

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/150664>

How to cite:

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

How to Guide a Non-Cooperative Learner to Cooperate: Exploiting No-Regret Algorithms in System Design

Extended Abstract

Nicholas Bishop
University of Southampton
United Kingdom
nb8g13@soton.ac.uk

Le Cong Dinh
University of Southampton
United Kingdom
l.c.dinh@soton.ac.uk

Long Tran-Thanh
University of Warwick
United Kingdom
long.tran-thanh@warwick.ac.uk

ABSTRACT

We investigate a repeated two-player game setting where the column player is also a designer of the system, and has full control over payoff matrices. In addition, we assume that the row player uses a no-regret algorithm to efficiently learn how to adapt their strategy to the column player’s behaviour over time. The goal of the column player is to guide her opponent into picking a mixed strategy which is preferred by the system designer. Therefore, she needs to: (i) design appropriate payoffs for both players; and (ii) strategically interact with the row player during a sequence of plays in order to guide her opponent to converge to the desired mixed strategy. To design appropriate payoffs, we propose a novel zero-sum game construction whose unique minimax solution contains the desired behaviour. We also propose another construction in which only the minimax strategy of the row player is unique. Finally, we propose a new game playing algorithm for the system designer and show that it can guide the row player to its minimax strategy, under the assumption that the row player adopts a *stable* no-regret algorithm.

KEYWORDS

System Design, Unique Nash Equilibrium, Last Round Convergence

ACM Reference Format:

Nicholas Bishop, Le Cong Dinh, and Long Tran-Thanh. 2021. How to Guide a Non-Cooperative Learner to Cooperate: Exploiting No-Regret Algorithms in System Design: Extended Abstract. In *Proc. of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021)*, Online, May 3–7, 2021, IFAAMAS, 3 pages.

1 INTRODUCTION

We consider a repeated two-player game setting, in which one player (say the column player) is also a designer of the system (i.e., she can design the payoffs for both players), and her opponent (the row player) is a strategic utility maximiser, who can efficiently learn to adapt their strategy to the column player’s behaviour over time in order to achieve good total payoff. The goal of the column player is to guide her opponent into selecting a mixed strategy which is favourable for the system designer. In particular, she needs to achieve this by: (i) designing appropriate payoffs for both players; and (ii) strategically interacting with the row player during a sequence of plays in order to ensure her opponent converges to the

desired behaviour. In this paper, we propose an approach for solving this problem, which consists of two components corresponding to (i) and (ii) respectively:

Games with a unique minimax solution. To begin, the system designer must decide upon the payoffs for both players. For two-player zero-sum games, it is well known that, in the full information setting, the best payoff rational players can attain is their payoff at any minimax equilibrium. [2]. Thus, a natural idea is to construct a zero-sum game, A , in which the *only* minimax strategy available to the row player is the desired behaviour. In doing so, the system designer can hope that any reasonable player will eventually begin to play their minimax strategy, and thus adopt the desired behaviour. Construction of games with unique equilibrium solutions is a long running problem within the literature, beginning with the seminal work of Shapley, Karlin, and Bohnenblust [1]. In this work, we aim to provide zero-sum game constructions which offer the system designer more flexibility in terms of the payoff matrix they choose. Such flexibility may be useful when enforcing system payoffs correspond to costly real world actions.

Last round convergence in zero-sum games. Once a zero-game has been chosen, the system designer must strategically incentivise the row player to converge to the desired behaviour through repeated play. Recall that, if both players commit to a no-regret algorithm at each phase of play, then payoffs will converge in expectation to the value of the game. Additionally, during this process, no exchange of information takes place, as both players require only the payoffs they observe in order to update their strategies [2, 6]. As a result, one may hope that, by adopting a no-regret algorithm, the system designer can naturally guide the row player to their minimax strategy. Unfortunately, convergence of average payoffs, in general, does not imply convergence in strategies [7]. This issue is known as the last round convergence problem in the online learning literature, and has recently attracted a good amount of attention [3–5]. In what follows, we propose a novel no-regret algorithm which leverages the information advantage of the system designer to guide the row player to its minimax strategy over time, under a mild set of assumptions.

For the remainder, we will describe each component of our approach in more detail, beginning with the construction of zero-sum games with unique minimax equilibria.

