

Solving Two-person Zero-sum Stochastic Games with Incomplete Information using Learning Automata with Artificial Barriers

Anis Yazidi, *Senior Member, IEEE*, Daniel Silvestre, and B. John Oommen, *Life Fellow, IEEE*, ,

Abstract—Learning Automata (LA) with artificially absorbing barriers was a completely new horizon of research in the 1980s [8]. These new machines yielded properties which were previously unknown. More recently absorbing barriers have been introduced in continuous estimator algorithms so that the proofs could follow a martingale property, as opposed to monotonicity [18], [19], [20]. However, the applications of LA with artificial barriers are almost non-existent. In that regard, this paper is pioneering in that it provides effective and accurate solutions to an extremely complex application domain, namely that of solving two-person zero-sum stochastic games which are provided with incomplete information. LA have been previously used [13] to design algorithms capable of converging to the game’s Nash equilibrium under limited information. Those algorithms have focused on the case where the Saddle Point of the game exists in a pure strategy. However, the majority of the LA algorithms used for games are absorbing in the probability simplex space, and thus they converge to an exclusive choice of a single action. These LA are thus unable to converge to other *mixed* Nash equilibria when the game possesses no Saddle Point for a pure strategy. The pioneering contribution of this paper¹ is that we propose a LA solution that is able to converge to an optimal mixed Nash equilibrium even though there may be no Saddle Point when a pure strategy is invoked. The scheme, being of the Linear Reward-Inaction (L_{R-I}) paradigm, is in and of itself, absorbing. However, by incorporating artificial barriers, we prevent it from being “stuck” or getting absorbed in pure strategies. Unlike the Linear Reward- ϵ Penalty ($L_{R-\epsilon P}$) scheme proposed by Lakshminarayanan and Narendra [1] almost four decades ago, our new scheme achieves the same goal with much less parameter tuning, and in a more elegant manner. The paper includes the non-trivial proofs of the theoretical results characterizing our scheme, and also contains experimental verification that confirms our

theoretical findings.

Index Terms—Learning Automata (LA), Games with Incomplete Information, LA with Artificial Barriers.

I. INTRODUCTION

Learning Automata: The term Learning Automata (LA) denotes a whole sub-field of research within adaptive systems with several books being dedicated to its study [2], [5], [6], [12], [14]. The work on LA dates to the the Soviet Union in the 1960s, when the mathematical giant, Tsetlin, [15] devised the so-called Tsetlin Machine which is a learning mechanism with finite memory. Tsetlin’s learning machines were demonstrated to give birth to self-organizing behavior through collective learning. In his work, Tsetlin pioneered the Goore game which is a distributed coordination game with limited feedback that has many practical applications, as shown by Tung and Kleinrock [16]. The early works in the field of LA, such as the Tsetlin Machine, fall under the family of Fixed Structure Learning Automata. The main stream of current LA research concerns the family of Variable Structure LA (VSLA) which, loosely speaking, differs from Fixed Structure LA in the fact that they operate with a probability vector that is updated dynamically over time. In Fixed Structure LA, the choice is governed by a transition matrix whose transitions do not depend on time and that describes how the internal states of the LA are updated based on the Environment’s feedback. The term Learning Automata was coined for the first time by Narendra and Thathachar in [6].

Markovian Representations of LA: LA can also be characterized by their Markovian representations. They thus fall into one of two families, being either ergodic or those that possess absorbing barriers. Such a characterization is crucial to the tenets of this paper. Absorbing automata have underlying Markov Chains that get absorbed or locked into a barrier state. Sometimes this can occur even after a relatively small, *finite* number of iterations. The classic references [2], [5], [6], [12], [14] report numerous LA families that contain such absorbing barriers. On the other hand, as these same references explain, the literature has also reported scores of ergodic automata, which converge in distribution. In these cases, the asymptotic *distribution* of the action probability vector converges to a value that is independent of its initial vector. Absorbing LA are usually designed to operate in stationary Environments. As opposed to these, ergodic LA are preferred for non-stationary Environments, namely those that possess

Anis Yazidi, Author’s status: *Professor*. This author can be contacted at: Oslo Metropolitan University, Department of Computer Science, Pilestredet 35, Oslo, Norway. This author is also an *Adjunct Professor* with University of Science and Technology, Trondheim, Norway and a Senior Researcher at Oslo University Hospital, Oslo, Norway. E-mail: anis.yazidi@oslomet.no.

Daniel Silvestre, Author’s status: *Assistant Professor* with the Lusófona University and also with the Institute for Systems and Robotics, Instituto Superior Técnico, University of Lisbon, 1049-001 Lisbon, Portugal. dsilvestre@isr.tecnico.ulisboa.pt. This author’s work was partially supported by the Portuguese Fundação para a Ciência e a Tecnologia (FCT) through Institute for Systems and Robotics (ISR), under Laboratory for Robotics and Engineering Systems (LARSyS) project UIDB/50009/2020, through project PCIF/MPG/0156/2019 FirePuma and through COPELABS, University Lusófona project UIDB/04111/2020.

B. John Oommen, *Chancellor’s Professor ; Life Fellow: IEEE and Fellow: IAPR*. The work of this author was partially supported by NSERC, the Natural Sciences and Engineering Council of Canada. This author can be contacted at: School of Computer Science, Carleton University, Ottawa, Canada : K1S 5B6. This author is also an *Adjunct Professor* with the University of Agder in Grimstad, Norway. E-mail address: oommen@scs.carleton.ca.

¹We are very grateful for the feedback from the anonymous Referees of the *original* submission. Their input significantly improved the quality of this final version.

time-dependent reward probabilities. These characterizations, and their corresponding implications for game playing, will be explained presently.

Continuous or Discretized VSLA: VSLA can also be characterized as being Continuous or Discretized. This depends on the values that the action probabilities can take. Continuous LA allow the action probabilities to assume any value in the interval $[0, 1]$. Such algorithms have a relatively slow rate of convergence. The problem with continuous LA is that they *approach* a goal but never *reach* there. This was mitigated in the 1980s by introducing the concept of discretization, where if an action probability was close enough to zero or unity, it could jump to that end point in a single step. This also rendered the LA to have a faster convergence, because one could increase their speeds of convergence, by incorporating this phenomenon [3], [4], [9]. This is implemented by constraining the action selection probability to be one of a finite number of values in the interval $[0, 1]$. By incorporating discretization, almost all of the reported VSSA of the continuous type have been also discretized [9], [10], [19].

LA with Artificially Absorbing Barriers: LA with artificially-introduced absorbing barriers was a novelty in the 1980s. These yielded machines which had properties that were previously unknown. This was due to the fact that a discretized machine, even though it was ergodic, could be rendered absorbing by forcing the machine to stay at one of the absorbing barriers [8]. Ironically, this simple step introduced families of new LA, with properties that were previously unknown. For example, the ADL_{R-P} and ADL_{I-P} are absorbing versions of their corresponding ergodic counterparts, but have been proven to be ϵ -optimal in all random environments. This phenomenon, of including artificially absorbing barriers, has been recently applied to the family of Pursuit LA [18].

Estimator LA with Artificial Barriers: The concept of introducing absorbing barriers is also central to the proofs of estimator algorithms. For three decades, these pursuit algorithms were “proven” to be ϵ -optimal by virtue of the monotonicity property. However, recently, these proofs have been shown to be flawed. To remedy this, absorbing barriers have been introduced in continuous estimator algorithms so that the proofs could follow a martingale property, as opposed to monotonicity. Consequently, Zhang and others have shown that one can invoke this weaker property, namely, the martingale property, by artificially providing such an absorbing barrier. Thus, whenever an action probability is close enough to unity, the LA is forced to jump to this absorbing barrier [18], [19], [20].

Applications of LA: LA have boasted scores of applications. These include theoretical problems like the graph partitioning problem. They have been used in controlling intelligent vehicles. When it concerns neural networks and hidden Markov models, Meybodi *et al.* have used them in adapting the former, and others have applied them in training the latter. Network call admission, traffic control and quality of service routing have been resolved using LA, while others They have also found applications in tackling problems involving network and communications issues. Apart from these, the entire field of LA and stochastic learning, has had a myriad of applications

listed in the reference books [2], [5], [6], [14]. In the interest of the page-limit constraints, the citations to these applications are not included. But they can be easily found by executing a simple search, and many are included in the above benchmark references.

Game Playing with LA: While artificially-introduced barriers have been shown to have powerful theoretical and design implications, the applications of them are few. This is where the present paper finds its place – it presents one such application. LA have also been used to resolve stochastic games with incomplete information. This paper pioneers a merge of the above two issues. First of all, we present a mechanism by which LA can be augmented with artificial barriers, but unlike the state-of-the-art, these barriers are non-absorbing. We then proceed to use these to play zero sum games with incomplete information. Games of this type were studied four decades ago for scenarios when the game matrix had a Saddle Point using traditional L_{R-I} and L_{R-P} LA. Our results generalize those when the game does not possess a Nash equilibrium. Rather, we propose the non-trivial use of LA with artificial non-absorbing barriers to resolve such games. The paper contains the theoretical results and those from simulations using corresponding benchmark games.

