## Paths to Equilibrium in Normal-Form Games

**Bora Yongacoglu** University of Toronto

Gürdal Arslan University of Hawaii at Manoa

Lacra Pavel University of Toronto

Serdar Yüksel Queen's University BORA.YONGACOGLU@UTORONTO.CA

GURDAL@HAWAII.EDU

PAVEL@CONTROL.TORONTO.EDU

YUKSEL@QUEENSU.CA

## Abstract

In multi-agent reinforcement learning (MARL), agents repeatedly interact across time and revise their strategies as new data arrives, producing a sequence of strategy profiles. This paper studies sequences of strategies satisfying a pairwise constraint inspired by policy updating in reinforcement learning, where an agent who is best responding in period t does not switch its strategy in the next period t + 1. This constraint merely requires that optimizing agents do not switch strategies, but does not constrain the other non-optimizing agents in any way, and thus allows for exploration. Sequences with this property are called satisficing paths, and arise naturally in many MARL algorithms. A fundamental question about strategic dynamics is such: for a given game and initial strategy profile, is it always possible to construct a satisficing path that terminates at an equilibrium strategy? The resolution of this question has implications about the capabilities or limitations of a class of MARL algorithms. We answer this question in the affirmative for mixed extensions of finite normal-form games.

## 1. Introduction

Game theory is a mathematical framework for studying strategic interaction between multiple selfinterested agents, called players. Game theoretic models are pervasive in machine learning, appearing in multi-agent systems (Zhang et al., 2021; Gronauer and Diepold, 2022), single-agent reinforcement learning with multiple objectives (Hayes et al., 2022), in adversarial model training (Goodfellow et al., 2014, 2020; Bose et al., 2020), and throughout online learning theory in the form of worst-case performance guarantees (Cesa-Bianchi and Lugosi, 2006).

In an *n*-player game, each player  $i = 1, \dots, n$ , selects a strategy  $x^i \in \mathcal{X}^i$  and receives a reward  $R^i(x^1, \dots, x^n)$ , which depends on the collective *strategy profile*  $\mathbf{x} = (x^1, \dots, x^n) =: (x^i, \mathbf{x}^{-i})$ . Player *i*'s optimization problem is to *best respond* to the strategy  $\mathbf{x}^{-i}$  of its counterparts, choosing  $x^i \in \mathcal{X}^i$  to maximize  $R^i(x^i, \mathbf{x}^{-i})$ . A strategy profile  $\mathbf{x}_* = (x^i_*)_{i=1}^n$  is called a *Nash equilibrium* if all players are simultaneously best responding to one another:

$$R^{i}(x_{*}^{i}, \mathbf{x}_{*}^{-i}) \geq R^{i}(y^{i}, \mathbf{x}_{*}^{-i}), \quad \forall y^{i} \in \mathcal{X}^{i}, \forall i = 1, \cdots, n.$$

The related tasks of computing, approximating, and learning Nash equilibrium have attracted enduring attention in the broader machine learning community (Singh et al., 2000; Jafari et al., 2001;

© B. Yongacoglu, G. Arslan, L. Pavel & S. Yüksel.

Hu and Wellman, 2003; Daskalakis et al., 2010; Nowé et al., 2012; Bravo et al., 2018; Flokas et al., 2020; Hsieh et al., 2021; Lu, 2023).

In the field of multi-agent reinforcement learning (MARL), players use learning algorithms to iteratively revise their strategies in response to the observed history of play, producing a sequence  $\{\widehat{\mathbf{x}}_t\}_{t\geq 1}$  in the set of strategy profiles  $\mathbf{X} := \mathcal{X}^1 \times \cdots \times \mathcal{X}^n$ . In the multi-agent setting, each player's learning problem is complicated by numerous inherent challenges. First, there is a non-stationarity issue, as a given individual's reward function changes whenever other agents revise their strategies, which subsequently prompts the individual to revise its own strategy. Second, key information about the game or about the strategies of other players may only be partially observable to an individual, necessitating estimation of various quantities. Analyzing the convergence properties of MARL algorithms can therefore be difficult, and the development of theoretical tools for such analysis is an important aspect of multi-agent learning theory.

A number of MARL algorithms are designed to approximate dynamical systems  $\{\mathbf{x}_t\}_{t\geq 1}$  on the set of strategy profiles  $\mathbf{X}$  in which the next strategy for player *i* is selected as  $x_{t+1}^i = f^i(\mathbf{x}_t)$ , where  $\mathbf{x}_t = (x_t^1, \ldots, x_t^n)$  is the strategy profile in period *t*. A representative sample of MARL algorithms of this type is offered in the next section. This approach facilitates analysis of the convergence behavior of the MARL algorithm, as it allows the analyst to separately consider the convergence of the idealized dynamical process  $\{\mathbf{x}_t\}_{t\geq 1}$  induced by the update functions  $\{f^i\}_{i=1}^n$  on one hand and then consider the approximation of  $\{\mathbf{x}_t\}_{t\geq 1}$  by the MARL algorithm's iterates  $\{\hat{\mathbf{x}}_t\}_{t\geq 1}$  on the other.

Our primary interest in this work centers on update functions satisfying a quasi-rationality condition called *satisficing*: if  $x^i$  is a best response to  $\mathbf{x}^{-i}$ , then  $f^i(x^i, \mathbf{x}^{-i}) = x^i$ . That is, when an agent is already best responding, the update rule instructs the agent to continue using this strategy. This quasi-rationality constraint generalizes the well-studied best response update and is desirable for stability of the resulting dynamics, as it guarantees that Nash equilibrium strategy profiles are invariant under the dynamics. Moreover, the satisficing condition is only quasi-rational, in that it imposes no constraint on strategy updates when an agent is not best responding, and so allows for exploratory strategy updates. Such update rules are common in the literature on multi-agent learning theory (Blume, 1993; Marden and Shamma, 2012; Chasparis et al., 2013).

The basic motivation of this paper is to better understand the capabilities of MARL algorithms that operate using the satisficing principle in selecting successive strategies, potentially augmented with random exploration when an agent is not best responding. Rather than studying a particular collection of strategy update functions, we study the problem on the level of sequences in  $\mathbf{X}$ , which allows us to implicitly account for experimental strategy updates. A sequence  $(\mathbf{x}_t)_{t\geq 1}$  of strategy profiles is called a *satisficing path* if, for each player *i* and time *t*, one has that  $x_{t+1}^i = x_t^i$  whenever  $x_t^i$  is a best response to  $\mathbf{x}_t^{-i}$ . The central research question of this paper is such:

# For a normal-form game $\Gamma$ and an initial strategy profile $\mathbf{x}_1$ , is it always possible to construct a satisficing path from $\mathbf{x}_1$ to a Nash equilibrium of the game $\Gamma$ ?

The resolution of this question has implications for the possible effectiveness of a class of MARL algorithms designed to seek Nash equilibrium. Indeed, the question has been answered in the affirmative for two-player normal-form games by Foster and Young (2006) and for n-player symmetric Markov games by Yongacoglu et al. (2023), and in both classes of games this has directly lead to MARL algorithms with convergence guarantees for approximating equilibria. A positive resolution of this question would remove a theoretical obstacle and establish that uncoordinated,

distributed random search can effectively assist Nash-seeking algorithms achieve last-iterate convergence guarantees in a more general class of games than previously possible.

**Our Contributions.** In this work, we give a positive answer to the preceding research question for (the mixed extension of) any *n*-player normal-form game with finite action sets. That is, we show that for a finite *n*-player game  $\Gamma$  and any initial strategy profile  $\mathbf{x}_1$ , there exists a satisficing path of finite length beginning at  $\mathbf{x}_1$  and ending at a Nash equilibrium of the game  $\Gamma$ . This partially answers an open question posed by Yongacoglu et al. (2023).

We prove our main result by analytically constructing a satisficing path from an arbitrary initial strategy profile to a Nash equilibrium. Our approach is somewhat counterintuitive, in that it does not attempt to seek Nash equilibrium by improving the performance of unsatisfied players (players who are not best responding at a given strategy profile), but rather by updating strategies in a way that *increases* the number of unsatisfied players at each round. This tactic heavily leverages the freedom afforded to unsatisfied players to explore their strategy space and avoids the well-observed challenges of cyclical strategy revision that occurs when agents attempt to best respond to their counterparts.

