

Learning from Induced Changes in Opponent (Re)Actions in Multi-Agent Games

P.J. 't Hoen
hoen@cw.nl

S.M. Bohte
sbohte@cw.nl

J.A. La Poutré
hlp@cw.nl

CWI, The Netherlands Centre for Mathematics and Computer Science
Kruislaan 413, NL-1098 SJ Amsterdam, The Netherlands

ABSTRACT

Multi-agent learning is a growing area of research. An important topic is to formulate how an agent can learn a good policy in the face of adaptive, competitive opponents. Most research has focused on extensions of single agent learning techniques originally designed for agents in more static environments. These techniques however fail to incorporate a notion of the effect of own previous actions on the development of the policy of the other agents in the system. We argue that incorporation of this property is beneficial in competitive settings. In this paper, we present a novel algorithm to capture this notion, and present experimental results to validate our claims.

General Terms

Multi-Agent Learning

Keywords

Multi-Agent RL, Opponent Modeling, Learning Reactions to Actions

1. INTRODUCTION

Acting intelligently in a dynamic environment shared with other competing agents is a hard task, as daily life as well as a host of work on learning in multi-agent settings demonstrates. Without the presence of a teacher, and with only the observed actions and rewards to learn from, the framework of Reinforcement Learning (RL) (Sutton & Barto, 1998) is an obvious choice. Conceptually, the focus moves from the realm of Markov Decision Problems to that of Game Theory (GT) and stochastic games (Shoham, Powers, & Grenager, 2004).

Recent work has proposed a number of extensions from the single agent setting to the multi-agent domain. State-of-the-art Multi-Agent Reinforcement Learning (MARL) algorithms improve on single-agent RL approaches by incorporating models of the *current* opponent agent's behavior.

This allows an agent to model the current "environment" including opponents, against which it has to optimize its own behavior by playing a best-response. In a multi-agent setting, the policy that maximizes payoff will depend on the *changing* actions of adversaries. Importantly, the adaptations of adversaries are likely to be reactions to *ones own* actions. This point is largely ignored in current (MARL) algorithms. Such strategic reasoning can be paramount when an agent is repeatedly interacting with the same opponents, and earlier actions influence the future behavior of these learning adversaries.

The Prisoner's Dilemma (PD) is illustrative in that in repeated play the optimal policy differs from blindly playing the myopic best response, and that it can be profitable to have a good estimation of the opponents reactions to ones own play. In one single game, two competing agents can each choose from actions cooperate or defect ($\{C,D\}$), and the maximum *joint* payoff is achieved when both agents choose to Cooperate. However, when the other agents plays Cooperate, an agent obtains an even larger payoff from playing Defect. From a game theoretical perspective then, for the single shot game $\{D,D\}$ is the dominant strategy and a Nash-Equilibrium (see Section 2).

In iterated play of the PD (iPD) game, there is no longer a clear dominant strategy, as both agents can achieve a higher aggregated *and* individual reward by cooperating and playing the joint action $\{C,C\}$, *provided* there is some strategic incentive to rarely unilaterally defect. That is: defecting as a strategy is discouraged by the likely *negative future* reaction of the opponent (see also the famous Folk Theorem discussed in Section 2).

It is important to notice that just modeling the *current* behavior of the opponent in iPD can lead to $\{D,D\}$ outcome in repeated play when the reaction of the opponent to ones own play of defection is not taken into account. An agent that does not take into account the changing behavior of the opponent may estimate that a cooperating opponent may be exploited, and ignores the fact that defection will result in an opponent that also defects. When repeatedly playing a game against an adversary, the question is which moves one should play to maximize the individual reward given the *dependent* adaptive behavior of the adversary.

We argue that current state-of-the-art MARL algorithms (M. H. Bowling & Veloso, 2002; M. Bowling, 2004; Tesauro, 2003; Conitzer & Sandholm, 2003; Greenwald & Hall, 2003; Banerjee & Peng, 2004; Weinberg & Rosenschein, 2004) (see also Shoham et al. (2004) for a more in depth discussion) do not incorporate a sufficient notion of the longer term im-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

AAMAS'06 May 8–12 2006, Hakodate, Hokkaido, Japan.
Copyright 2006 ACM 1-59593-303-4/06/0005 ...\$5.00.

part of their actions on the dynamics of the environment. This environment includes other adaptive agents that react to moves of their opponents as well. For example, current MARL algorithms that play the iPD will universally converge to the worst possible outcome of $\{D,D\}$ in self-play¹ (Crandall & Goodrich, 2005). The one exception to the full defection converge outcome in repeated play of the iPD we are aware of is a recent paper by Crandall and Goodrich (2005), which we will discuss further below.

In this paper, we present a novel MARL framework called Strategic Opponent Policy Modeling (StrOPM) that goes beyond current Myopic Opponent Policy Modeling (MOPM) MARL approaches. StrOPM, as most other MARL’s, takes into account the current (estimated) policy of the opponents and the rewards for actions. Importantly, our algorithm also estimates how the policies of the opponents change over time due to actions played by StrOPM. The policy of the agent using the StrOPM algorithm is adjusted to increase the expected future reward by taking into consideration 1) the immediate and future rewards of actions 2) the estimated policy of the opponent, and 3), most importantly, the impact of the chosen actions on the future policy of the opponent. The last point is the novel contribution of our work.

Our algorithm is a policy-gradient like method with roots in Q-learning. Q-like values for states are calculated based on estimated values of actions and estimated opponent policies. Policies are updated along the gradient of expected increased rewards. To more accurately model the (future) changes in the environment due to an agent’s own actions, the agent computes a continuously updated estimate of the *change* in the choice of opponents actions due to its own actions. StrOPM then adjusts the reward gradient by this estimate to account for the opponent’s likely reaction.

We present our results for the well studied domain of iterated-play two-player two-action matrix games. We show that our algorithm quickly arrives at the desired Nash-Equilibrium for games that are unproblematic. These games are relatively straightforward in the sense that optimal play in a single-shot play of the game is also a good solution to iterated play of the game.