Landscape of our Present Work: In this paper, we propose an algorithm addressing zero sum games, which can be generalized to non-zero sum games in a manner similar to the principle by which the method in [1] was generalized in [17]. In the latter, Xing and Chandramouli proved that the Linear Reward – ϵ Penalty ($L_{R-\epsilon P}$) algorithm, devised in [1], is able to work in non-zero sum games. Thus, without further elaborating on this², our results are generalizable to non-zero sum games.

Since the game is zero sum, the outcomes are either a loss for player A , with reward -1 , and the corresponding win for player B with value $+1$, or the converse for the case of a win for player A . We emphasize that this is a limited information game where each player is unaware of both the mixed strategy and the selected action of the other player. The available information to each player is whether its action resulted in a win or a loss. The reader should note that either/both players might not even be aware of the *existence* of another player, and be working with the assumption that he is playing against Nature, as in the classical multi-armed bandit algorithms. However, if both players learn using our algorithm based on the assumption that they are operating with an adversarial environment, we show that they will both converge to the desired equilibrium. Our proposed scheme has players adjusting their strategy whenever it obtains a “win” for that round. This conforms to the Linear Reward-Inaction, L_{R-I} , paradigm, described in detail, presently. It is thus, unarguably, radically different from the mechanism proposed by Lakshminarayanan [1], where the probability updates are performed upon receiving both reward and penalty responses, and which thus render changes to occur at every time instant.

Objective and Contribution of this paper: Based on the

²Some preliminary unpublished work is being conducted for extending this work to non-zero sum games.

above discussion, one can summarize the objective of this paper as to study the behavior of a non-absorbing barrier-based L_{R-I} mechanism in a stochastic zero-sum game played by two players A and B , with two actions each, as earlier done in [1]. Each player uses an LA to decide his strategy, where the only received feedback from the environment is the reward of the joint actions of both players. The game is played iteratively and the players are able to revise their mixed strategies.

Applications of the proposed method: Learning within the context of games has a natural application in the realm of Game Theory. However, in the context of Multi-Agent Systems (MAS) this has been shown to be suitable for the cooperative control of robotic systems [21]. In such a design, it is assumed that the mission can be fully described as a potential game, where the utility function measures how well the nodes in the network are complying with the objectives. Nevertheless, having robots converging to pure strategies means that the network designer is favoring exploitation and disregarding exploration. If the environment changes and causes a different payoff matrix, the agent would be locked into repeatedly playing the same strategy. Moreover, this assumes that the utility must be known and deterministic. Therefore, instead of designing application-specific algorithms, the proposed learning algorithm can be used to address problems in cooperative control such as the so-called “rendezvous” problem for a fleet of robots [26], [27], the desynchronization of the use of a shared medium [23], [22], a consensus algorithm to have the agents agree on a common value [24], or to solve distributed computation such as the PageRank [25], by only considering the current stochastic payoff.

It is also pertinent to mention that the mechanism that we propose here, can be used by the agents to learn how to act if the payoff corresponds to how successful they are in following the objectives of the “mission”. Much can be said about this, but we terminate these discussions here in the interest of brevity and due to space limitations. However, with respect to future research, it is wise to mention that the question of whether they can be applied to synchronization, as in the analysis of the family of so-called “Firefly” algorithms, is yet open.

A. The Notation Used

Most of the notation that we use, is well-established from the theory of matrices and in the field of LA [2], [6], and stating them would trivialize the paper. However, we mention that apart from the well-established notation used in these areas, we will use the notation that the conditional expectation of some variable v with respect to w is written as $E[v|w]$, and the partial derivative of a variable $v(t)$ with respect to time t is denoted by $\frac{\partial v(t)}{\partial t}$.

II. THE GAME MODEL

To initiate discussions, we formalize the game model that is being investigated. Let $P(t) = [p_1(t) \ p_2(t)]^T$ denote the mixed strategy of player A at time instant t , where $p_1(t)$ accounts for the probability of adopting strategy 1 and, conversely, $p_2(t)$ stands for the probability of adopting strategy

2. Thus, $P(t)$ describes the distribution over the strategies of player A . Similarly, we can define the mixed strategy of player B at time t as $Q(t) = [q_1(t) \ q_2(t)]^T$. The extension to more than two actions per player is straightforward following the method analogous to what was used by Papavassilopoulos [11], which extended the work of Lakshmivarahan and Narendra [1].

Let $\alpha_A(t) \in \{1, 2\}$ be the action chosen by player A at time instant t and $\alpha_B(t) \in \{1, 2\}$ be the one chosen by player B , following the probability distributions $P(t)$ and $Q(t)$, respectively. The pair $(\alpha_A(t), \alpha_B(t))$ constitutes the joint action at time t , and are pure strategies. Specifically, if $(\alpha_A(t), \alpha_B(t)) = (i, j)$, the probability of gain for player A is determined by d_{ij} , as formalized in [1]. We thus construct a matrix with the set of probabilities $D = [d_{ij}], 1 \leq i \leq 2$, which is the so-called payoff matrix associated with the game.

The matrix D is given by:

$$D = \begin{pmatrix} d_{11} & d_{12} \\ d_{21} & d_{22} \end{pmatrix}, \quad (1)$$

where all the entries are probabilities.

Clearly, the actual game matrix G is given by $g_{ij} = 2d_{ij} - 1$, with entries in the interval $[-1, 1]$. Without loss of generality, player A corresponds to the row player while B is the column player. Further, when referring to a “gain” we are seeing this from the perspective of player A .

In zero-sum games, Nash equilibria are equivalently called the “Saddle Points” for the game. Since the outcome for a given joint action is stochastic, the game is the stochastic form of a zero-sum game. The “zero-sum” property implies that at any time t , there is only one winning player³.

In the interest of completeness, we present the original scheme proposed in [1] based on the $L_{R-\epsilon P}$ rule. It uses two parameters θ_R and θ_P as the learning rates associated with the reward and penalty responses, respectively. When player A gains at time instant t by playing action i , he updates his mixed strategy as:

$$\begin{aligned} p_i(t+1) &= p_i(t) + \theta_R(1 - p_i(t)) \\ p_s(t+1) &= p_s(t) - \theta_R p_s(t) \quad \text{for } s \neq i. \end{aligned}$$

However, if player A loses after using action i , his mixed strategy is updated by the following:

$$\begin{aligned} p_i(t+1) &= p_i(t) - \theta_P p_i(t) \\ p_s(t+1) &= p_s(t) + \theta_P(1 - p_s(t)) \quad \text{for } s \neq i. \end{aligned}$$

The exact update mechanism for player B is obtained by replacing the corresponding $p(t)$ by $q(t)$, and by recalling that a gain for A maps onto a loss scenario for player B . We now introduce our novel solution that is proposed to learn a new mixed strategy.

³The results inferred from this paper can be extended to non-zero sum games. However, for the sake of simplicity we only consider the case of zero-sum games.

III. LA ALGORITHM BASED ON THE L_{R-I} WITH ARTIFICIAL BARRIERS

A. Non-Absorbing Artificial Barriers

We have earlier seen that an ergodic LA can be made absorbing by artificially rendering the end states to become absorbing. This was briefly addressed above. But what has not been discussed in the literature is a strategy by which a scheme which is, in and of itself, absorbing, can be rendered to be ergodic. In other words, the LA is allowed to move within the probability simplex by utilizing an absorbing scheme. But when it enters an absorbing barrier, the scheme is forced to go back into the simplex in order to render it to be ergodic. No such scheme has ever been reported in the literature, and the advantage of having such a scheme is that one does not get locked into a sub-optimal absorbing barrier. Rather, we can permit it to move around so that it can migrate stochastically towards an optimal mixed strategy. This is, precisely, what we shall do.

B. Non-Absorbing Game Playing

We now present our strategic LA-based game algorithm together with a formal analysis that demonstrates the convergence to the Saddle Points of the game even if the Saddle Point corresponds to a *mixed* Nash equilibrium. Our LA solution is based on the L_{R-I} scheme, but as alluded to earlier, it has been modified in order to non-trivially provide non-absorbing barriers. The proof of convergence is based on Norman's theory for learning processes characterized by small learning steps [6], [7].

Considering that p_{max} denotes an artificial barrier, we use the notation that $p_{min} = 1 - p_{max}$. We further constrain the probability for each action by restricting it, by design, to belong to the interval $[p_{min}, p_{max}]$ if $p_1(0)$ and $q_1(0)$ are initially chosen to belong to the same interval. If the outcome from the environment is a gain at a time t for action $i \in \{1, 2\}$, the update rule is given by:

$$\begin{aligned} p_i(t+1) &= p_i(t) + \theta(p_{max} - p_i(t)) \\ p_s(t+1) &= p_s(t) + \theta(p_{min} - p_s(t)) \quad \text{for } s \neq i. \end{aligned} \quad (2)$$

The reader will observe that this update mechanism is identical to the well-established linear schemes, except that p_{min} and p_{max} replace the values zero and unity respectively. When the player receives a loss, the probabilities are not updated, which translates into:

$$\begin{aligned} p_i(t+1) &= p_i(t) \\ p_s(t+1) &= p_s(t) \quad \text{for } s \neq i. \end{aligned} \quad (3)$$

The update rules for the mixed strategy $q(t+1)$ are defined in a similar fashion by recalling the dichotomy that whenever player A gains, it corresponds to a loss for player B , and vice-versa. Analogous to the L_{R-I} paradigm, mixed strategies are not changed in the case of a loss.