Notation. We let  $\mathbb{P}$  and  $\mathbb{E}$  denote probability and expectation, respectively. For a finite set A, we let  $\mathbb{R}^A$  denote the vector space  $\mathbb{R}^{|A|}$  with components indexed by elements of A. We let  $\Delta_A$  denote the set of probability measures over a set A. For  $n \in \mathbb{N}$ , we let  $[n] := \{1, 2, ..., n\}$ . For a point x, the Dirac measure centered at x is denoted  $\delta_x$ . Agent indices are typically superscripts, while time/iteration indices are typically subscripts. Boldface characters are reserved for multi-agent quantities. When discussing a fixed agent i, the remaining collection of agents are called i's counterparts or counterplayers.

#### **Related Work**

Beginning with fictitious play (Brown, 1951), a vast number of MARL algorithms have been proposed for iterative strategy adjustment while playing a game under various assumptions on observability of counterplayer strategies. The most widely studied class of algorithms of this type involve each player running a no-regret algorithm on its own stream of rewards. Fictitious play and its descendants, such as stochastic fictitious play (Hofbauer and Sandholm, 2002) and generalized weakened fictitious play (Leslie and Collins, 2006) are special cases of this class and have been extensively studied. Although the convergence behavior of fictitious play and its variants has been studied in several game models, including normal-form games and Markov games, convergence results are typically available only for games exhibiting special structure, such as zero-sum rewards or other advantageous properties amenable to analysis (Hofbauer and Sandholm, 2002; Baudin and Laraki, 2022; Sayin et al., 2022a,b).

The formalization of the *fictitious play property* is especially relevant to this paper. A game is said to have the fictitious play property if the empirical frequencies of strategies converge to a Nash equilibrium from any initialization. Examples of games with and without the fictitious play property date to (Robinson, 1951) and (Shapley, 1964), respectively, and the identification of games with and without this property was viewed as an important research question (Monderer and Sela, 1996; Monderer and Shapley, 1996a).

The convergence properties of the fictitious play algorithm are intimately connected to those of best response dynamics, a full information dynamical system evolving in continuous time where the evolution rule for player i's strategy is governed by its best response multi-function. By harnessing such connections, convergence results for fictitious play and a number of other MARL algorithms have been obtained by analyzing the dynamical systems induced by specific update rules (Benaïm et al., 2005; Leslie and Collins, 2005; Swenson et al., 2018b).

A second research direction relevant to the present work involves (im)possibility results for strategic dynamics defined by strategy update functions, taking the form  $x_{t+1}^i = f^i(\mathbf{x}_t)$  in discrete time or an analogous form in continuous time. The question of when such dynamics converge to Nash equilibrium has received persistent attention. In the case of deterministic strategy updates, Hart and Mas-Colell (2003) studied strategic dynamics in continuous time and showed that if the strategy update functions, analogous to  $f^i$  above, satisfy regularity conditions as well as a desirable property called uncoupledness, by which  $f^i$  cannot depend on the reward functions of *i*'s counterplayers, then the resulting dynamics are not Nash convergent in general. These results were recently generalized by Milionis et al. (2023), who obtained impossibility results in both continuous and discrete time while requiring only continuity of the deterministic dynamical system. Additional possibility and impossibility results were presented by Babichenko (2012), who studied strategic dynamics in a different setting, where players observe only their own actions and not the actions of their counterplayers.

Passing from deterministic strategic dynamics to stochastic strategic dynamics, a number of positive results were obtained by incorporating exogenous randomness into one's strategy update, along with finite recall of recent play (Hart and Mas-Colell, 2006; Foster and Young, 2006; Germano and Lugosi, 2007).

In the regret testing algorithm of Foster and Young (2006), players revise their strategies according to whether or not their most recent strategy met a satisfaction criterion: if  $x_t^i$  performed within  $\epsilon$  of the optimal performance against  $\mathbf{x}_t^{-i}$ , player *i* continues using it and picks  $x_{t+1}^i = x_t^i$ . Otherwise, player *i* experiments and selects  $x_{t+1}^i$  according to a uniformly positive probability distribution over  $\mathcal{X}^i$ . Conditional strategy updates similar to this have appeared in several other works, such as (Chien and Sinclair, 2011; Candogan et al., 2013; Chasparis et al., 2013), and the regret testing algorithm has been extended in several ways (Germano and Lugosi, 2007; Arslan and Yüksel, 2017).

A game is said to have the *satisficing paths property* if it admits a finite length satisficing path ending at equilibrium with an arbitrary initial strategy profile. As we discuss in the next section, satisficing paths can be interpreted as a natural generalization of best response paths. Consequently, the problem of identifying games that have the satisficing paths property is a theoretically relevant question analogous to characterizing exact potential games (Monderer and Shapley, 1996b) or determining when a game has the fictitious play property.

The concept of satisficing paths was first formalized by Yongacoglu et al. (2023) in the context of multi-state Markov games, where it was shown that *n*-player symmetric Markov games have the satisficing paths property and this fact could be used to produce a convergent MARL algorithm. However, the core idea of satisficing paths appeared earlier, before this formalization: in the convergence analysis of the regret testing algorithm in (Foster and Young, 2006), it was shown that two-player normal-form games have the satisficing paths property, though this terminology was not used.

## 2. Model

A finite, *n*-player normal-form game  $\Gamma$  is described by a list

$$\Gamma = (n, \mathbf{A}, \mathbf{r}),$$

where *n* is the number of players,  $\mathbf{A} = \mathbb{A}^1 \times \cdots \times \mathbb{A}^n$  is a finite set of action profiles, and  $\mathbf{r} = (r^i)_{i \in [n]}$  is a collection of reward functions, where  $r^i : \mathbf{A} \to \mathbb{R}$  describes the reward of player *i* as a function of the action profile. The *i*<sup>th</sup> component of  $\mathbf{A}$  is player *i*'s action set  $\mathbb{A}^i$ .

**Description of play.** Each player  $i \in [n]$  selects its action  $a^i$  according to a probability vector  $x^i \in \Delta_{\mathbb{A}^i}$ . That is,  $a^i \sim x^i$ . The vector  $x^i$  is called player *i*'s mixed strategy, and we denote player *i*'s set of mixed strategies by  $\mathcal{X}^i := \Delta_{\mathbb{A}^i}$ . Players are assumed to select their actions simultaneously, or at least without observing one another's actions, and the collection of actions  $\{a^i : i \in [n]\}$  is assumed to be mutually independent. The set of mixed strategy profiles is denoted  $\mathbf{X} := \mathcal{X}^1 \times \cdots \mathcal{X}^n$  and corresponds to the set of product measures on  $\mathbf{A}$ . After the action profile is selected, each player *i* receives reward  $r^i(a^1, \ldots, a^n)$ .

Player *i*'s performance criterion is its expected reward, which depends jointly on its mixed strategy and the strategies of *i*'s counterplayers. For each strategy profile  $x \in X$ , player *i*'s expected reward is

$$R^{i}(x^{i}, \mathbf{x}^{-i}) = \mathbb{E}_{\mathbf{a} \sim \mathbf{x}} \left[ r^{i}(a^{1}, \dots, a^{n}) \right],$$

where  $\mathbb{E}_{\mathbf{a}\sim\mathbf{x}}$  signifies that  $a^j \sim x^j$  for each player  $j \in [n]$  and we have used the convention that  $\mathbf{x} = (x^i, \mathbf{x}^{-i})$  and  $\mathbf{x}^{-i} = (x^1, \dots, x^{i-1}, x^{i+1}, \dots, x_n)$ . Since player *i*'s objective depends on the strategies of its counterplayers, the relevant optimality notion is that of  $(\epsilon$ -) best responding.

**Definition 1** A mixed strategy  $x_*^i \in \mathcal{X}^i$  is called an  $\epsilon$ -best response to the strategy  $\mathbf{x}^{-i} \in \mathbf{X}^{-i}$  if

$$R^{i}(x_{*}^{i}, \mathbf{x}^{-i}) \ge R^{i}(x^{i}, \mathbf{x}^{-i}) - \epsilon \quad \forall x^{i} \in \mathcal{X}^{i}.$$

The standard solution concept for *n*-player normal form games is that of ( $\epsilon$ -) Nash equilibrium, which entails a situation in which all players are simultaneously ( $\epsilon$ -) best responding to one another.

**Definition 2** For  $\epsilon \ge 0$ , a strategy profile  $\mathbf{x}_* = (x_*^i, \mathbf{x}_*^{-i}) \in \mathbf{X}$  is called an  $\epsilon$ -Nash equilibrium if, for every  $i \in [n]$ ,

$$R^{i}(x_{*}^{i}, \mathbf{x}_{*}^{-i}) \geq R^{i}(x^{i}, \mathbf{x}_{*}^{-i}) - \epsilon \quad \forall x^{i} \in \mathcal{X}^{i}.$$

Putting  $\epsilon = 0$  in the definitions above, one recovers the classical definitions of *best responding* and *Nash equilibrium*. For any  $\epsilon \ge 0$ , the set of  $\epsilon$ -best responses to a specified strategy  $\mathbf{x}^{-i}$  is denoted  $BR^i_{\epsilon}(\mathbf{x}^{-i}) \subseteq \mathcal{X}^i$ .