Importantly, we show that the proposed StrOPM algorithm achieves cooperation ($\{C,C\}$) in repeated play of the PD. After an initial exploration period, StrOPM learns that defection leads to reciprocal defection, and an overall lower value of future play than a high level of cooperation. The agents learn to optimize their own interest as a function of the played game and the reciprocal behavior of the evolving opponent. The (possible) cooperation strategy is not a precoded option in the algorithm, but is achieved by optimizing beyond the immediate best response to the current opponent play. The StrOPM agents in self-play exhibit a learned Tit-For-Tat strategy (Axelrod, 1984). The Tit-For-Tat algorithm plays C until the opponent defects upon which a D is played followed by again C until the next defection to encourage cooperation. The StrOPM agents learn this strategy and cooperates to reach a good equilibrium, and yet also evolves a threat state where defection is returned in order to guard against exploitation by the opponent.

We also compare our approach to the recent algorithm M-Qubed of Crandall and Goodrich (2005) for the studied matrix games. M-Qubed achieves the (C,C) equilibrium in

the iPD by a meta step. The M-Qubed algorithm is designed to balance play of a best-response strategy with the precoded strategy of playing one action consistently. In particular, the algorithm during play continuously considers the value of switching to the pure, precoded strategy of playing C . Using a precoded strategy may lead to being exploited, or be unsuitable for more complex settings, as we will show.

We present results for StrOPM versus M-Qubed and show that the two type of agents playing against each other can together achieve cooperation, but from entirely different principles. The StrOPM algorithm in self-play finds good strategies significantly faster than M-Qubed for the iPD game. When playing against each other, this high speed of learning leads to an initial period where StrOPM exploits M-Qubed. Furthermore, StrOPM is shown to exploit the M-Qubed algorithm in a simple matrix game where the pre-coded strategy of playing one action predictably can result in inferior strategies in repeated play. Using precoded strategies has limited application. For example, precoded strategies cannot be used when the number of possible precoded strategies must be large, the precoded strategies cannot be determined, or do not generalize to novel, or more complex settings. Fundamentally, a MARL should remain flexible in novel settings.

In summary, we have presented three important contributions. First of all, we have signaled the need for Multi-Agent learning algorithms to take into account the impact of their actions on the learned policy of the opponents. This concept is as yet lacking in the MARL community. Secondly, we presented a novel framework that allows agents to apply the above observation to achieve desired outcomes for representative iterated matrix games, and, more importantly, also for the iPD. The latter is known to be difficult for MARL approaches; the presented StrOPM framework solves the iPD game by learning and recognizing “threat states”, where it will punish an opponent’s uncooperative behavior with a Tit-for-Tat type strategy (and also learns the reciprocal reaction). Lastly, we demonstrate the danger of incorporating pre-coded fixed strategies in learning algorithms as these are vulnerable to exploitation or may not be suited to novel or more complex settings.

Structure. The remainder of this document is structured as follows. In Section 2, we present iterated matrix games, and the iterated Prisoners Dilemma. In Section 3, we present our general framework. We present a Multi-Agent RL framework, StrOPM, that in the limit plays best-response against stationary opponents. Importantly, StrOPM can reason strategically in games where it is relevant in how an opponent will react to the history of play. In Section 4, we present experimental results. In Section 5, we discuss and conclude.

2. MODEL AND PROBLEM SETTING

In this section, we give some introductory definitions and notation from Game Theory, Reinforcement Learning, and the iterated matrix games we study as problem domain. We discuss the iterated Prisoners Dilemma as a game of special interest.

2.1 Agents and Matrix games

We consider multi-agent Reinforcement Learning for iterated play of games, more specifically **matrix games**, and introduce some well-known concepts from Game Theory (GT).

¹An agent playing against an identical opponent.

In general, let S denote the set of states in the game and let A_i denote the set of actions that agent/player i may select in each state $s \in S$. Let $a = (a_1, a_2, \dots, a_n)$, where $a_i \in A_i$ be a joint action for n agents, and let $A = A_1 \times \dots \times A_n$ be the set of possible joint actions. **Zero-sum games** are games where the rewards of the agents for each joint action sum to zero. **General sum games** allow for any sum of values for the reward of a joint action.

A **strategy (or policy)** for agent i is a probability distribution $\pi(\cdot)$ over its actions set A_i . Let $\pi(S)$ denote a strategy over all states $s \in S$ and let $\pi(s)$ (or π_i) denote a strategy in a single state s . A strategy may be a **pure strategy** (an agent selects an action deterministically) or according to a **mixed strategy** (a strategy that plays a random action, according a probability distribution). A **joint strategy** played by n agents is denoted by $\pi = (\pi_1, \dots, \pi_n)$. Let a_{-i} and π_{-i} refer to the joint action and strategy of all agents except agent i .

We focus on the more restricted **matrix game**, defined by a set of matrices $R = \{R_1, \dots, R_n\}$. Let $R(\pi) = (R_1(\pi), \dots, R_n(\pi))$ be a vector of expected payoffs when the joint strategy π is played. Also, let $R_i(\pi_i, \pi_{-i})$ be the expected payoff to agent i when it plays strategy π_i and the other agents play π_{-i} . A strategy is **dominant** if, regardless of what any other players do, the strategy earns a player a larger payoff than any other strategy. Also, let $R_i\left(\begin{smallmatrix} a_i \\ a_{-i} \end{smallmatrix}\right)$ be the payoff for agent i playing action a_i while the other agents play action a_{-i} .

A **stage game** is a single iteration of a matrix game, and a **repeated game** is the indefinite repetition of the stage game between the same agents. While matrix games do not have state, agents can encode the previous w joint actions taken by the agents as state information, as for example illustrated by Sandholm and Crites (1995).