We now proceed to analyze the convergence properties of the proposed algorithm. To aid in the analysis, we identify the Nash equilibrium of the game by the pair (p_{opt}, q_{opt}) . To render the presentation to be less cumbersome, we divide the analysis into two cases.

a) *Case 1: Only One Mixed Nash Equilibrium Case (No Saddle Point in pure strategies):* The first case depicts the situation where no Saddle Point exists in pure strategies. In other words, the only Nash equilibrium is a mixed one. Based on the fundamentals of Game Theory, the optimal mixed strategies can be easily shown to be the following:

$$p_{opt} = \frac{d_{22} - d_{21}}{L}, \quad q_{opt} = \frac{d_{22} - d_{12}}{L},$$

where $L = (d_{11} + d_{22}) - (d_{12} + d_{21})$. Without loss of generality, we assume that:

$$d_{11} > \max\{d_{12}, d_{21}\} \quad \text{and} \quad d_{22} > \max\{d_{12}, d_{21}\}. \quad (4)$$

Notice that the above inequalities are not restrictive, as games not satisfying them can be mapped in a symmetric manner by re-indexing the actions of the players and/or the indices of the players.

b) *Case 2: There is a Saddle Point in pure strategies:* The case where the game matrix has Saddle Points in pure strategies corresponds to either:

- $d_{11} > d_{12}$, $d_{12} < d_{21}$, $d_{21} > d_{22}$ and $d_{22} < d_{11}$;
- Or in the symmetric case, where $d_{11} < d_{12}$, $d_{12} > d_{21}$, $d_{21} < d_{22}$ and $d_{22} > d_{11}$.

Since the other cases can be proven in identical manners, in the interest of brevity, we consider only the case where:

$$d_{21} < d_{11} < d_{12}. \quad (5)$$

In this case, $p_{opt} = 1$ and $q_{opt} = 1$. The other sub-cases within Case 2 can be obtained by re-indexing the actions of the players and/or the indices of the players, as in Case 1.

Let the vector $X(t) = [p_1(t) \quad q_1(t)]^T$. We introduce the notation that $\Delta X(t) = X(t+1) - X(t)$. We also represent the conditional expected value operator by $\mathbb{E}[\cdot]$. Using these, we claim the next theorem.

Theorem 1. *Consider a zero sum game with a payoff matrix as in Eq. (1) and a learning algorithm defined by equations Eq. (2) and Eq. (3) for both players A and B , with learning rate θ . Then, $E[\Delta X(t)|X(t)] = \theta W(x)$ and for every $\epsilon > 0$, there exists a unique stationary point $X^* = [p_1^* \quad q_1^*]^T$ satisfying:*

- 1) $W(X^*) = 0$;
- 2) $|X^* - X_{opt}| < \epsilon$.

Proof. Let us first compute the conditional expected value⁴ of the increment $\Delta X(t)$:

$$\begin{aligned} E[\Delta X(t)|X(t)] &= E[X(t+1) - X(t)|X(t)] \\ &= \begin{bmatrix} E[p_1(t+1) - p_1(t)|X(t)] \\ E[q_1(t+1) - q_1(t)|X(t)] \end{bmatrix} \\ &= \theta \begin{bmatrix} W_1(X(t)) \\ W_2(X(t)) \end{bmatrix} \\ &= \theta W(X(t)), \end{aligned}$$

⁴Computing the ‘‘expected value of the increment’’ is a standard procedure in the theory of LA. This is because the increment, in and of itself, is a random variable, which is sometimes positive and sometimes negative. Quantifying the latter is not possible due to the randomness of the updating rule. However, the *conditional* expected value of the increment can be determined, whence (by invoking the ‘‘Law of the Unconscious Statistician’’), one can determine the expected value of the increment itself.

where the above format is possible since all possible updates share the form $\Delta X(t) = \theta W(t)$, for some $W(t)$, as given in Eq. (2).

For ease of notation, we drop the dependence on t with the implicit assumption that all occurrences of X , p_1 and q_1 represent $X(t)$, $p_1(t)$ and $q_1(t)$ respectively. $W_1(x)$ is then:

$$\begin{aligned} W_1(X) &= p_1 q_1 d_{11} (p_{max} - p_1) + p_1 (1 - q_1) d_{12} (p_{max} - p_1) \\ &\quad + (1 - p_1) q_1 d_{21} (p_{min} - p_1) \\ &\quad + (1 - p_1) (1 - q_1) d_{22} (p_{min} - p_1) \\ &= p_1 [q_1 d_{11} + (1 - q_1) d_{12}] (p_{max} - p_1) \\ &\quad + (1 - p_1) [q_1 d_{21} + (1 - q_1) d_{22}] (p_{min} - p_1) \\ &= p_1 (p_{max} - p_1) D_1^A(q_1) + (1 - p_1) (p_{min} - p_1) D_2^A(q_1), \end{aligned} \quad (6)$$

where,

$$D_1^A(q_1) = q_1 d_{11} + (1 - q_1) d_{12} \quad (7)$$

$$D_2^A(q_1) = q_1 d_{21} + (1 - q_1) d_{22}. \quad (8)$$

By replacing $p_{max} = 1 - p_{min}$ and rearranging the expression we get:

$$\begin{aligned} W_1(X) &= p_1 (1 - p_1) D_1^A(q_1) - p_1 p_{min} D_1^A(q_1) \\ &\quad + (1 - p_1) p_{min} D_2^A(q_1) - p_1 (1 - p_1) D_2^A(q_1) \\ &= p_1 (1 - p_1) [D_1^A(q_1) - D_2^A(q_1)] \\ &\quad - p_{min} [p_1 D_1^A(q_1) - (1 - p_1) D_2^A(q_1)]. \end{aligned}$$

Similarly, we can get

$$\begin{aligned} W_2(X) &= q_1 p_1 (1 - d_{11}) (p_{max} - q_1) + \\ &\quad q_1 (1 - p_1) (1 - d_{12}) (p_{max} - q_1) \\ &\quad + (1 - q_1) p_1 (1 - d_{21}) (p_{min} - q_1) + \\ &\quad (1 - q_1) (1 - p_1) (1 - d_{22}) (p_{min} - q_1) \\ &= q_1 [p_1 (1 - d_{11}) + (1 - p_1) (1 - d_{12})] (p_{max} - q_1) \\ &\quad + (1 - q_1) [p_1 (1 - d_{21}) + (1 - p_1) (1 - d_{22})] (p_{min} - q_1) \\ &= q_1 (p_{max} - q_1) [1 - D_1^B(p_1)] + \\ &\quad (1 - q_1) (p_{min} - q_1) [1 - D_2^B(p_1)] \end{aligned} \quad (9)$$

where

$$D_1^B(p_1) = p_1 d_{11} + (1 - p_1) d_{21} \quad (10)$$

$$D_2^B(p_1) = p_1 d_{12} + (1 - p_1) d_{22}. \quad (11)$$

By replacing $p_{max} = 1 - p_{min}$ and rearranging the expression we get:

$$\begin{aligned} W_2(X) &= q_1 (1 - q_1) (1 - D_1^B(p_1)) - q_1 p_{min} (1 - D_1^B(p_1)) \\ &\quad + (1 - q_1) p_{min} (1 - D_2^B(p_1)) - q_1 (1 - q_1) (1 - D_2^B(p_1)) \\ &= -q_1 (1 - q_1) [D_1^B(p_1) - D_2^B(p_1)] \\ &\quad + p_{min} [-q_1 (1 - D_1^B(p_1)) + (1 - q_1) (1 - D_2^B(p_1))] \\ &= -q_1 (1 - q_1) [D_1^B(p_1) - D_2^B(p_1)] + \\ &\quad p_{min} [q_1 D_1^B(p_1) - (1 - q_1) D_2^B(p_1) + (1 - 2q_1)]. \end{aligned} \quad (12)$$

We need to address the two identified cases. Consider Case 1), where there is only a single mixed equilibrium. According to Eq. (4), we get:

$$\begin{aligned} D_{12}^A(q_1) &= D_1^A(q_1) - D_2^A(q_1) \\ &= (d_{12} - d_{22}) + L q_1. \end{aligned} \quad (13)$$

Given that $L > 0$, since $d_{11} > d_{12}$ and $d_{22} > d_{21}$, $D_{12}^A(q_1)$ is an increasing function of q_1 and

$$\begin{cases} D_{12}^A(q_1) < 0, & \text{if } q_1 < q_{opt}, \\ D_{12}^A(q_1) = 0, & \text{if } q_1 = q_{opt}, \\ D_{12}^A(q_1) > 0, & \text{if } q_1 > q_{opt}. \end{cases} \quad (14)$$

