

Learning through Gradient-Type Reinforcement for N-person Repeated Constrained Games: Counter-Coalition Space Approach

Alexander S. Poznyak^a, Member, IEEE, Martín Godoy-Alcántar^b and Eduardo Gómez-Ramírez^c

^a CINEVESTAV-IPN, Departamento de Control Automático, A.P. 14-740, CP 07000 México D.F., México, e-mail:

apoznyak@ctrl.cinvestav.mx.

^b CINEVESTAV-IPN, Departamento de Control Automático, A.P. 14-740, CP 07000 México D.F., México, and Instituto Mexicano del

Petróleo, e-mail: mgodoy@ctrl.cinvestav.mx

^c Laboratorio de Investigación y Desarrollo de Tecnología Avanzada, Universidad La Salle,

Benjamín Franklin No. 47 Col. Condesa, CP 06140, México, D. F., México, e-mail: egomez@ci.ulsal.mx

Abstract— The paper tackles the design and analysis of a learning gradient-type strategy for N -person averaged constrained game with incomplete information. Each player is modelled by a stochastic variable-structure learning automaton (a simplest single state Markov Chain). Using the "joint payoff function", the considered game problem is formulated in terms of, so-called, counter-coalition variables. A special δ -regularization is introduced. Such approach does not require the "diagonal concavity" conditions to guarantee the uniqueness of the Nash equilibrium. The asymptotic convergence of the suggested learning procedure is analyzed.

Keywords— Reinforcement learning, Repeated game, Nash equilibrium

I. INTRODUCTION

Games theory dealing with data of stochastic nature (or, with imperfect information) has seen a tremendous growth in the last decades. Two models of players are more frequently considered in the literature: *dynamic and static models*. The first ones are given by a Markov chain characterized by sets (may be, finite or infinite) of states and actions (several actions are allowed to be applied at each state) with the corresponding transition (in general, probabilistic) rules. The second ones (*static models*) represent the simplest case of a finite Markov chain with a single state, that is, there are not any internal dynamic effects and the only actions selected by the participants are responsible for the further behavior of each player; last models are considered within the *Learning Automata paradigm* [8], [6], [3] and [4]. In this paper the simplest static player models will be considered.

A. Motivation

A few papers and rigorous results, dealing with learning (for static or dynamic players model) in constrained repeated stochastic games, are actually known. The fundamental paper of Rosen [9] concerns the uniqueness problem for Nash equilibrium but with no stochastic in the model. Several games with learning were considered by Lakshmi-varahan [2] and by Narendra and Thathachar [6], but there were no studies of constraint game situations. The spread-

ing of the Rosen's ideas to the class of strictly diagonal concave games with reinforcement automata learning was realized by Nazin and Poznyak in [7] (see, also [8]). The "*strict convexity condition*" (see also [9]) is also required that restricts significantly the areas of possible applications. Recently, the results dealing with learning to such sort of game have been published by Poznyak and Najim [5] where the Bush-Mosteller reinforcement technique was applied. But this paper also requires the restricting "*strict convexity condition*".

The rigorous description of the class of the repeated games having non unique equilibrium points as well as the analysis of the convergence and rate of learning for different reinforcement schemes still practically stay open that serves a lot for the motivation of this study.

B. Discussion of basic assumptions and restrictions

The most important assumptions (and, as the result, restrictions) concerning the considered game are as follows:

- this is the multi-person repeated game where the behavior of each player is the static model given by the simplest single-state Markov Chain, or, a *finite state automaton* with "time-variable structure" (at each stage of the game the distribution of the selected control actions can vary according to a reinforcement procedure);
- if some controls are realized by players at a current stage, then each participant immediately get the information on the realized values of his payoff and constraint functions which are individual for each player. The necessary information (realizations of both payoffs and constraints) is obtained during the course of the game. These realized payoffs and constraints (below called "*realizations*") are random variables (non obligatory binary) having constant first (conditional) moments which are assumed to be *a priori* unknown as well as uniformly bounded second conditional moment. In other words, the game is one of incomplete information.
- each player has *his own set of constraints* (constraint functions) which may be dependent on the selected control actions of other participants (the strategic interdepen-

dence);

- all participants accept some sort of agreement to use the same reinforcement learning procedure (in this case the stochastic pseudo-gradient scheme) and not to change it to another one during all long-range time of the game. Any other "cooperations" during the game are prohibited. Only the parameters of this procedure can be modified by player during the game, but not the *fixed structure* of the reinforcement procedure.

C. Main Contribution

Briefly, the basic contributions of this paper can be summarized as follows:

- using the "joint payoff function", the given game problem is formulated in terms of, so-called, "counter-coalition variables" that significantly simplifies the problem description making it linear with respect to these variables;
- for each player the suggested pseudo-gradient reinforcement procedure, based on Lagrange multipliers and an appropriate regularization, is shown to be oriented to obtain the *optimal response* (in average sense) within the given constraints;
- this learning tactics is proven to *lead to the unique Nash equilibrium point* corresponding to the "extended adjusted vector" with a minimal norm;
- any "Strictly Diagonal Conditions", commonly supposed in the finite game theory, do not required;
- based on the stochastic approximation technique (the Robbins- Sigmund theorem), the *convergence* (with probability one and in mean square sense) of the considered learning procedure to the Nash equilibrium is stated and the *rate of learning* is also estimated.

II. N-PERSON REPEATED CONSTRAINT GAME

Let (Ω, F, P) be a probability space, where Ω is the sample space, F a minimal σ -algebra of subsets of Ω , and P a probability measure on (Ω, F) . The symbol ω denotes the canonical point (event) in Ω . All subsequent random variables will be defined in this space.