Each individual matrix game has certain classic game theoretic values. The **minimax** value for player i is $m_i = \max_{\pi_i} \min_{a_{-i}} R_i(\pi_i, a_{-i})$, i.e. the least reward that can be achieved if the game is known and the game is only played once. A **Best-Response (BR)** to the opponents strategy π_{-i} is defined by $BR = \pi^* = \max_{\pi} R_i(\pi, \pi_{-i})$. This corresponds to the (expected) most reward that can be gained from playing, under the assumption that the game is known; the game is only played once; and the opponent strategy is known.

A **Nash Equilibrium (NE)** is a joint strategy such that no agent may unilaterally change its strategy without lowering its expected payoff in the one shot play of the game. Nash (1951) showed that every n -player matrix game has at least one such NE. A **Pareto optimal** solution of the game is a joint strategy such that no agent may unilaterally increase its expected payoff without making another agent worse off. A (joint) strategy π_1 is said to **Pareto dominate** a strategy π_2 if the expected payoff for π_1 is at least as high as for π_2 and higher for at least one of the agents. A joint strategy is **Pareto deficient** if it is not Pareto optimal.

In the studied matrix games, we assume that an agent can observe its own payoffs as well as the actions taken by all agents in each stage game, but only after the fact. All agents concurrently choose their actions. A possible adaptation of the agents' policy, i.e. learning as a result of observed opponent behavior, only takes effect in the next stage game. Repeated games are modeled as a series of stage games with

	R=0.35	T=0.5
R=0.35		S=0
T=0.5	S=0	P=0.1

Table 1: Example payoffs for the (symmetrical) PD

the same opponent(s). Each agent then aims to maximize its reward from iterated play of the same matrix game.

We restrict our investigation to two-player, two-action games as these are well classified (Rapoport, Guyer, & Gordon, 1976). The arguments and formalism presented in the rest of the paper are however applicable to more general settings. The next section discusses the Prisoners Dilemma (PD), a special two-player two-action game.

2.2 The iterated Prisoner's Dilemma

We present the iterated **Prisoners Dilemma (iPD)** as a special matrix game. In this game, each player has a choice of two operations: either **cooperate (C)** with the other player or **defect (D)**. The payoff matrix for joint actions is shown in Table 1. If both players cooperate, they both receive a given payoff R . However, if one player plays Cooperate, and the other plays Defect, the defector receives a payoff $T > R$, and the cooperator receives a much lower payoff $S < R$. If both players play Defect, they receive the low payoff P , that is however better than the payoff S ($P > S$). For each individual, the incentive is thus present to defect, hoping that the other player plays "cooperate". Myopic play by even one of the players already quickly leads both players to arrive at the suboptimal outcome of both players receiving a low reward P . Furthermore, the PD has as constraint that $R > \frac{(T+S)}{2}$; full cooperation is better than playing $\begin{bmatrix} C \\ D \end{bmatrix}$ and then $\begin{bmatrix} D \\ C \end{bmatrix}$ for the agents over two stage games.

The PD game is of particular interest as it has specific properties that make it unique among the matrix games. The PD is the only game of the taxonomy of two-player, two-action games listed by Rapoport et al. (1976) for which the natural outcome is stable (a NE), but Pareto-deficient: there are outcomes that Pareto-dominate the NE. This unique property also makes it non-trivial for competitive learning algorithms in repeated play; how can learners consider the possible future gains that are non-existent in the single shot game?

Most multi-agent learning algorithms to date have focused on an individual agent learning a (myopic) Best Response to the *current* strategies of the other agents. Play between such agents using this approach often converge, and have as goal to converge, to a one-shot NE. However, a famous result from game theory (the **folk theorem**) (Myerson, 1991) suggests that the goal of reaching a one-shot NE may be inappropriate in repeated games.

The folk theorem implies that, in many games, there exists NEs for repeated games, repeated Nash-Equilibria (rNEs), that yield higher individual payoffs to all agents than do one-shot NEs, i.e. the rNE Pareto dominates the NE. Hence, in repeated games, a successful set of agents should learn to play profitable rNEs. However, since many repeated games have an infinite number of rNEs, the folk theorem does little

to indicate which one the agents should play. Littman and Stone (2003) present an algorithm for computing rNEs that satisfies a set of desiderata, but how to learn these strategy online is unknown. Additionally, an agent may have preferences between rNEs and play one above the other, if allowed by its opponents.

It is however not given that the mainstream, or state of the art RL algorithms, will (approximately) learn these rNE equilibria in self play or against various classes of opponents. As discussed in Section 1, and by Crandall and Goodrich (2005), the claim is that current MARL algorithms will not converge to good rNE for iPD-like games.

3. THE StrOPM FRAMEWORK

Here, we refine the StrOPM framework. We show (very successful) experimental results with an instantiation of the framework for matrix games in Section 4.

We can formulate the goal of an agent i as wanting to maximize the total reward it obtains from playing T times against an adversary:

$$R_i^{tot} = \sum_{t=0}^{t=T} val_i(\pi_{i,t}, \pi_{-i,t}). \quad (1)$$

where π_t is the policy at time t and val_i the expected reward for a agent i using policy $\pi_{i,t}$ and opponent policies $\pi_{-i,t}$. Thus, we want to specify or learn the time-varying policy $\pi_{i,t}$ that maximizes the reward obtained.

We make several observations: first, we have to start playing from some policy $\pi_{i,t=0}$, and we may not know the payoff of (joint) actions yet. We may also not know the policy of the adversary yet, nor to what extent it is adaptive.

We assume that we can define a state representation for an agent i that includes sufficient history to make the evolution of both its own policy and the opponent’s policy Markovian: $\pi_{i,t+1} = f(s(t), \pi_{i,t})$, where $f(s(t), \pi_{i,t})$ is the state dependent function that adapts the agent’s policy based on the agent’s current policy $\pi_{i,t}$ and the current state $s(t)$. Likewise, the agent assumes that the opponent behaves similarly Markovian: $\pi_{-i,t+1} = g(s(t), \pi_{-i,t})$, where $g(s(t), \pi_{-i,t})$ is an unknown function that adapts the adversarie’s policy. We assume that part of the state information in the $g(s(t))$ term depends on the actions taken by the agent i .