For a given q_1 , $W_1(X)$ is quadratic in p_1 . Also, we have:

$$\begin{aligned} W_1 \left(\begin{bmatrix} 0 \\ q_1 \end{bmatrix} \right) &= p_{min} D_2^A(q_1) > 0 \\ W_1 \left(\begin{bmatrix} 1 \\ q_1 \end{bmatrix} \right) &= -p_{min} D_1^A(q_1) < 0. \end{aligned} \quad (15)$$

Since $W_1(X)$ is quadratic with a negative second derivative with respect to p_1 , and since the inequalities in Eq. (15) are strict, it admits a single root p_1 for $p_1 \in [0, 1]$. Moreover, we have $W_1(X) = 0$ for some p_1 such that:

$$\begin{cases} p_1 < \frac{1}{2}, & \text{if } q_1 < q_{opt}, \\ p_1 = \frac{1}{2}, & \text{if } q_1 = q_{opt}, \\ p_1 > \frac{1}{2}, & \text{if } q_1 > q_{opt}. \end{cases} \quad (16)$$

Using a similar argument, we can see that there exists a single solution for each p_1 , and as $p_{min} \rightarrow 0$, we conclude that $W_1(X) = 0$ whenever $p_1 \in \{0, p_{opt}, 1\}$. Arguing in a similar manner we see that $W_2(X) = 0$ when:

$$X \in \left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} p_{opt} \\ q_{opt} \end{bmatrix} \right\}.$$

Thus, there exists a small enough value for p_{min} such that $X^* = [p^*, q^*]^T$ satisfies $W_2(X^*) = 0$, proving Case 1).

In the proof of Case 1), we have utilized the fact that for small enough p_{min} , the learning algorithm admits a stationary point, and also identified the corresponding possible values for this point. It is thus always possible to select a small enough $p_{min} > 0$ such that X^* approaches X_{opt} , concluding the proof for Case 1.)

Case 2) can be derived in a similar manner, and the details are omitted to avoid repetition. \square

In the next theorem, we show that the expected value of $\Delta X(t)$ has a negative definite gradient.

Theorem 2. *The matrix of partial derivatives, $\frac{\partial W(X^*)}{\partial x}$ is negative definite.*

Proof. We start the proof by writing the explicit format for $\frac{\partial W(X)}{\partial X} = \begin{bmatrix} \frac{\partial W_1(X)}{\partial p_1} & \frac{\partial W_1(X)}{\partial q_1} \\ \frac{\partial W_2(X)}{\partial p_1} & \frac{\partial W_2(X)}{\partial q_1} \end{bmatrix}$ and then computing each of the entries as below:

$$\begin{aligned} \frac{\partial W_1(X)}{\partial p_1} &= (1 - 2p_1) (D_1^A(q_1) - D_2^A(q_1)) - \\ &\quad p_{min} (D_1^A(q_1) + D_2^A(q_1)) \\ &= (1 - 2p_1) D_{12}^A(q_1) - \\ &\quad p_{min} (D_1^A(q_1) + D_2^A(q_1)). \end{aligned}$$

$$\begin{aligned} \frac{\partial W_1(X)}{\partial q_1} &= p_1 (1 - p_1) L - p_{min} (p_1 (d_{11} - d_{12}) + \\ &\quad (1 - p_1) (d_{22} - d_{21})). \end{aligned}$$

$$\frac{\partial W_2(X)}{\partial p_1} = -q_1(1-q_1)L + p_{min}((q_1(d_{11}-d_{21}) - (1-q_1)(d_{12}-d_{22})).$$

$$\frac{\partial W_2(X)}{\partial q_1} = -(1-2q_1)(D_1^B(p_1) - D_2^B(p_1)) + p_{min}(D_1^B(p_1) + D_2^B(p_1) - 2).$$

As seen in Theorem 1, for a small enough value for p_{min} , we can ignore the terms that are weighted by p_{min} , and we will thus have $\frac{\partial W(X^*)}{\partial X} \approx \frac{\partial W(X_{opt})}{\partial X}$. We now subdivide the analysis in the two cases identified as above, which are equivalent to:

- Case 1: No Saddle Point in pure strategies
- Case 2: There is a Saddle point in pure strategies.

c) *Case 1: No Saddle Point in pure strategies:* In this case, we have:

$$D_1^A(q_{opt}) = D_2^A(q_{opt}) \quad \text{and} \quad D_1^B(p_{opt}) = D_2^B(p_{opt})$$

which makes

$$\frac{\partial W_1(X_{opt})}{\partial p_1} = -2p_{min}D_1^A(q_{opt}). \quad (17)$$

Similarly, we can compute

$$\frac{\partial W_1(X_{opt})}{\partial q_1} = (1-2p_{min})p_{opt}(1-p_{opt})L. \quad (18)$$

The entry $\frac{\partial W_2(X_{opt})}{\partial p_1}$ can be simplified to:

$$\frac{\partial W_2(X_{opt})}{\partial p_1} = -(1-2p_{min})q_{opt}(1-q_{opt})L \quad (19)$$

and

$$\frac{\partial W_2(X_{opt})}{\partial q_1} = -2p_{min}(1-D_1^B(p_{opt})) \quad (20)$$

resulting in:

$$\frac{\partial W(X_{opt})}{\partial X} = \begin{bmatrix} -2p_{min}D_1^A(q_{opt}) & (1-2p_{min})p_{opt}(1-p_{opt})L \\ -(1-2p_{min})q_{opt}(1-q_{opt})L & -2p_{min}(1-D_1^B(p_{opt})) \end{bmatrix}. \quad (21)$$

The matrix given in Eq. (21) satisfies:

$$\det\left(\frac{\partial W(X_{opt})}{\partial x}\right) > 0, \quad \text{trace}\left(\frac{\partial W(X_{opt})}{\partial x}\right) < 0, \quad (22)$$

which implies the 2×2 matrix is negative definite.

d) *Case 2: There is a Saddle Point in pure strategies:*

In Theorem 1, Case 2 reduces to considering $q_{opt} = 1$ and $p_{opt} = 1$.

Computing the entries of the matrix for this case yields:

$$\frac{\partial W_1(X_{opt})}{\partial p_1} = -(d_{11}-d_{21}) - p_{min}(d_{11}+d_{21}), \quad (23)$$

and

$$\frac{\partial W_1(X_{opt})}{\partial q_1} = -p_{min}(d_{11}-d_{12}). \quad (24)$$

The entry $\frac{\partial W_2(X_{opt})}{\partial p_1}$ can be simplified to:

$$\frac{\partial W_2(X_{opt})}{\partial p_1} = p_{min}(d_{11}-d_{21}) \quad (25)$$

and

$$\frac{\partial W_2(X_{opt})}{\partial q_1} = (d_{11}-d_{12}) - p_{min}(2-d_{11}-d_{12}) \quad (26)$$

resulting in:

$$\frac{\partial W(X_{opt})}{\partial X} = \begin{bmatrix} -(d_{11}-d_{21}) - p_{min}(d_{11}+d_{21}) & -p_{min}(d_{11}-d_{12}) \\ p_{min}(d_{11}-d_{21}) & (d_{11}-d_{12}) - p_{min}(2-d_{11}-d_{12}) \end{bmatrix}. \quad (27)$$

The matrix in (27) satisfies:

$$\det\left(\frac{\partial W(X_{opt})}{\partial X}\right) > 0, \quad \text{trace}\left(\frac{\partial W(X_{opt})}{\partial X}\right) < 0 \quad (28)$$

for a sufficiently small value of p_{min} , which again implies that the 2×2 matrix is negative definite. \square

Theorem 3. *Let V be the von Neumann value of the game given by matrix D . Let $\mathbf{p}(t) = [p_1, p_2]$ and $\mathbf{q}(t) = [q_1, q_2]$. For a sufficiently small p_{min} approaching 0, $\eta(t)$ converges to V as $\theta \rightarrow 0$ where:*

$$\eta(t) \triangleq E[\mathbf{p}(t)]DE[\mathbf{q}^T(t)] \quad (23)$$

Proof. The proof of this results requires a classic result due to Norman [7], given in the Appendix A, in the interest of completeness.

The convergence of $[E(p_1(t)) \ E(q_1(t))]$ to $[p_{opt}^* \ q_{opt}^*]$ is a consequence of this theorem. Interestingly enough, this theorem is a classical fundamental result that has been used to prove many of the convergence results in LA. It has, for example, been used by the seminal paper by Lakshmivarahan and Narendra to derive similar convergence properties of the $L_{R-\epsilon P}$ [1], applicable for the same game settings as ours. Indeed, it is easy to verify that Assumptions (1)-(6) required for Norman's result are satisfied. Thus, by further invoking Theorem 1 and Theorem 2, the result follows. \square

We conclude this section by mentioning that like all LA algorithms, the computational complexity of our scheme is linear in the size of the action probability vector. This is because, at the most, all the action probabilities are updated at every time instant.