#### 2.1. Satisficing Paths

This section contains several definitions useful for the of study strategic dynamics in normal-form games and presents satisficing paths as a generalization of best response paths.

**Definition 3** A sequence of strategy profiles  $(\mathbf{x}_t)_{t\geq 1}$  in  $\mathbf{X}$  is called a best response path if, for every  $t \geq 1$  and every player  $i \in [n]$ , we have

$$x_{t+1}^{i} = \begin{cases} x_{t}^{i}, & \text{if } x_{t}^{i} \in BR_{0}^{i}(\mathbf{x}_{t}^{-i}), \\ \text{some } x_{\star}^{i} \in BR_{0}^{i}(\mathbf{x}_{t}^{-i}), & \text{else.} \end{cases}$$

The preceding definition of best response paths can be relaxed in several ways, and such relaxations are often desirable to avoid non-convergent cycling behavior. A common relaxation involves synchronizing players or incorporating inertia, so that only a subset of players switch their strategies at a given time, which can be help achieve coordination in cooperative settings (Marden and Shamma, 2012; Swenson et al., 2018a; Yongacoglu et al., 2022).

Beyond cooperative settings, the use of best response dynamics to seek Nash equilibrium may not be justified. In purely adversarial settings, for instance, best response paths cycle and fail to converge (Balcan et al., 2023), and some alternative strategic dynamics are needed to drive play to equilibrium. Consider the following generalization of the best response update:

$$x_{t+1}^i = \begin{cases} x_t^i, & \text{if } x_t^i \in BR_0^i(\mathbf{x}_t^{-i}), \\ f^i(x_t^i, \mathbf{x}_t^{-i}) & \text{else.} \end{cases}$$

The update defined above is characterized by a "win stay, lose shift" principle, which only constrains the player to continue using a strategy when it is optimal. On the other hand, the player is not forced to use a best response when  $x_t^i \notin BR_0^i(\mathbf{x}_t^{-i})$ , and may experiment with suboptimal responses according to a function  $f^i : \mathbf{X} \to \mathcal{X}^i$ .<sup>1</sup> Allowing the function  $f^i$  to be any function from  $\mathbf{X}$  to  $\mathcal{X}^i$ , one generalizes best response updates and obtains a much larger set of sequences  $(\mathbf{x}_t)_{t\geq 1}$  and greater flexibility to approach equilibrium from new directions. This motivates the following definition of satisficing paths.

**Definition 4** A sequence of strategy profiles  $(\mathbf{x}_t)_{t=1}^T$ , where  $T \in \mathbb{N} \cup \{\infty\}$ , is called a satisficing path if it satisfies the following pairwise satisfaction constraint for any player  $i \in [n]$  and any t:

$$x_t^i \in \mathrm{BR}_0^i(\mathbf{x}_t^{-i}) \Rightarrow x_{t+1}^i = x_t^i.$$
(1)

The intuition behind satisficing paths is that they are the result of an iterative search process in which players settle upon finding an optimal strategy (i.e. a best response to the strategies of counterplayers) but are free to explore different strategies when they are not already behaving optimally. Note, however, that the definition above is merely a formal property of sequences of strategy profiles in  $\mathbf{X}$  and is agnostic to how a satisficing path is produced. The latter point will be important in the coming sections, where we analytically obtain a particular satisficing path as part of an existence proof.

We note that Condition (1) constrains only optimizing players. It does not mandate a particular update for the so-called unsatisfied player *i*, for whom  $x_t^i \notin BR_0^i(\mathbf{x}_t^{-i})$ . In particular,  $x_{t+1}^i$  can be any strategy without restriction, and  $x_{t+1}^i \notin BR_0^i(\mathbf{x}_t^{-i})$  is allowed. In addition to best response paths, constant sequences  $(\mathbf{x}_t)_{t\geq 1}$  with  $\mathbf{x}_t \equiv \mathbf{x}$  are always satisficing paths, even when  $\mathbf{x}$  is not a Nash equilibrium. Moreover, since arbitrary strategy revisions are allowed when a player is unsatisfied, if  $\mathbf{x}_1 \in \mathbf{X}$  is a strategy profile for which all players are unsatisfied, then  $(\mathbf{x}_1, \mathbf{x}_2)$  is a satisficing path for any  $\mathbf{x}_2 \in \mathbf{X}$ .

<sup>1.</sup> As a special case,  $f^i$  may simply be a best response selector, recovering the best response update.

**Definition 5** The game  $\Gamma$  has the satisficing paths property if for any  $\mathbf{x}_1 \in \mathbf{X}$ , there exists a satisficing path  $(\mathbf{x}_1, \mathbf{x}_2, ...)$  such that, for some finite  $T = T(\mathbf{x}_1)$ , the strategy profile  $\mathbf{x}_T$  is a Nash equilibrium.<sup>2</sup>

Satisficing paths were initially formalized by Yongacoglu et al. (2023), who proved that twoplayer games and n-player symmetric games have the satisficing paths property. However, whether general-sum n-player games have the satisficing paths property was left as an open question.

#### 3. Main result

**Theorem 6** Any finite normal-form game  $\Gamma$  has the satisficing paths property.

**Proof sketch.** Before presenting the formal proof, we describe the intuition of its main argument. In the proof of Theorem 6, we construct a satisficing path from an arbitrary initial strategy  $x_1$  to a Nash equilibrium by repeatedly switching the strategies of unsatisfied players in a way that grows the set of *unsatisfied* players after the update. Once the set of unsatisfied players is maximal, we argue that a Nash equilibrium can be reached in one step by switching the strategies of the unsatisfied players. The final point represents the main technical challenge in the proof, as switching the strategies of unsatisfied players. We address this challenge by showing the existence of a Nash equilibrium on the boundary of a strategy subset in which previously satisfied players remain satisfied.

To give the complete proof, we will require some additional notation, detailed in the next subsection, and some supporting results, detailed in Appendix A and Appendix B.

#### 3.1. Additional notation

In order to describe the construction of our satisficing path from  $x_1$  to a Nash equilibrium, we require the following sets, defined for any  $x \in X$ :

$$\operatorname{Sat}(\mathbf{x}) := \left\{ i \in [n] : x^i \in \operatorname{BR}_0^i(\mathbf{x}^{-i}) \right\}, \text{ and } \operatorname{UnSat}(\mathbf{x}) := [n] \setminus \operatorname{Sat}(\mathbf{x}).$$

A player in  $Sat(\mathbf{x}) \subseteq [n]$  is called *satisfied* (at  $\mathbf{x}$ ), and a player in  $UnSat(\mathbf{x})$  is called *unsatisfied* (at  $\mathbf{x}$ ). For  $\mathbf{x} \in \mathbf{X}$ , we also define

$$\operatorname{Access}(\mathbf{x}) := \left\{ \mathbf{y} \in \mathbf{X} : y^i = x^i, \, \forall i \in \operatorname{Sat}(\mathbf{x}) \right\}.$$

 $Access(\mathbf{x})$  is the subset of strategies that are accessible from strategy  $\mathbf{x}$ , to mean one can obtain strategy  $\mathbf{y} \in Access(\mathbf{x}) \subseteq \mathbf{X}$  from  $\mathbf{x}$  by switching (at most) the strategies of players who were unsatisfied at  $\mathbf{x}$ . We define a subset  $NoBetter(\mathbf{x}) \subseteq Access(\mathbf{x})$  as

NoBetter(
$$\mathbf{x}$$
) := { $\mathbf{y} \in \operatorname{Access}(\mathbf{x}) : \operatorname{UnSat}(\mathbf{x}) \subseteq \operatorname{UnSat}(\mathbf{y})$ }  
= { $\mathbf{y} \in \operatorname{Access}(\mathbf{x}) | \forall i \in \operatorname{UnSat}(\mathbf{x}), i \in \operatorname{UnSat}(\mathbf{y})$ },

The set NoBetter(x) consists of strategies y that are accessible from x and also fail to improve the status of players who were previously unsatisfied. The set name NoBetter(x) is chosen to

<sup>2.</sup> A more general definition, involving  $\epsilon \ge 0$  best responding and strategy subsets was studied in (Yongacoglu et al., 2023). In this paper, we consider true optimality and no strategic constraints, which additionally aids clarity.

suggest that the players unsatisfied at  $\mathbf{x}$  are not better off at  $\mathbf{y} \in \text{NoBetter}(\mathbf{x})$ , since they are unsatisfied at both  $\mathbf{x}$  and  $\mathbf{y}$ . We observe  $\mathbf{x} \in \text{NoBetter}(\mathbf{x})$ , hence  $\text{NoBetter}(\mathbf{x})$  is non-empty.