A. Player model as Learning Automaton

Each player is modelled by a *stochastic variable-structure learning automaton* which consists of a simple Markov chain containing only one state (memoryless or static system) [3] and [4]. A stochastic automaton operating in a random environment (medium) is an adaptive discrete machine described by a collection $\{\Xi, U^k, \{\xi_n^k\}, \{u_n^k\}, \{p_n^k\}\}$ where Ξ is the automaton input bounded set, U^k denotes the set $\{u^k(1), u^k(2), \dots, u^k(N_k)\}$ ($k = \overline{1, N}$) of actions of the player (automaton), $\{\xi_n^k\}, \{\eta_n^{k, l_k}\}$ is a sequence of the k -player (automaton) inputs (payoffs and constraints $(\xi_n^k; \eta_n^{k, l_k}) \in \Xi$, $l_k = \overline{1, M_k}$) provided by the game mechanism in a binary (P -model environment) or continuous (S -model environment) form, $\{u_n^k\}$ is a sequence of the k -player (automaton) outputs (actions) and $p_n^k = [p_n^k(1), p_n^k(2), \dots, p_n^k(N_k)]^T$ is the conditional

probability distribution at time n , that is, $p_n^k(i) = P\{\omega \in \Omega : u_n^k = u^k(i) / F_{n-1}\}$, $\sum_{i=1}^{N_k} p_n^k(i) = 1$. Here $F_n = \sigma(u_1^k, p_1^k, \xi_1^k, \eta_1^{k, l_k}; \dots; u_n^k, p_n^k, \xi_n^k, \eta_n^{k, l_k})$ is the minimal σ -algebra generated by the corresponding events ($F_n \subseteq F$).

B. How game is played

This game is played in the following way. According to the probability distributions (mixed strategy) p_n^k , at each stage n , simultaneously and independently (we consider a non-cooperative game in which each player has only its own payoff realization $\xi_n^k = \xi_n^k(\omega)$ and the realizations $\eta_n^{k, l_k} = \eta_n^{k, l_k}(\omega)$ of the constraints), each player chooses randomly an action, for instance $u^k(i_k)$ ($i_k = \overline{1, N_k}, k = \overline{1, N}$):

$$\begin{cases} \xi_n^k = \xi^k(\omega; u_n^k = u^k(i_k)), & k = \overline{1, N} \\ \eta_n^k = (\eta_n^{k, 1}, \dots, \eta_n^{k, M_k})^T, & \eta_n^{k, l_k} = \eta^{k, l_k}(\omega; u_n^k = u^k(i_k)) \end{cases}$$

The k^{th} random component ξ_n^k represents the payoff (Borel function) earned by the k^{th} player, and η_n^{k, l_k} is the random realization of the l_k -th constraint for this player. To obtain reasonable generality without excessive complexity, the following assumptions will be in force throughout this paper: (A1) The conditional expectations of ξ_n^k and η_n^{k, l_k} are independent of the history of the game (the game rules are assumed to be unchangeable):

$$\begin{aligned} E\{\xi_n^k | u_n^k = u^k(i_k) \wedge F_{n-1}, k = \overline{1, N}\} &\stackrel{a.s.}{=} v_{i_1, \dots, i_N}^k \\ E\{\eta_n^{k, l_k} | u_n^k = u^k(i_k) \wedge F_{n-1}, k = \overline{1, N}\} &\stackrel{a.s.}{=} w_{i_1, \dots, i_N}^{k, l_k} \end{aligned} \quad (1)$$

(A2) For any realized actions u_n^k at current stage n of the play, the second absolute moments of ξ_n^k and η_n^{k, l_k} are uniformly (on ω and n) bounded with probability one, that is,

$$\begin{aligned} E\left\{\left|\xi^k(\omega; u_n^k = u^k(i_k), k = \overline{1, N})\right|^2 | F_{n-1}, \right. \\ \left. k = \overline{1, N}\right\} &\stackrel{a.s.}{\leq} \sigma_{\xi, k}^+ < \infty \\ E\left\{\left|\eta^{k, l_k}(\omega; u_n^k = u^k(i_k), k = \overline{1, N})\right|^2 | F_{n-1}, \right. \\ \left. k = \overline{1, N}\right\} &\stackrel{a.s.}{\leq} \sigma_{\eta, k, l_k}^+ < \infty \end{aligned} \quad (2)$$

where $|\cdot|$ stands for the absolute value. For each player the average payoff and constraints form the collection of $M_k + 1$ tensors

$$\begin{cases} V^k = [v_{i_1, \dots, i_N}^k] & (i_k = \overline{1, N_k}, k = \overline{1, N}) \\ W^{k, l_k} = [w_{i_1, \dots, i_N}^{k, l_k}] & (l_k = \overline{1, M_k}) \end{cases}$$

which are assumed to be *a priori* unknown. Then each player changes his mixed strategy according to the accepted reinforcement $T_n^k : \left(p_n^k \xrightarrow{T_n^k} p_{n+1}^k\right)$ trying for large n to maximize his individual averaged payoff $\liminf_{n \rightarrow \infty} \Phi_n^k$ where $\Phi_n^k := n^{-1} \sum_{t=1}^n \xi_t^k$ maintaining (with probability one) the constraints $\liminf_{n \rightarrow \infty} \Psi_n^{k, l_k} \leq b^{k, l_k}$, $\Psi_n^{k, l_k} := n^{-1} \sum_{t=1}^n \eta_t^{k, l_k}$.