For the benefit of future researchers, we believe that it will be profitable to record the hurdles we encountered in this research. The break-through came when we were able to devise/design LA systems which possessed no-absorbing barriers. In others words, it involved the concept of forcing the LA back into the probability space when it was close enough to the absorbing barriers. This was a phenomenon which we had not earlier seen in the literature. The consequent problem was the analysis. The underlying Markov process could not be easily analyzed using the properties of absorbing chains. Neither could it be trivially modelled as an ergodic chain converging to an equilibrium distribution. The analysis that we presented here came as a "brain-wave", and once the building blocks were established, everything naturally seemed to fall in place. These few sentences, requested by an Anonymous Referee, should clarify the difficulties encountered in this research, in order to show that the present research is pioneering, and that is not a trivial extension of existing methodologies.

Table I: Error for different values of θ and p_{max} , when $p_{opt} = 0.5789$ and $q_{opt} = 0.7368$ for the game specified by the D matrix given by Eq. (29). The point that you have raised is pertinent.

p_{max}	$\theta = 0.001$	$\theta = 0.0001$
0.999	2.1621×10^{-3}	1.6820×10^{-3}
0.998	2.5456×10^{-3}	1.7059×10^{-3}
0.997	3.7380×10^{-3}	2.2332×10^{-3}
0.996	3.4007×10^{-3}	2.0155×10^{-3}
0.995	5.4371×10^{-3}	3.7888×10^{-3}
0.994	5.5962×10^{-3}	4.2018×10^{-3}
0.993	7.3416×10^{-3}	5.4064×10^{-3}
0.992	7.7319×10^{-3}	7.8230×10^{-3}
0.991	9.6127×10^{-3}	6.7476×10^{-3}
0.990	9.3467×10^{-3}	9.6713×10^{-3}

IV. SIMULATIONS

In this section, we present simulations to confirm the above-mentioned theoretical properties of the proposed learning algorithm. In the interest of maintaining benchmarks, we adopt the same examples as those reported in [1]. Also, by using different instances of the payoff matrix D , we are able to experimentally cover the two cases referred to in Section III. Again, we refer to those cases as Case 1 and Case 2, as done in [1].

A. Convergence in Case 1

We consider an instance of the game where only one mixed Nash equilibrium exists, i.e., there is no Saddle Point in pure strategies. We adopt the same game matrix D as in [1] given by:

$$D = \begin{pmatrix} 0.6 & 0.2 \\ 0.35 & 0.9 \end{pmatrix} \quad (29)$$

which admits $p_{opt} = 0.5789$ and $q_{opt} = 0.7368$.

In order to eliminate the Monte Carlo error, we ran our scheme for 5×10^6 iterations, and report the error in Table I for different values of p_{max} and θ as the difference between X_{opt} and the mean over time of $X(t)$ after convergence⁵. An important remark is that the error decreases as p_{max} approaches 1 (i.e., when $p_{min} \rightarrow 0$). Please observe that in this case, we have particularly chosen to not let p_{max} be unity. If we allow it to be precisely unity, it would mean that we would not require an *artificial* barrier close to unity (for example, between 0.990 and 0.999 as in Table 1). In fact, for $p_{max} = 0.999$ and $\theta = 0.001$, the method achieves an error of 2.1621×10^{-3} , and further reducing $\theta = 0.0001$ leads to an error of 1.6820×10^{-3} .

To better visualize the scheme, Figure 1 depicts the evolution over time of the mixed strategies for both players (given by $X(t)$) for an ensemble of 1,000 runs using $\theta = 0.01$ and $p_{max} = 0.999$.

The trajectory of the ensemble allows us to perceive the mean evolution of the mixed strategies. The spiral pattern is caused by one of the players adapting to the strategy being used by the other, before the former learns by over-correcting

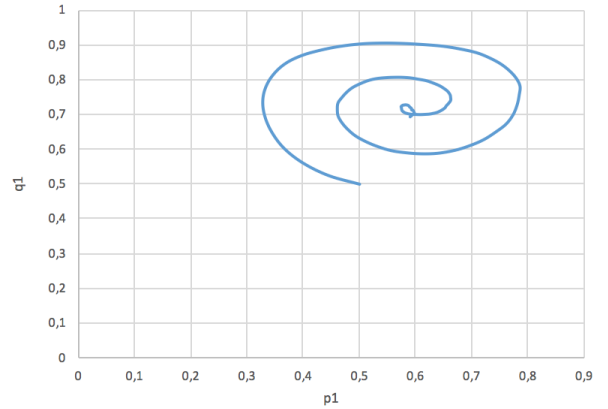


Figure 1: Time evolution of $[p_1(t), q_1(t)]^\top$ for the same settings as in Figure 2.

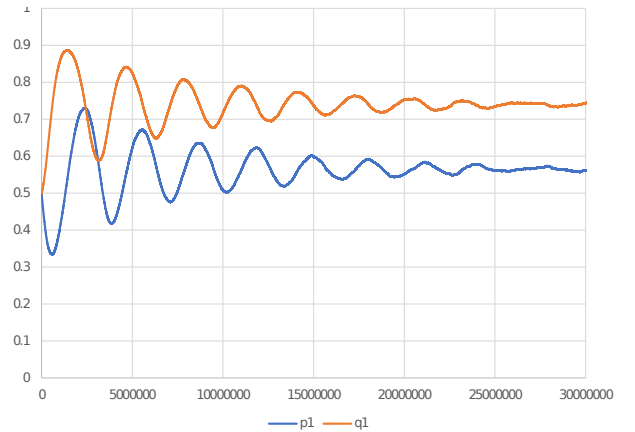


Figure 2: Trajectory of $X(t)$ for the case of the D matrix given by Eq. (29) with $p_{opt} = 0.5789$ and $q_{opt} = 0.7368$, and using $p_{max} = 0.99$ and $\theta = 0.00001$.

its strategy. The procedure is continued leading to smaller corrections until the players reach the Nash equilibrium.

The above-mentioned behavior can also be visualized in Figure 2 that presents the trajectory for a single experiment with $p_{max} = 0.99$ and $\theta = 0.00001$ over 3×10^7 steps. The described oscillatory behavior is attenuated as the players play for more iterations. The reader should particularly observe that a larger value of θ will cause more steady state error (as specified in Theorem 1), but it will also perturb this behavior as the nodes take larger updates whenever they win. On the other hand, further decreasing θ results in a smaller error of the stationary point of the method, but also decreases the convergence speed. This well-established inherent trade-off between the steady state error and rate of convergence can be better visualized by comparing Figure 1 with $\theta = 0.001$ against Figure 3 for a smaller value of $\theta = 10^{-5}$.

Further, in order to clearly emphasize the necessity of using an artificial barrier, we have specifically repeated the same experiment except that we have included an absorbing barrier instead, i.e., set $p_{max} = 1$. The result is illustrated in Figure 4. In this case, we expect that the scheme enters an absorbing barrier. Since it is impossible for the human eye

⁵The mean is taken over the last 10% of the total number of iterations.

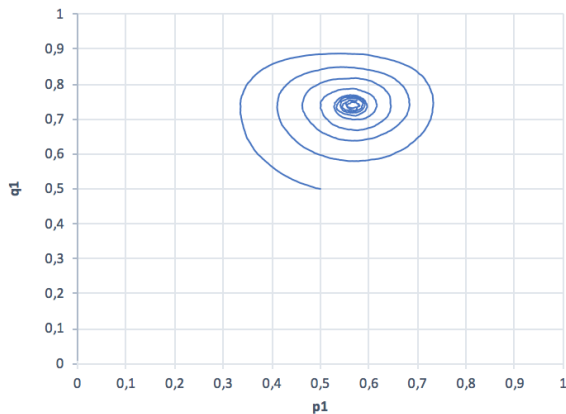


Figure 3: Time evolution of $X(t)$ where $p_{opt} = 0.5789$ and $q_{opt} = 0.7368$, using $p_{max} = 0.99$ and $\theta = 0.00001$.

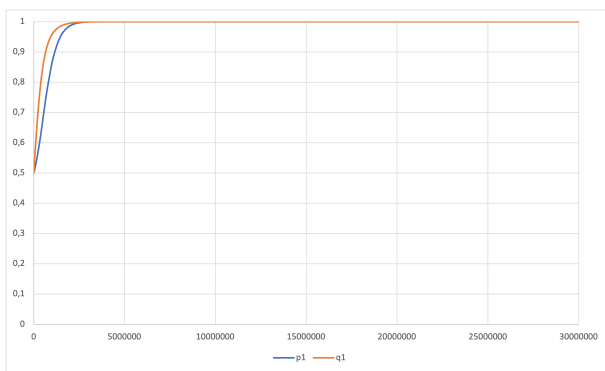


Figure 4: Trajectory of $X(t)$ for the case of the D matrix given by Eq. (29) and using an absorbing barrier $p_{max} = 1$ and $\theta = 0.00001$.

to detect whether or not we entered an absorbing barrier by merely examining the graph, we also manually checked the log of the experiment and verified that the probabilities became exactly unity after around 6,798,000 iterations. Although the theoretical convergence should have occurred in the limit and not after a finite number of iterations, the machine limited accuracy rounded the probabilities to unity after this juncture.