2 DESIGNING GAMES WITH A UNIQUE MINIMAX SOLUTION

Without loss of generality, suppose there exists a preferred strategy, y^* , that the column player would like to play whilst at a minimax

equilibrium. That is, the column player wishes not only to ensure that the row player's minimax strategy is \mathbf{x}^* , but also that their own minimax strategy is \mathbf{y}^* . Under this assumption, we observe that the strategy pair $(\mathbf{x}^*, \mathbf{y}^*)$ must form a minimax equilibrium of A . Moreover, as previously mentioned, \mathbf{x}^* must be the unique minimax strategy for the row player.

In order to satisfy both conditions, the support of \mathbf{y}^* must be greater than equal to the support of \mathbf{x}^* . For if this is not the case, then it is provably impossible to construct a matrix A with minimax equilibrium $(\mathbf{x}^*, \mathbf{y}^*)$ whilst guaranteeing the uniqueness of the minimax strategy \mathbf{x}^* [1]. From now on, without loss of generality, we shall assume that for any strategy with support k , that the first k entries are nonzero.

If the support of \mathbf{y}^* is greater than the support of \mathbf{x}^* , then we construct A according to Theorem 1. The other case, in which the support of \mathbf{y}^* is equal to the support of \mathbf{x}^* is dealt with in the full version of the paper. Note that, in both theorems, \mathbf{y}^* is not necessarily the unique minimax strategy for the column player, whilst \mathbf{x}^* is the unique minimax strategy for the row player.

THEOREM 1. *Let $\mathbf{x} \in \Delta_n, \mathbf{y} \in \Delta_m$ such that $k = \text{support}(\mathbf{x}) < l = \text{support}(\mathbf{y})$. Let the matrix A be of the form*

$$A = \begin{bmatrix} a_1 & a_2 & \dots & a_k & \beta_1 & \dots & \beta_l \\ \alpha_1 & \alpha_2 & \dots & \alpha_k & \beta_2 & \dots & \beta_l \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \alpha_1 & \alpha_2 & \dots & \alpha_k & \beta_k & \dots & \beta_l \\ \alpha_1 - z & \alpha_2 - z & \dots & \alpha_k - z & v & \dots & v \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \alpha_1 - z & \alpha_2 - z & \dots & \alpha_k - z & v & \dots & v \end{bmatrix}$$

where the parameters of A satisfy

$$0 < v_1 < v\bar{y}, \quad \bar{y} = \sum_{i=k+1}^l y_i, \quad z = \frac{v\bar{y} - v_1}{\sum_{i=1}^k y_i},$$

$$\beta_i = v, \quad \alpha_i = v + \frac{x_i(v\bar{y} - v_1)}{y_i}, \quad a_i = \alpha_i - \frac{v\bar{y} - v_1}{y_i} \quad \forall i \in [k].$$

then \mathbf{x} is the unique minimax strategy for the row player in the zero-sum game described by A .

In Theorem 1, the parameters a_i and z ensure that \mathbf{x} is the unique minimax strategy for the row player, even in the case where the support of \mathbf{x} is less than n . Meanwhile, the parameters β_i ensure that \mathbf{y} is a minimax strategy for the column player. Lastly, the parameter y ensures that all entries of A are nonnegative (although this is not strictly required).

3 LAST ROUND CONVERGENCE IN TWO-PLAYER ZERO-SUM GAMES

Next, given a matrix A , as constructed in the previous section, we investigate how the system designer can guide the row player to its minimax strategy, under the assumption that the row player uses a no-regret algorithm.

Firstly, we show that a naïve approach, namely to repeatedly playing \mathbf{y}^* , will not always lead to the desired last round convergence (i.e., the row player will converge to \mathbf{x}^*):

CLAIM 2. *If $\text{support}(\mathbf{x}) > 1$, then there is no guarantee that if the column player repeatedly plays \mathbf{y}^* , the row player will eventually converge to \mathbf{x}^* .*

Given this result, we need to design a different game playing policy for the column player. More specifically, we require a policy which actively exploits the information advantage possessed by the column player. The LRCA algorithm, originally proposed by Dinh *et al.* [5], was designed with this goal in mind. On odd rounds, LRCA selects the minimax strategy in order to stabilise the trajectory of both players. Meanwhile, in even rounds, the LRCA algorithm exploits any weaknesses in the row player's strategy by moving in the direction of the highest payoff strategy in the previous round, where the rate of moving depends on how far the row player's strategy from the Nash equilibrium (i.e., $f(\mathbf{x}_{t-1}) - v$). The LRCA algorithm guarantees last round convergence (for both players) against a number of popular no-regret algorithms including the multiplicative weight update algorithm, online mirror descent, and the linear multiplicative weight update algorithm.