Definition 1: At stage n , for the k^{th} player of the considered game, and for any F_{n-1} -measurable conditional probability distribution (p_n^1, \dots, p_n^N) , the Borel function

$$V^k(p_n^1, \dots, p_n^N) := \sum_{i_1=1}^{N_1} \cdot \sum_{i_N=1}^{N_N} v_{i_1, \dots, i_N}^k \prod_{s=1}^N p_n^s(i_s)$$

are said to be the **expected payoff**, and

$$W^{k, l_k}(p_n^1, \dots, p_n^N) := \sum_{i_1=1}^{N_1} \cdot \sum_{i_N=1}^{N_N} w_{i_1, \dots, i_N}^{k, l_k} \prod_{s=1}^N p_n^s(i_s)$$

corresponds to the **expected constraints**.

Define also the *averaged expected payoff* and *constraints* as follows:

$$V_n^k := n^{-1} \sum_{t=1}^n V^k(p_t^1, \dots, p_t^N)$$

$$W_n^{k, l_k} := n^{-1} \sum_{t=1}^n W^{k, l_k}(p_t^1, \dots, p_t^N)$$

The following lemma (see [5]) states the asymptotic equivalence between $(\Phi_n^k, \Psi_n^{l_k})$ and (V_n^k, W_n^{k, l_k}) .

Lemma 1: Under the condition **A1** for any conditional distributions $\{p_n^1, \dots, p_n^N\}$ it follows

$$\Phi_n^k \stackrel{a.s.}{=} V_n^k + o_\omega\left(n^{-\frac{1}{2}}\right), \Psi_n^{l_k} \stackrel{a.s.}{=} W_n^{k, l_k} + o_\omega\left(n^{-\frac{1}{2}}\right)$$

III. MIXED STRATEGIES AND NASH EQUILIBRIUM

A. Randomized strategies

The *randomized (or mixed) strategy* of the k^{th} player is any sequence D^k of vectors $\{p_n^k\}$ ($k = \overline{1, N}$) with F_{n-1} -measurable components belonging to the simplex $S_{\varepsilon=0}^{N_k}$, that is,

$$p_n^k \in S_{\varepsilon^k=0}^{N_k} := \left\{ p_n^k \in R^{N_k} : p_n^k(i) \geq \varepsilon^k \geq 0, \sum_{i=1}^{N_k} p_n^k(i) = 1 \right\}$$

Definition 2: The strategies $\overline{D}^1, \dots, \overline{D}^N$ are said to be the **non-cooperative equilibrium strategies** (in the Nash sense) if

1) they are admissible, that is, $\min_{p_1^s, s=1, N} \liminf_{n \rightarrow \infty} W_n^{k, l_k} := \widetilde{W}(\overline{D}^1, \dots, \overline{D}^N) \leq b^{k, l_k}$ for any $k = \overline{1, N}$, $l_k = \overline{1, M_k}$, where $b^{k, l_k} \in R$ are *a priori* given;

2) for any integer k and any admissible strategy D^k

$$\begin{aligned} \widetilde{V}^k(\overline{D}^1, \dots, \overline{D}^N) &:= \min_{p_1^s, s=1, N} \liminf_{n \rightarrow \infty} V_n^k \\ &\stackrel{a.s.}{\geq} \widetilde{V}(\overline{D}^1, \dots, \overline{D}^{k-1}, D^k, \overline{D}^{k+1}, \dots, \overline{D}^N) \end{aligned} \quad (3)$$

where the minimization is done over all the initial probability distributions $p_1^s, s = \overline{1, N}$.

B. Equilibrium within the subclass of stationary mixed strategies

Consider now the subclass of stationary mixed strategies $D^k = \{p^k\}$ ($k = \overline{1, N}$).

Definition 3: The point $(\overline{p}^1, \dots, \overline{p}^N)$ is said to be an **equilibrium point** [10] of the given N -person game within the class of stationary strategies if

$$V^k(\overline{p}^1, \dots, \overline{p}^N) = \max_{p^k \in R^k} V^k(\overline{p}^1, \dots, p^k, \dots, \overline{p}^N)$$

$$R^k := S_0^{N_k} \bigcap_{l=1}^{M_k} \{p^k : W^{k, l_k}(\overline{p}^1, \dots, \overline{p}^N) \leq b^{k, l_k}\}$$

for each $k = \overline{1, N}$

At this point no player can increase his payoff by a unilateral change in his strategy.

Remark 1: The set of all equilibrium strategies \overline{D}^k ($k = \overline{1, N}$) contains (see [5]) the subset of **admissible stationary strategies** $\{\overline{p}^k\}$ realizing the inequalities $V^k(\overline{p}^1, \dots, \overline{p}^N) \geq V^k(\overline{p}^1, \dots, \overline{p}^{k-1}, p^k, \overline{p}^{k+1}, \dots, \overline{p}^N)$, $W^{k, l_k}(\overline{p}^1, \dots, \overline{p}^N) \leq b^{k, l_k}$ ($k = \overline{1, N}, l_k = \overline{1, M_k}$) for any $\overline{p}^k \in S_{\varepsilon^k=0}^{N_k}$.

IV. GAME AS REINFORCEMENT: PROBLEM SETTING

Now we are ready to formulate the N -person repeated constrained game problem with *a priori* unknown average payoffs and constraints: *based on current information, generate randomized (mixed) admissible strategies* $\{p_n^k\}$ ($k = \overline{1, N}$) *in order to achieve a Nash equilibrium realizable within the subclass of stationary strategies*.