B. Pure equilibrium

In order to assess the performance of the proposed learning algorithm on cases with a pure equilibrium, we consider two instances of games falling in the category of Case 2 with $p_{opt} = 1$ and $q_{opt} = 1$. The payoff matrices D_1 and D_2 for the two games are given by:

$$D_1 = \begin{pmatrix} 0.6 & 0.8 \\ 0.35 & 0.9 \end{pmatrix}$$

$$D_2 = \begin{pmatrix} 0.7 & 0.9 \\ 0.6 & 0.8 \end{pmatrix}$$

We first show the convergence errors of our method for both games D_1 and D_2 in Table II and Table III, respectively. As in the previous simulation for Case 1, the errors are on the order

Table II: Error for different values of θ and p_{max} for D_1 .

p_{max}	$\theta = 0.001$	$\theta = 0.0001$
0.999	2.1073×10^{-3}	2.0971×10^{-3}
0.998	4.2753×10^{-3}	4.4573×10^{-3}
0.997	6.6147×10^{-3}	6.9025×10^{-3}
0.996	8.7588×10^{-3}	8.9192×10^{-3}
0.995	1.0815×10^{-2}	1.1044×10^{-2}
0.994	1.3424×10^{-2}	1.2894×10^{-2}
0.993	1.5005×10^{-2}	1.5415×10^{-2}
0.992	1.7347×10^{-2}	1.7805×10^{-2}
0.990	1.9772×10^{-2}	1.9670×10^{-2}
0.99	2.2516×10^{-2}	2.2548×10^{-2}

Table III: Error for different values of θ and p_{max} for D_2 .

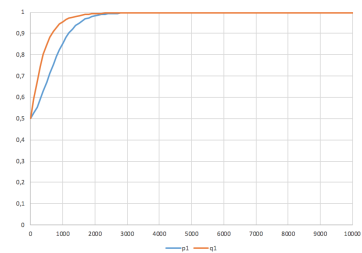
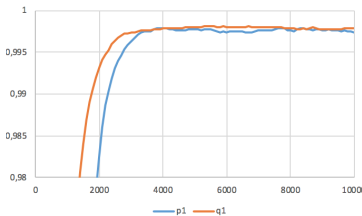
p_{max}	$\theta = 0.001$	$\theta = 0.0001$
0.999	6.2439×10^{-3}	4.8359×10^{-3}
0.998	9.8350×10^{-3}	9.9923×10^{-3}
0.997	1.5470×10^{-2}	1.3583×10^{-2}
0.996	1.8769×10^{-2}	2.1704×10^{-2}
0.995	2.5573×10^{-2}	2.4587×10^{-2}
0.994	3.1426×10^{-2}	2.8006×10^{-2}
0.993	3.6112×10^{-2}	3.5181×10^{-2}
0.992	3.7508×10^{-2}	4.0789×10^{-2}
0.991	4.6255×10^{-2}	4.2545×10^{-2}
0.99	4.8299×10^{-2}	4.5069×10^{-2}

to 10^{-3} for larger values of p_{max} . However, given that our algorithm uses artificial barriers to prevent absorbing states, the error is lower bounded by p_{min} . A similar issue is present in game D_2 . We have also included this simulation since it is a more challenging game to learn with our method for a larger steady-state error, even for very small values of θ .

In Figure 5, we depict the time evolution of the two components of the vector $X(t)$ using the proposed algorithm for an ensemble of 1,000 runs. In the case of having a Pure Nash equilibrium, there is no oscillatory behavior as when a player assigns more probability to an action, since the other player reinforces the strategy. However, Figure 5a could lead make one believe that the LA method has converged to a pure strategy. Figure 5b zooms around the point where the strategies have converged to showcase that their maximum value is limited by p_{max} , as per the design of our updating rule. This mechanism is particularly favourable to prevent players from converging to absorbing states for games with time-varying payoff matrices. However, the study of such a scenario is left for future research, namely that of determining how to design p_{max} and θ that represent a good trade-off between learning the game and adapting to a change in the payoffs.

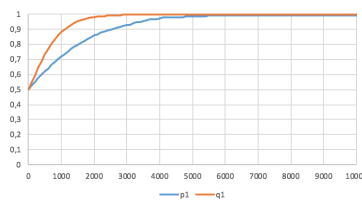
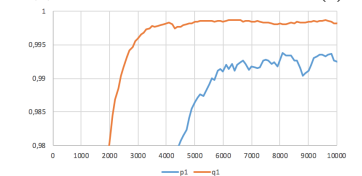
Game D_2 presents a harder challenge for our method as we can see from its larger steady state error. Figure 6 depicts the time evolution of the probabilities for each player when the algorithm is applied to D_2 with $\theta = 0.01$, $p_{max} = 0.999$ and for an ensemble with 1,000 runs.

The main remark regarding the results presented in Figure 6a is that the convergence is much slower when compared to game D_1 . This behavior is governed by the fact that the entries in matrix D_2 are closer to each other, unlike in D_1 where there is a clear disadvantage for player A when selecting action 2. There will, thus, be much fewer updates for player A where

(a) Evolution over time of $X(t)$.

(b) Zoomed version around convergence.

Figure 5: The figure shows a) the evolution over time of $X(t)$ for $\theta = 0.01$ and $P_{max} = 0.999$ when applied to game with payoffs D_1 , and b) is a zoomed version around the steady-state value.

(a) Evolution over time of $X(t)$.

(b) Zoomed version when the mixed strategies are approaching convergence.

Figure 6: The figure shows a) the evolution over time of $X(t)$ for $\theta = 0.01$ and $P_{max} = 0.999$ when applied to game with payoffs D_2 , and b) is a zoomed version around the steady-state value.

it increases the probability of action 2 in game D_1 – which is not pertinent in game D_2 . Figure 6b further emphasizes this remark by displaying a zoom, and depicting a sharper change in the probabilities in comparison with the smooth behavior in game D_1 .

C. Comparisons with Related Works

Now that we have explained our new techniques and established its theoretical basis, we continue this discussion with a brief comparison with some of the prior art⁶.

⁶We are thankful to the anonymous Referee who requested this comprehensive sub-section. It significantly adds to the quality of the paper.

First of all, one possible alternative when the payoff matrix is known can be to consider the problem as that of designing a local controller for each of the agents. One alternative is to explore the results in [28] and further investigated in [29]. However, this is only possible when D is known, which is not the scenario that we have assumed in this paper.

It is not out of place to review some of the relevant works in game theory that are not necessarily solved using LA. However, in the interest of space and brevity, we will not aim at submitting an extensive review of the field of game theory. Rather, we shall cite some pertinent works inasmuch as our main contribution in this article centers on advancing the field of *LA-based* solutions, and more specifically those dealing with the special case of games with “incomplete information”.

There are different variants of zero-sum stochastic games in the literature. In [30], Flesch *et al.* have proven the general result that every positive zero-sum stochastic game with countable state and action spaces, admits a value if at least one player has a finite action space. A similar value-existence result was obtained for a zero-sum stochastic game [31] with a continuous-time Markov chain, where the players have also the possibility of stopping the game. Ziliotto [32] considered weighted-average stochastic games, that is, stochastic games where Player 1 maximizes (in expectation) a fixed weighted average of the sequence of rewards. A so-called pumping algorithm was proposed in [33] for two-person zero-sum undiscounted stochastic games. Other approaches map the game onto a dynamic programming problem, and solve it based on Bellman’s optimal principle using concepts from the theory of optimal control [34].

The research on game theoretical learning with *incomplete information* [13] is scarce in the literature. *Incomplete information* is a taxonomy used within the field of LA games to denote the case where the players do not observe the action of the opponent players, and where each player does not know his own payoff function but only observes outcomes in the form of a reward or a penalty. The informed reader would observe that the games we deal with in this paper falls under this class of games characterized by such incomplete information.

The case of incomplete information is not usually treated by the main stream of literature in game theory. Indeed, the main game learning algorithms available in the literature, such as fictitious play [35], best response dynamics, and gradient-based learning approaches, deal with the complete knowledge case, where the players know their own payoff function, and observe the history of the choices of other players. Fictitious play is one of the few algorithms that can converge to a mixed strategy equilibrium by maintaining various frequency-based beliefs over the action of the opponent players, and using those beliefs, for deciding the next action to be played. However, the fictitious play algorithm can not solve our settings of incomplete information.

When it comes to games with incomplete information, different algorithms have been suggested which are based on the Bush-Mosteller learning paradigm. Notable examples include the ones reported in [36], [37], [38], [39]. All those algorithms share a similar structure to our proposed LA, in particular, and to Variable Structure LA in general, in the

sense, that the action probabilities are updated iteratively based on feedback, and using some learning parameter. In this context, one should note that many LA models can be seen as extensions of Bush-Mosteller learning. However, the difference with our work is the fact that all the aforementioned algorithms have absorbing barriers. The theoretical analyses of the convergence to pure equilibria for this family of algorithms rely usually on the theory replicator dynamics.