Finally, we define a set  $\operatorname{Worse}(\mathbf{x}) \subseteq \operatorname{NoBetter}(\mathbf{x})$  as

$$Worse(\mathbf{x}) := \{ \mathbf{y} \in NoBetter(\mathbf{x}) : UnSat(\mathbf{x}) \subsetneq UnSat(\mathbf{y}) \}$$
$$= \{ \mathbf{y} \in NoBetter(\mathbf{x}) | \exists i \in Sat(\mathbf{x}) : i \in UnSat(\mathbf{y}) \}$$

The set  $Worse(\mathbf{x})$  consists of strategies that are accessible from  $\mathbf{x}$ , that leave all previously unsatisfied players unsatisfied, and flip at least one previously satisfied player to being unsatisfied. In particular, if  $\mathbf{y} \in Worse(\mathbf{x})$ , this means  $|UnSat(\mathbf{y})| \ge |UnSat(\mathbf{x})| + 1$ . We observe that  $Worse(\mathbf{x})$  may be empty, and  $Worse(\mathbf{x}) \subseteq NoBetter(\mathbf{x}) \subseteq Access(\mathbf{x})$ .

#### 3.2. Proof of Theorem 6

**Remark 7** In the proof below, we analytically construct a satisficing path from  $\mathbf{x}_1$  to a Nash equilibrium. The process of selecting successive strategies  $\mathbf{x}_1, \mathbf{x}_2, \cdots$  and switching the component strategy of each player is done centrally, by the analyst, and should not be mistaken for some type of distributed learning algorithm.

**Proof** Let  $\mathbf{x}_1 \in \mathbf{X}$  be any initial strategy profile. We must produce a satisficing path of finite length terminating at a Nash equilibrium. Equivalently, we must produce a sequence  $\mathbf{x}_1, \ldots, \mathbf{x}_T$  with  $\mathbf{x}_{t+1} \in \operatorname{Access}(\mathbf{x}_t)$  for each t and  $\mathbf{x}_T$  a Nash equilibrium, where the length T may depend on  $\mathbf{x}_1$ . In the trivial case that  $\mathbf{x}_1$  is a Nash equilibrium, we put T = 1. The remainder of this proof focuses on the non-trivial case, where  $\mathbf{x}_1$  is not a Nash equilibrium.

To begin, we produce a satisficing path  $\mathbf{x}_1, \ldots, \mathbf{x}_k$  as follows. We put t = 1, and while both  $\operatorname{Sat}(\mathbf{x}_t) \neq \emptyset$  and  $\operatorname{Worse}(\mathbf{x}_t) \neq \emptyset$ , we arbitrarily fix  $\mathbf{x}_{t+1} \in \operatorname{Worse}(\mathbf{x}_t)$  and increment  $t \leftarrow t+1$ . By construction, we have

$$\emptyset \neq \text{UnSat}(\mathbf{x}_1) \subsetneq \cdots \subsetneq \text{UnSat}(\mathbf{x}_t) \subsetneq \text{UnSat}(\mathbf{x}_{t+1})$$

for each non-terminal iteration t, where the inequality holds because  $\mathbf{x}_1$  is not a Nash equilibrium. Thus, the number of unsatisfied players is strictly increasing along this satisficing path. Since the number of unsatisfied players is bounded above by n, and since we have assumed  $|\text{UnSat}(\mathbf{x}_1)| \ge 1$ , this process terminates in at most n-1 steps. Letting k denote the terminal index of this process, we have  $k \le n-1$ .

By the construction of the path  $(\mathbf{x}_1, \ldots, \mathbf{x}_k)$ , (at least) one of the following holds at index k: either  $\operatorname{Sat}(\mathbf{x}_k) = \emptyset$  or  $\operatorname{Worse}(\mathbf{x}_k) = \emptyset$ . In other words, either no player is satisfied at  $\mathbf{x}_k$ , or there is no accessible strategy that grows the subset of unsatisfied players.

**Case 1:**  $\operatorname{Sat}(\mathbf{x}_k) = \emptyset$ , and all players are unsatisfied at  $\mathbf{x}_k$ . In this case, we may switch the strategy of each player  $i \in [n]$  to any successor strategy. That is,  $\operatorname{Access}(\mathbf{x}_k) = \mathbf{X}$ . We fix an arbitrary Nash equilibrium  $\mathbf{z}_{\star}$ , put  $\mathbf{x}_{k+1} = \mathbf{z}_{\star}$ , and let T = k + 1. Then,  $(\mathbf{x}_1, \ldots, \mathbf{x}_T)$  is a satisficing path terminating at equilibrium.

**Case 2:**  $\operatorname{Sat}(\mathbf{x}_k) \neq \emptyset$  and  $\operatorname{Worse}(\mathbf{x}_k) = \emptyset$ . In this case, there are no accessible strategies that strictly grow the set of unsatisfied players, which is also non-empty.

Since  $Worse(\mathbf{x}_k) = \emptyset$ , the following holds: for any strategy  $\mathbf{y} \in NoBetter(\mathbf{x}_k)$  and any satisfied player  $i \in Sat(\mathbf{x}_k)$ , we have that  $i \in Sat(\mathbf{y})$ . (Otherwise, if  $i \in UnSat(\mathbf{y})$ , then  $\mathbf{y} \in Worse(\mathbf{x}_k)$ , since it flipped a satisfied player. But this contradicts the emptiness of  $Worse(\mathbf{x}_k)$ .)

We now argue that there exists a strategy profile  $\mathbf{x}_{\star}$  accessible from  $\mathbf{x}_k$  such that all players unsatisfied at  $\mathbf{x}_k$  are satisfied at  $\mathbf{x}_{\star}$ . That is, there exists an accessible strategy  $\mathbf{x}_{\star} \in \operatorname{Access}(\mathbf{x}_k)$ such that

$$\mathrm{UnSat}(\mathbf{x}_k) \subset \mathrm{Sat}(\mathbf{x}_\star). \tag{2}$$

To see that such a strategy  $\mathbf{x}_{\star}$  exists, note that fixing the strategies of the *m* players satisfied at  $\mathbf{x}_k$  defines a new game, say  $\tilde{\Gamma}$ , with n - m players, and the new game  $\tilde{\Gamma}$  admits a Nash equilibrium  $\tilde{\mathbf{x}}_{\star} = (\tilde{x}^i_{\star})_{i \in \text{UnSat}(\mathbf{x}_k)}$ . We extend  $\tilde{\mathbf{x}}_{\star}$  to be a strategy profile in the larger game  $\Gamma$  by putting  $x^i_{\star} = x^i_k$  for players  $i \in \text{Sat}(\mathbf{x}_k)$  while putting  $x^j_{\star} = \tilde{x}^j_{\star}$  for players  $j \in \text{UnSat}(\mathbf{x}_k)$ . By construction, we have that  $x^j_{\star} \in \text{BR}^0_0(\mathbf{x}^{-j}_{\star})$  for each  $j \in \text{UnSat}(\mathbf{x}_k)$ , so (2) holds.

From the set containment in (2), it is clear that  $\mathbf{x}_{\star} \notin \text{NoBetter}(\mathbf{x}_k)$ , since  $\text{NoBetter}(\mathbf{x}_k)$  consists of strategies accessible from  $\mathbf{x}_k$  in which unsatisfied agents remain unsatisfied, while the previously unsatisfied agents are satisfied at  $\mathbf{x}_{\star}$ . We now state a key technical lemma, which asserts that although  $\mathbf{x}_{\star}$  does not belong to  $\text{NoBetter}(\mathbf{x}_k)$ , it does lie on its boundary.

**Lemma 8** There exists a sequence  $\{\mathbf{y}\}_{t=1}^{\infty}$ , with  $\mathbf{y}_t \in \text{NoBetter}(\mathbf{x}_k)$  for each t, such that

$$\lim_{t\to\infty}\mathbf{y}_t = \mathbf{x}_\star$$

A proof of Lemma 8 given in Appendix B.