To achieve this objective, first emphasize the following fact. According to Nash theorem [10], the set of admissible stationary distributions $(\overline{p}^1, \dots, \overline{p}^N)$, may contain more than one element. As it is shown in [5] that if an admissible mixed strategy D^k ($k = \overline{1, N}$) converges to such probability distribution $(\overline{p}^1, \dots, \overline{p}^N)$, that is, $\limsup_{n \rightarrow \infty} n^\tau E \left\{ \sum_{k=1}^N \|p_n^k - \overline{p}^k\|^2 \right\} < \infty$ for some $\tau > 0$, or,

in another words, the strategy D^k ($k = \overline{1, N}$) is asymptotically stationary realizing a Nash equilibrium, then the associated random functions Φ_n^k and Ψ_n^{k, l_k} also converge (with the same rate) to the corresponding average values $V^k(\overline{p}^1, \dots, \overline{p}^N)$ and $W^{k, l_k}(\overline{p}^1, \dots, \overline{p}^N)$ respectively, that is, $\limsup_{n \rightarrow \infty} n^\tau E \left\{ \sum_{k=1}^N |\Phi_n^k - V^k(\overline{p}^1, \dots, \overline{p}^N)|^2 \right\} < \infty$

and $\limsup_{n \rightarrow \infty} n^\tau E \left\{ \sum_{l_k=1}^{M_k} |\Psi_n^{k, l_k} - W^{k, l_k}(\overline{p}^1, \dots, \overline{p}^N)|^2 \right\} < \infty$.

So, if we construct an asymptotically stationary mixed strategy $\{p_n^k\}$ converging to a stationary distribution $(\overline{p}^1, \dots, \overline{p}^N)$ realizing a Nash equilibrium, we will be able to attain the main aim of the game. But to do it in a rigorous manner, first, the questions related to the *existence* and the *uniqueness* of the Nash equilibrium (within stationary strategies) should be resolved. Talking on the uniqueness

of the finite matrix games, one can consider the uniqueness of the optimal policy for each player, and the uniqueness of the Nash equilibrium point.

As it is shown by Rosen [9], the non-uniqueness of the Nash equilibrium points cannot be offset by a small regularization (perturbation) term. The condition for the uniqueness of equilibria are known as "strict diagonal concavity". As it was shown in [8] and [5], if the given (nonregularized) matrix game is "diagonal concave", then the corresponding regularized game turns out to be strictly diagonal concave that, by Rosen's theorem (the theorem 2 in [9]), implies the uniqueness of the equilibria policy. But in this paper we suggest another approach which does not requires the "diagonal concavity condition", but demands the representation of the initial problem in new (counter-coalition) variables with some additional regularization.

V. EQUIVALENT COUNTER-COALITION VARIABLES REPRESENTATION

A. Joint Payoff Function and Nash Equilibrium

Following the approach presented in [9], let us introduce the joint payoff function $\rho(p, q)$ defined as

$$\rho(p, q) := \sum_{k=1}^N V^k(p^1, \dots, p^{k-1}, q^k, p^{k+1}, \dots, p^N) = \sum_{k=1}^N V^k(p^{\hat{k}}, q^k) \quad (4)$$

for any $(p, q) \in (S \times Q) \times (S \times Q)$ where:

$$V^k(p^{\hat{k}}, q^k) = \sum_{j_1=1}^{N_1} \sum_{j_2=1}^{N_2} \dots \sum_{j_N=1}^{N_N} V_{j_1, \dots, j_N}^k p_{(\cdot)}^{\hat{k}} q_{j_k}^k$$

with

$$p_{(\cdot)}^{\hat{k}} = p_{j_1, \dots, j_{k-1}, j_{k+1}, \dots, j_N}^{\hat{k}} = p_{j_{\hat{k}}}^{\hat{k}} := \prod_{s=1, s \neq k}^N p_{j_s}^s \in S^{N_{\hat{k}}}$$

where $N_{\hat{k}} := \prod_{s=1, s \neq k}^N N_s$ being referred to as the *extended counter coalition vector*.

Example 1: For a three player game with two action for each one we have for player one

$$p^{\hat{1}} := (p_1^2 p_1^3, p_1^2 p_2^3, p_2^2 p_1^3, p_2^2 p_2^3)^T$$

Using Kakutani's fixed point theorem, Rosen proved the following important result [9], that any N -person matrix game with constraints has an equilibrium point $p^* \in R^{\bar{N}}$ ($R^{\bar{N}} := R^{N_1} \times R^{N_2} \times \dots \times R^{N_N}$), in Nash sense (see Definition 3), satisfying

$$\rho(p^*, p^*) = \max_{q \in S \times Q} \rho(p^*, q) \quad (5)$$

where $Q \subset R^{\bar{N}}$ defined as

$$\begin{aligned} Q &:= Q^1 \times \dots \times Q^N; \quad Q^k := Q^{k,1} \times \dots \times Q^{k,M_k} \\ Q^{k,l_k} &:= \left\{ p \in R^{\bar{N}} \mid k = \overline{1, N}, \quad l_k = \overline{1, M_k}, \quad W^{k,l_k}(p) \right. \\ &\quad \left. := \sum_{j_1=1}^{N_1} \dots \sum_{j_N=1}^{N_N} W_{j_1, \dots, j_N}^{k,l_k} \prod_{i=1}^N p_{j_i}^i \leq b^{k,l_k} \right\} \end{aligned} \quad (6)$$

The constraints (6) can be written as

$$p^{\hat{k}, \top} \overline{W}^{k,l_k} p^k - b^{k,l_k} \leq 0, \quad (k = \overline{1, N}, \quad l_k = \overline{1, M_k})$$

where the matrix \overline{W}^{k,l_k} is given by $\overline{W}^{k,l_k} = [w_{j_{\hat{k}}, j_k}^{k,l_k}]$. The index j_k of the actions of player k grows in the columns and the combination of indexes $j_{\hat{k}}$ of the counter coalition grows in the rows of \overline{W}^{k,l_k} , or, in an equivalent form,

$$\widetilde{W}^{k,l_k}(p^k, p^{\hat{k}}) = p^{k, \top} \left((\overline{W}^{k,l_k})^T - b^{k,l_k} \mathbf{1}^k \right) p^{\hat{k}}(p^k) \leq 0$$

where $\mathbf{1}^k$ is a matrix of ones of size N_k by $N_{\hat{k}}$ where $N_{\hat{k}} := \prod_{i=1, i \neq k}^N N_i$.