At its core, the StrOPM algorithm attempts to estimate this unknown adaptive function $g(s(t), \pi_{-i,t})$, and then computes forward all possible future action sequences resulting from the thus evolving policies $\pi_{i,t}$ and $\pi_{-i,t}$. It then updates the policy $\pi_{i,(t+1)}$ in the direction of those future action sequences that promise the most reward.

In the StrOPM framework, an agent applies reinforcement learning to a state-based policy as described in Algorithm 1. At each epoch of learning, the agent adapts its policy along the gradient of increasing reward. The gradient of reward is calculated including the expected changing behavior of the opponent, as a reaction of the agent’s own actions. The StrOPM algorithm tracks the changes in observed opponent policy, on a state by state basis, over time. It is assumed that these changes in the opponent behavior, at least in part, reflect reaction to actions chosen by the StrOPM algorithm. The StrOPM algorithm not only tracks the changing behavior of the opponents, but also how these changes are potentially caused by the agent’s own action. The policy of the StrOPM is then optimized with expected future reac-

tions taken into account.

We describe the StrOPM algorithm from the perspectives of an agent i and its opponents, agents $-i$. We use this notation with subscripts to indicate policy, states, action, etc . . . of the two types of agents. E.g., a_i and a_{-i} are actions of agents i and $-i$ respectively.

States. The set of states that an agent i can visit, S_i , fulfills the unichain assumption (Puterman, 1994): One set of “recurrent” class of states. Starting from any state in the class, the probability of visiting all the states in the class is 1. We introduce a **transition function** $T : S_i \times A_i \times A_{-i} \rightarrow S_i$ to return the next state of agent i upon playing a (joint) action from the current state.

Policies. For a state $s \in S_i$, $\pi_i(s)$ is the state-based policy of agent i , and $\pi_{-i}(s)$ is the estimate of the opponent $-i$ policies when agent i is in state s . The opponent policy is estimated online using Exponential Moving Average (EMA). The estimate of $\pi_{-i}(s)$ after observing action a_{-i} is adjusted according, for action a_i :

$$\pi_{-i,t+1}(s)(a_{-i}) = (1 - \alpha_{EMAI})\pi_{-i,t}(s)(a_{-i}) + \alpha_{EMAI}. \quad (2)$$

After this update, the policy $\pi_{-i,t+1}(s)$ is normalized to retain $\pi_{-i}(s)$ as a probability distribution by equally reducing the probabilities of the other actions to retain $\pi_{-i}(s)$ as a probability distribution.

Actions to update the policy. We introduce **policy update actions**. A policy update action $pua_i(s)$ dictates whether the likelihood of an action a_i should be increased for state s . Additionally, we introduce the *null* policy update action $pua_{null}(s)$ to indicate that the policy should not be changed.

Let $\pi_i(s)^{pua_i}$ be the policy achieved by applying the policy update action $pua_i(s)$ to $\pi(s)_i$. The policy $\pi(s)_i$ of agent i given $pua_i(s)$ not equal to the null action is updated according to:

$$\pi_{i,t+1}(s)^{pua_i} = (1 - \alpha_{LEARN})\pi_{i,t}(s)(a_i) + \alpha_{LEARN}, \quad (3)$$

where α_{LEARN} is the learning rate. The probabilities $\pi_i(s)(\cdot)$ for actions $a_j \neq a_i$ are then normalized to retain $\pi_i(s)$ as a probability distribution.

We wish to select policy update action $pua_i(s)^*$ in (3) that maximizes a measure of expected future payoff of the changed policy. Below, we explain how to compute this.

Estimation of impact actions. We introduce $\xi(\pi_{-i}(s), a_i)$ to estimate the impact of an agent i playing an action a_i on the development of the policies $\pi_{-i}(s)$ of the opponent. This function predicts the change in policy of the opponents upon playing action a_i ; i.e.

$$\pi_{-i,t+1}(s) = \xi(\pi_{-i,t}(s), a_i). \quad (4)$$

As a first implementation of this function, we take the approach that the changes in the opponent policy are, at least in part, caused by an agent’s own actions. Our StrOPM implementation estimates the change in policy as a continuation of the change in policy from estimation of the opponent policy N epochs in the past, i.e. how was the opponent policy at time $t - N$? This is done for each state: for state s reached after playing action joint action the change in policy for this new state is estimated as a linear extrapolation of

the change in the estimated opponent policy:

$$\xi(\pi_{-i,t}(s), a_i) = \frac{\pi_{-i,t}(s) - \pi_{-i,t-N}(s)}{N} + \pi_{-i,t}(s), \quad (5)$$

where we limit ξ to $[0, 1]$.

Thus estimating the change in opponent policies, the question becomes how to compute the value of proposed policy updates. Here, we make use of the unichain assumption that says we only need to consider loops starting from and returning to the current state to compute this value. The algorithm then updates the policy for the possible loops to increase the expected average reward.

Loops. From a state s_0 , a single loop $L(s_0)$ is defined as:

$$L(s_0) = s_0, \begin{bmatrix} a_{i,0} \\ a_{-i,0} \end{bmatrix}, s_1, \begin{bmatrix} a_{i,1} \\ a_{-i,1} \end{bmatrix}, s_2, \dots, s_{(n-1)}, \begin{bmatrix} a_{i,(n-1)} \\ a_{-i,(n-1)} \end{bmatrix}, s_0, \quad (6)$$

for a sequence of joint moves for agent i starting in s_0 , going through s_1, s_2, \dots to s_n and ending again in s_0 . Each next state reached is through a joint move; $T(s_j, \begin{bmatrix} a_{i,j} \\ a_{-i,j} \end{bmatrix}) = s_{j+1}$.