To conclude the proof, we will argue that  $\mathbf{x}_{\star}$  is in fact a Nash equilibrium for the game  $\Gamma$ . That is, in addition to (2), we also have  $\operatorname{Sat}(\mathbf{x}_k) \subset \operatorname{Sat}(\mathbf{x}_{\star})$ . To do so, we introduce functions  $F^i : \mathbf{X} \to \mathbb{R}$  for each player  $i \in [n]$ , given by

$$F^{i}(x^{i}, \mathbf{x}^{-i}) = \max_{a^{i} \in \mathbb{A}^{i}} R^{i}(\delta_{a^{i}}, \mathbf{x}^{-i}) - R^{i}(x^{i}, \mathbf{x}^{-i}),$$

for each  $\mathbf{x} = (x^i, \mathbf{x}^{-i}) \in \mathbf{X}$ . The functions  $\{F^i\}_{i=1}^n$  have the following useful properties, which are well known (Maschler et al., 2020), and are summarized in Appendix A. For each player  $i \in [n]$ ,

- a.  $F^i$  is continuous on **X**,
- b.  $F^i(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \mathbf{X}$ , and
- c. For any  $\mathbf{x}^{-i} \in \mathbf{X}^{-i}$ , a strategy  $x^i$  is a best response to  $\mathbf{x}^{-i}$  if and only if  $F^i(x^i, \mathbf{x}^{-i}) = 0$ .

Let  $(\mathbf{y}_t)_{t=1}^{\infty}$  be a sequence in NoBetter $(\mathbf{x}_k)$  converging to  $\mathbf{x}_{\star}$ , which exists by Lemma 8. For any previously satisfied player  $i \in \text{Sat}(\mathbf{x}_k)$ , since  $\text{Worse}(\mathbf{x}_k) = \emptyset$  and  $\mathbf{y}_t \in \text{NoBetter}(\mathbf{x}_k)$ , from a previous observation, we have that  $i \in \text{Sat}(\mathbf{y}_t)$ . Equivalently,  $x_k^i \in \text{BR}_0^i(\mathbf{y}_t^{-i})$ . Re-writing this using the function  $F^i$  and the notation  $y_t^i = x_k^i$  for satisfied players  $i \in \text{Sat}(\mathbf{x}_k)$ , we have

$$F^i(y_t^i, \mathbf{y}_t^{-i}) = 0, \quad \forall t \in \mathbb{N}$$

for any  $i \in \text{Sat}(\mathbf{x}_k)$ . By continuity of  $F^i$ , we have

$$0 = \lim_{t \to \infty} F^{i}(\mathbf{y}_{t}) = F^{i}\left(\lim_{t \to \infty} \mathbf{y}_{t}\right) = F^{i}(\mathbf{x}_{\star}),$$

establishing that player *i* is satisfied at  $\mathbf{x}_{\star}$ . Since  $i \in \text{Sat}(\mathbf{x}_k)$  was generic, we have  $\text{Sat}(\mathbf{x}_k) \subset \text{Sat}(\mathbf{x}_{\star})$ . Then, by (2), we also had  $\text{UnSat}(\mathbf{x}_k) \subset \text{Sat}(\mathbf{x}_{\star})$ , so indeed  $\text{Sat}(\mathbf{x}_{\star}) = [n]$ , and  $\mathbf{x}_{\star}$  is a Nash equilibrium accessible from  $\mathbf{x}_k$ .

We put T = k + 1 and  $\mathbf{x}_T = \mathbf{x}_{\star}$ , which completes the proof, since  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is a satisficing path terminating at a Nash equilibrium.

## 4. Discussion

#### **On decentralized learning**

Multi-agent reinforcement learning algorithms based on the "win stay, lose shift" principle characteristic of satisficing paths are especially well suited to decentralized applications, since players are often able to estimate the performance of their current strategy as well as the performance of an optimal strategy, even under partial information. In decentralized problems, coordinated search of the set  $\mathbf{X}$  of strategy profiles for a Nash equilibrium is typically infeasible, and players must select successor strategies in a way the depends only on quantities that can be locally accessed or estimated.

For instance, consider a trivial coordinated search method, where player *i* selects  $x_{t+1}^i$  uniformly at random from  $\mathcal{X}^i$  whenever  $\mathbf{x}_t$  was not a Nash equilibrium and selects  $x_{t+1}^i = x_t^i$  only when  $\mathbf{x}_t$  is a Nash equilibrium. This process is clearly ill-suited to decentralized applications, because player *i*'s strategy update depends on both a locally estimable condition (whether player *i* is best responding to  $\mathbf{x}_t^{-i}$ ) as well as a condition that cannot be locally estimated (whether another player  $j \neq i$  is best responding to  $\mathbf{x}_t^{-j}$ .) The satisfaction ("win stay") constraint plays a key role as a *local* stopping condition for satisficing paths, and rules out coordinated search of the set  $\mathbf{X}$  such as the trivial update outlined above. Examples of decentralized or partially decentralized learning algorithms leveraging satisficing paths in their analysis include (Foster and Young, 2006; Marden et al., 2009; Arslan and Yüksel, 2017; Yongacoglu et al., 2023). The analytic results of this paper suggest that algorithms such as these can be extended to wider classes of games and enjoy equilibrium guarantees under different informational constraints on the players.

#### On complexity and dynamics

In Theorem 6, we showed that for any finite *n*-player normal-form game  $\Gamma$  and any initial strategy profile  $\mathbf{x}_1 \in \mathbf{X}$ , there exists a satisficing path  $\mathbf{x}_1, \ldots, \mathbf{x}_T$  of finite length  $T = T(\mathbf{x}_1)$  terminating at a Nash equilibrium  $\mathbf{x}_T$ . From the proof of Theorem 6, one makes the following observations.

- 1. The length of such a path can be uniformly bounded above as  $T(\mathbf{x}_1) \leq n$ .
- 2. There exists a collection of strategy update functions  $\{f_{\Gamma}^i : \mathbf{X} \to \mathcal{X}^i | i \in [n]\}$  whose joint orbit is the satisficing path described by the proof of Theorem 6. That is,  $f_{\Gamma}^i(\mathbf{x}_t) = x_{t+1}^i$  for each player  $i \in [n]$ , every  $0 \le t \le T 1$ , and every  $\mathbf{x}_1 \in \mathbf{X}$ , where  $x_t^i$  is player i's component of  $\mathbf{x}_t$  in the satisficing path initialized at  $\mathbf{x}_1$ .

The proof of Theorem 6 is semi-constructive. At each step along the path, we describe how the next strategy profile should be picked (e.g.  $\mathbf{x}_{t+1} \in \text{Worse}(\mathbf{x}_t)$ ), but we do not suggest an algorithm for computing it. In at least one place, namely Case 1 where we put  $\mathbf{x}_T := \mathbf{z}_*$ , the path construction involves moving jointly to a Nash equilibrium in one step. The computational complexity of such a step is prohibitive (Daskalakis et al., 2009), underscoring that ours is an analytical existence result rather than a computational prescription.

Although we have shown that there exists a discrete-time dynamical system on **X** that converges to Nash equilibrium in *n* steps and can be characterized by update functions  $\{f_{\Gamma}^i\}_{i=1}^n$ , we note that our possibility result does not contradict the impossibility results of (Hart and Mas-Colell, 2003; Babichenko, 2012) or (Milionis et al., 2023). In particular, the functions  $\{f_{\Gamma}^i\}_{i=1}^n$  need not be (and usually will not be) continuous, violating the regularity conditions of Hart and Mas-Colell (2003) and Milionis et al. (2023), and furthermore the functions  $\{f_{\Gamma}^i\}_{i=1}^n$  depend crucially on the game  $\Gamma$  in a way that violates the uncoupledness conditions of (Hart and Mas-Colell, 2003) and (Babichenko, 2012).

#### **Open questions and future directions**

Several interesting questions about satisficing paths remain open. We now briefly describe some that we find especially practical or theoretically relevant.

While this paper dealt with satisficing paths defined using a best responding constraint, the original definition was stated using an  $\epsilon$ -best responding constraint, according to which a player who was  $\epsilon$ -best responding was not allowed to switch its strategy. Putting  $\epsilon = 0$ , one recovers the definition used here, but one may also select  $\epsilon > 0$ , which can be desirable to accommodate for estimation error in multi-agent reinforcement learning applications. The added constraint reduces freedom to switch strategies, and thus makes it more challenging to construct paths starting from a given strategy profile. On the other hand, the collection of Nash equilibria is a strict subset of the set of  $\epsilon$ -Nash equilibria, and one can attempt to guide the process to a different terminal point in a larger set. At this time, it is not clear to us whether the main result of this paper holds for small  $\epsilon > 0$ . It is clear, however, that the proof technique used here will have to be modified, since we have relied on Lemma 8, whose proof involved an indifference condition and invoked the fundamental theorem of algebra, and relaxing to  $\epsilon > 0$  would render such an argument ineffective.

