)$$

$$a_n^k := \left( 2 \frac{(\sigma + 2\lambda_n^+)}{\varepsilon_n} + \frac{N_k K_k}{N_k K_k - 1} \delta_n \right)^{-1}$$

$$b_n^k := a_n^k (\sigma + 2\lambda_n^+)$$

*Step 2.* Calculate the elements  $c_{n+1}^{k,il}$  and  $\lambda_{n+1}^{k,j}$  as

$$c_{n+1}^{k,il} = \begin{cases} c_n^{k, \alpha_k \gamma_k} + \gamma_n^{k,c} (1 - c_n^{k, \alpha_k \gamma_k} - \zeta_n^k) \\ \quad (i = \alpha_k, j = \gamma_k) \wedge l = \beta_k \\ c_n^{k,il} - \gamma_n^{k,c} \left( c_n^{k,il} - \frac{\zeta_n^k}{N_k K_k - 1} \right) \\ \quad (i \neq \alpha_k, j \neq \gamma_k) \vee l \neq \beta_k \end{cases}$$

$\gamma_n^{k,c} \in [0, 1]$ , and the Lagrange multipliers are adjusted as follows:

$$\lambda_{n+1}^{k,j} = \left[ \lambda_n^{k,j} + \gamma_n^\lambda \psi_n^{k,j} \right]_{-\lambda_{n+1}^+}^{\lambda_{n+1}^+}$$

$$\psi_n^{k,j} = (-1)^k \left[ \sum_{l=1}^N c_n^{k,jl} - \chi(x^k(j) = x_{n+1}^k) + \delta_n \lambda_n^{k,j} \right]$$

$$[y]_{-\lambda_{n+1}^+}^{\lambda_{n+1}^+} := \begin{cases} y & \text{if } y \in [-\lambda_{n+1}^+, \lambda_{n+1}^+] \\ \lambda_{n+1}^+ & \text{if } y > \lambda_{n+1}^+ \\ -\lambda_{n+1}^+ & \text{if } y < -\lambda_{n+1}^+ \end{cases} \quad (4)$$

*Step 3.* Construct the stochastic matrices

$$d_{n+1}^{k,il} = c_{n+1}^{k,il} \left( \sum_{r=1}^{N_k} c_{n+1}^{k,ir} \right)^{-1} \quad (i = \overline{1, K_k}, l = \overline{1, N_k})$$

and according to

$$\Pr \{ u_{n+1}^k = u^k(l) \mid x_{n+1}^k = x^k(\gamma_k) \wedge \mathcal{F}_n \} = d_{n+1}^{k,\gamma_l}$$

generate randomly new discrete random variables  $u_{n+1}^k$  as in learning stochastic automata implementation [7], and get the new observation (realization)  $\eta_{n+1}$  which corresponds to the transition to states  $x_{n+1}^k$ .

*Step 4* return to Step 1.

The positive sequences  $\{\varepsilon_n\}$ ,  $\{\delta_n\}$ ,  $\{\lambda_n^+\}$ ,  $\{\gamma_n^{k,c}\}$  and  $\{\gamma_n^\lambda\}$  will be defined in what follows.

The normalization procedure (3) is an affine transformation. It is a kind of mapping or projection scheme for obtaining a new variable, namely  $\zeta_n^k$  which belongs to the unit segment in order to preserve the probability measure. In view of Lemma 2 in Chapter 2 [9] it follows that if  $\gamma_n^{k,c} \in [0, 1]$ ,  $\mathbf{c}_n^k \in S_{\varepsilon_n}^{K_k N_k}$ , i.e.,  $c_n^{\alpha\beta} \geq \varepsilon_n$  and  $\delta_n \downarrow 0$ , then  $\zeta_n^k \in [\zeta_n^{k,-}, \zeta_n^{k,+}] \subset [0, 1]$  where

$$\begin{cases} \zeta_n^{1,-} = 1 - a_n^1 \delta_n, \quad \zeta_n^{1,+} = a_n^1 \delta_n / (N_1 K_1 - 1) \\ \zeta_n^{2,-} = a_n^2 \delta_n, \quad \zeta_n^{2,+} = 1 - a_n^2 \delta_n / (N_2 K_2 - 1) \\ \mathbf{c}_{n+1}^k \in S_{\varepsilon_{n+1}}^{K_k N_k} \end{cases}$$