All intermediate states reached in the loop are unique, and not equal to s_0 . For brevity, we denote the j -th state s_j in the sequence by $L(s_0)^j$, the j -th action pair $\begin{bmatrix} a_{i,j} \\ a_{-i,j} \end{bmatrix}$ by $L(s_0)_j$, and $L(s_0)_{j+}$ and $L(s_0)_{j-}$ the components of the j -th action pair: $a_{i,j}$ and $a_{-i,j}$. The length of the sequence $L(s_0)$, $|L(s_0)|$, is equal to n ; the number of joint actions played before the considered state is reached again.

Sets of loops. Let $Bag(s, n) = \{L(s) \mid |L(s)| \leq n\}$, i.e. $Bag(s, n)$ are all the loops starting in state s with length of at most n . The probability of a particular loop $L(s)$ occurring, denoted by $Pr(L(s))$ is:

$$Pr(L(s)) = \prod_{0 \leq j < |L|} Pr(L(s)_{j+} \mid \pi_{i,j}(L(s)^j)) \times Pr(L(s)_{j-} \mid \pi_{-i,j}(L(s)^j)), \quad (7)$$

where $L(s)_j$ and $L(s)^j$ are the respective elements in the loop as defined above, $\pi_{i,j}$ and $\pi_{-i,j}$ are the respective policies at time j , evolving according to Equation (4), and $Pr(\cdot)$ denotes the probability of a specific transition along the sequence given the respective policies $\pi(s)$.

The expected reward over the possible loops can then be expressed as the weighted expected value of individual loops:

$$E(Bag(s, n)) = \sum_{L(s) \in Bag(s, n)} Pr(L(s)) \times E(L(s)), \quad (8)$$

where $E(L(s))$ is the expected average reward for each joint action of a loop starting in state s :

$$E(L(s)) = \frac{\sum_{0 \leq j < |L(s)|} V(L(s)_j)}{|L(s)|}, \quad (9)$$

where $V(L(s)_j)$ denotes the estimated the value of the j -th joint action in the loop, $L(s)_j$, to agent i .

Learning the value of joint actions. The values $V_i : A_i \times A_{-i} \rightarrow \mathfrak{R}$ track the estimates the value of a single joint action to agent i and through V_i , agent i learns its part of

the rewards of the game, dependent upon the actions of the opponent. The function V_i is updated using EMA:

$$V_{i,t+1} = (1 - \alpha_{EMA2})V_{i,t} + \alpha_{EMA2} \times r_t, \quad (10)$$

where r_t is the reward received by agent i in epoch t , and $0 < \alpha_{EMA2} \leq 1$ is the learning rate for the update.

Strategic Updating. Let the opponent changes in policy be estimated by $\xi(\pi_{-i}(s), a_i)$, we can then compute the policy update with highest expected payoff:

$$pua_i(s)^* = \max_i E(Bag(s, n))_{pua_i(s)}^\xi, \quad (11)$$

where $E(Bag(s, n))_{pua_i(s)}^\xi$ denotes the expected average reward for the possible loops of at most length n , starting in state s , after effecting policy update action $pua_i(s)$ and taking $\xi(\pi_{-i}(s), a_i)$ as the estimated opponent adaptation. Our StrOPM algorithm thus obtained is outlined in Algorithm 1.

Algorithm 1 StrOPM

- 1: Initialize $\pi_i(s)$, $\pi_{-i}(s)$ for all $s \in S_i$ and V for all all joint actions. Set the initial state.
 - 2: Do in each epoch sequentially for each agent i in state s :
 - 3: loop
 - 4: calculate the highest valued $pua_i(s)^*$ using $E(Bag(s, n))_{pua_i(s)}^\xi$ as value for the individual policy update actions.
 - 5: Update policy $\pi_i(s)$ using the highest valued policy update $pua_i(s)^*$ action.
 - 6: Play action $a_i \in A_i$ based on the chosen policy update action $pua_i(s)^*$. StrOPM plays action a_i for chosen policy update action $pua_i(s) = pua_i(s)^*$ if pua_i^* is not the null policy update action. Otherwise choose the action according to $\pi_i(s)$.
 - 7: Receive reward $R_i(\begin{bmatrix} a_i \\ a_{-i} \end{bmatrix})$ for the joint action determined by agents $-i$.
 - 8: Update the estimate of the reward for the joint action $V(\begin{bmatrix} a_i \\ a_{-i} \end{bmatrix})$ and the opponent policy π_{-i} .
 - 9: set the current state to $T(s, \begin{bmatrix} a_i \\ a_{-i} \end{bmatrix})$.
 - 10: end loop
-

4. EXPERIMENTS

We first investigate the StrOPM algorithm for various well-known two-player, two-action games.

In Figures 1b,d,f, we show the results for StrOPM in a variety of games. Figure 1b shows the results for self-play for the game of ‘‘matching pennies’’, Figure 1d for a variant of a coordination game², and Figure 1f for the game of Chicken (see also Rapoport et al. (1976)). The corresponding payoff matrices are plotted above the respective graphs in Figures 1a,c,e. Experiments are averaged over 50 runs. A learning

²In the simple coordination game, the two agents must learn to either play $\begin{bmatrix} C \\ D \end{bmatrix}$ or $\begin{bmatrix} D \\ C \end{bmatrix}$ to achieve the highest reward.