LRCA is described in detail by Algorithm 1 below:

Algorithm 1: Last Round Convergence with Asymmetry (LRCA) algorithm

Input: Current iteration t , past feedback $\mathbf{x}_{t-1}^\top A$ of the row player, minimax strategy \mathbf{y}^* and value v of the game.

Output: Strategy \mathbf{y}_t for the column player

if $t = 2k - 1, k \in \mathbb{N}$ **then**

$\mathbf{y}_t = \mathbf{y}^*$

if $t = 2k, k \in \mathbb{N}$ **then**

$\mathbf{e}_t := \text{argmax}_{e \in \{e_1, e_2, \dots, e_m\}} \mathbf{x}_{t-1}^\top A e$
 $f(\mathbf{x}_{t-1}) := \max_{\mathbf{y} \in \Delta_m} \mathbf{x}_{t-1}^\top A \mathbf{y}; \quad \alpha_t := \frac{f(\mathbf{x}_{t-1}) - v}{\max(\frac{v}{4}, 2)}$
 $\mathbf{y}_t := (1 - \alpha_t) \mathbf{y}^* + \alpha_t \mathbf{e}_t$

More generally, we prove that LRCA also guarantees last round convergence when paired with *any* no-regret algorithm, as long as this no-regret algorithm possesses the "stability" property, as defined below:

Definition 3. A no-regret algorithm is *stable* if $\forall t : \mathbf{y}_t = \mathbf{y}^* \implies \mathbf{x}_{t+1} = \mathbf{x}_t$.

Note that a wide range of no-regret algorithms adhere to this property. For example, the class of Follow the Regularised Leader algorithms are stable. A proof of this fact is given in the full paper.

THEOREM 4. *Assume that the row player follows a stable no-regret algorithm and n is the dimension of the row player's strategy. Then, by following LRCA, for any $\epsilon > 0$, there exists $l \in \mathbb{N}$ such that $\frac{\mathcal{R}_l}{l} = O(\frac{\epsilon^2}{n})$ and $f(\mathbf{x}_l) - v \leq \epsilon$.*

Note that $(\mathbf{x}_l, \mathbf{y}^*)$ are ϵ -Nash equilibria. For no-regret algorithms with the optimal regret bound $\mathcal{R}_l = O(l^{\frac{1}{2}})$, Theorem 4 guarantees that the row player will reach an ϵ -Nash equilibrium in at most $O(\epsilon^{-4})$ rounds. Thus, LRCA is successful in guiding the row player towards the desired behaviour, as long as the row player adopts a stable no-regret algorithm.

REFERENCES

- [1] HF Bohnenblust, S Karlin, and LS Shapley. 1950. Solutions of discrete, two-person games. *Contributions to the Theory of Games* 1 (1950), 51–72.
- [2] Nicolo Cesa-Bianchi and Gábor Lugosi. 2006. *Prediction, learning, and games*. Cambridge university press.
- [3] Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. 2017. Training gans with optimism. *arXiv preprint arXiv:1711.00141* (2017).
- [4] C. Daskalakis and I. Panageas. 2018a. Last-iterate convergence: Zero-sum games and constrained min-max optimization. *arXiv preprint arXiv:1807.04252* (2018a).
- [5] Le Cong Dinh, Long Tran-Thanh, Tri-Dung Nguyen, and Alain B Zemkoho. 2020. Last Round Convergence and No-Instant Regret in Repeated Games with Asymmetric Information. *arXiv preprint arXiv:2003.11727* (2020).
- [6] Yoav Freund and Robert E Schapire. 1999. Adaptive game playing using multiplicative weights. *Games and Economic Behavior* 29, 1-2 (1999), 79–103.
- [7] Panayotis Mertikopoulos, Christos Papadimitriou, and Georgios Piliouras. 2018. Cycles in adversarial regularized learning. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2703–2717.