## 5 Convergence Analysis

Introduce the following Lyapunov function

$$W_n := \sum_{k=1}^2 \left( \|\mathbf{c}_n^k - \mathbf{c}_{\delta_n}^{k,*}\|^2 + \|\boldsymbol{\lambda}_n^k - \boldsymbol{\lambda}_{\delta_n}^{k,*}\|^2 \right)$$

starting from

$$n \geq \inf_{l \geq 1} \left\{ t : \max_k \|\boldsymbol{\lambda}_{\delta_n}^{k,*}\| \leq \lambda_n^+ \right\}$$

**Theorem 5** *If the suggested adaptive control is used with the design parameters  $\varepsilon_n$ ,  $\delta_n$ ,  $\lambda_n^+$ ,  $\gamma_n^{k,c}$  and  $\gamma_n^\lambda$  satisfying the following conditions*

$$0 < \delta_n \downarrow 0, \quad 0 < \lambda_n^+ \uparrow \infty, \quad \gamma_n^{k,c} \in (0, 1)$$

$$\gamma_n^\lambda = \frac{N_k K_k}{N_k K_k - 1} \gamma_n^{k,c} a_n^k, \quad \sum_{n=1}^{\infty} \gamma_n^{k,c} \varepsilon_n \delta_n (\lambda_n^+)^{-1} = \infty$$

then

1) if

$$\sum_{n=1}^{\infty} \mu_n < \infty$$

where

$$\begin{aligned} \mu_n &:= (\delta_n - \delta_{n+1})^2 \lambda_n^+ \left( \min_k \gamma_n^{k,c} \varepsilon_n \delta_n \right)^{-1} + \max_k (\gamma_n^{k,c})^2 \\ &\quad + |\delta_n - \delta_{n+1}| \max_k \gamma_n^{k,c} \left( 1 + (1 + \delta_n \lambda_n^+)^2 \varepsilon_n / \lambda_n^+ \right) \end{aligned}$$

then, for any fixed distribution of the initial states of the Markov chains associated to the players, the control policy converges with probability one to the equilibrium point of the game, that is, with probability one

$$W_n \rightarrow 0$$

2) if

$$\frac{\mu_n \lambda_n^+}{\varepsilon_n \delta_n \min_k \gamma_n^{k,c}} \rightarrow 0$$

then, we obtain the convergence in the mean squares sense:

$$E \{ W_n \} \rightarrow 0$$

**Proof:** The regularized Lagrange function (??), is strictly concave with respect to  $\mathbf{c}^1$  and  $\boldsymbol{\lambda}^2$  and strictly convex with respect to  $\mathbf{c}^2$  and  $\boldsymbol{\lambda}^1$ , hence (see Lemma 1 Chapter 4 in [8])

$$\begin{aligned} & (\boldsymbol{\lambda}^1 - \boldsymbol{\lambda}_{\delta}^{1,*})^T \nabla_{\boldsymbol{\lambda}^1} L_{\delta}(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) \\ & - (\mathbf{c}^1 - \mathbf{c}_{\delta}^{1,*})^T \nabla_{\mathbf{c}^1} L_{\delta}(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) \\ & + (\mathbf{c}^2 - \mathbf{c}_{\delta}^{2,*})^T \nabla_{\mathbf{c}^2} L_{\delta}(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) \\ & - (\boldsymbol{\lambda}^2 - \boldsymbol{\lambda}_{\delta}^{2,*})^T \nabla_{\boldsymbol{\lambda}^2} L_{\delta}(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) \end{aligned} \quad (5)$$

$$\geq (\delta/2) \sum_{k=1}^2 \left( \|\mathbf{c}^k - \mathbf{c}_{\delta}^{k,*}\|^2 + \|\boldsymbol{\lambda}^k - \boldsymbol{\lambda}_{\delta}^{k,*}\|^2 \right)$$