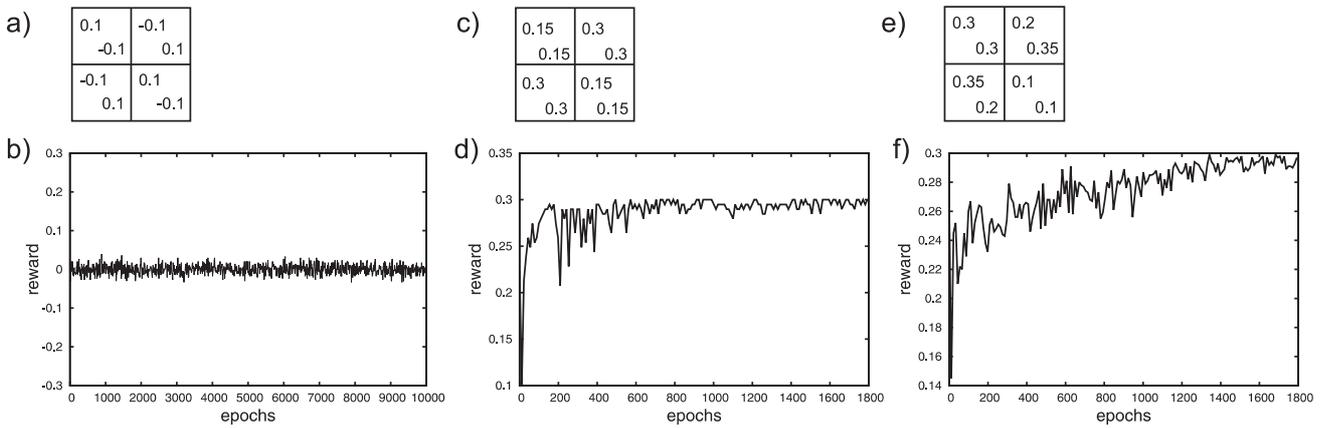


Figure 1: a) payoff matrix for “Matching Pennies”, reward representation as in Table 1. b) two StrOPM players playing Matching Pennies. c,d) payoff matrix and two StrOPM players playing a Coordination Game. e,f) payoff matrix and two StrOPM players playing a game of Chicken.

rate of 0.01 was used for all the EMA equations of Section 3. The StrOPM algorithm looks back $N = 10$ epochs in Equation 5. Additionally, we added an $\epsilon = 0.01$ probability of the StrOPM agent taking an exploration move in each epoch to ensure all states are sufficiently sampled in play. Note that we reuse the payoff labels C and D for ease of exposition.

In the games shown in Figure 1, the StrOPM algorithm finds the optimal Nash-Equilibria rapidly, as was to be expected since these games are unproblematic in the sense that the Nash-Equilibria of the games is also a good solution for iterated play of the game. The results show that the StrOPM algorithm finds good solutions for these diverse games quickly, as it should since it is not designed with a bias for any particular matrix game.

StrOPM in the iPD. Of greater interest than the previous games is the iterated Prisoner’s Dilemma, as it is a much harder, open problem for MARLs and is also the motivating example for the work of Crandall and Goodrich (2005).

In Figure 2a, we present results of the StrOPM algorithm for self-play in the iPD game of Figure 2a (using the payoff matrix of Table 1). To highlight the difference between a forward-thinking algorithm like StrOPM, and more reactive MOPM algorithms, we compare in Figure 1 the results for the StrOPM algorithm in self-play with that of self-play for an MOPM player in the form of a reduced variant of StrOPM that does not anticipate reactions to its actions. With the latter, we mean that the used ξ functions do not predict a change in policy of the opponent, i.e. $\xi(\pi_{-i}, a_i) = \pi_{-i}$. As such, the StrOPM implementation reverts to an opponent modeler type player that simply moves its policy to a best-response to the recently observed play of the opponent.

The StrOPM learner in self-play (Figure 2a, solid line) converges to playing the (optimal) Cooperate-Cooperate ($\{C, C\}$) strategy with reward 0.35. The full cooperation equilibrium is reached as the StrOPM learner estimates that leaving the full cooperation state leads to states where more and more defection is expected. Additionally, the full cooperation state is expected to lead to more cooperation. In contrast, the MOPM variant quickly learns cooperation is

risky and moves to full defection (Figure 2a, dashed line).

Furthermore, a closer study of the state-based policy of the StrOPM players reveals that the agents in self-play repeatedly exhibited a learned Tit-For-Tat strategy (Axelrod, 1984). The Tit-For-Tat algorithm plays C until the opponent defects upon which a D is played followed by again C until the next defection to encourage cooperation. For a four state agent, with each state encoding the last joint action of play, the policy is to play C in state $\begin{bmatrix} C \\ C \end{bmatrix}$, play D from state $\begin{bmatrix} C \\ D \end{bmatrix}$, and C from $\begin{bmatrix} D \\ D \end{bmatrix}$. The StrOPM agents learn this strategy and play C to cooperate and reach a good equilibrium, except for occasional exploratory moves. At the same time the StrOPM players also evolve a “threat” state where defection is retaliated in order to guard against exploitation by the opponent. The need, and the possibility, for learning a threat state as in the Tit-For-Tat play of the iPD is an interesting venue of research for iterated play of games. For example, a StrOPM player playing against a hard-coded Tit-For-Tat player quickly learned to cooperate (not shown), but did not evolve such a threat state.

As a negative result, a StrOPM agent playing against a MOPM agent converged to mutual defection (not shown) as the MOPM player learned to play defect. The MOPM agent playing against a hard coded Tit-For-Tat player was however able to achieve full cooperation as the MOPM player was consistently taught by the Tit-For-Tat agent that defection was not worthwhile. This strategy was however not yet available to the StrOPM agent as it still had to be learned. This indicates that players in the iPD must both be reasonably savvy, be that through their learning algorithm or a-priori imparted strategies, for mutual cooperation to emerge.

Ongoing preliminary results (not shown here) however indicate that generalised StrOPM-type agents are able to achieve high levels of cooperation in the generalised n-player iPD, the nIPD (Yao & Darwen, 1994). For example, for three agents, full cooperation is reached in the 3IPD game. In (’t Hoen & La Poutré, 2005), we have shown that agents participating in a sequence of auctions for individual items with as goal to win a bundle of these items can face a nIPD-like problem. Each of the agents has an individual incentive

to bid strategically. The prices paid for the items are however much higher if all the agents follow this strategy. A venue of research is to investigate whether the concepts of the StrOPM algorithm can be applied to this more complex setting to arrive at more beneficial outcomes from the perspective of the buyers.