Another family of methods that can operate with limited information include the Erev-Roth algorithm [40] and the Arthur algorithm [41], which in turn, can be seen as a variant of the Erev-Roth algorithm. The Erev-Roth algorithm is alternatively called the Erev-Roth payoff matching algorithm, and relies on updating the so-called “propensity” for each action, which is, loosely speaking, the cumulative payoff for that action. Thereafter, each action is played in a manner proportional to its corresponding relative propensity. The Erev-Roth algorithm is one of the few examples of limited-information game learning approaches that converge to unique mixed strategy equilibria. However, the Erev-Roth algorithm requires storing the *entire* history of rewards and penalties for each action. Furthermore, we have not been able to locate any research study that reports the analysis of the Erev-Roth algorithm for the case of our stochastic zero-sum game. We therefore opted to implement it for our game. Experimental results (not reported here, in the interest of not distracting from the main contribution of this paper) show that it neither converges to the desired equilibrium, nor does it possess consistent convergence results.

D. Real-life Application Scenarios

One Referee had requested a brief explanation of a complex environment, or different scenarios in a game, by which we could utilize our newly-proposed solution. We agree that providing an insightful discussion could be insightful for interested readers and active researchers. This, of course, can be open-ended, but to satisfy the Referee, we present the following brief example.

Our learning algorithm admits potential applications in many security games as well as in communication problems. The intersection between game theory and security is an emerging field of research. Algorithms that can converge to mixed equilibria are of great interest to the security community, because mixed equilibria are usually preferred over pure ones. In fact, randomization gives less predictive ability to the attacker to guess the deployed strategy of the defender [42]. For instance, let us take a repetitive game involving a jammer and a transmitter, which, in turn, constitute our players [43]. The jammer aims to disturb and block a communication between a transmitter and its associated receiver. The transmitter can choose the channel over which his message is communicated, while the jammer chooses a channel to attack. We suppose that the outcome is stochastic depending on the choice of the attacker (jammer) and defender (transmitter), and the stochastic characteristics of the channel. Both the jammer and transmitter can observe whether the attack was successful or not, and for instance, this common observation can be due to the receiver acknowledging the correct reception of the

message over a wireless channel that both the attacker and jammer can overhear. Thus, the game is stochastic zero-sum.

V. CONCLUSION

The theoretical applications Learning Automata (LA) with artificially *absorbing* barriers have been reported since the 1980s [8], and more recently, in Estimator LA [18], [19], [20]. This paper pioneers the study of LA with artificial non-absorbing barriers. LA have been previously used [13] to design algorithms capable of converging to the game’s Nash equilibrium under limited information. The majority of the LA algorithms used for games are absorbing in the probability simplex space, and they converge to an exclusive choice of a single action. These LA are, thus, unable to converge to other *mixed* Nash equilibria when the game possesses no Saddle Point for a pure strategy. As opposed to these, we propose a LA solution that is able to converge to an optimal mixed Nash equilibrium even though there may be no Saddle Point when a pure strategy is invoked. The scheme is inherently of the absorbing Linear Reward-Inaction (L_{R-I}) paradigm. However, by introducing reflecting barriers, we prevent it from being “stuck” or getting absorbed in pure strategies. Unlike the Linear Reward- ϵ Penalty ($L_{R-\epsilon P}$) scheme proposed in [1], our new scheme achieves the same goal with much less parameter tuning, and in a more elegant manner.

As far as know, our method is only the second reported algorithm in the literature capable of finding mixed strategies whenever no Saddle Point exists for pure strategies. If a Saddle Point exists for pure strategies, the scheme converges to a near-optimal solution close to the pure strategies in the probability simplex. The paper includes the non-trial proofs of the theoretical results characterizing the convergence and stability of the algorithm. These are presented and illustrated through simulations for benchmark games presented in the literature.

With regard to future work, we believe that it will be useful in real-life applications that can be modeled using such game-like behavior.

APPENDIX

Norman theorem

Theorem 4. *Let $X(t)$ be a stationary Markov process dependent on a constant parameter $\theta \in [0, 1]$. Each $X(t) \in I$, where I is a subset of the real line. Let $\Delta X(t) = X(t+1) - X(t)$. The following are assumed to hold:*

- 1) I is compact.
- 2) $E[\Delta X(t)|X(t) = y] = \theta w(y) + O(\theta^2)$
- 3) $Var[\Delta X(t)|X(t) = y] = \theta^2 s(y) + o(\theta^2)$
- 4) $E[\Delta X(t)^3|X(t) = y] = O(\theta^3)$ where $\sup_{y \in I} \frac{O(\theta^k)}{\theta^k} < \infty$ for $K = 2, 3$ and $\sup_{y \in I} \frac{o(\theta^2)}{\theta^2} \rightarrow 0$ as $\theta \rightarrow 0$.
- 5) $w(y)$ has a Lipschitz derivative in I .
- 6) $s(y)$ is Lipschitz I .

If Assumptions (1)-(6) hold, $w(y)$ has a unique root y^ in I and $\left. \frac{dw}{dy} \right|_{y=y^*} \leq 0$ then*

- 1) $\text{var}[\Delta X(t)|X(0) = x] = 0(\theta)$ uniformly for all $x \in I$ and $t \geq 0$. For any $x \in I$, the differential equation $\frac{dy(\tau)}{d\tau} = w(y(\tau))$ has a unique solution $y(\tau) = y(\tau, x)$ with $y(0) = x$ and $E[\delta X(t)|X(0) = x] = y(t\theta) + O(\theta)$ uniformly for all $x \in I$ and $t \geq 0$.
- 2) $\frac{X(t) - y(t\theta)}{\sqrt{\theta}}$ has a normal distribution with zero mean and finite variance as $\theta \rightarrow 0$ and $t\theta \rightarrow \infty$.