Making use of the identity

$$E \{ Z \mid x_n = x(\alpha) \wedge \mathbf{F}_{n-1} \} = \\ E \left\{ Z \frac{\chi(x_n=x(\alpha))}{p(\alpha)} \mid \mathbf{F}_{n-1} \right\}, p(\alpha) > 0$$

it follows

$$E \left\{ \zeta_n^1 \mathbf{e}(x_n^1 \wedge u_n^1) \mid \mathbf{F}_{n-1} \right\} \stackrel{a.s.}{=} \\ 1 - a_n^1 \frac{\partial}{\partial \mathbf{c}^1} L_\delta(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) - b_n^1 \mathbf{e}^{N_1 K_1} \\ E \left\{ \zeta_n^2 \mathbf{e}(x_n^2 \wedge u_n^2) \mid \mathbf{F}_{n-1} \right\} \stackrel{a.s.}{=} \\ a_n^2 \frac{\partial}{\partial \mathbf{c}^2} L_\delta(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) + b_n^2 \mathbf{e}^{N_2 K_2} \\ E \left\{ \boldsymbol{\psi}_n^k \mid \mathbf{F}_{n-1} \right\} \stackrel{a.s.}{=} (-1)^k \frac{\partial}{\partial \boldsymbol{\lambda}^k} L_\delta(\mathbf{c}^1, \mathbf{c}^2, \boldsymbol{\lambda}^1, \boldsymbol{\lambda}^2) \\ \boldsymbol{\psi}_n^k := \left( \psi_n^{k,1}, \dots, \psi_n^{k,K_k} \right)^T$$

where the  $N_k K_k$ -dimensional vector  $\mathbf{e}(x_n \wedge u_n)$  is defined as follows

$$\mathbf{e}(x_n^k \wedge u_n^k) := \begin{bmatrix} \chi(x_n^k = x^k(1), u_n^k = u^k(1)) \\ \dots \\ \chi(x_n^k = x^k(K_k), u_n^k = u^k(1)) \\ \dots \\ \chi(x_n^k = x^k(1), u_n^k = u^k(l)) \\ \dots \\ \chi(x_n^k = x^k(K_k), u_n^k = u^k(l)) \\ \dots \\ \chi(x_n^k = x^k(K_k), u_n^k = u^k(N_k)) \end{bmatrix}$$

To make the analysis easy, let us rewrite the algorithm in a vector form

$$\mathbf{c}_{n+1}^k = \mathbf{c}_n^k + \gamma_n^{k,c} \left( \mathbf{e}(x_n^k \wedge u_n^k) - \mathbf{c}_n^k \right) \\ + \zeta_n^k \frac{\mathbf{e}^{N_k K_k} - N_k K_k \mathbf{e}(x_n^k \wedge u_n^k)}{N_k K_k - 1} \\ \boldsymbol{\lambda}_{n+1}^k = \left[ \boldsymbol{\lambda}_n^k + \gamma_n^\lambda \boldsymbol{\psi}_n^k \right]_{-\lambda_{n+1}^+}^{\lambda_{n+1}^+} \\ \boldsymbol{\psi}_n^k = (-1)^k \left[ \begin{pmatrix} \sum_{l=1}^{N_k} c_n^{k,1l} \\ \vdots \\ \sum_{l=1}^{N_k} c_n^{k,K_k l} \end{pmatrix} + \delta_n \boldsymbol{\lambda}_n^k - \mathbf{e}(x_{n+1}^k) \right]$$

where  $\mathbf{e}(x_{n+1}^k)$  is defined as follows: for  $x_{n+1}^k = x^k(i_k)$

$$\mathbf{e}(x_{n+1}^k) := \left[ \underbrace{0, \dots, 0}_{i_k}, 1, 0, \dots, 0 \right]^T \in R^{K_k}$$

Therefore, using the property of the projection operator (4)

$$\left\| \left[ \boldsymbol{\lambda}_n^k + \gamma_n^{k,\lambda} \boldsymbol{\psi}_n^k \right]_{-\lambda_{n+1}^+}^{\lambda_{n+1}^+} - \boldsymbol{\lambda}_{\delta_{n+1}}^{k,*} \right\| \\ \leq \left\| \boldsymbol{\lambda}_n^k + \gamma_n^\lambda \boldsymbol{\psi}_n^k - \boldsymbol{\lambda}_{\delta_{n+1}}^{k,*} \right\|$$