StrOPM versus M-Qubed. In Figure 2b we show results for an agent using the StrOPM algorithm (using the same settings as above) playing against an agent using the M-Qubed algorithm by Crandall and Goodrich (2005). We choose M-Qubed as a state-of-the-art MARL algorithm to compare StrOPM to, as Crandall & Goodrich claim it performs on par with other state-of-the-art MARLs on simple matrix games, and is additionally capable of successfully solving the iPD game.

M-Qubed as introduced in Crandall and Goodrich (2005) operates as a mix of a Q-learning type algorithm and a pure strategy player. The algorithm learns Q-like values that encode the received reward and the discounted future reward. The algorithm investigates the policy of playing according to the Q-like values various exploration strategies. At the same time, the algorithm considers whether to play any of its actions consistently to promote possibly good interactions with the opponent. A parameter $\beta \in [0, 1]$ is used in M-Qubed and adapted in play to learn the best overall strategy; M-Qubed plays a game using the Q-values with probability $(1 - \beta)$ or it plays the highest valued action as a pure strategy with probability β . We used for M-Qubed the same settings as Crandall & Goodrich note in their work³.

When a StrOPM player plays against a player using M-Qubed, the full average reward of 0.35 is not reached as alternating Cooperate/Defect followed by Defect/Cooperate patterns often emerged, giving an average reward of 0.25 per epoch. This is not the optimal C, C outcome, but still much better than the typical D, D outcome other MARLs arrive at. This result suggests that the switching policy in M-Qubed is repeatedly “tripped”. The non-linear switching policy threshold is hard to approximate in StrOPM’s linear opponent change model. Including more history of non-linear update terms in future extensions may solve this issue.

The trouble with precoded strategies. The M-Qubed algorithm has a clear potential weakness since it has to decide between the pure (or at least precoded) strategy of playing full cooperation, and a myopic Q-Learning type strategy. We would expect it to perform badly when the optimal strategy involves a more complex strategy that is not precoded. Our StrOPM algorithm on the other hand is not biased to playing precoded strategies.

We tested this hypothesis in the game with the payoff matrix of Figure 3a, dubbed “Switching Bait”, which we specifically designed so that the highest (average) joint reward is achieved when agents alternate Cooperate/Defect and Defect/Cooperate (or at least play these two joint actions equally often). Playing the Cooperate/Defect, Defect/Cooperate mixed strategy yields an average reward of 0.45 per epoch for both agents.

In experiments playing the Switching Bait game, M-Qubed

³These settings allowed us to successfully replicate the results in Crandall and Goodrich (2005).

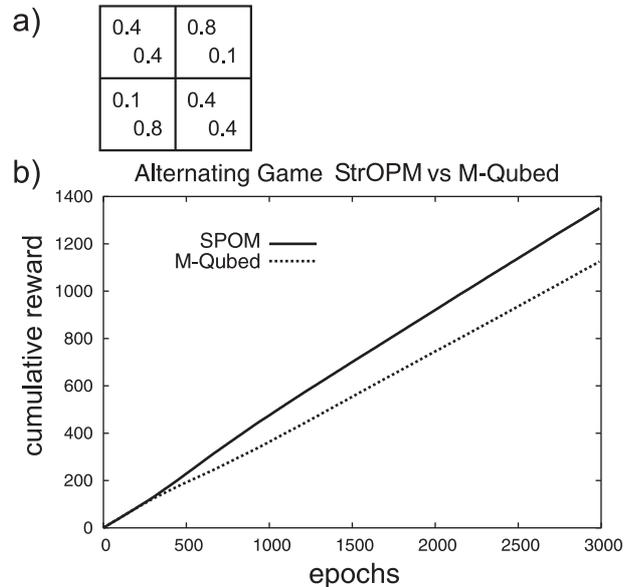


Figure 3: a) payoff matrix for “Switching Bait” game b) respective cumulative reward obtained by a StrOPM player playing against an M-Qubed player.

was frequently not able to find the optimal mixed strategy in selfplay, as it predominantly converged to playing C or D and achieves a reward of only 0.4 per epoch. Higher rewards of $\frac{0.8+0.1}{2}$ were achieved by one of the M-Qubed learners by playing C and D alternately, while the other M-Qubed player plays C or D continuously for an average reward of $\frac{0.1+0.4}{2}$. For the Switching Bait new game, M-Qubed often decides to play the safer strategy of full cooperation or defection.

The StrOPM algorithm, unlike M-Qubed, is not biased to playing pure strategies. For the Switching Bait game, it often finds the alternating Cooperate/Defect pattern and matches it to the opponent’s Defect/Cooperate pattern in self play. Importantly, a StrOPM player playing against an M-Qubed player was found to either achieve this mutually beneficial equilibrium or was repeatedly able to exploit M-Qubed if the opponent decided to play a (precoded) pure strategy. An example of this latter behavior is shown in Figure 3b, where the cumulative payoff achieved during play is plotted. The StrOPM algorithm achieves a higher payoff as M-Qubed opts for the safer, pure strategy. This clearly demonstrates the weakness of using precoded fixed strategies (as in M-Qubed).