Assumption A3: We will assume that the Slater's condition hold, that is, there exist pairs $(\hat{p}^k, \hat{p}^{\hat{k}}) \in S^{N_k} \times S^{N_{\hat{k}}}$ such that

$$\widetilde{W}^{k,l_k}(\hat{p}^k, \hat{p}^{\hat{k}}) < 0$$

for all $k = \overline{1, N}$.

B. Criterion of Nash-equilibrium using counter-coalition variables

Define the matrices \overline{V}^k given by $\overline{V}^k = [V_{j_{\hat{k}}, j_k}^k]$ where the index j_k of the actions of player k grows in the columns and the combination of indexes $j_{\hat{k}}$ of the counter coalition grows in the rows of \overline{V}^k and

$$\widetilde{W}^k(\beta^k) = \sum_{l_k=1}^{M_k} \beta_{l_k}^k \widetilde{W}^{k,l_k}; \quad \widetilde{W}^{k,l_k} = \overline{W}^{k,l_k} - b^{k,l_k} \mathbf{1}^{k, \top}$$

The following theorem represents the criterion of a Nash-equilibrium (see [1]).

Theorem 1: A necessary and sufficient condition that \hat{p}^* be a Nash equilibrium point is that the collection $(\hat{p}^*, \alpha^*, \beta^*, t^*)$ is the solution of the following bilinear programming problem

$$\begin{aligned} \overline{F}(t, \alpha) &= \sum_{k=1}^N t^k - \sum_{k=1}^N \alpha^k \rightarrow \max_{\substack{\hat{p}^k \in R^{N_{\hat{k}}}, \alpha^k, t^k \in R, \beta^k \in R^{M_k}}} \\ &\text{subject to} \end{aligned}$$

$$\begin{aligned} &(\overline{V}^{k, \top} - \widetilde{W}^{k, \top}(\beta^k)) p^{\hat{k}} - \alpha^k e^k \leq 0 \\ &-(\overline{V}^{k, \top} - \widetilde{W}^{k, \top}(\beta^k)) p^{\hat{k}} + t^k e^k \leq 0 \\ &\overline{M}^{k, \hat{s}} p^{\hat{s}} = \overline{M}^{k, \hat{r}} p^{\hat{r}}; \quad p^{\hat{k}} \in S^{N_{\hat{k}}} \\ &k = \overline{1, N}, \quad l_k = \overline{1, M_k}, \quad s = \overline{1, N}, \quad s \neq k, \quad s \neq r \end{aligned} \quad (7)$$

The following property holds: $\overline{F}(t^*, \alpha^*) = 0$.

VI. LEARNING ALGORITHM

A. Regularization, Lagrange multipliers and projection gradient procedure: the complete information case

Using the Lagrange approach define the δ -regularized Lagrangian as:

$$\begin{aligned} L_\delta(\hat{p}, \alpha, t, \beta, \lambda_\alpha, \lambda_t, \lambda_{\hat{p}}) &:= \sum_{k=1}^N t^k - \sum_{k=1}^N \alpha^k \\ &- \sum_{k=1}^N \lambda_\alpha^{k,T} \left[\left(\bar{V}^{k,\top} - \bar{W}^{k,\top} \left(\beta^k \right) \right) p^{\hat{k}} - \alpha^k e^k \right] \\ &- \sum_{k=1}^N \lambda_t^{k,T} \left[- \left(\bar{V}^{k,\top} - \bar{W}^{k,\top} \left(\beta^k \right) \right) p^{\hat{k}} + t^k e^k \right] \\ &- \sum_{k=1}^N \sum_{r=1, r \neq k}^N \sum_{s=1, s \neq r, k}^N \lambda_{\hat{p}, \hat{r}, \hat{s}}^{k,T} \left[M^{k, \hat{s}} p^{\hat{s}} - M^{k, \hat{r}} p^{\hat{r}} \right] \quad (8) \\ &+ \frac{\delta}{2} \sum_{k=1}^N \left(- \|p^{\hat{k}}\|^2 - (\alpha^k)^2 - (t^k)^2 - \|\beta^k\|^2 + \right. \\ &\quad \left. \|\lambda_\alpha^k\|^2 + \|\lambda_t^k\|^2 + \|\lambda_{\hat{p}}^k\|^2 \right) \\ &\rightarrow \max_{p^{\hat{k}} \in S_{\varepsilon_n^k}^{N_{\hat{k}}}, \alpha^k, t^k \in R, \beta^k \in R^{M_k}} \min_{\lambda_\alpha^k, \lambda_t^k, \lambda_{\hat{p}, \hat{r}, \hat{s}}^k \in R^{N_k}} \end{aligned}$$

In the case when the complete information on the expected payoffs and constraints is available, then the *gradient-like* technique may be applied to attain the equilibrium point:

$$\begin{aligned} p_{n+1}^{\hat{k}} &= \pi_{S_{\varepsilon_n^k}^{N_{\hat{k}}}} \left\{ p_n^{\hat{k}} + \gamma_n^{\hat{p}} \nabla_{p^{\hat{k}}} L_{\delta_n}^k(\cdot) \right\} \\ \alpha_{n+1}^k &= \alpha_n^k + \gamma_n^\alpha \nabla_{\alpha^k} L_{\delta_n}^k(\cdot) \\ t_{n+1}^k &= t_n^k + \gamma_n^t \nabla_{t^k} L_{\delta_n}^k(\cdot) \\ \beta_{n+1}^k &= \left[\beta_n^k + \gamma_n^\beta \nabla_{\beta^k} L_{\delta_n}^k(\cdot) \right]_+ \\ \lambda_{\alpha, n+1}^k &= \left[\lambda_{\alpha, n}^k - \gamma_n^{\lambda_\alpha} \nabla_{\lambda_\alpha^k} L_{\delta_n}^k(\cdot) \right]_+ \\ \lambda_{t, n+1}^k &= \left[\lambda_{t, n}^k - \gamma_n^{\lambda_t} \nabla_{\lambda_t^k} L_{\delta_n}^k(\cdot) \right]_+ \\ \lambda_{\hat{p}, \hat{r}, \hat{s}, n+1}^k &= \left[\lambda_{\hat{p}, \hat{r}, \hat{s}, n}^k - \gamma_n^{\lambda_{\hat{p}}} \nabla_{\lambda_{\hat{p}}^k} L_{\delta_n}^k(\cdot) \right]_+ \\ k &= 1, \bar{N}, n = 1, 2, \dots \end{aligned} \quad (9)$$

where

$$L_{\delta_n}^k(\cdot) = L_{\delta_n}^k(\hat{p}_n, \alpha_n, t_n, \beta_n, \lambda_{\alpha, n}, \lambda_{t, n}, \lambda_{\hat{p}, n})$$

and $\pi_{S_{\varepsilon_n^k}^{N_{\hat{k}}}}\{\cdot\}$ is the projection operator to the simplex set $S_{\varepsilon_n^k}^{N_{\hat{k}}}$ and $[\cdot]_+$ is the "positive part" operator. If the parameters of this procedure satisfy

$$0 < \gamma_n^{\hat{p}}, \delta_n, \varepsilon_n^k \rightarrow 0$$

$$\sum_{n=1}^{\infty} \gamma_n^{\hat{p}} \delta_n = \infty, \sum_{n=1}^{\infty} \kappa_n < \infty$$

(κ_n is defined by (12)), it provides the convergence of the estimates $p_n^{\hat{k}}$ to the corresponding unique equilibrium strategy $\bar{p}^{\hat{1}}(\vartheta), \dots, \bar{p}^{\hat{N}}(\vartheta)$ having the minimal norm over all possible equilibrium points.

B. Stochastic approximation approach and the "measure keeping problem": incomplete information case

When we deal with incomplete information case, that is, the only current realizations of payoff ξ_n^k and constraints η_n^{k, l_k} are available, the "stochastic approximation" version of the procedure (9) can be applied (see, for example, [4]) where instead of $\nabla_{p^{\hat{k}}} L_{\delta_n}^k(\cdot)$, $\nabla_{\alpha_n^k} L_{\delta_n}^k(\cdot)$, $\nabla_{t_n^k} L_{\delta_n}^k(\cdot)$, $\nabla_{\beta_n^k} L_{\delta_n}^k(\cdot)$, $\nabla_{\lambda_\alpha^k} L_{\delta_n}^k(\cdot)$, $\nabla_{\lambda_t^k} L_{\delta_n}^k(\cdot)$, and $\nabla_{\lambda_{\hat{p}}^k} L_{\delta_n}^k(\cdot)$ their estimates using the current realizations are implemented. Such procedure is known as the "reinforcement learning algorithm" and can be realized by different ways (with different estimates in the use). Below, we present the concrete learning procedure which is based on Learning Automata concept [3], [4].

C. Pseudo-gradient reinforcement with Lagrange multipliers adjustment

The pseudo-gradient reinforcement scheme presented in will be applied hereafter to design a new learning algorithm for N-person constrained repeated games in terms of the complemented variables with unknown expected payoff and constraints. In fact, we assume that after each stage, the payoff to each player as well as the constraints are random variables. No information concerning the distribution of the payoff and constraints is *a priori* available. The necessary information is obtained during the course of the game.

Learning control is an iteration process involving an adaptation at each stage (time step). We now present the "four-step" recursive algorithm.

Step 1. On the basis of the available data

$$u_n^k = u^k(i_k), \xi_n^k, \eta_n^{k, l_k}, p_n^{\hat{k}}(p_n^k(i_k) \geq \varepsilon_n^k > 0) \\ \lambda_{t, n}^k, \lambda_{\alpha, n}^k, \lambda_{\hat{p}, \hat{r}, \hat{s}, n}^k$$

built the following functions for $k=1, \bar{N}$, $n=1, 2, \dots$

$$\begin{aligned} R\left(\nabla_{p^{\hat{k}}} L_{\delta_n}^k(\cdot)\right) &:= \frac{\left(\xi_n^k - \sum_{l_k=1}^{M_k} \beta_{l_k}^k \eta_n^{k, l_k}\right)}{\prod_{s=1, s \neq k}^N p^{\hat{s}}(u_n^{\hat{s}})} I^{\hat{k}, k} \left(u_n^{\hat{k}}, u_n^k\right) \times \\ &\left(\lambda_t^k - \lambda_\alpha^k\right) + \sum_{i=1, i \neq k}^N \sum_{r=1, r \neq i}^N \sum_{s=1, s \neq r, i}^N (M^{i, \hat{r}})^T \lambda_{\hat{p}, \hat{r}, \hat{s}}^i - \delta_n p_n^{\hat{k}} \\ R\left(\nabla_{\alpha^k} L_{\delta_n}^k(\cdot)\right) &:= -1 + \lambda_{\alpha, n}^{k, T} e^k - \delta_n \alpha_n^k \\ R\left(\nabla_{\beta^k} L_{\delta_n}^k(\cdot)\right)_{l_k} &:= \left(\lambda_t^k - \lambda_\alpha^k\right)^T \\ &\times \frac{\eta_n^{k, l_k}}{\prod_{s=1, s \neq k}^N p^{\hat{s}}(u_n^{\hat{s}})} \left(I^{\hat{k}, k} \left(u_n^{\hat{k}}, u_n^k\right)\right)^T p_n^{\hat{k}} - \delta_n \beta_n^k \end{aligned}$$