and taking into account that

$$\left\| \boldsymbol{\psi}_n^k \right\| \leq \text{Const} \gamma_n^\lambda (1 + \delta_n \lambda_n^+)$$

it follows

$$W_{n+1} \leq W_n + \text{Const} \cdot \mu_{1,n} \sqrt{W_n} + \text{Const} \cdot \mu_{2,n} + \varphi_n \\ \mu_{1,n} := |\delta_n - \delta_{n+1}| \\ \mu_{2,n} := \sum_{k=1}^2 (\gamma_n^{k,c})^2 + (\delta_n - \delta_{n+1})^2 + \\ + (1 + \delta_n \lambda_n^+)^2 (\gamma_n^\lambda)^2 \\ + |\delta_n - \delta_{n+1}| \left( \sum_{k=1}^2 \gamma_n^{k,c} + \gamma_n^\lambda (1 + \delta_n \lambda_n^+) \right) \\ \varphi_n := 2 \sum_{k=1}^2 \gamma_n^{k,c} \left( \mathbf{c}_n^k - \mathbf{c}_{\delta_n}^{k,*} \right)^T \left( \mathbf{e}(x_n^k \wedge u_n^k) - \mathbf{c}_n^k \right) \\ + \zeta_n^k \frac{\mathbf{e}^{N_k K_k} - N_k K_k \mathbf{e}(x_n^k \wedge u_n^k)}{N_k K_k - 1} \\ + 2 \gamma_n^\lambda \sum_{k=1}^2 \left( \boldsymbol{\lambda}_n^k - \boldsymbol{\lambda}_{\delta_n}^{k,*} \right)^T \boldsymbol{\psi}_n^k$$

Notice that

$$\left( \mathbf{c}_n^k - \mathbf{c}_{\delta_n}^{k,*} \right)^T \frac{\mathbf{e}^{N_k K_k}}{N_k K_k - 1} = 0$$

The assumptions of this theorem, and the strict convex property (5) lead to

$$E \{ \varphi_n \mid \mathbf{F}_{n-1} \} \stackrel{a.s.}{\leq} -\gamma_n^\lambda \delta_n W_n$$

The last inequality implies

$$E \{ W_{n+1} \mid \mathbf{F}_{n-1} \} \stackrel{a.s.}{\leq} (1 - \gamma_n^\lambda \delta_n) W_n \\ + \text{Const} (\mu_{1,n} \sqrt{W_n} + \mu_{2,n})$$

In view of

$$2\mu_{1,n} \sqrt{W_n} \leq \mu_{1,n}^2 \rho_n^{-1} + W_n \rho_n$$

(which is valid for any  $\rho_n > 0$ ) for  $\rho_n := \gamma_n^\lambda \delta_n$  one establishes

$$2\mu_{1,n} \sqrt{W_n} \leq (\mu_{1,n})^2 (\gamma_n^\lambda \delta_n)^{-1} + W_n \gamma_n^\lambda \delta_n$$

In view of the following estimation

$$(\mu_{1,n})^2 (\gamma_n^\lambda \delta_n)^{-1} + \mu_{2,n} \leq \text{Const} \cdot \mu_n$$

from the previous inequality we finally, obtain

$$E \{ W_{n+1} \mid \mathbf{F}_{n-1} \} \stackrel{a.s.}{\leq} (1 - \gamma_n^\lambda \delta_n) W_n + \text{Const} \cdot \mu_n \quad (6)$$

Observe that  $a_n^k = O(\varepsilon_n / \lambda_n^+)$ . It follows from (6) that  $\{W_n, \mathbf{F}_n\}$  is a nonnegative quasimartingale. Appealing to the assumptions of this theorem, and Robbins-Siegmund theorem [9]) the convergence with probability one is obtained.  $\blacksquare$