5. DISCUSSION AND CONCLUSIONS

We have presented a number of important contributions in this paper. First, we argued the need for Multi-Agent learning algorithms to take into account the impact of their actions on the adaptive policy of the opponents. We then presented a novel framework that allows agents to apply the above observation to achieve desired outcomes for representative iterated matrix games, and, importantly, also for the iterated Prisoner’s Dilemma. The latter is notoriously difficult to solve for current state-of-the-art MARL approaches. Our MARL framework allows an agent to learn profitable strategies for long term behavior. Finally, we demonstrated

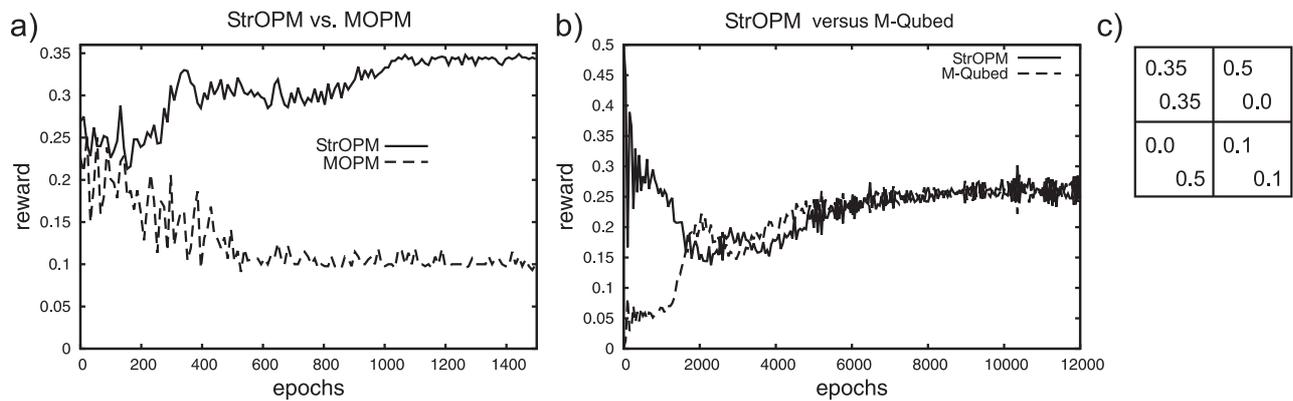


Figure 2: a) StrOPM players playing the iPD (solid line) and the MOPM-variant (dashed line). b) A StrOPM player playing vs M-Qubed for the iPD. c) iPD payoff matrix

that incorporating fixed precoded strategies in MARL algorithms makes agents vulnerable to exploitation, especially when mixed strategies can be optimal (of which there is a continuous spectrum). Additionally, fixed precoded strategies may not be suited to novel settings.

The important novel component of this work is the concept of learning the changes in opponents policy due ones own action. This is an idea that can in fact be incorporated into the quickly growing literature on MARL. One can foresee more and more complex nested opponent models (Hu & Wellman, 1998) to extract every bit of reward from complex games like iPD. Although perfectly learning about an opponent while at the same time perfectly learning to adjust oneself is problematic (Nachbar & Zame, 1996), there is still much scope to be “smarter” than your opponents.

Our future challenge is to apply our new concept to large, complex state spaces. One interesting direction may be to integrate our framework with new ideas about state space representations like Predictive State Representations (Wolfe, James, & Singh, 2005).

References

- Axelrod, R. (1984). *The evolution of cooperation*. Basic Books, New York, NY.
- Banerjee, B., & Peng, J. (2004). The role of reactivity in multiagent learning. In *Proc. 3rd AAMAS* (p. 538-545).
- Bowling, M. (2004). Convergence and no-regret in multiagent learning. In *Proc. NIPS-17* (pp. 209-216).
- Bowling, M. H., & Veloso, M. M. (2002). Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136(2), 215-250.
- Conitzer, V., & Sandholm, T. (2003). AWESOME: A General Multiagent Learning Algorithm that Converges in Self-Play and Learns a Best Response Against Stationary Opponents. In *Proc. 20th ICML* (p. 83-90).
- Crandall, J. W., & Goodrich, M. A. (2005). Learning to compete, compromise, and cooperate in repeated general-sum games. In *Proc. 22nd ICML*.
- Greenwald, A. R., & Hall, K. (2003). Correlated Q-learning. In *Proc. 20th ICML* (p. 242-249).
- Hu, J., & Wellman, M. (1998). Online learning about other agents in a dynamic multiagent system. In *Proc ACM Conf. on Autonomous Agents* (p. 239-246).
- Littman, M. L., & Stone, P. (2003). A polynomial-time nash equilibrium algorithm for repeated games. In *Proc. 4th ACM Conf. on Electronic Commerce* (p. 48-54).
- Myerson, R. B. (1991). *Game theory. analysis of conflict*. Harvard University Press.
- Nachbar, J. H., & Zame, W. R. (1996). Non-computable strategies and discounted repeated games. *Economic Theory*, 8, 103-122.
- Nash, J. (1951). Non-cooperative games. *Annals of Mathematics*, 54, 286-295.
- Puterman, M. L. (1994). *Markov decision process*. John Wiley and Sons, Inc., New York.
- Rapoport, A., Guyer, M., & Gordon, D. (1976). *The 2x2 game*. MI: University of Michigan Press.
- Sandholm, T., & Crites, R. (1995). Multiagent reinforcement learning in the iterated prisoner’s dilemma. *Biosystems*, 37, 147-166.
- Shoham, Y., Powers, R., & Grenager, T. (2004). Multi-agent reinforcement learning: a critical survey. In *AAAI Fall Symposium on Artificial Multi-Agent Learning*.
- Sutton, R., & Barto, A. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press.
- ’t Hoen, P., & La Poutré, J. (2005). Repeated auctions with complementarities. In *AMEC VII*. (to appear, LNCS)
- Tesauro, G. (2003). Extending Q-learning to general adaptive multi-agent systems. In *Nips-16* (pp. 871-878).
- Weinberg, M., & Rosenschein, J. S. (2004). Best-response multiagent learning in non-stationary environments. In *Proc. 3rd AAMAS* (p. 506-513). New York.
- Wolfe, B., James, M. R., & Singh, S. (2005). Learning predictive state representations in dynamical systems without reset. In *Proc. 22nd ICML*.
- Yao, X., & Darwen, P. J. (1994). An experimental study of n-person iterated prisoner’s dilemma games. *Informatica*, 18, 435-450.