$$\begin{aligned}
R\left(\nabla_{\lambda_{\alpha}^k} L_{\delta_n}^k(\cdot)\right) &:= \frac{\left(\xi_n^k - \sum_{l_k=1}^{M_k} \beta_{l_k}^k \eta_n^{k,l_k}\right)}{\prod_{s=1, s \neq k}^N p^{\hat{s}}(u_n^{\hat{s}})} \times \\
&\quad \left(I^{\hat{k},k}\left(u_n^{\hat{k}}, u_n^k\right)\right)^T p_n^{\hat{k}} + \alpha_n^k e^k + \delta_n \lambda_{\alpha,n}^{k,T} \\
R\left(\nabla_{\lambda_t^k} L_{\delta_n}^k(\cdot)\right) &:= 1 - \lambda_{t,n}^{k,T} e^k - \delta_n t_n^k \\
&\quad \frac{\left(\xi_n^k - \sum_{l_k=1}^{M_k} \beta_{l_k}^k \eta_n^{k,l_k}\right)}{\prod_{s=1, s \neq k}^N p^{\hat{s}}(u_n^{\hat{s}})} \times \\
&\quad \left(I^{\hat{k},k}\left(u_n^{\hat{k}}, u_n^k\right)\right)^T p_n^{\hat{k}} - t_n^k e^k + \delta_n \lambda_{\alpha,n}^{k,T} \\
R\left(\nabla_{\lambda_{\hat{p}, \hat{r}, \hat{s}}^k} L_{\delta_n}^k(\cdot)\right) &:= M^{k,\hat{s}} p^{\hat{s}} - M^{k,\hat{r}} p^{\hat{r}} + \delta_n \lambda_{\hat{p}, \hat{r}, \hat{s}}^{k,T}
\end{aligned} \tag{10}$$

Step 2. Update the probability distributions p_{n+1}^k and the Lagrange multipliers λ_{n+1} using the following iterative schemes:

$$\begin{aligned}
p_{n+1}^{\hat{k}} &= \pi_{S_{\varepsilon_n^k}} \left\{ p_n^{\hat{k}} + \gamma_n^{\hat{p}} R\left(\nabla_{p^{\hat{k}}} L_{\delta_n}^k(\cdot)\right) \right\} \\
\alpha_{n+1}^k &= \alpha_n^k + \gamma_n^{\alpha} R\left(\nabla_{\alpha^k} L_{\delta_n}^k(\cdot)\right) \\
t_{n+1}^k &= t_n^k + \gamma_n^t R\left(\nabla_{t^k} L_{\delta_n}^k(\cdot)\right) \\
\beta_{n+1}^k &= \beta_n^k + \gamma_n^{\beta} R\left(\nabla_{\beta^k} L_{\delta_n}^k(\cdot)\right) \\
\lambda_{\alpha}^k(n+1) &= \left[\lambda_{\alpha}^k(n) - \gamma_n^{\alpha} R\left(\nabla_{\lambda_{\alpha}^k} L_{\delta_n}^k(\cdot)\right) \right]_0^{\lambda_{\alpha,n+1}^+} \\
\lambda_t^k(n+1) &= \left[\lambda_t^k(n) - \gamma_n^t R\left(\nabla_{\lambda_t^k} L_{\delta_n}^k(\cdot)\right) \right]_0^{\lambda_{t,n+1}^+} \\
\lambda_{\hat{p}, \hat{r}, \hat{s}}^k(n+1) &= \left[\lambda_{\hat{p}, \hat{r}, \hat{s}}^k(n) - \gamma_n^{\lambda_{\hat{p}, \hat{r}, \hat{s}}} R\left(\nabla_{\lambda_{\hat{p}, \hat{r}, \hat{s}}^k} L_{\delta_n}^k(\cdot)\right) \right]_0^{\lambda_{\hat{p}, \hat{r}, \hat{s}, n+1}^+} \\
k &= \overline{1, N}, n = 1, 2, \dots
\end{aligned} \tag{11}$$

where if $u_n^k = u^k(j_k)$ and $u_n^{\hat{k}} = u^{\hat{k}}(j_{\hat{k}})$ then $I^{\hat{k},k}\left(u_n^{\hat{k}}, u_n^k\right) = \left[\delta_{(\hat{i}_k, i_k)}\right] \in R^{N_{\hat{k}} \times N_k}$ with $\delta_{(\hat{i}_k, i_k)} = 1$ if $\hat{i}_k = j_{\hat{k}}$ and $i_k = j_k$ and $\delta_{(\hat{i}_k, i_k)} = 0$ otherwise. The operator $[y]_0^{\lambda_{n+1}^+}$ is defined as follows:

$$[y]_0^{\lambda_{n+1}^+} = \begin{cases} y & \text{if } y \in [0, \lambda_{n+1}^+] \\ \lambda_{n+1}^+ & \text{if } y > \lambda_{n+1}^+ \\ 0 & \text{if } y < 0 \end{cases}$$

Step 3. According to $\Pr\{u_{n+1}^k = u^k(i) \mid F_n\} = p_{n+1}^k(i)$ generate randomly (for each player) new discrete random variables u_{n+1}^k as in learning stochastic automata implementation, and get a new observations (realizations) ξ_{n+1}^k and η_{n+1}^{k,l_k} that corresponds to the environment vector-reactions.

Step 4. Return to Step 1.