A second interesting question for future work is whether multi-state Markov games with n > 2 players have the satisficing paths property. The case with n = 2 was resolved by Yongacoglu et al. (2023), but the proof technique used there did not generalize to  $n \ge 3$ . By contrast, our proof technique readily accommodates any number of players, but is designed for stateless normal-form games. Our proof used multi-linearity of the expected reward functions  $\{R^i\}_{i=1}^n$ , which does not generally hold in the multi-state setting.

In this work, satisficing paths were defined in a way that allowed an unsatisfied player *i* to change its strategy to any strategy in its set  $\mathcal{X}^i$ , without constraint. This is interesting in many problems where the set of strategies can be explicitly and directly parameterized, but may be unrealistic in games where the set of strategies is poorly understood or in which a player can effectively represent only a subset of its strategies  $\mathcal{Y}^i \subsetneq \mathcal{X}^i$ . In such games, the question more relevant for algorithm design is whether the game admits satisficing paths to equilibrium within the restricted subset  $\mathcal{Y}^1 \times \cdots \times \mathcal{Y}^n$ . This point was implicitly appreciated by both Foster and Young (2006) and Germano and Lugosi (2007) and explicitly noted in Yongacoglu et al. (2023).

## 5. Conclusion

Satisficing paths can be interpreted as a natural generalization of best response paths in which players may experimentally select their next strategy in periods when they fail to best respond to their counterplayers. While (inertial) best response dynamics drive play to equilibrium in certain well-structured classes of games, such as potential games and weakly acyclic games (Fabrikant et al., 2010), the constraint of best responding limits the efficacy of these dynamics in games with cycles in the best response graph (Pangallo et al., 2019). In such games, best response paths leading to equilibrium do not exist, and multi-agent reinforcement learning algorithms designed to produce such paths will not lead to equilibrium.

In this paper, we have shown that every finite normal-form game enjoys the satisficing paths property. By relaxing the best response constraint for unsatisfied players, one ensures that paths to equilibrium exist from any initial strategy profile. Multi-agent reinforcement learning algorithms designed to produce satisficing paths, rather than best response paths, thus do not face the same fundamental obstacle of algorithms based on best responding. While algorithms based on satisficing have previously been developed for two-player games normal-form games, symmetric Markov games, and several other subclasses of games, the findings of this paper suggest that similar algorithms can be devised for the wider class of *n*-player general-sum normal-form games.

## References

- Gürdal Arslan and Serdar Yüksel. Decentralized Q-learning for stochastic teams and games. *IEEE Transactions on Automatic Control*, 62(4):1545–1558, 2017.
- Yakov Babichenko. Completely uncoupled dynamics and Nash equilibria. *Games and Economic Behavior*, 76(1):1–14, 2012.
- Maria-Florina Balcan, Rattana Pukdee, Pradeep Ravikumar, and Hongyang Zhang. Nash equilibria and pitfalls of adversarial training in adversarial robustness games. In *International Conference* on Artificial Intelligence and Statistics, pages 9607–9636. PMLR, 2023.
- Lucas Baudin and Rida Laraki. Fictitious play and best-response dynamics in identical interest and zero-sum stochastic games. In *International Conference on Machine Learning*, pages 1664–1690. PMLR, 2022.
- Michel Benaïm, Josef Hofbauer, and Sylvain Sorin. Stochastic approximations and differential inclusions. *SIAM Journal on Control and Optimization*, 44(1):328–348, 2005.
- Lawrence E Blume. The statistical mechanics of strategic interaction. *Games and Economic Behavior*, 5(3):387–424, 1993.
- Joey Bose, Gauthier Gidel, Hugo Berard, Andre Cianflone, Pascal Vincent, Simon Lacoste-Julien, and Will Hamilton. Adversarial example games. *Advances in Neural Information Processing Systems*, 33:8921–8934, 2020.
- Mario Bravo, David Leslie, and Panayotis Mertikopoulos. Bandit learning in concave n-person games. Advances in Neural Information Processing Systems, 31, 2018.

- George W Brown. Iterative solution of games by fictitious play. *Act. Anal. Prod Allocation*, 13(1): 374, 1951.
- Ozan Candogan, Asuman Ozdaglar, and Pablo A Parrilo. Near-potential games: Geometry and dynamics. ACM Transactions on Economics and Computation (TEAC), 1(2):1–32, 2013.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- Georgios C Chasparis, Ari Arapostathis, and Jeff S Shamma. Aspiration learning in coordination games. *SIAM Journal on Control and Optimization*, 51(1):465–490, 2013.
- Steve Chien and Alistair Sinclair. Convergence to approximate Nash equilibria in congestion games. *Games and Economic Behavior*, 71(2):315–327, 2011.
- Constantinos Daskalakis, Paul W Goldberg, and Christos H Papadimitriou. The complexity of computing a Nash equilibrium. *Communications of the ACM*, 52(2):89–97, 2009.
- Constantinos Daskalakis, Rafael Frongillo, Christos H Papadimitriou, George Pierrakos, and Gregory Valiant. On learning algorithms for Nash equilibria. In Algorithmic Game Theory: Third International Symposium, SAGT 2010, Athens, Greece, October 18-20, 2010. Proceedings 3, pages 114–125. Springer, 2010.
- Alex Fabrikant, Aaron D Jaggard, and Michael Schapira. On the structure of weakly acyclic games. In Algorithmic Game Theory: Third International Symposium, SAGT 2010, Athens, Greece, October 18-20, 2010. Proceedings 3, pages 126–137. Springer, 2010.
- Lampros Flokas, Emmanouil-Vasileios Vlatakis-Gkaragkounis, Thanasis Lianeas, Panayotis Mertikopoulos, and Georgios Piliouras. No-regret learning and mixed Nash equilibria: They do not mix. Advances in Neural Information Processing Systems, 33:1380–1391, 2020.
- Dean Foster and Hobart Peyton Young. Regret testing: Learning to play Nash equilibrium without knowing you have an opponent. *Theoretical Economics*, 1(3):341–367, 2006.
- Fabrizio Germano and Gabor Lugosi. Global Nash convergence of Foster and Young's regret testing. *Games and Economic Behavior*, 60(1):135–154, 2007.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: a survey. Artificial Intelligence Review, pages 1–49, 2022.
- Sergiu Hart and Andreu Mas-Colell. Uncoupled dynamics do not lead to Nash equilibrium. American Economic Review, 93(5):1830–1836, 2003.

- Sergiu Hart and Andreu Mas-Colell. Stochastic uncoupled dynamics and Nash equilibrium. *Games* and economic behavior, 57(2):286–303, 2006.
- Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane, Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz, et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous Agents and Multi-Agent Systems*, 36(1):26, 2022.
- Josef Hofbauer and William H Sandholm. On the global convergence of stochastic fictitious play. *Econometrica*, 70(6):2265–2294, 2002.
- Yu-Guan Hsieh, Kimon Antonakopoulos, and Panayotis Mertikopoulos. Adaptive learning in continuous games: Optimal regret bounds and convergence to Nash equilibrium. In *Conference on Learning Theory*, pages 2388–2422. PMLR, 2021.
- Junling Hu and Michael P Wellman. Nash Q-learning for general-sum stochastic games. Journal of Machine Learning Research, 4(Nov):1039–1069, 2003.
- Amir Jafari, Amy Greenwald, David Gondek, and Gunes Ercal. On no-regret learning, fictitious play, and Nash equilibrium. In *ICML*, volume 1, pages 226–233, 2001.
- David S Leslie and Edmund J Collins. Individual Q-learning in normal form games. *SIAM Journal on Control and Optimization*, 44(2):495–514, 2005.
- David S Leslie and Edmund J Collins. Generalised weakened fictitious play. Games and Economic Behavior, 56(2):285–298, 2006.
- Yulong Lu. Two-scale gradient descent ascent dynamics finds mixed Nash equilibria of continuous games: A mean-field perspective. In *International Conference on Machine Learning*, pages 22790–22811. PMLR, 2023.
- Jason R Marden and Jeff S Shamma. Revisiting log-linear learning: Asynchrony, completeness and payoff-based implementation. *Games and Economic Behavior*, 75(2):788–808, 2012.
- Jason R Marden, H Peyton Young, Gürdal Arslan, and Jeff S Shamma. Payoff-based dynamics for multiplayer weakly acyclic games. SIAM Journal on Control and Optimization, 48(1):373–396, 2009.
- Michael Maschler, Shmuel Zamir, and Eilon Solan. *Game theory*. Cambridge University Press, 2020.
- Jason Milionis, Christos Papadimitriou, Georgios Piliouras, and Kelly Spendlove. An impossibility theorem in game dynamics. *Proceedings of the National Academy of Sciences*, 120(41): e2305349120, 2023.
- Dov Monderer and Aner Sela. A 2× 2game without the fictitious play property. *Games and Economic Behavior*, 14(1):144–148, 1996.
- Dov Monderer and Lloyd S Shapley. Fictitious play property for games with identical interests. *Journal of Economic Theory*, 68(1):258–265, 1996a.