REFERENCES

- [1] S. Lakshmivarahan and K. S. Narendra, "Learning algorithms for two-person zero-sum stochastic games with incomplete information: A unified approach," *SIAM Journal on Control and Optimization*, vol. 20, no. 4, pp. 541–552, 1982.
- [2] S. Lakshmivarahan, *Learning Algorithms Theory and Applications: Theory and Applications*. Springer Science & Business Media, 2012.
- [3] Lanctot, J.K., Oommen, B.J.: On discretizing estimator-based learning algorithms. *IEEE Trans. on Systems, Man, and Cybernetics, Part B: Cybernetics* 2, 1417–1422 (1991)
- [4] Lanctot, J.K., Oommen, B.J.: Discretized estimator learning automata. *IEEE Trans. on Systems, Man, and Cybernetics, Part B: Cybernetics* 22(6), 1473–1483 (1992)
- [5] K. Najim and A. S. Poznyak, *Learning automata: theory and applications*. Elsevier, 2014.
- [6] K. S. Narendra and M. A. Thathachar, *Learning automata: an introduction*. Courier corporation, 2012.
- [7] M. F. Norman, *Markov processes and learning models*. Academic Press New York, 1972, vol. 84.
- [8] Oommen, B.J.: Absorbing and ergodic discretized two-action learning automata. *IEEE Transactions on Systems, Man, and Cybernetics* 16, 282–296 (1986)
- [9] Oommen, B.J., Lanctot, J.K.: Discretized pursuit learning automata. *IEEE Transactions on Systems, Man, and Cybernetics* 20, 931–938 (1990)
- [10] Oommen, B.J., Agache, M.: Continuous and discretized pursuit learning schemes: various algorithms and their comparison. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics* 31(3), 277–287 (2001)
- [11] G. Papavassilopoulos, "Learning algorithms for repeated bimatrix nash games with incomplete information," *Journal of optimization theory and applications*, vol. 62, no. 3, pp. 467–488, 1989.
- [12] A. Rezvanian, A. M. Saghiri, S. M. Vahidipour, M. Esnaashari, and M. R. Meybodi, *Recent Advances in Learning Automata*. Springer, 2018, vol. 754.
- [13] P. Sastry, V. Phansalkar, and M. Thathachar, "Decentralized learning of nash equilibria in multi-person stochastic games with incomplete information," *IEEE Transactions on systems, man, and cybernetics*, vol. 24, no. 5, pp. 769–777, 1994.
- [14] M. A. L. Thathachar and P. S. Sastry, *Networks of Learning Automata: Techniques for Online Stochastic Optimization*. Boston: Kluwer Academic, 2003.
- [15] M. L. Tsetlin et al., *Automaton theory and modeling of biological systems*. Academic Press New York, 1973.
- [16] B. Tung and L. Kleinrock, "Using finite state automata to produce self-optimization and self-control," *IEEE Trans. Parallel Distrib. Syst.*, vol. 7, no. 4, pp. 439–448, 1996.
- [17] Y. Xing and R. Chandramouli, "Stochastic learning solution for distributed discrete power control game in wireless data networks," *IEEE/ACM Transactions on networking*, vol. 16, no. 4, pp. 932–944, 2008.
- [18] Zhang, X., Granmo, O.C., Oommen, B.J., Jiao, L.: A formal proof of the ϵ -optimality of absorbing continuous pursuit algorithms using the theory of regular functions. *Applied Intelligence* 41, 974–985 (2014)
- [19] Zhang, X., Oommen, B.J., Granmo, O.C., Jiao, L.: A formal proof of the ϵ -optimality of discretized pursuit algorithms. *Applied Intelligence*, DOI 10.1007/s10489-015-0670-1 (2015).
- [20] Zhang, X., Oommen, B.J., Granmo, O.C.: The Design of Absorbing Bayesian Pursuit Algorithms and the Formal Analyses of their ϵ -Optimality. *Pattern Analysis and Applications* 20(3), (2015)
- [21] J. R. Marden, G. Arslan and J. S. Shamma: Cooperative Control and Potential Games. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 39, no. 6, pp. 1393-1407, Dec. 2009, doi: 10.1109/TSMCB.2009.2017273.
- [22] D. Silvestre, J. P. Hespanha and C. Silvestre, "Resilient Desynchronization for Decentralized Medium Access Control," in *IEEE Control Systems Letters*, vol. 5, no. 3, pp. 803-808, July 2021, doi: 10.1109/LC-SYS.2020.3005819.
- [23] D. Silvestre, J. Hespanha and C. Silvestre, "Desynchronization for Decentralized Medium Access Control based on Gauss-Seidel Iterations," 2019 American Control Conference (ACC), Philadelphia, PA, USA, 2019, pp. 4049-4054, doi: 10.23919/ACC.2019.8814471.
- [24] D. Silvestre, J. P. Hespanha and C. Silvestre, "Broadcast and Gossip Stochastic Average Consensus Algorithms in Directed Topologies," in *IEEE Transactions on Control of Network Systems*, vol. 6, no. 2, pp. 474–486, June 2019, doi: 10.1109/TCNS.2018.2839341.
- [25] D. Silvestre, J. Hespanha and C. Silvestre, "A PageRank Algorithm based on Asynchronous Gauss-Seidel Iterations," 2018 Annual American Control Conference (ACC), Milwaukee, WI, 2018, pp. 484-489, doi: 10.23919/ACC.2018.8431212.
- [26] R. Ribeiro, D. Silvestre and C. Silvestre, "A Rendezvous Algorithm for Multi-agent Systems in Disconnected Network Topologies," 2020 28th Mediterranean Conference on Control and Automation (MED), Saint-Raphaël, France, 2020, pp. 592-597, doi: 10.1109/MED48518.2020.9183093.
- [27] Ribeiro R., Silvestre D., Silvestre C. (2021) Decentralized Control for Multi-agent Missions Based on Flocking Rules. In: Gonçalves J.A., Braz-César M., Coelho J.P. (eds) CONTROLO 2020. CONTROLO 2020. Lecture Notes in Electrical Engineering, vol 695. Springer, Cham. https://doi.org/10.1007/978-3-030-58653-9_43.
- [28] Y. Wu and R. Lu, "Output Synchronization and L_2 -Gain Analysis for Network Systems," in *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 48, no. 12, pp. 2105-2114, Dec. 2018, doi: 10.1109/TSMC.2017.2754544.
- [29] Y. Wu, A. Isidori, R. Lu and H. K. Khalil, "Performance Recovery of Dynamic Feedback-Linearization Methods for Multivariable Nonlinear Systems," in *IEEE Transactions on Automatic Control*, vol. 65, no. 4, pp. 1365-1380, April 2020, doi: 10.1109/TAC.2019.2924176.
- [30] J. Flesch, A. Predtetchinski, and W. Sudderth, "Positive zero-sum stochastic games with countable state and action spaces," *Applied Mathematics & Optimization*, pp. 1–18, 2018.
- [31] C. Pal and S. Saha, "Continuous-time zero-sum stochastic game with stopping and control," *Operations Research Letters*, vol. 48, no. 6, pp. 715–719, 2020.
- [32] B. Ziliotto, "A tauberian theorem for nonexpansive operators and applications to zero-sum stochastic games," *Mathematics of Operations Research*, vol. 41, no. 4, pp. 1522–1534, 2016.
- [33] E. Boros, K. Elbassioni, V. Gurvich, and K. Makino, "A potential reduction algorithm for two-person zero-sum mean payoff stochastic games," *Dynamic Games and Applications*, vol. 8, no. 1, pp. 22–41, 2018.
- [34] K. Du, R. Song, Q. Wei, and B. Zhao, "A solution of two-person zero sum differential games with incomplete state information," in *International Symposium on Neural Networks*. Springer, 2019, pp. 434–443.
- [35] J. Hofbauer and W. H. Sandholm, "On the global convergence of stochastic fictitious play," *Econometrica*, vol. 70, no. 6, pp. 2265–2294, 2002.
- [36] T. Börgers and R. Sarin, "Learning through reinforcement and replicator dynamics," *Journal of economic theory*, vol. 77, no. 1, pp. 1–14, 1997.
- [37] L. R. Izquierdo, S. S. Izquierdo, N. M. Gotts, and J. G. Polhill, "Transient and asymptotic dynamics of reinforcement learning in games," *Games and Economic Behavior*, vol. 61, no. 2, pp. 259–276, 2007.
- [38] A. S. Poznyak and K. Najim, "Bush-mosteller learning for a zero-sum repeated game with random pay-offs," *International Journal of Systems Science*, vol. 32, no. 10, pp. 1251–1260, 2001.
- [39] Q. Zhu, H. Tembine, and T. Başar, "Heterogeneous learning in zero-sum stochastic games with incomplete information," in *49th IEEE conference on decision and control (CDC)*. IEEE, 2010, pp. 219–224.
- [40] I. Erev and A. E. Roth, "Multi-agent learning and the descriptive value of simple models," *Artificial Intelligence*, vol. 171, no. 7, pp. 423–428, 2007.
- [41] W. B. Arthur, "On designing economic agents that behave like human agents," *Journal of Evolutionary Economics*, vol. 3, no. 1, pp. 1–22, 1993.
- [42] M. H. Manshaei, Q. Zhu, T. Alpcan, T. Başçar, and J.-P. Hubaux, "Game theory meets network security and privacy," *ACM Computing Surveys (CSUR)*, vol. 45, no. 3, pp. 1–39, 2013.
- [43] V. Vadori, M. Scalabrin, A. V. Guglielmi, and L. Badia, "Jamming in underwater sensor networks as a bayesian zero-sum game with

position uncertainty,” in *2015 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2015, pp. 1–6.



Anis Yazidi received the M.Sc. and Ph.D. degrees from the University of Agder, Grimstad, Norway, in 2008 and 2012, respectively. He was a Researcher with Teknova AS, Grimstad, Norway. From 2014 til 2019 he was an associate professor with the Department of Computer Science, Oslo Metropolitan University, Oslo, Norway. He is currently a Full Professor with the same department where he is leading the research group in Applied Artificial Intelligence. He is also Professor II with the Norwegian University of Science and Technology (NTNU),

Trondheim, Norway. His current research interests include machine learning, learning automata, stochastic optimization, and autonomous computing.



Daniel Silvestre received his B.Sc. in Computer Networks in 2008 from the Instituto Superior Técnico (IST), Lisbon, Portugal, and an M.Sc. in Advanced Computing in 2009 from the Imperial College London, London, United Kingdom. In 2017, he got his Ph.D. (with the highest honors) in Electrical and Computer Engineering from the former university, and was a visiting scholar at the University of California at Santa Barbara. Currently, Dr. Silvestre is with the Institute for Systems and Robotics, at the Instituto Superior Técnico in Lisbon

(PT), and he holds a research assistant appointment with the University of Macau. His research interests span the fields of fault detection and isolation, distributed systems, network control systems, computer networks, set-valued estimation and control methods, and randomized algorithms.



B. John Oommen was born in India in 1953. He received his Bachelor of Technology in Electrical Engineering at the Indian Institute of Technology in Madras, India in 1975. He then pursued his Master of Engineering degree at the Indian Institute of Science in Bangalore, India receiving his degree in 1977. At both these institutions, he won the medal for being the Best Graduating Student. He received a Master of Science degree in 1979, and a Ph.D. in Electrical Engineering in 1982, both from Purdue University, Indiana, USA.

In 2003, Dr. Oommen was nominated as a Fellow of the Institute of Electrical and Electronic Engineers (IEEE) for research in a subfield of Artificial Intelligence, namely in Learning Automata. He is currently a Life Fellow of the IEEE. He was also nominated as a Fellow of the International Association of Pattern Recognition (IAPR) in August 2006 for contributions to fundamental and applied problems in Syntactic and Statistical Pattern Recognition. He has served on the editorial board of the journals *IEEE Transactions on Systems, Man and Cybernetics*, and *Pattern Recognition*. Dr. Oommen has been teaching in the School of Computer Science at Carleton University since 1981, and was elevated to be a Chancellor’s Professor at Carleton University in 2006. He has published more than 485 refereed publications, many of which have been award-winning. He has also won Carleton University’s Research Achievement Award four times, in 1995, 2001, 2007 and 2015.