VII. CONVERGENCE ANALYSIS

A. Main theorem on the convergence with probability one

Theorem 2 (convergence with probab one) Suppose that assumptions **A1**- **A2** hold for the learning reinforcement procedure (10)-(11) and the given game is diagonal convex. In addition, assume that:

1. there exist the following nonnegative sequences $\{\gamma_n^{\hat{p}}\}$, $\{\gamma_n^{\alpha}\}$, $\{\gamma_n^t\}$, $\{\gamma_n^{\beta}\}$, $\{\gamma_n^{\lambda_{\alpha}}\}$, $\{\gamma_n^{\lambda_t}\}$, $\{\gamma_n^{\lambda_{\hat{p}}}\}$, $\{\delta_n\}$ and $\{\varepsilon_n^k\}$ such that

$$\begin{aligned}
&\gamma_n^{\hat{p}}, \gamma_n^{\alpha}, \gamma_n^t, \gamma_n^{\beta}, \gamma_n^{\lambda_{\alpha}}, \gamma_n^{\lambda_t}, \gamma_n^{\lambda_{\hat{p}}} \downarrow 0, \delta_n \in (0, \delta^+), \delta_n \downarrow 0 \\
&\gamma_n^{\hat{p}} = c_{\alpha} \gamma_n^{\alpha} = c_t \gamma_n^t = c_{\beta} \gamma_n^{\beta} = c_{\lambda_{\alpha}} \gamma_n^{\lambda_{\alpha}} = c_{\lambda_t} \gamma_n^{\lambda_t} = c_{\lambda_{\hat{p}}} \gamma_n^{\lambda_{\hat{p}}} \\
&\varepsilon_n^k \in \left(0, \frac{1}{N_{\hat{k}}}\right), \varepsilon_n^k \downarrow 0, \limsup_n \frac{\varepsilon_n^k}{\delta_n} < \infty
\end{aligned}$$

2. the updating factors $\gamma_n^{\hat{p}}, \gamma_n^{\alpha}, \gamma_n^t, \gamma_n^{\beta}, \gamma_n^{\lambda_{\alpha}}, \gamma_n^{\lambda_t}, \gamma_n^{\lambda_{\hat{p}}}$ satisfy $\sum_{n=1}^{\infty} \gamma_n^{\hat{p}} \delta_n = \infty$ and $\sum_{n=1}^{\infty} \left[\phi_n + \kappa_n^2 (\gamma_n^{\hat{p}} \delta_n)^{-1}\right] < \infty$ where

$$\begin{aligned}
\kappa_n &:= C_1 \left\| \varepsilon_{n+1}^k - \varepsilon_n^k \right\| + C_2 |\delta_{n+1} - \delta_n| \\
&\quad + C_3 \left\| \varepsilon_{n+1}^k \delta_{n+1}^{-1} - \varepsilon_n^k \delta_n^{-1} \right\|
\end{aligned} \tag{12}$$

$$\begin{aligned}
\phi_n &:= \kappa_n^2 + C_A^2 \sum_{k=1}^N (\gamma_n^{\hat{p}})^2 + 2C_A \kappa_n \sum_{k=1}^N \gamma_n^{\hat{p}} \\
&\quad + 2M (\gamma_n^{\hat{p}})^2 \left[(\delta_n \lambda_n^+) + (\sigma_n^+) \right] \\
&\quad + 2\sqrt{M} \gamma_n^{\hat{p}} \kappa_n (\delta_n \lambda_n^+ + \sigma_n^+)
\end{aligned}$$

Then the mixed strategies of the players ensure the convergence of the game to the equilibrium point with probability one, that is,

$$\begin{aligned}
&\sum_{k=1}^N r_k \left(\left\| p_{n+1}^{\hat{k}} - \bar{p}^{\hat{k}}(\varepsilon_{n+1}^k, \delta_{n+1}) \right\|^2 \right. \\
&\quad \left. + \|\lambda_{n+1} - \lambda(\varepsilon_{n+1}, \delta_{n+1})\|^2 \right) \xrightarrow{a.s.} 0
\end{aligned}$$

VIII. Conclusion

This paper was concerned with the development of a new learning algorithm for N-person constrained repeated game with unknown expected payoff and constraints. Based on Lyapunov approach and martingale theory, the convergence with probability one has been stated. Simulations results justifying this approach are available.

REFERENCES

- [1] M. Godoy-Alcántar, E. Gómez-Ramírez, A. Poznyak and K. Najim, Alternate Linear Programming Approach for Noncooperative Constrained Finite Games, *International Journal Systems Science*, submitted.
- [2] S. Lakshmivarahan, Learning Algorithms Theory and Applications, *Springer-Verlag*, Berlin, 1981.
- [3] K. Najim and A. S. Poznyak, *Learning Automata: Theory and Applications*. Oxford: Pergamon Press, 1994.
- [4] A. S. Poznyak and K. Najim, *Learning Automata and Stochastic Optimization*. Berlin: Springer-Verlag, 1997.
- [5] A. Poznyak and K. Najim, "Learning through reinforcement for N-person repeated constraint games", *IEEE Trans. on Systems, Man, and Cybernetics: Part B Cybernetics*, vol. 32, No.6, 2002.
- [6] K. S. Narendra and M. A. L. Thathachar, *Learning Automata: an Introduction*, *Prentice-Hall*, Englewood Cliffs, New Jersey, 1989.
- [7] A. V. Nazin and A. S. Poznyak, "Matrix N-person game with incomplete information," *Economics and Mathematical Methods*, vol. 14, pp. 958-968, 1978.
- [8] Nazin, A. V. and A. S. Poznyak, *Adaptive Choice of Variants* (in Russian). Moscow: Nauka, 1986.
- [9] J. B. Rosen, Existence and uniqueness of equilibrium points for concave N-persons games, *Econometrica*, vol. 33, pp. 520-534, 1965.
- [10] J. Nash, "Equilibrium points in n-person games," *Proc. Nat. Acad. U.S.A.*, vol.36, pp. 48-49, 1950.