- Dov Monderer and Lloyd S Shapley. Potential games. *Games and Economic Behavior*, 14(1): 124–143, 1996b.
- Ann Nowé, Peter Vrancx, and Yann-Michaël De Hauwere. Game theory and multi-agent reinforcement learning. *Reinforcement Learning: State-of-the-Art*, pages 441–470, 2012.
- Marco Pangallo, Torsten Heinrich, and J Doyne Farmer. Best reply structure and equilibrium convergence in generic games. *Science Advances*, 5(2):eaat1328, 2019.
- Julia Robinson. An iterative method of solving a game. *Annals of mathematics*, pages 296–301, 1951.
- Muhammed O Sayin, Francesca Parise, and Asuman Ozdaglar. Fictitious play in zero-sum stochastic games. *SIAM Journal on Control and Optimization*, 60(4):2095–2114, 2022a.
- Muhammed O Sayin, Kaiqing Zhang, and Asuman Ozdaglar. Fictitious play in Markov games with single controller. In *Proceedings of the 23rd ACM Conference on Economics and Computation*, pages 919–936, 2022b.
- Lloyd Shapley. Some topics in two-person games. Advances in Game Theory, 52:1–29, 1964.
- Satinder Singh, Michael J Kearns, and Yishay Mansour. Nash convergence of gradient dynamics in general-sum games. In UAI, pages 541–548, 2000.
- Brian Swenson, Ceyhun Eksin, Soummya Kar, and Alejandro Ribeiro. Distributed inertial bestresponse dynamics. *IEEE Transactions on Automatic Control*, 63(12):4294–4300, 2018a.
- Brian Swenson, Ryan Murray, and Soummya Kar. On best-response dynamics in potential games. *SIAM Journal on Control and Optimization*, 56(4):2734–2767, 2018b.
- Bora Yongacoglu, Gürdal Arslan, and Serdar Yüksel. Decentralized learning for optimality in stochastic dynamic teams and games with local control and global state information. *IEEE Transactions on Automatic Control*, 67(10):5230–5245, 2022.
- Bora Yongacoglu, Gürdal Arslan, and Serdar Yüksel. Satisficing paths and independent multiagent reinforcement learning in stochastic games. *SIAM Journal on Mathematics of Data Science*, 5 (3):745–773, 2023.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pages 321–384, 2021.

## **Proofs of technical lemmas**

We now discuss the properties of the auxiliary functions  $\{F^i : i \in [n]\}$  that were used in the proof of Theorem 6, and we prove Lemma 8.

We remark that for each player  $i \in [n]$ , we identify their set of mixed strategies  $\mathcal{X}^i = \Delta_{\mathbb{A}^i}$  with the probability simplex in  $\mathbb{R}^{\mathbb{A}^i}$ . Thus,  $\mathcal{X}^i$  inherits the Euclidean metric from  $\mathbb{R}^{|\mathbb{A}^i|}$ . Neighborhoods and limits in  $\mathcal{X}^i$  (or its subsets) are defined with respect to this metric. Similarly, we inherit a Euclidean metric for **X**. For  $\zeta > 0$ , we let  $N_{\zeta}(\mathbf{x})$  denote the  $\zeta$ -neighborhood of the strategy profile  $\mathbf{x} \in \mathbf{X}$ .

## Appendix A. Properties of the auxiliary functions

We begin by discussing the properties of the auxiliary functions  $\{F^i : i \in [n]\}$ , as they are relevant to characterizing best responses. The facts below are well-known. For a reference, see the text of Maschler et al. (2020).

Recall that for each player  $i \in [n]$ , the function  $F^i : \mathbf{X} \to \mathbb{R}$  is defined as

$$F^{i}(x^{i}, \mathbf{x}^{-i}) = \max_{a^{i} \in \mathbb{A}^{i}} R^{i}(\delta_{a^{i}}, \mathbf{x}^{-i}) - R^{i}(x^{i}, \mathbf{x}^{-i}), \quad \forall \mathbf{x} \in \mathbf{X}.$$

We now show that for any  $i \in [n]$ , the following hold:

- a.  $F^i$  is continuous on **X**,
- b.  $F^i(\mathbf{x}) \ge 0$  for all  $\mathbf{x} \in \mathbf{X}$ , and
- c. For any  $\mathbf{x}^{-i} \in \mathbf{X}^{-i}$ , a strategy  $x^i$  is a best response to  $\mathbf{x}^{-i}$  if and only if  $F^i(x^i, \mathbf{x}^{-i}) = 0$ .

The expected reward function  $R^i(\mathbf{x}) = \mathbb{E}_{\mathbf{a} \sim \mathbf{x}} \left[ r^i(\mathbf{a}) \right]$  can be expressed as a sum of products:

$$R^{i}(\mathbf{x}) = \sum_{\tilde{\mathbf{a}} \in \mathbf{A}} r^{i}(\mathbf{a}) \mathbb{P}_{\mathbf{a} \sim \mathbf{x}} \left( \mathbf{a} = \tilde{\mathbf{a}} \right) = \sum_{\tilde{\mathbf{a}} \in \mathbf{A}} r^{i}(\tilde{a}^{1}, \dots, \tilde{a}^{n}) \prod_{j=1}^{n} x^{j}(\tilde{a}^{j}), \quad \forall \mathbf{x} \in \mathbf{X}.$$

From this form, it is immediate that  $R^i$  is continuous on **X**. Moreover, it can easily be shown that  $R^i$  is multi-linear in **x**. That is, for any  $j \in [n]$ , fixing  $\mathbf{x}^{-j}$ , we have that  $x^j \mapsto R^i(x^j, \mathbf{x}^{-j})$  is linear.<sup>3</sup>

Since  $R^i$  is continuous on X and  $\mathbb{A}^i$  is a finite set, one has that the pointwise maximum of finitely many continuous functions is continuous. Thus, the function

$$\mathbf{x}^{-i} \mapsto \max_{a^i \in \mathbb{A}^i} R^i \left( \delta_{a^i}, \mathbf{x}^{-i} \right)$$

is continuous on  $\mathbf{X}^{-i}$ . Since  $F^i(x^i, \mathbf{x}^{-i}) = \max_{a^i \in \mathbb{A}^i} R^i \left( \delta_{a^i}, \mathbf{x}^{-i} \right) - R^i(x^i, \mathbf{x}^{-i})$  is the difference of continuous functions,  $F^i$  is also continuous. This proves item a.

From the multi-linearity of  $R^i$ , we have that, for fixed  $\mathbf{x}^{-i} \in \mathbf{X}^{-i}$ , the optimization problem  $\sup_{x^i \in \mathcal{X}^i} R^i(x^i, \mathbf{x}^{-i})$  is equivalent to a linear program

$$\sup_{x^i \in \mathbb{R}^{\mathbb{A}^i}} w_{\mathbf{x}^{-i}}^\top x^i, \quad \text{subject to } \begin{cases} 1^\top x^i = 1, \\ x^i \ge 0 \end{cases}$$

,

<sup>3.</sup> Of course, scaling inputs of  $R^i$  means the resulting argument is no longer a probability vector. However, one can simply linearly extend  $R^i$  to be a function on  $\mathbb{R}^d$ , where  $d = \sum_{j=1}^n |\mathbb{A}^j|$ .

where  $w_{\mathbf{x}^{-i}} \in \mathbb{R}^{\mathbb{A}^i}$  is a vector defined by  $w_{\mathbf{x}^{-i}}(a^i) := R^i(\delta_{a^i}, \mathbf{x}^{-i}).$ 

The vertices of the feasible set for the latter linear program are precisely the points  $\{\delta_{a^i} : a^i \in \mathbb{A}^i\}$ . This implies that  $\max_{a^i} R^i(\delta_{a^i}, \mathbf{x}^{-i}) \ge R^i(x^i, \mathbf{x}^{-i})$  for any  $x^i, \mathbf{x}^{-i}$ . Items b and c follow. From this formulation, one can also see that a player  $i \in [n]$  is satisfied at  $\mathbf{x} \in \mathbf{X}$  if and only if its strategy  $x^i$  is supported on the set of maximizers  $\operatorname{argmax}_{a^i \in \mathbb{A}^i} \{R^i(\delta_{a^i}, \mathbf{x}^{-i})\}$ .