**Corollary 6** *If the design parameters of the adaptive algorithm belong to the following class of parameters*

$$\left( \varepsilon_0 \in \left[ 0, \left( \max_k N_k K_k \right)^{-1} \right), \varepsilon \geq 0 \right)$$

$$\delta_n := \frac{\delta_0}{n^\delta} \quad (\delta_0, \delta > 0)$$

$$\lambda_n^+ := \lambda_0^+ (1 + n^\lambda \ln n)^{-1} \quad (\lambda_0^+ > 0, \lambda \geq 0)$$

$$\gamma_n^{k,c} := \frac{\gamma_0^k}{n^\gamma} \quad (\gamma_0^k \in (0, 1), \gamma \geq 0)$$

with

$$\gamma + \varepsilon + \lambda + \delta \leq 1, \delta \leq \varepsilon$$

then

1. *the convergence with probability one is achieved if*

$$1 - \varepsilon - \lambda + \delta > \gamma > 1/2$$

2. *the mean squares convergence is ensured if*

$$\min \{ 2 - 2(\gamma + \varepsilon + \lambda), \gamma - \varepsilon - \lambda - \delta, 1 + 2(\delta - \lambda) \} > 0$$

**Theorem 7** *Assume that the conditions of the previous theorems are fulfilled, then for the class of parameters defined in the previous corollary, there exists  $\nu > 0$  such that*

$$0 < \nu < \min \left\{ 1 - \gamma - \varepsilon - \lambda + \delta; 2\gamma - 1; 2\delta \right\}$$

$$:= \nu^*(\gamma, \varepsilon, \delta, \lambda)$$

The proof of this theorem is essentially that of Theorem 2 of chapter 2 in [8].

The optimal convergence rate is given by the next corollary.

**Corollary 8** *The following optimal parameters  $\varepsilon^*, \delta^*, \lambda^*$  and  $\gamma^*$*

$$\varepsilon = \varepsilon^* = \delta = \delta^* = \frac{1}{6}, \lambda = \lambda^* = 0, \gamma = \gamma^* = \frac{2}{3}$$

*lead to the maximum convergence rate*

$$\nu^*(\gamma^*, \varepsilon^*, \delta^*, \lambda^*) = 2\gamma^* - 1 = \frac{1}{3}$$

## References

- [1] Gardner R. *Games for Business and Economics*. John Wiley & Sons, Inc. 1995.
- [2] Filar, J. and Vrieze, K. (1979). *Competitive Markov Decision Processes*. Springer-Verlag, Berlin.
- [3] Flesch, J., Thuijsman, F. and Vrieze, O. J. (1997). Markov strategies are better than stationary strategies. *Report M 97-09*, Maastricht University, Mathematics, The Netherlands.
- [4] Krausz, A. and Rieder, U. (1997). Markov Games with Incomplete Information. *Mathematical Methods of Operations Research*, 46, 263-279.
- [5] Mertens, J. F. and Neyman, A. (1981). Stochastic games. *International Journal of Game Theory*, 10, 53-66.
- [6] Mertens, J. F., Sorin, S. and Zamir, S. (1994). On stochastic games. Core Discussion Paper 9420, 9421, 9422, Université Catholique de Louvain.
- [7] Poznyak, A. S. and Najim, K. (1997). *Learning automata and stochastic optimization*. Springer-Verlag, Berlin.
- [8] Poznyak, A. S. and Najim, K. (1999). Adaptive control of constrained finite Markov chains, *Automatica*, 35, May.
- [9] Poznyak, A. S., Najim, K. and Gomez, E. (1999). *Self-learning control of finite Markov chains*. Marcel Dekker, New York.
- [10] Raghavan, T. E. S. and Filar, J. A. (1991). Algorithms for stochastic games - A survey. *ZOR - Methods and Models for Operations Research*, 35, 437-472.
- [11] Raghavan, T. E. S., Ferguson, T. S., Parthasarathy, T. and Vrieze, O. J. (1991). *Stochastic games and related topics*. Kluwer Academic Publishers, London.
- [12] Thuijsman, F. and Vrieze O. J. (1993). Stationary  $\varepsilon$ -optimal strategies in stochastic games. *OR Spectrum*, 15, 9-15.
- [13] Thuijsman, F. (1997). A survey on optimality and equilibria in stochastic games. *Ten years LNMB*, Klein et al. (eds.), CWI-Tract n<sup>o</sup>. 122, (Center for Mathematics and Computer Science, Amsterdam).
- [14] Vrieze O. J. and Thuijsman, F. (1987). Stochastic games and optimal stationary strategies. *Methods of Operations Research*, eds, Domschke, W., Krabs, W., Lehn, J. and Sppelluci, 513-529.