#### Appendix B. Proof of Lemma 8

Recall that in the proof of Theorem 6,  $\mathbf{x}_{\star}$  was defined to be some strategy accessible from  $\mathbf{x}_{k} \in \mathbf{X}$  such that all players unsatisfied at  $\mathbf{x}_{k}$  were satisfied at  $\mathbf{x}_{\star}$ . The statement of Lemma 8 was the following.

**Lemma 8** There exists a sequence  $\{\mathbf{y}\}_{t=1}^{\infty}$ , with  $\mathbf{y}_t \in \text{NoBetter}(\mathbf{x}_k)$  for each t, such that  $\lim_{t\to\infty} \mathbf{y}_t = \mathbf{x}_{\star}$ .

**Proof** Suppose, to the contrary, that no such sequence exists. Then, there exists some  $\zeta > 0$  such that for every  $\mathbf{z} \in \operatorname{Access}(\mathbf{x}_k) \cap N_{\zeta}(\mathbf{x}_{\star})$ , one has  $\mathbf{z} \notin \operatorname{NoBetter}(\mathbf{x}_k)$ . That is, some player unsatisfied at  $\mathbf{x}_k$  is satisfied at  $\mathbf{z}$ . Equivalently, for some  $i \in \operatorname{UnSat}(\mathbf{x}_k)$ , we have  $z^i \in \operatorname{BR}_0^i(\mathbf{z}^{-i})$ . This implies that for that player *i*, that value of  $\zeta$ , and the strategy profile  $(z^i, \mathbf{z}^{-i}) \in N_{\zeta}(\mathbf{x}_{\star}), z^i$  is supported on the set  $\operatorname{argmax}_{a^i \in \mathbb{A}^i} \{R^i(\delta_{a^i}, \mathbf{z}^{-i})\}$ .

For each  $\xi \ge 0$ , we define a strategy profile  $\mathbf{w}_{\xi} \in \mathbf{X}$  as follows:

$$w_{\xi}^{i} := \begin{cases} (1-\xi)x_{k}^{i} + \xi \operatorname{Uniform}(\mathbb{A}^{i}), & \text{if } i \in \operatorname{UnSat}(\mathbf{x}_{k}) \\ x_{k}^{i}, & \text{else.} \end{cases}$$

Fixing  $\xi > 0$  at a sufficiently small value, by the continuity of the functions  $\{F^i\}_{i \in [n]}$ , we have that  $\mathbf{w}_{\xi} \in \text{NoBetter}(\mathbf{x}_k)$ . By the earlier discussion, we have that  $\mathbf{w}_{\xi} \notin N_{\zeta}(\mathbf{x}_{\star})$ .

A very important aspect of this construction is that  $w_{\xi}^i(a^i) > 0$  for each  $i \in \text{UnSat}(\mathbf{x}_k)$  and action  $a^i \in \mathbb{A}^i$ , so that  $w_{\xi}^i$  is fully mixed for each player who was unsatisfied at  $\mathbf{x}_k$ .

Next, for each  $\lambda \in [0, 1]$  and player  $i \in \text{UnSat}(\mathbf{x}_k)$ , we define

$$z_{\lambda}^{i} = (1 - \lambda)x_{\star}^{i} + \lambda w_{\xi}^{i}.$$

We also define  $z_{\lambda}^{i} = x_{k}^{i}$  for players  $i \in \text{Sat}(\mathbf{x}_{k})$ . For sufficiently small values of  $\lambda$ , say  $\lambda \leq \overline{\lambda}$ , we have that  $\mathbf{z}_{\lambda} \in N_{\zeta}(\mathbf{x}_{\star})$ , which implies  $\mathbf{z}_{\lambda} \notin \text{NoBetter}(\mathbf{x}_{k})$ .

This implies that there exists a player  $i^{\dagger} \in \text{UnSat}(\mathbf{x}_k)$  for whom

$$z_{\lambda}^{i^{\dagger}} \in \mathrm{BR}_{0}^{i^{\dagger}}\left(\mathbf{z}_{\lambda}^{-i^{\dagger}}\right)$$
, for infinitely many  $\lambda \in \left(0, \bar{\lambda}\right]$ .

(The existence of such a player is perhaps not obvious. As we previously noted, for  $\lambda < \overline{\lambda}$ , we have  $\mathbf{z}_{\lambda} \notin \text{NoBetter}(\mathbf{x}_k)$ , which means there exists *some* player  $i^{\dagger}(\lambda)$  that was unsatisfied at  $\mathbf{x}_k$  and is satisfied at  $\mathbf{z}_{\lambda}$ . The identity of this player may change with  $\lambda$ . To see that some particular individual must satisfy this best response condition infinitely often, one can apply the pigeonhole principle to the set  $\{\overline{\lambda}, \overline{\lambda}/2, \dots, \overline{\lambda}/m\}$  for arbitrarily large m.)

By our definition of  $z_{\lambda}^{i^{\dagger}}$  as a convex combination involving  $\text{Uniform}(\mathbb{A}^{i^{\dagger}})$ , we have that  $z_{\lambda}^{i^{\dagger}}$  is fully mixed and puts positive probability on each action in  $\mathbb{A}^{i^{\dagger}}$ . Using the characterization involving  $F^{i^{\dagger}}$ , the fact that  $z_{\lambda}^{i^{\dagger}} \in \text{BR}_{0}^{i^{\dagger}}(\mathbf{z}_{\lambda}^{-i^{\dagger}})$  and the fact that  $z_{\lambda}^{i^{\dagger}}$  is fully mixed together imply that  $R^{i^{\dagger}}(\delta_{a}, \mathbf{z}_{\lambda}^{-i^{\dagger}}) = R^{i^{\dagger}}(\delta_{a'}, \mathbf{z}_{\lambda}^{-i^{\dagger}})$ , for any  $a, a' \in \mathbb{A}^{i^{\dagger}}$ . This can be equivalently re-written as

$$\sum_{\mathbf{a}^{-i^{\dagger}}} r^{i^{\dagger}}(a, \mathbf{a}^{-i^{\dagger}}) \prod_{j \neq i^{\dagger}} \left\{ (1 - \lambda) x_{\star}^{j}(a^{j}) + \lambda w_{\xi}^{j}(a^{j}) \right\}$$
$$= \sum_{\mathbf{a}^{-i^{\dagger}}} r^{i^{\dagger}}(a', \mathbf{a}^{-i^{\dagger}}) \prod_{j \neq i^{\dagger}} \left\{ (1 - \lambda) x_{\star}^{j}(a^{j}) + \lambda w_{\xi}^{j}(a^{j}) \right\}$$
$$\iff \sum_{\mathbf{a}^{-i^{\dagger}}} \left[ r^{i^{\dagger}}(a, \mathbf{a}^{-i^{\dagger}}) - r^{i^{\dagger}}(a', \mathbf{a}^{-i^{\dagger}}) \right] \prod_{j \neq i^{\dagger}} \left\{ (1 - \lambda) x_{\star}^{j}(a^{j}) + \lambda w_{\xi}^{j}(a^{j}) \right\} = 0$$
(3)

for any  $a, a' \in \mathbb{A}^{i^{\dagger}}$ .

The lefthand side of the final equality (3) is a polynomial in  $\lambda$  of finite degree, but admits infinitely many solutions (from our choice of  $i^{\dagger}$ ). This implies that it is the zero polynomial. In turn, this implies that the left side of (3) holds for any  $\lambda \in [0, 1]$ , and in particular for  $\lambda = 1$ . This means  $z_1^{i^{\dagger}} \in BR_0^{i^{\dagger}}(\mathbf{z}_1^{-i^{\dagger}})$ , meaning  $\mathbf{z}_1 \notin NoBetter(\mathbf{x}_k)$ . On the other hand, we have  $\mathbf{z}_1 = \mathbf{w}_{\xi} \in NoBetter(\mathbf{x}_k)$ , a contradiction.

Thus, we see that there exists a sequence  $\{\mathbf{y}_t\}_{t=1}^{\infty}$ , with  $\mathbf{y}_t \in \text{NoBetter}(\mathbf{x}_k)$  for all t, such that  $\lim_{t\to\infty} \mathbf{y}_t = \mathbf{x}_*$ .